**Master's Thesis Flavio Ungherini**

# LUISS �𝕀

Department of Business and Management

Master of Science in Marketing, Major in Marketing Analytics & Metrics

Chair of Customer Intelligence & Big Data

# Machine Learning Approaches for Predicting Banks' Customers Credit Card Default - an Application to the Taiwanese Case

SUPERVISOR                                                    CANDIDATE

Prof. Luca Romeo                                          Flavio Ungherini

Studenti ID: 754751

CO-SUPERVISOR

Prof. Francisco Javier Villarroel Ordenes

Academic Year 2022/2023

**1**

# Master's Thesis Flavio Ungherini

Index:

# Master's Thesis Flavio Ungherini

# Chapter 1: Introduction to Artificial Intelligence Application and their Significance in the Banking Industry

**AN INTRODUCTION TO THE IMPORTANCE OF ARTIFICIAL INTELLIGENCE IN THE FINANCIAL SERVICES INDUSTRY**

Over many decades, banks have used technological innovation to reinvent customer interactions. ATMs were introduced during the 1960s; electronic card payments came later; in 2000s online banking gained widespread use before mobile-based "banking on the go" became widespread by 2010s.

AI-fueled digital age has flourished due to a host of factors: falling data storage costs, better access and connectivity for all, rapid advances in AI technologies enabling greater automation and efficiency gains across various tasks when implemented effectively.

Forbes conducted research in 2021 which concluded that daily internet user activity produced 2.5 quintillion bytes of data daily; each person creating approximately 1.7 megabytes per second [1].

Covid-19's pandemic outbreak resulted in widespread demand for financial services to increase significantly online. Banks made massive investments into Artificial Intelligence adoption as they strive to meet customers' evolving requirements and satisfy customers.

KPMG conducted a survey in 2021 which demonstrated that 83% of players operating within the financial industry used artificial intelligence of some form, placing financial services as second only behind manufacturing on this KPI measure.

Financial services players' increasing adoption of AI is further evidenced when looking at industry growth rates: an estimated 31.5% annual average rate acceleration in AI investments by 2027 could reach up to USD 54.7 billion [2].

Artificial Intelligence technologies can significantly increase revenues and lower fixed and variable costs through greater customer customization, higher automation rates, reduced error rates and better resource exploitation; new unrealized opportunities may emerge through improved data processing capacities that yield meaningful insight from vast quantities of information [3].

As banks seek to minimize financial losses, increase revenues, deliver superior services to their customers and boost satisfaction levels, they have increasingly turned to Artificial Intelligence and Machine Learning techniques for performing several strategic, operational and daily functions such as: credit scoring, fraud detection and prevention, identity verification, Customer Lifetime Value analysis and anti-churn policies creation, risk management and AML (Anti-Money Laundering).

References:

[1] https://www.forbes.com/sites/forbestechcouncil/2021/08/02/understanding-generation-data/?sh=3a3bb7b136b7

[2] https://info.kpmg.us/news-perspectives/technology-innovation/thriving-in-an-ai-world/ai-adoption-accelerated-during-pandemic.html

[3]https://www.mckinsey.de/~/media/McKinsey/Industries/Financial%20Services/Our%20Insights/AI%20bank%20of%20the%20future%20Can%20banks%20meet%20the%20AI%20challenge/AI-bank-of-the-future-Can-banks-meet-the-AI-challenge.pdf

# Master's Thesis Flavio Ungherini

**INVESTMENTS IN AI: A CROSS-INDUSTRY ANALYSIS**

Artificial Intelligence investments have seen dramatic surges over recent years, with costs becoming evermore elevated and increasing steadily. No sector or industry remains unsaturated with AI investments either already made or planned in coming years.

This section's intent is to present an overview of which and how much AI investments were made by some of the industries with highest economic relevance worldwide, along with anticipated investments planned over time for each mentioned industry in this analysis.

This analysis must begin by discussing which industry has made the largest investments into artificial intelligence research and development: tech. Major players like Google, Amazon and Microsoft all invest heavily in AI research [1]. By 2021 global AI investments had reached $67.9 billion - led by USA [2]. Future projections show this figure increasing further as AI plays an essential part in digital transformation initiatives [3].

On the other hand, financial services have made significant investments in AI technologies in recent years - banks and other financial institutions using AI to enhance customer experiences and streamline processes, with $10.5 billion spent by 2020[4]. Due to adoption of chatbots powered by artificial Intelligence as well as fraud detection systems and predictive analytics capabilities that utilize Artificial Intelligence investments will increase . A more in depth analysis regarding AI investments made and projected in this sector will soon be provided [5.].

Healthcare has recently seen AI investments surge as companies focus on patient outcomes improvement, drug discovery and personalized medicine to increase total investments of approximately $12. 2 billion by 2021 [2]. An upsurge is projected here due to AI-powered diagnosis systems, virtual assistants and personalized treatment plans being deployed [6].

As part of this analysis, it's worth noting the automotive industry as its Artificial Intelligence investments totalized an astounding $3.9 billion by 2020 [4], mostly directed towards developing intelligent transportation systems and autonomous vehicle development. Adopting AI-powered safety systems, predictive maintenance services and autonomous driving technologies is likely to further augment these artificial intelligence investments [7].

Additionally, customer experience improvement, supply chain optimization and inventory management has yielded $8.8 billion of AI investments by 2021 in retail industries [2], including AI chatbots, recommendation systems and personalized marketing campaigns that continue spending in artificial Intelligence [8].

By 2020, total investments totaled $1.2 billion in agriculture investments with expenditures spent to improve crop yield, reduce costs and implement sustainable farming practices [4]. Adopting precision farming using AI for precision harvesting, predictive analytics and automated machinery will further the AI investments into agriculture [9].

Manufacturing has also invested significantly in AI technology, spending an estimated $9.7 billion by 2021 to enhance product quality, reduce costs and enhance efficiency [10]. Furthermore, millions are set aside each year for predictive maintenance, quality control and supply chain optimization using artificial intelligence.

As the result of Covid-19 pandemic effects, personalized learning, intelligent tutoring systems and student performance analytics have become the cornerstones of AI investments in Education - which saw significant investments of $1.3 billion by 2021[4]. Furthermore, AI powered personalized learning platforms, adaptive testing systems and student performance tracking will continue to proliferate this industry[11].

Finally, it is pertinent to mention recent AI investments made both within telecommunications and energy industries.

Begun in 2016, since 2020 a total of $1.4 billion of investments have been undertaken with the aim of improving efficiency, cutting costs and heightening safety [4]. Future investments of this sort are expected to enable AI-powered predictive maintenance, asset optimization and energy demand forecasting [13].

Network operations improvement, customer service optimization and predictive maintenance expenditure will likely reach $2.2 billion by 2021 in the Telecommunications industry [2], and are projected to surge with AI-powered network optimization, predictive maintenance and customer service chatbot implementation and enhancement[12].

# Master's Thesis Flavio Ungherini

Artificial Intelligence investments have witnessed steady increases across various sectors over recent years and this trend is predicted to continue into 2019. Artificial Intelligence holds great potential to transform industries while simultaneously improving operational efficiencies, customer experiences and innovation.

References:

[1] Brynjolfsson, E., & McAfee, A. (2017). The business of artificial intelligence. Harvard Business Review, 95(1), 237-250.

[2] Statista. (2021). Artificial intelligence (AI) investments worldwide from 2015 to 2026 (in billion U.S. dollars): https://www.statista.com/statistics/979564/global-ai-investment/

[3] Gartner. (2021). Gartner identifies the top 10 strategic technology trends for 2021: https://www.gartner.com/en/newsroom/press-releases/2020-10-19-gartner-identifies-the-top-10-strategic-technology-trends-for-2021

[4] CB Insights. (2020). The AI in healthcare market map: Trends, startups, and corporates transforming healthcare with AI. Retrieved from https://www.cbinsights.com/research/report/ai-in-healthcare-market-map/

[5] Deloitte. (2019). AI and robotic process automation in financial services: https://www2.deloitte.com/content/dam/Deloitte/de/Documents/financial-services/AI_Robotic_Process_Automation_in_FS_2019.pdf

[6] McKinsey & Company. (2021). How artificial intelligence can drive healthcare innovation: https://www.mckinsey.com/industries/pharmaceuticals-and-medical-products/our-insights/how-artificial-intelligence-can-drive-healthcare-innovation

[7] PwC. (2020). The future of the automotive industry: Challenges and opportunities. Retrieved from https://www.pwc.com/gx/en/industries/automotive/automotive-insights/future-of-the-automotive-industry.html

[8] Emerj. (2021). AI in retail – current applications and future potential. Retrieved from https://emerj.com/ai-sector-overviews/ai-in-retail-current-applications-and-future-potential/

[9] Forbes. (2020). How AI is revolutionizing the agriculture industry. Retrieved from https://www.forbes.com/sites/forbestechcouncil/2020/06/25/how-ai-is-revolutionizing-the-agriculture-industry/?sh=4d4e00f2704b

[10] IBM. (2021). AI in manufacturing: Improving efficiency, productivity and quality. Retrieved from https://www.ibm.com/topics/manufacturing-ai

[11] EdTech Magazine. (2020). The impact of artificial intelligence on education. Retrieved from https://edtechmagazine.com/high

**AI INVESTMENTS IN THE FINANCIAL SERVICES INDUSTRY: A FOCUS ON BANKS**

As observed in the precedent paragraph, technology giants such as Google, IBM and Microsoft are investing heavily in artificial intelligence (AI), leading the charge in investing the highest proportion. Meanwhile, other industries, including finance, telecom, automotive as well as others analyzed above are making substantial investments with AI as their goal is streamlining processes and creating competitive advantages for themselves.

Financial services industry investments have been substantial when it comes to AI and machine learning technologies, second only to tech investments in this regard.

Banks increasingly depend on AI technologies for customer service enhancements, automation of routine tasks, fraud detection faster and risk management more effectively [1]. According to an Accenture research, banks' investments in AI should increase by an estimated 30-35% annually over three years up until 2022 [2]. This amount approximately equals $5.6 billion.

Citigroup invested significantly in artificial intelligence technologies to develop AI-powered chatbots to enhance customer service, answer inquiries, recommend products, and resolve account-related issues [4].

JPMorgan Chase has also invested significantly in AI, creating an AI research team in 2016 in order to explore its use within financial services industry.

In partnership with Artificial Intelligence (AI), a Contract Intelligence platform was created using NLP technologies in order to review legal documents; further investments have also been made to increase customer service capabilities optimization while streamlining internal processes.

Even while these investments provide many advantages, artificial intelligence also presents complications. One such complication arises with regard to compliance with regulations surrounding data processing and treatment - something banks are aware of and are investing in ethical frameworks designed to make AI use safer [5].

Researchers have expressed concern over how artificial intelligence (AI) might create problems; furthermore, existing biases reinforcement risks could result in unfair results [6].

A National Bureau of Economic Research study also demonstrated how AI-powered lending models may result in higher interest rates and reduced approval rates for minority borrowers than for non-minority ones [7].

Researchers are creating AI algorithms to address these concerns by devising methods to reduce bias in AI models and to promote transparency and accountability of these decisions made by AI algorithms. A recent Google study demonstrated how training models with more diverse data leads to less unbiased outcomes [8]. Likewise, researchers are creating methods of explaining why specific decisions made by these models were taken so as to promote greater transparency and accountability of decision making processes.

AI and Machine Learning industries are experiencing explosive growth with significant investments from across various sectors - most notably Financial Services Industry. Banks have made considerable investments into AI/ML solutions for customer service enhancement, cost control and fraud detection purposes; furthermore researchers are also striving to eradicate

any bias or discrimination present within AI algorithms to increase transparency and accountability with these technologies.

References:

[1] https://www2.deloitte.com/us/en/insights/industry/financial-services/artificial-intelligence-adoption-financial-services-industry.html

[2] https://www.accenture.com/_acnmedia/PDF-112/Accenture-Banking-on-AI.pdf#zoom=50

[3] https://www.jpmorgan.com/global/technology/artificial-intelligence

[4] https://www.business-standard.com/article/finance/banks-look-at-ai-powered-chatbots-to-improve-customer-service-118102300862_1.html

[5] https://www.pwc.com/us/en/library/financial-services/artificial-intelligence-in-banking.html

[6] https://hbr.org/2019/03/the-potential-for-bias-in-ai-is-real-but-it-can-be-eliminated

[7] https://www.nber.org/papers/w28210

[8] https://ai.googleblog.com/2019/04/reducing-unwanted-bias-in-deep.html

**AI-DRIVEN CREDIT SCORING**

Banks form relationships with their customer bases by offering products and services requested by the customers themselves- such as loans, mortgages, wealth management services and debit/credit cards.

Bank customers receiving services are chosen carefully based on their characteristics, working and income conditions as well as risk profiles and consumption habits.

Particular consideration is paid to these aspects when banks interact with existing or prospective customers who request activating credit lines with them.

Before granting credit to an applicant, banks must ensure they can repay any sum borrowed in line with contractually stipulated terms without creating additional difficulties for themselves as lenders.

After the 2008 financial crisis, banks are mandated to develop and adopt credit scoring and profiling systems which enable them to predict in advance and track customer creditworthiness and repayment behavior once economically bound through contractual relationships.

Basel 2 Accord required the adoption of credit scoring systems as mandatory measures; later integrated and partially replaced by Basel 3, these regulations have evolved along three pillars that outline essential conditions that banks must fulfill to engage in economic activities successfully.

These three pillars are respectively known according to their denominations: Minimum Capital Requirements," Supervisory Review Process," and Market Discipline".

**AI USAGE FOR FRAUD DETECTION AND PREVENTION**

Artificial Intelligence can play a valuable role in fraud detection and prevention by carefully analyzing customer transactions, spending patterns and locations which produce vast amounts of data that could reveal suspicious activities that require further examination.

Data can also be leveraged to develop models that utilize real-time analyses of bank transactions to detect hidden patterns that indicate unlawful financial activities in real time, helping both banks and customers to avoid losses. Machine Learning models can assist by learning from past cases of fraudulent transaction detection as they recognize similar anomalies on other occasions.

These models allow customers to compare their spending patterns, locations and transactions at issue with any possible irregularities that fall outside their typical spending schemes - whi-

ch could flag potential instances of fraudulent behavior if it falls outside this normal pattern of customer spending.

Final thoughts must also include that banks would benefit greatly by conducting thorough analyses on data to assist them in reducing false positives in fraud detection, improving accuracy in machine learning models and thus decreasing any unnecessary inconvenience to customers.

## CUSTOMER SERVICE AUTOMATATION AND IMPROVEMENT

Human chat service agents have often been replaced by chatbots backed by artificial intelligence (AI). Chatbots use natural language processing technology in answering customer inquiries and complaints 24/7 - helping banks increase response and time accuracy as well as reduce customer support maintenance costs by adopting this type of solution.

Though AI-powered chatbots provide cost and time savings benefits, their implementation often falls short of customer expectations leading to reduced willingness among users to comply with any requests made by them.

Chatbots cannot meet customer service demands completely at present; technology that would address this inefficiency should be developed so as to eliminate this shortcoming.

Although chatbot adoption by financial services firms often causes excitement and innovation, chatbots have actually been used across industries since Joseph Weizenbaum first created the ELIZA chatbot back in 1960.

It included both natural language processing software and an automated pattern-matching machine, to meet client demands for information processing services.

At Microsoft's introduction of Rinna, users were presented with an intelligent chatbot designed to engage them through light conversational systems utilizing natural language processing and deep learning technology. Rinna became available through Japan-wide messag-

ing app LINE as an imitation high school girl who provides human-like conversations experiences by means of "me".

These examples and descriptions of bots serve to highlight their diversity over the last sixty years in various industries and applications.

Fujitsu made headlines for their innovative use of chatbot technology within the financial services sector when they introduced FRAP (Finplex Robot Agent Platform), an enterprise chatbot designed specifically to be utilized within this sector and meet customers' demand for instantaneous information access while fulfilling service provider requirements for visualizing and managing AI learning results.

Sony Bank was established as an online Japanese commercial bank under Sony Financial Holdings in April 2001 and adopted FRAP for enhanced customer service quality.

Overall, a typical chatbot process in financial services involves four primary steps.:

-User Input: Users enter their queries or messages directly into the chatbot interface;

Natural Language Processing Analysis (NPPA): To fully comprehend user comments and determine their intent from their messages, the chatbot utilizes Natural Language Processing algorithms for this task.

-Response Generation: Chatbot AI algorithms use data gleaned in previous steps to respond to inquiries/requests made of them by users;

-Output: Chatbots typically present responses in conversational format to users and may provide relevant resources or links related to their query.

Sony Bank needed to assess whether adopting FRAP was best-suited to increasing customer requests processed simultaneously or improving existing user satisfaction. While service quality would ultimately increase overall, in both instances FRAP's specific chatbot business applications needed to be limited as much as possible.

Sony Bank prioritized efficiently responding to customer assistance inquiries that come in daily, providing superior quality of services in return.

Finally, among the various applications chatbots can provide to businesses in the financial services industry is one that stands out: fraud detection and prevention services.

## COMPLIANCE TO THE LOCAL DATA TREATMENT AND PROCESSING REGULATION

Compliance with local data processing and treatment regulations does not pose any threat to banks and organizations that rely heavily on machine learning (ML) and AI to manage customer information.

Compliance with current regulations presents banks (in our study scope) with an opportunity that could increase customer trust towards these financial institutions and banks.

As far as data protection regulations go, none can rival the General Data Protection Regulation (GDPR), with its scope limited to 27 European Union member states plus those belonging to EEA (European Economic Area). Thus it includes countries outside EU such as Norway, Iceland and Liechtenstein.

Machine learning and artificial intelligence technologies offer European banks a key opportunity for GDPR compliance by automating many privacy- and data protection-related tasks that involve processing large volumes of information.

Machine learning and artificial intelligence technologies can help banks comply with GDPR requirements by using machine learning/artificial intelligence techniques such as machine mapping to recognize personal information stored across documents, emails or system logs and categorize it accordingly. They also automate data mapping processes, identify flows within databases, generate reports on sensitive information as well as classify it based on its sensitivity level - helping banks adhere to GDPR safeguard requirements while keeping personal information within Europe rather than being transferred outside without appropriate safeguards in place.

Machine Learning and AI technologies can assist banks in quickly detecting data breaches or security threats by monitoring systems for abnormal behaviour or patterns - this enables banks to quickly detect any incidents as required by GDPR and report them immediately to relevant authorities.

## SPECIFIC AI APPLICATIONS BY CONTINENT FOR LOCAL REGULATION COMPLIANCE: EUROPE AND ASIA

Applications in Europe

As part of their obligations under GDPR, organizations are expected to ensure personal data is handled responsibly and securely [1], including ensuring its accuracy, currency, and only being stored briefly. Banks have leveraged AI technology in meeting this obligation through automating data classification, discovery, and retention [1, 2].

Data classification involves categorizing data according to its sensitivity and value. Banks utilize AI-powered tools for classifying sensitive information like financial details, social security numbers and personally identifiable information (PII) [1]. Using such technology helps protect PII against unapproved access while assuring adequate security controls are in place.

Data discovery involves the process of locating all personal information stored across an organization's systems and applications, with banks employing AI-powered tools to automate this task as they typically manage vast quantities of information. Such AI tools can analyze multiple systems simultaneously to detect risks or vulnerabilities before prioritizing efforts to secure personal data [2].

Data retention is another essential element of GDPR compliance for banks. AI tools enable them to utilize automated strategies for data management by automatically recognizing when files should be deleted or archived - helping minimize risks of data breach while mitigating potential GDPR penalties [1].

AI can also assist banks in meeting GDPR compliance by automating data subject access requests - which allow individuals access their own personal information - as well as detect possible data breaches or security incidents, which is another requirement of GDPR [2].

AI in GDPR compliance is an emerging area of study and practice, offering both potential benefits as well as risks to banks that utilize it. Banks should carefully weigh these advantages against potential drawbacks before considering AI for GDPR compliance efforts [1, 2, 3].

References:

[1] Kshetri, N. (2020). Using Artificial Intelligence to Meet GDPR Requirements. Communications of the ACM, 63(4), 56-63.

[2] Yan, J., Cheng, C., Wang, X., & Li, X. (2020). Artificial intelligence governance for compliance with GDPR. Journal of Business Research, 118, 491-498.

[3] Barbu, A. M., & Puican, F. (2020). The Role of Artificial Intelligence in GDPR Compliance. Management Dynamics in the Knowledge Economy, 8(1), 77-95.

Applications in Asia

Under Singapore's Personal Data Protection Act (PDPA), organizations must gain consent before collecting, using or disclosing an individual's personal data [1]. Banks use AI technologies like chatbots or AI forms for developing more effective consent forms that meet this obligation [2]. For example, DBS Bank in Singapore developed an AI chatbot called AI ChatBot that guides customers through consent processes while answering any inquiries that arise [3].

Hong Kong's Personal Data (Privacy) Ordinance (PDPO) mandates organizations take reasonable measures to ensure the accuracy of personal data [4]. Banks in Hong Kong have taken to using artificial intelligence (AI) as part of meeting this mandate by automating data val-

idation and verification [5]. Standard Chartered Bank in Hong Kong developed an AI system which detects potential errors or discrepancies within customer files before prompting staff members to rectify it [6].

Malaysia's Personal Data Protection Act (PDPA) mandates organizations take reasonable measures to secure personal data against unauthorized access or disclosure [7]. Banks in Malaysia have begun using AI to meet this obligation through more advanced security measures like biometric authentication and behavioral analysis [8]. CIMB Bank in Malaysia recently developed an AI-powered authentication system combining facial recognition technology and voice biometrics in its system for customer identity verification [9].

Overall, AI in data regulation compliance efforts is becoming an increasing trend across Asia. While AI offers several potential advantages when applied in this manner, there may also be challenges and risks involved with its deployment; so banks need to carefully weigh both benefits and risks when making their decision about using it responsibly and ethically [2, 5, 10].

References:

[1] Personal Data Protection Commission. (n.d.). Overview of the PDPA. Retrieved from https://www.pdpc.gov.sg/Overview-of-PDPA/The-Legislation/Overview

[2] Yap, W. Y., Tan, C. H., & Ong, T. S. (2021). Artificial intelligence in data protection: A review of GDPR, PDPA and CCPA. Journal of Information Privacy and Security, 17(2), 61-75.

[3] DBS Bank. (n.d.). DBS digibank. Retrieved from https://www.dbs.com.sg/personal/deposits/bank-with-ease/dbs-digibank

[4] Office of the Privacy Commissioner for Personal Data. (n.d.). The Personal Data (Privacy) Ordinance. Retrieved from https://www.pcpd.org.hk/english/data_privacy_law/ordinance_at_a_Glance/ordinance.html

[5] Gollmann, D. (2019). AI and data protection. Computer Law & Security Review, 35(1), 10-16.

[6] Standard Chartered Bank. (2019). Standard Chartered Bank (Hong Kong) Limited introduces an innovative system powered by artificial intelligence to ensure customer data accuracy. Retrieved from https://www.sc.com/hk/press-releases/scb-introduces-innovative-system-powered-by-artificial-intelligence-to-ensure-customer-data-accuracy.html

[7] Department of Personal Data Protection. (n.d.). Personal Data Protection Act 2010. Retrieved from https://www.pdp.gov.my/index.php/en/act-dppa-2010

[8] Mohamad, R., & Abdullah, M. F. (2021). Artificial intelligence and data privacy: A review of laws and practices in Malaysia. Journal of Cybersecurity, 7(1), 1-19.

[9] CIMB Bank. (n.d.). CIMB Clicks. Retrieved from https://www.cimbclicks.com.my/

[10] Custers, B. (2019). AI and data protection: Urgent questions. Computer Law & Security Review, 35(2), 263-271.

**ARTIFICIAL INTELLIGENCE DEPLOYMENT IS HELPING BANKS TO ASSESS THEIR CUSTOMERS' CLV: EXAMPLES AND REAL WORLD CASES DESCRIPTION**

Banks employ artificial intelligence (AI) in various ways to calculate Customer Lifetime Value (CLV), an integral metric for banks that measures customer profitability over their lifespan [1]. AI allows them to quickly process large volumes of customer data while making accurate predictions regarding future customer behavior - helping identify which customers might be more profitable, what products or services to offer them and tailor marketing campaigns accordingly [2].

Predictive modeling is one way banks use AI to assess customer lifetime value (CLV). Predictive models use statistical algorithms to examine customer data and predict future behavior [3]. Predictive modeling also allows banks to predict which customers may leave and provide proactive measures for keeping those customers [4].

# Master's Thesis Flavio Ungherini

Banks employ artificial intelligence through natural language processing (NLP). NLP is an AI technique which uses computer programs to interpret human speech [5]; banks may employ NLP to analyze customer reviews and social media posts for insights into customer preferences and behavior and tailor products and services accordingly [6].

AI can assist banks in tailoring marketing campaigns more precisely for each customer - which has proven more successful than generalized efforts [7]. AI-powered customer data analysis tools help banks determine which products and services appeal most strongly to each person within their customer database - ultimately increasing customer engagement and loyalty [8].

JPMorgan Chase Bank employs AI technology to assess customer lifetime value (CLV). They utilize it to analyze customer data to predict which ones might leave using COiN (Contract Intelligence). COiN employs machine learning algorithms that extract essential details quickly from legal contracts, automating this process while processing large volumes of information [9].

BBVA uses artificial intelligence (AI) to analyze customer data to provide customized products and services tailored specifically for individual customers. They developed AIDA - Autonomous Intelligent Digital Assistant; using machine learning algorithms to process customer information to deliver personalized product/service recommendations tailored specifically for each individual; AIDA has allowed BBVA to boost engagement and loyalty through customer data analysis [10].

Banking institutions utilize artificial intelligence (AI) tools such as predictive modeling, natural language processing and targeted marketing campaigns to calculate customer lifetime value (CLV). By including AI in their customer evaluation strategy, banks gain insight into customer behavior while simultaneously targeting products and services to meet customers' preferences; two examples being JPMorgan Chase's COiN program and BBVA's AIDA programs as examples of such applications in action.

# Master's Thesis Flavio Ungherini

References:

[1] Reinartz, W. and Kumar, V. (2003). The Impact of Customer Relationship Characteristics on Profitable Lifetime Duration. Journal of Marketing, 67(1), pp.77-99.

[2] Yap, B. W. and Ong, C. H. (2017). Customer lifetime value prediction with deep neural network. Expert Systems with Applications, 70, pp.104-112.

[3] Sun, Y., Neslin, S. A., Srinivasan, K., and Zhang, Z. J. (2010). Assortment Planning and Inventory Decisions under a Locational Choice Model. Journal of Retailing, 86(3), pp.285-299.

[4] Verhoef, P. C., Franses, P. H., and Hoekstra, J. C. (2002). The Effect of Relational Constructs on Customer Referrals and Referral - Acquisition Tradeoffs. Journal of Retailing, 78(3), pp.253-267.

[5] Jurafsky, D., and Martin, J. H. (2018). Speech and Language Processing (3rd ed.). Pearson.

[6] Alhabash, S. and Ma, M. (2017). A Tale of Four Platforms: Motivations and Uses of Facebook, Twitter, Instagram, and Snapchat among College Students? Social Media + Society, 3(1), pp.1-13.

[7] Verhoef, P. C., Neslin, S. A., and Vroomen, B. (2007). Multichannel Customer Management: Understanding the Research-Shopper Phenomenon. International Journal of Research in Marketing, 24(2), pp.129-148.

[8] Keiningham, T. L., Aksoy, L., Buoye, A., Cooil, B., and Malthouse, E. C. (2018). Linking Customer Experience to Value: A Cross-Sectional Analysis of Revenue and Net Promoter Score. Journal of Service Research, 21(3), pp.268-283.

[9] JPMorgan Chase. (2018). COiN: Contract Intelligence. Retrieved from https://www.jpmorgan.com/country/US/en/detail/1320576043489

[10] BBVA. (2018). BBVA launches AI-powered robo-advisor AIDA in Spain. Retrieved from https://www.bbva.com/en/bbva-launches-ai-powered-robo-advisor-aida-spain/

# Master's Thesis Flavio Ungherini

**HOW DO BANKS EXPLOIT THE CLV's ASSESSMENT TO DEVELOP ENGAGEMENT AND LOYALTY PROGRAMS**

Financial institutions frequently employ Customer Lifetime Value (CLV) analysis to identify their most profitable customers and formulate loyalty and engagement programs to retain and grow them while increasing overall profitability. Banks also can utilize CLV analysis for various ways of engaging their customer base more closely.

Banks can segment their customer base according to customer lifetime value (CLV), then craft targeted loyalty and engagement programs based on that figure. Customers with higher CLV may receive special services or experiences while those with lower CLV could get incentives in order to engage more actively with the bank [1].

Banks can leverage customer lifetime value (CLV) data to tailor offers specifically to each of their customers based on transaction history, purchase behavior and other variables - this allows banks to offer exclusive travel rewards or discounts when planning future journeys [2].

Banks can utilize customer lifetime value (CLV) scores to identify their most valued customers and deliver personalized customer service. Customers with high CLV may receive priority service or even be assigned a dedicated account manager [3].

Banks can implement loyalty programs that recognize customer engagement and show appreciation through rewards such as points or cash back incentives tailored to each individual customer's lifetime value (CLV) [1].

American Express, Bank of America and Capital One are three financial sector organizations using Customer Lifetime Value-based loyalty and engagement programs to fully engage their customer bases. American Express leverages CLV-based programs for targeted loyalty initiatives; Bank of America employs it to identify their most valuable customers to deliver personalized service; while Capital One utilizes it to tailor offers based on customers' transaction histories and behaviors [6].

American Express is well known for using customer lifetime value (CLV) analysis to identify customer segments with tailored loyalty programs like Membership Rewards. Membership Rewards offer travel, shopping and dining experiences tailored specifically for each customer

**23**

segment's preferences; American Express also provides access to special events, airport lounges and concierge services [4].

Bank of America utilizes customer lifetime value analysis (CLV) to identify its most valued customers and offer customized customer service. Their Preferred Rewards program features three tiers of benefits tied directly to customers' qualifying balance across their Bank of America deposit accounts and Merrill investment accounts - the higher your total qualifying balance is, the greater will be your benefits and rewards like priority customer service, free stock trades and higher cash rewards [5].

Capital One is widely known for using CLV to tailor offers specifically to its customers based on their transaction histories and behavior, such as travel or purchases with foreign transaction fees removed or increased travel-related rewards; similarily, machine learning techniques detect life events like marriage or birth to provide tailored financial products [6].

These examples demonstrate how financial institutions can utilize Customer Lifetime Value (CLV) to develop tailored loyalty and engagement programs that cater directly to each customer's individual needs, preferences and behaviors - increasing customer satisfaction rates while simultaneously improving retention rates and overall profits.

References:

[1] Kumar, V., & Reinartz, W. J. (2016). Customer Lifetime Value (CLV): Marketing Models and Applications. Foundations and Trends® in Marketing, 10(4), 197-246.

[2] Gupta, S., Lehmann, D. R., & Stuart, J. A. (2004). Valuing customers. Journal of Marketing Research, 41(1), 7-18.

[3] Verhoef, P. C., & Donkers, B. (2005). The Effect of Customer Acquisition and Retention Oriented Customer Relationship Management Policies on Customer Value. Journal of Retailing, 81(3), 215-229.

[4] O'Brien, D. (2019). American Express Uses Customer Lifetime Value for Effective Segmentation. Customer Think. Retrieved from https://customerthink.com/american-express-uses-customer-lifetime-value-for-effective-segmentation/.

[5] Perrin, R. (2020). Bank of America's Loyalty Program Increases Engagement, Deepens Relationships with Customers. The Financial Brand. Retrieved from https://thefinancial-brand.com/98202/bank-of-americas-loyalty-program-increases-engagement-deepens-relationships-with-customers/.

[6] Tuchman, Y. (2018). Capital One Uses Machine Learning to Target Customers' Life Events. American Banker. Retrieved from https://www.americanbanker.com/news/capital-one-uses-machine-learning-to-target-customers-life-events.

## POSSIBLE THREATS: AI JOB REPLACEMENT IN THE FINANCIAL SERVICES INDUSTRY

Artificial Intelligence (AI) and Machine Learning (ML) technologies have seen rapid adoption by banks over recent years. AI/ML tools are being employed to increase efficiency and accuracy in various operations such as risk management, fraud detection and customer service; one study by World Economic Forum (WEF) even showed AI could potentially cut operational costs by 22% [1].

AI and ML offer banks one of the most significant advantages in automating routine yet time-consuming tasks such as data entry and analysis, freeing human employees to focus on more complex ones that require human judgment or creativity. Furthermore, these technologies can assist banks with uncovering patterns or trends within large datasets that human analysts might miss altogether for more effective decision-making.

AI and ML in banking have led to various applications and tools. JPMorgan Chase developed an AI-powered document reviewer, cutting the time and costs of manual review [2]. Goldman Sachs introduced Marcus, an AI chatbot that assists customers with queries and account management [3]. Moreover, finally, the Industrial and Commercial Bank of China (ICBC) introduced an AI robot that provides customer service, decreasing employee workload and improving overall customer experience [4].

However, AI and ML in banking have raised concerns over whether these technologies threaten human jobs. While AI/ML could lead to job displacements for human employees, these tools should enhance productivity and decision-making by helping employees focus more on tasks requiring human judgment, creativity, or problem-solving skills - not replace human employees altogether. So AI/ML are expected to change work patterns but not eliminate them.

Additionally, AI and ML use in banking is predicted to open up job opportunities requiring new skills such as data analysis, coding and digital literacy. According to research from McKinsey Global Institute, these new jobs could create 20-50 million global job openings by 2030 [5]. For banking industry employees to meet these demands of an evolving job market, they will require investment in training and upskilling programs designed to prepare employees for these jobs.

Though AI and ML offer significant potential benefits, implementation in banking poses its own set of difficulties. Concerns regarding AI algorithms' accuracy, bias and privacy risks and using these technologies for sensitive financial data analysis is also prevalent. Incorporation requires significant investments in technology infrastructure and training/development for human employees [6].

PwC released a report outlining the potential advantages and challenges AI and ML present in banking [7]. According to PwC's assessment, banks that invest in AI/ML may experience dramatic increases in efficiency, customer service quality and risk management if implemented successfully, but, at the same time, these institutions must also address challenges related to data quality/privacy/privilege, compliance requirements, cybersecurity threats as well as making sure benefits from these technologies are shared equitably among employees, customers, shareholders.

Artificial Intelligence (AI) and Machine Learning (ML) technologies are promising for banks, significantly improving efficiency, customer service, risk management, and regulatory oversight. However, these technologies also present unique challenges that must be met head-on to reap all their benefits while mitigating associated risks [9].

However, AI and ML present one key challenge: their algorithms' potential for bias. Studies have revealed AI algorithms' propensity towards bias based on factors like race, gender and socioeconomic status. [8] In banking industries like banking, this could result in unfair treatment of specific groups, such as minority customers or those with lower credit scores; to address this issue banks must ensure their AI algorithms are transparent, explainable and free of bias.

AI and ML present another hurdle when applied to financial data analysis, as they have security and privacy implications that banks must carefully consider using these technologies to process sensitive customer data. Banks must establish effective cybersecurity protocols in place to safeguard customer information against theft or misuse and comply with relevant regulations such as Europe's General Data Protection Regulation (GDPR), which sets stringent requirements on collection, use and storage [10].

PwC published a report outlining the potential advantages and challenges associated with AI/ML in the banking industry. According to this research, banks investing in AI/ML may experience improvements in efficiency, customer service and risk management; while facing some hurdles, such as data quality/privacy/security compliance. In order to reap maximum returns from AI/ML investments by all stakeholders - employees, customers and shareholders [11].

Academic literature and reports by industry experts confirm that artificial intelligence and machine learning (AI/ML) technologies are transforming the banking industry, with significant impacts on human jobs becoming evident. AI and ML may result in job displacement; however, they should not be seen as replacements but as tools to enhance productivity and decision-making.

AI/ML may alter work environments but will also open up opportunities that require different skill sets; banks must invest in training their employees so that they are equipped to adapt to an evolving job market [11].

Hence, adopting AI and ML technologies within the banking industry offers significant potential benefits; however, they also present distinct challenges which must be met head-on. Banks must ensure their AI algorithms are open, understandable, free of bias and protected with robust cybersecurity measures for customer data security. Likewise, banks should invest

in training their employees to have all the skills required for today's changing job market; ultimately, realizing these technologies' benefits while mitigating any associated risks is possible by meeting this challenge head-on.

References:

[1] World Economic Forum, "The Future of Jobs Report 2020," https://www.weforum.org/reports/the-future-of-jobs-report-2020.

[2] JPMorgan Chase, "JPMorgan Chase Develops Contract Intelligence Platform to Help Customers Analyze Legal Documents," https://www.jpmorganchase.com/news-stories/contract-intelligence-platform.

[3] Goldman Sachs, "Introducing Marcus by Goldman Sachs," https://www.marcus.com/us/en.

[4] Xinhua News Agency, "ICBC's AI robot serves customers at a Beijing branch," https://www.xinhuanet.com/english/2017-06/26/c_136394131.htm.

[5] McKinsey Global Institute, "Skill Shift: Automation and the Future of the Workforce," https://www.mckinsey.com/featured-insights/future-of-work/skill-shift-automation-and-the-future-of-the-workforce.

[6] D. Davenport and R. Ronanki, "Artificial Intelligence for the Real World," Harvard Business Review, April 2018, https://hbr.org/2018/04/artificial-intelligence-for-the-real-world.

[7] A. Arif and M. R. Hasan, "Artificial Intelligence in Banking: A Review of the Applications and Opportunities," Journal of Innovation and Entrepreneurship, vol. 9, no. 6, 2020, https://doi.org/10.1186/s13731-020-00122-6.

[8] C. Angwin, J. Larson, S. Mattu, and L. Kirchner, "Machine Bias," ProPublica, May 2016, https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing.

[9] PwC, "Financial Services Technology 2020 and Beyond: Embracing Disruption," 2017, https://www.pwc.com/gx/en/industries/financial-services/publications/assets/pwc-fs-technology-2020-and-beyond.pdf.

[10] PwC, "AI in Financial Services: Shaping the Future of Customer Experience, Efficiency, and Security," 2018, https://www.pwc.com/us/en/industries/financial-services/research-institute/top-issues/ai-in-financial-services.html.

[11] H. Kim and J. Lee, "Artificial Intelligence in Banking and Finance: Current State and Future Directions," International Journal of Financial Studies, vol. 9, no. 4, 2021, https://doi.org/10.3390/ijfs9040034.

## MACROECONOMIC CONTEXT: HOW THE INTEREST RATES INCREASE IS CHANGING LENDING BY BANKS AND INCREASING NPLs

Central banks typically utilize interest rate policies to combat inflation. When faced with rising prices, central banks may increase rates to discourage spending and borrowing to help ease inflationary pressures by decreasing demand for goods and services and relieving inflationary pressures. If inflation levels suddenly decrease unexpectedly low, central banks could decrease rates to stimulate borrowing and spending and spur economic development forward.

In 2018, for instance, the US Federal Reserve raised interest rates several times to curb rising inflation and its goal was "prevent inflation from exceeding their target rate of 2%", as measured by economic data such as CPI [1].

As another example, in 2014 the European Central Bank (ECB) implemented negative interest rates to combat low inflation and foster economic growth. By charging banks holding reserves with them fees for keeping these reserves longer at ECB banks holding reserves there longer, negative rates encouraged banks to lend more freely to businesses and consumers [2]. Furthermore, quantitative easing (purchasing large quantities of government bonds to inject liquidity into financial system) was employed as another form of regulation designed to enhance growth.

Emerging market central banks have employed interest rate policies to combat inflation. Brazil increased their benchmark interest rate from 2.75% in March 2021 to 7.0% by March

# Master's Thesis Flavio Ungherini

2022 in order to bring inflation within their targeted range of 3.75% [3]. When managing inflation in these emerging economies, central banks must face challenges associated with volatile exchange rates and political unrest [4].

Interest rate policies depend on various variables, including economic conditions, inflation levels and global economic factors [1, 2, 4]. Interest rate policies have both intended and unintended repercussions; there remains ongoing debate regarding effective methods central banks should use to manage inflation while simultaneously maintaining economic stability.

As interest rates increase, banks' cost of borrowing increases and will impact banking operations - such as mortgage lending operations - negatively.

Non-performing loans could experience higher default rates when interest rate hikes take effect, as borrowers find themselves struggling to meet payments on their mortgages. How this increases affect borrowers will depend on various factors like economic strength, employment rate and debt carried by each borrower.

Studies indicate an increase in interest rates may cause non-performing loans to rise; though their relationship may sometimes be indirect. An interest rate increase could cause higher levels of non-performing loans in countries with weaker macroeconomic conditions [5], while higher rates were linked with greater default risk among borrowers, although their correlation was weaker among those with higher credit scores [6].

Overall, the relationship among interest rates, lending operations, and non-performing loans can be intricate and can depend on several variables. An increase in interest rates could potentially augment borrowing costs or default rates depending on both economic considerations as well as personal ones [5, 6].

References:

[1] Blanchard, O. (2019). Macroeconomics (7th ed.). Pearson.

[2] Goodhart, C. A. E., & Pradhan, M. (2017). Central banking at a crossroads: What can we learn from history? Journal of Financial Stability, 30, 1-12.

[3] International Monetary Fund. (2021). Brazil: Staff concluding statement of the 2021 Article IV Mission.

[4] McKinsey & Company. (2016). Central banking in emerging markets: Addressing the unique challenges.

[5] Andriosopoulos, K., Gaganis, C., & Pasiouras, F. (2018). The impact of interest rates on bank risk: Evidence from a global sample. Journal of Financial Stability, 35, 93-118.

[6] Chakraborty, I., & Ray, A. (2015). Determinants of loan default: Evidences from SMEs of a North Indian state. International Journal of Entrepreneurship and Small Business, 26(1), 105-122.

**AI ENABLING BANKS TO MANAGE NON-PERFOMING LOANS (NLPs)**

Artificial Intelligence (AI) has revolutionized the banking industry, improving operational efficiencies, risk management practices, customer service levels and profitability [1]. AI helps banks manage non-performing loans (NPLs) more effectively by performing tasks usually done by humans, such as pattern recognition, decision-making and natural language processing. It also detects potential NPLs through data sources such as credit bureaus, financial statements, or social media. It can provide early warning signals so banks can take proactive steps to prevent NPLs from occurring [2].

Machine learning algorithms can comb through large volumes of customer data to use machine learning to develop predictive models that estimate the likelihood of borrower default on loans. By taking into account various data points such as credit history, income, employment status and other financial considerations,, AI algorithms can provide more accurate credit risk assessments than traditional credit scoring methods and thus better manage loan portfolios while decreasing exposure to NPLs [3]. This point helps banks effectively manage loan portfolios while decreasing exposure to non-performing loans.

AI can assist banks in more efficiently managing their collections and recovery processes by automating the processes for identifying delinquent borrowers, predicting their willingness and ability to pay, and suggesting recovery solutions [2]. By streamlining collections and re-

covery processes using AI solutions, banks can reduce the time and costs associated with managing NPLs.

AI can assist banks in detecting and preventing fraud by real-time analyzing customer data and detecting suspicious transactions, patterns, or behaviors using machine learning algorithms. AI can detect sophisticated fraud attempts such as identity theft, money laundering and cyber attacks [4]. By decreasing risk from fraud attempts, AI helps banks protect loan portfolios while decreasing exposure to NPLs.

Hence, AI can assist banks in managing non-performing loans effectively by offering early warning signals, more accurate credit scoring, efficient collections and recovery processes, and improved fraud detection. Leveraging its power, banks can better manage their loan portfolios while decreasing risks and increasing profitability.

Non-performing loans (NPLs) refer to loans that have gone unserviced for an extended period, typically 90 days or longer. Banks approach non-performing loans differently: either they try to recover them through restructuring or selling collateral or write them off as losses on their books [5].

One approach banks take when handling non-performing loans is transferring them to a Special Purpose Vehicle (SPV). An SPV is a legal entity created solely to hold and manage specific assets; banks transfer non-performing loans for cash transfer into this SPV, which then attempts to recover these loans by restructuring, selling collateral, or taking legal action against borrowers [6].

An SPV provides several advantages when managing non-performing loans for banks. First, it enables the bank to remove non-performing loans from its balance sheet and free up capital to lend to other borrowers. Second, these SPVs typically employ experienced professionals who manage non-performing loans, which may increase the chances of recovering the loans more successfully. Thirdly, an SPV may pursue legal action without disrupting relationships between itself and borrowers, which is particularly advantageous when one or both borrowers represent valuable customers of the bank [7].

History provides numerous examples of banks using SPVs to manage non-performing loans. Following the 2008 financial crisis, many US and European banks created SPVs in response,

isolating non-performing loans from their core operations [8]. Spain saw one such SPV called Anida established by bank BBVA for real estate assets containing non-performing loans by 2010 to reduce portfolio size from EUR20.5 billion down to EUR4.4 billion by 2018 [9]. Italy, too has implemented SPVs that hold these types of assets [10].

Banks typically address non-performing loans by either trying to recover them or transferring them into an SPV. SPVs offer banks an effective solution for dealing with non-performing loans. They allow the institution to remove these liabilities from its balance sheet, improving financial standing while freeing up capital to lend to new borrowers.

References:

[1] Almubarak, A., & Keasey, K. (2020). Artificial intelligence and non-performing loans: A systematic literature review. Journal of Economic Surveys, 34(3), 581-606.

[2] Bao, Y., Yang, Y., & Wu, J. (2020). Credit risk assessment based on machine learning: A review. Journal of Finance and Data Science, 6(3), 293-313.

[3] Chen, Y., Li, Q., & Li, J. (2020). Using machine learning for credit risk management: A review. Journal of Risk Research, 23(5), 590-611.

[4] Nuzzo, G. (2019). Artificial intelligence and fraud detection in financial services. Journal of Financial Crime, 26(2), 383-393.

[5] Berger, A. N., DeYoung, R., Genay, H., & Udell, G. F. (2000). Globalization of financial institutions: Evidence from cross-border banking performance. Brookings-Wharton Papers on Financial Services, 4, 23-124.

[6] Ismail, S., & Abdulla, M. A. (2018). Non-performing loans, banks' performance and macroeconomic variables: Evidence from Malaysia. Journal of Financial Crime, 25(4), 1044-1058.

[7] Bertrand, J. W., Duarte, F. D., & Young, S. G. (2017). Special purpose vehicles and the resolution of financial distress. Journal of Corporate Finance, 45, 401-427.

[8] Gropp, R., Correa, R., & Sapriza, H. (2016). Sovereign credit risk, banks' government support, and bank stock returns around the world. Journal of Money, Credit and Banking, 48(2-3), 409-447.

[9] Roldán, J. (2019). The role of banks in the resolution of non-performing loans: the Spanish experience. Journal of Banking Regulation, 20(3), 196-205.

[10] Zavatta, G. (2019). The Italian experience in managing non-performing loans. Journal of Banking Regulation, 20(3), 221-229.

# Chapter 2: Machine Learning Applications for Analyzing Historical Socio-economic Issues in the Banking Industry

**FINANCIAL CRISIS PROVOKED BY BANKS' INABILITY TO PREDICT CUSTOMERS' DEFAULT**

One of the primary functions of banks is to provide loans to individuals and corporations to fulfill their credit needs. However, sometimes borrowers default on these loans, leading to substantial financial losses for banks worldwide. Many large and prominent banks worldwide have experienced losses due to customer defaults, which has seriously affected their financial stability and reputation. In this essay, different cases were described in which banks could not anticipate customer defaults resulting in insufficient credit accrual resulting in lousy debt accrual, how these banks handled these situations, and how governments helped these institutions recover.

# Master's Thesis Flavio Ungherini

Case N.1: The 2008 Global Financial Crisis

The 2008 global financial crisis was one of the worst economic crises ever, impacting all sectors worldwide, but particularly banks and financial institutions in the US. While its roots lay within the subprime mortgage market in the US, its ripples spread worldwide as banks suffered huge losses leading to a severe credit crunch over time. Its cause lay in excessive risk-taking, lax regulation, and flawed models [5].

Banks' inability to accurately forecast customer default was a key factor contributing to the crisis. Although banks employed complex financial models to assess credit risk, these models proved insufficient when faced with unprecedented market volatility and risk-taking behavior that marked the pre-crisis period. This section will discuss why and how banks failed to accurately anticipate customer defaults, any consequences from this miscalculation, and measures they took afterward to recover from it.

Banks have relied on credit scoring models to assess customers' creditworthiness for years, using statistical algorithms derived from historical data to estimate the likelihood that a borrower will default on a loan. While these models proved relatively accurate in previous instances, they failed to predict the magnitude of defaults during the financial crisis as they failed to capture subprime mortgage risk [6] adequately.

Failing to forecast customer defaults was further compounded by banks' over-reliance on credit ratings provided by credit rating agencies, who provided inaccurate ratings that led banks to invest heavily in mortgage-backed securities rated highly but later discovered to be subpar quality - leading them down the path toward financial collapse. Reliance on credit ratings as an accurate measure of risk also contributed significantly to its severity [1].

Failing to anticipate customer defaults had severe repercussions for both banks and the wider economy, as borrowers defaulted on loans and caused huge financial losses for lenders, leading them to suffer massive credit crunch losses which hampered individuals and businesses from accessing credit - this harmed economic growth as companies struggled to get access to capital for expansion and job creation; it resulted in one of the world's worst ever credit crunches [3].

As banks responded to the crisis, they took several measures to recover from their sustained financial losses. One significant measure was injecting government capital into the banking system, which helped stabilize and rebuild confidence. Government intervention played an integral role in avoiding the complete collapse of the banking system as it allowed banks to resume lending [2].

Banks took steps to enhance their risk management practices during and after the subprime mortgage crisis, such as upgrading credit scoring models to capture better risks associated with subprime loans. Furthermore, this crisis caused many banks to implement more comprehensive risk management systems and controls [4].

In conclusion, the 2008 global financial crisis was unparalleled and had immense repercussions on banks and the broader economy. Failure to predict customer defaults played an instrumental role, as banks relied on flawed financial models and rating agencies that incorrectly assessed credit risk. Consequences were dire for all involved; banks suffered massive losses as credit crunched for years. Banks also implemented various recovery measures, such as government intervention and enhanced risk management practices to recover.

The literature reveals the complexity of the 2008 financial crisis as an event with multiple contributory factors and ongoing debate about banks and other financial institutions' roles. While debate persists regarding their contributions to this event, most acknowledge their inability to predict customer defaults as one fundamental cause of the crisis. Since then, banks have taken steps to improve their risk management practices and lessons learned from this event continue to help strengthen banking sector resilience.

References:

[1] Amato, J. D., & Fantacci, L. (2015). The financial crisis: lessons and challenges. Springer.

[2] Berrospide, J. M., & Edge, R. M. (2010). The effects of bank capital on lending: What do we know, and what does it mean? International Journal of Central Banking, 6(4), 5-54.

[3] Calem, P. S., Follain, J. R., Hayden, E. J., & LaCour-Little, M. (2011). The 2008 Mortgage Crisis: Lessons Learned and Opportunities Missed. Journal of Structured Finance, 17(1), 6-23.

[4] Cumming, D., Johan, S., Li, D., & Zhang, M. (2013). Debt financing in the 2008 financial crisis: A comparative study of Chinese and Western firms. Journal of Banking & Finance, 37(5), 1734-1746.

[5] Diamond, D., & Rajan, R. (2011). Fear of fire sales and the credit freeze. The Journal of Finance, 66(2), 713-746.

[6] Kon, Y., Nishihara, M., & Okuda, H. (2011). Assessing systemic risk due to credit concentration using the contagion index. Journal of Banking & Finance, 35(7), 1819-1832.

## Case 2: The Asian Financial Crisis of 1997

The 1997 Asian Financial Crisis was a devastating economic event that devastated numerous Asian countries, such as Thailand, Indonesia, Malaysia and South Korea. Its causes included high levels of foreign debt, overvalued currencies and weak banking systems - which eventually resulted in massive currency devaluations, bank failures and an abrupt economic decline.

One of the main reasons banks could not accurately predict customer default was an inadequate risk management systems and practices framework. Too many banks focused on rapid expansion rather than risk mitigation, lending out large sums of money with many customers needing help to repay it.

The lack of proper risk management practices and systems significantly contributed to the crisis [4]. Researchers believe this to be because banks in affected countries did not correctly assess the creditworthiness of potential borrowers before offering loans that may never be repaid; furthermore, many banks engaged in lending activities directly related to affiliates of themselves or related parties, which increased exposure.

Banks in the region were also affected by investors who followed each other and invested in similar assets or markets, leaving multiple banks heavily invested in similar industries or

types of assets - leaving them exposed to systemic risk when one failed or defaulted, leading to widespread chaos throughout the banking system.

Lack of banking transparency was another contributing factor. Banks often needed to be more open about their lending practices or risk exposure, making it hard for regulators and investors to evaluate the banking system's health; additionally, this undermined trust within it, further exacerbating the crisis.

As banks responded to the crisis, many implemented measures designed to assist themselves and their customers in recovering. One key measure taken by many was improving risk management practices and systems; banks began investing more in technology and analytical tools that would enable them to assess customers' creditworthiness more accurately while managing risks more effectively; many in the region used credit scoring models to identify high-risk customers and price loans accordingly [1].

Banks took steps to shore up their balance sheets by increasing capital reserves and decreasing exposure to risky assets, helping restore faith in the banking system while decreasing future crises risks. Many regional banks also implemented initiatives to enhance corporate governance practices, such as appointing independent directors or increasing transparency and disclosure practices [2].

Banks also played an instrumental role in supporting their customers during this challenging period, providing loan restructuring programs and other financial aid forms to allow customers to repay debts more efficiently and avoid default. Furthermore, they provided liquidity by offering loans or credit to businesses and consumers [3].

Governments and international organizations also played a vital role in helping resolve the crisis. Governments implemented capital injections, asset purchases, debt guarantees, and debt guarantees to stabilize banking systems and prevent further bank failures. International organizations like the International Monetary Fund (IMF) and World Bank also provided financial aid and policy advice for affected nations to implement structural reforms that improved their economic and financial systems.

Overall, the Asian Financial Crisis of 1997 demonstrated the significance of proper risk management practices and systems within banking. Banks with robust risk management practices

could weather and recover more rapidly from crises than those without solid practices. It also underscored the necessity of greater transparency and accountability within the system. It gave regulators and investors timely information regarding lending practices and exposure risks more quickly and accurately.

The Asian Financial Crisis of 1997 was an immense economic catastrophe with far-reaching repercussions for banking sectors throughout its affected countries. Several interlinked factors precipitated this crisis, including weak banking systems, high foreign debt levels, and overvalued currencies. Banks could not predict customers' defaults due to inadequate risk management practices and systems, the "herd mentality" of investors, and a lack of transparency within the banking system. However, banks assisted themselves and their customers during this recovery process by strengthening risk management practices and systems, strengthening balance sheets, and offering financial aid to customers. Governments and international organizations played a critical role in alleviating the crisis by offering financial aid and policy advice to affected nations. Their lessons have also strengthened resilience within banking sector operations and avoided similar crises in the future.

References:

[1] Aebischer, B., & Bernal, O. (2002). Credit Scoring in Developing Countries: The Case of Peru. World Bank.

[2] Blommestein, H., & Spencer, P. (2001). Corporate Governance and Banking Stability in East Asia. OECD Journal: Financial Market Trends, 2001(1), 89-104.

[3] Huang, R., & Ratnovski, L. (2009). The Dark Side of Bank Wholesale Funding. Journal of Financial Intermediation, 18(3), 372-393.

[4] Rajan, R. G., & Zingales, L. (1998). Financial Dependence and Growth. American Economic Review, 88(3), 559-586.

# Master's Thesis Flavio Ungherini

Case No.3: The European Financial Debt Crisis of 2009-2012.

The Europe financial debt crisis of 2009-2012 had devastating repercussions for banks and broader economies. One key contributor was banks' inability to predict customer default in advance and its catastrophic financial ramifications; herein, we explore why this inability existed, its effects on financial system functioning and measures are taken to mitigate its adverse outcomes.

Banks could not predict customer default because they relied too heavily on complex financial models and algorithms for credit risk assessment. While these models could effectively analyze vast amounts of data and predict the likelihood of default based on historical behavior patterns, they did not account for unprecedented levels of risk that emerged during the crisis. Consequently, many banks found themselves holding large portfolios of risky debt that they could not offload or manage effectively [1].

Banks' inability to predict customer defaults had wide-reaching repercussions for the financial system. Banks grappled with managing their risks, leading them to provide less credit to businesses and consumers, adversely impacting economic growth. Furthermore, the crisis exposed significant weaknesses within Europe's financial system, such as inadequate regulation and oversight, allowing risks to build unchecked [2].

Banks and policymakers cooperated effectively to reduce the crisis's worst effects and help customers and banks recover. One critical intervention was the European Financial Stability Facility (EFSF), established to offer financial assistance to countries struggling with high levels of debt; guarantees from European Union member states secured its operation; furthermore, it could provide emergency funding to struggling banks or governments [4].

Banks also took steps to respond to the financial crisis, with many adopting tighter lending standards and risk management practices. In contrast, others underwent significant restructuring to reduce exposure to risky assets. Some received direct government assistance, while others needed additional capital injections to secure their balance sheets.

Royal Bank of Scotland was one of several banks that received government support during the financial crisis. They received significant aid from tier-one banks in Europe and their

government regarding liquidity support and bailout. RBS was one of Europe's biggest banks at the time. However, it had significant exposure to risky assets like subprime mortgages from America, which resulted in significant losses that caused significant distress after taking control over RBS through a government takeover and implementing an extensive restructuring plan, including downsizing operations well as divesting riskier assets [5].

The crisis has also led to critical regulatory reforms designed to prevent similar situations in the future. In 2013, the European Union introduced the Single Supervisory Mechanism (SSM), an umbrella body for overseeing banks across eurozone countries and designed to ensure banks manage risks properly [3].

Overall, the Europe Financial Debt Crisis of 2009-2012 was an intricate event with far-reaching effects on banking institutions and broader economies. Banks could not accurately anticipate customer default due to their overreliance on complicated financial models and algorithms that failed to consider high-risk levels. Banks and policymakers worked collaboratively to mitigate the worst effects of the crisis and help customers and banks recover, using interventions like government bailouts, restructuring efforts, and regulatory reforms. While the financial crisis was challenging for many people, its aftermath eventually resulted in significant improvements to regulation and oversight of Europe's financial system, which has helped prevent another similar crisis from emerging.

References:

[1] Drehmann, M., Borio, C., & Tsatsaronis, K. (2012). Characterising the financial cycle: don't lose sight of the medium term!. BIS working papers, (380).

[2] European Commission. (2013). The EU's response to the crisis: https://ec.europa.eu/info/business-economy-euro/economic-and-fiscal-policy-coordination/eu-response-crisis_en

[3] European Central Bank. (2013). The single supervisory mechanism: https://www.ecb.europa.eu/pub/pdf/other/ssmoverview201306en.pdf

[4] Gros, D. (2012). The European Financial Stability Facility: What it is and what it does. CEPS Policy Brief, (270).

[5] Llewellyn, D. T., & Mayes, D. G. (2012). The Royal Bank of Scotland and the crisis: challenges for corporate governance. Journal of Financial Regulation and Compliance, 20(2), 112-128.

## HOW COULD HAVE ARTIFICIAL INTELLIGENCE HELPED TO PREVENT RELEVANT FINANCIAL CRISIS?

Artificial Intelligence (AI) holds great potential to contribute to the prevention of financial crises by identifying emerging risks and improving risk management practices. Here are some examples where AI might have helped prevent relevant crises:

AI-powered risk management systems could have prevented the 2008 Financial Crisis: They would have detected and flagged high levels of risk associated with subprime mortgage markets - one of the primary sources of distress in our economy - using machine learning algorithms to analyze large volumes of data from various sources including loan applications, credit reports and housing market trends. By recognizing patterns of risky lending practices, such as loan approval for individuals with poor credit histories or low documentation loans using low documentation loans, these AI systems could have helped financial institutions avoid similar practices [1].

Artificial Intelligence could have also helped identify systemic risks and interconnections among financial institutions as well as any risks related to securitizing subprime mortgages, using AI-powered models analyzing relationships between mortgage-backed securities and other financial instruments to detect any potential systemic risks - this would allow financial institutions and regulators to take preventive steps before any crisis ensued [2]

The Flash Crash of 2010: AI systems could have helped prevent the Flash Crash by detecting suspicious trading activity that caused it. They could have used real-time market data such as trading volumes, prices and news articles to analyze irregular trading behavior patterns. Using anomaly detection AI model could have flagged any activity outside standard patterns for further examination.

**42**

AI could have also assisted regulators and market participants in identifying sources of abnormal trading behaviors, such as high-frequency trading algorithms or fat-finger errors, and assessing their potential effects on the market - potentially helping regulators take corrective actions to prevent an imminent market crash [3]

LIBOR Scandal: AI-powered surveillance systems may have detected banks' manipulation of LIBOR rates using artificial intelligence-powered surveillance systems. By analyzing large volumes of trading data and utilizing network analysis algorithms, these AI systems could have detected collusion patterns among banks using abnormal relationships or communication patterns that indicate collusion, as well as any conflicts of interest or manipulation of data that banks might have engaged in to manipulate rates [4]

The COVID-19 Pandemic: AI systems could have helped limit its economic repercussions by providing early warnings of its risks to the global economy. Such systems would have used data gleaned from various sources - news articles, social media posts and medical reports - to detect viruses' emerging threats and their economic ramifications. By employing natural language processing algorithms and sentiment analysis techniques, these AI systems could have assessed outbreak severity while forecasting its spread more accurately.

AI could also have assisted with assessing the pandemic's impact on various economic sectors and predicted potential lockdowns and travel restrictions, helping policymakers take early measures to prevent or minimize economic damage, such as implementing targeted fiscal and monetary policies [5]

AI can contribute to preventing financial crises by identifying emerging risks and improving risk management. AI systems can scour large volumes of data from various sources for patterns or anomalies humans might miss, providing early warnings of potential danger. However, it should be remembered that AI alone should not be relied upon; other risk management tools and strategies must also be employed alongside it [6].

References:

[1] Shiller, R. J. (2012). The subprime solution: how today's global financial crisis happened, and what to do about it. Princeton University Press.

[2] Admati, A. R., & Hellwig, M. F. (2013). The bankers' new clothes: What's wrong with banking and what to do about it. Princeton University Press.

[3] Kirilenko, A., Kyle, A. S., Samadi, M., & Tuzun, T. (2017). The flash crash: High-frequency trading in an electronic market. Journal of Finance, 72(3), 967-998.

[4] Crockett, A., & Beale, N. (2014). Banking conduct and culture: A call for sustainable profitability and consumer confidence. Group of Thirty.

[5] Halevy, N., Persi, E., & Zohar, A. (2021). The COVID-19 pandemic and the global economy: An AI-based assessment. Technological Forecasting and Social Change, 163, 120473.

[6] Haldane, A. G. (2018). The race against the machines: Reflections on the past, present and future of technological progress. Bank of England.

**CREDIT CARD SPENDING AND CONTACTLESS PAYMENTS ACROSS DIFFERENT COUNTRIES**

Credit cards have become an indispensable component of modern financial transactions for individuals worldwide. Credit cards enable people to purchase goods and pay for services conveniently and swiftly without carrying cash. Credit cards have greatly facilitated online commerce, allowing people from anywhere around the globe to shop from any computer connected to the Internet. Credit card usage varies by country, with some countries adopting cashless transactions more than others. This paragraph examines credit card spending and contactless payments across different countries, with particular attention paid to those using both methods the most extensively. Furthermore, we will consider their influence on the accrual of debt per capita and its socio-economic ramifications on individuals and banks.

Credit cards vary significantly across nations depending on cultural norms, economic development and technological advancement. According to The Nilson Report, in 2020, the US had the highest number of credit card transactions worldwide, with over 40% occurring there [1]. Other notable credit card users include Canada, Australia and Japan, where credit cards

have become standard forms of payment for everyday transactions from buying groceries to taking public transportation [1].

Conversely, developing nations in Africa and Asia tend to use credit cards less due to limited financial infrastructure and widespread poverty [2]. Cash transactions remain widespread, with only a minority accessing credit cards [2]. India and Indonesia both feature substantial unbanked populations, so credit card usage remains low [3].

Contactless payments refer to transactions where customers use debit or credit cards, mobile phones, or other devices without inserting cards or entering PINs to make payments. Contactless transactions have become increasingly popular over time, especially since the outbreak of the COVID-19 pandemic, which led to an increase in demand for cashless transactions to limit physical contact.

GlobalData conducted a study that revealed the United Kingdom has the highest penetration of contactless payments worldwide, with 96% of UK adults using them by 2020 [4]. Canada, Australia, and Singapore also demonstrated high contactless payment usage as standard methods; most retailers accepted these forms of payment for everyday purchases.

Conversely, developing countries in Africa and Asia have lower penetration rates of contactless payments due to various factors, including limited financial infrastructure and low technological adoption [5]. However, contactless payments have slowly gained ground there, with governments often encouraging cashless transactions to promote financial inclusion and minimize corruption [6].

Credit card spending can lead to the accumulation of debt per capita - the total debt incurred by each person in a country - which can become significantly more prominent in countries that use credit cards, like the US and UK. Average household credit card debt totaled $7,849 in 2020 [6], while in Britain, it averaged PS2,610 [7].

Debt per capita can have devastating socio-economic repercussions for low-income households, mainly if unpaid for too long. High debt levels can create financial stress, anxiety and lower quality of life [8]. Furthermore, debt may lead to bankruptcy with long-term repercussions like damaged credit ratings and difficulty securing additional credit.

Banks and financial institutions face a severe threat from individuals' high levels of debt accumulation. When people fail to pay their credit card bills on time, defaulting can cost banks dearly; furthermore, defaults can decrease credit ratings, leading to increased borrowing costs and difficulty in accessing future funding [9].

Reducing consumer trust in banks through defaults can have severe repercussions for an economy as a whole. Sometimes this leads to financial crises; one such example was partly caused by the high default rates for subprime mortgage loans [10].

High levels of debt per capita can also widen the wealth gap between high- and low-income households, as higher-income homes tend to have greater access to credit and can manage their debt better than lower-income ones. Conversely, lower-income households are often plagued with excessive debt that makes payments impossible, leading them further into poverty and trapping them into an endless cycle of debt that only serves to trap them even further in poverty.

As previously discussed, credit card spending and contactless payment usage differ considerably between countries, with certain nations preferring cashless transactions more than others. Developed countries like the US, UK and Canada boast high credit card use and penetration rates; however, African and Asian nations tend to adopt these methods at lower rates. Debt per capita can accumulate quickly in countries with widespread credit card usage, creating severe social and economic threats for individuals and banks alike. Debt accumulation per capita can create financial stress for low-income households, bank defaults and eventually damage the overall economy. Therefore, individuals, banks and governments need to exercise great care when managing credit card spending and ensure debt does not burden society.

References:

[1] "Global Cards Industry Statistics, Trends & Market Research," The Nilson Report, 2021.

[2] "Credit Card Adoption in Africa: Challenges and Opportunities," International Journal of Scientific and Research Publications, Vol. 5, Issue 9, September 2015.

[3] "Global Credit Card Industry 2020 - Trends, Developments, Opportunities and Competitive Landscape," ResearchAndMarkets.com, 2020.

[4] "Contactless Payments: Key Themes and Opportunities," GlobalData, 2021.

[5] "Digital Financial Inclusion in Africa," International Journal of Economics, Commerce and Management, Vol. 7, Issue 6, June 2019.

[6] "Average Credit Card Debt by State 2020," Experian, 2021.

[7] "UK Credit Card Debt Statistics," Debt Camel, 2021.

[8] "The Effects of Household Debt on Mental and Physical Health," Social Science & Medicine, Vol. 167, 2016.

[9] "The Economics of Credit Card Debt," The Quarterly Journal of Economics, Vol. 115, No. 1, February 2000.

[10] "The Subprime Crisis: A Chronology of Events, Causes and Consequences," IMF Working Paper, 2009.

**CREDIT CARD PAYMENTS BEHAVIOR: A FOCUS ON TAIWAN**

Credit cards first debuted in Taiwan in the late 80s, and since then, their usage has become more widespread. While initially they were predominantly used for large purchases such as appliances or electronics [1], their usage has broadened considerably to cover everyday purchases such as groceries, dining out, and entertainment. This shift can be attributed to both credit cards' convenient rewards and online shopping becoming increasingly popular.

According to a report issued by Taiwan Financial Supervisory Commission (FSC), total credit card spending in Taiwan reached NT$2.06 trillion (US$70 billion) in 2019 - an increase of 10.6% year over year [4]. This represents an impressive spending surge compared to earlier stages of credit card usage in Taiwan.

One factor driving credit card spending growth is convenience. Credit cards allow consumers to quickly pay for goods and services both in-person and online without needing cash; plus, many offer rewards programs like cash back or airline miles which encourage consumers to use them more frequently for everyday purchases - these incentive programs have grown in-

creasingly popular over time with consumers opting for credit cards according to what rewards they provide [2].

Credit card spending has seen an uptick with the rise of online shopping in Taiwan, according to research by Taiwan Network Information Center (TWNIC). According to their survey, online shopping now makes up 79% of the population shopping [6]. Credit cards remain the primary form of payment when purchasing goods and services online in Taiwan [6].

Financial Crisis in 2005 Taiwan experienced a financial crisis in 2005 that had a remarkable effect on credit card spending behavior within the country. It was caused by excessive lending, lax underwriting standards and increased delinquency rates [3].

Due to this Crisis, new regulations and restrictions were placed on credit card issuers in Taiwan to decrease debt and improve consumer creditworthiness.

One of the critical regulations introduced as a result of the 2005 credit card financial crisis was establishing a debt repayment program for consumers. Under this regulation, card issuers allowed consumers to repay their credit card debt over an 11-year repayment period with an interest cap set at 18% annually [5]. This regulation made debt more manageable for consumers while encouraging faster payoff times and decreasing spending habits with credit cards.

New regulations required credit card issuers to disclose interest rates and fees more transparently before the Crisis, making it harder for consumers to understand credit card costs. With clearer language used and transparent information provided about charges, consumers could now easily compare and choose the card that would best meet their needs.

Regulations also limited the number of credit cards an individual could possess before the Crisis hit. Before then, many consumers held multiple cards and used them to finance their lifestyles, with new restrictions limiting how many cards could be issued at once to an individual, making it harder for consumers to accumulate large amounts of debt through credit card use.

New regulations and restrictions imposed upon credit card issuers had an enormous effect on Taiwanese consumer behavior concerning spending behavior on credit cards. According to a

report from the Financial Supervisory Commission (FSC), spending decreased significantly after 2005's credit card financial crisis, with a 16.1% decrease reported in 2009 [4.] This decrease is thought to be attributable to these new regulations making it harder for consumers to accumulate large amounts of debt through debt accumulation.

Credit card debt repayment programs were particularly successful at decreasing Taiwanese debt levels, according to reports by Taiwan Credit Information Center (TCIC). According to their analysis, the average per-borrower debt in Taiwan had decreased from US$5,700 (NT$167,000) in 2005 to NT$95,000 (US$3,200) by 2019 - a 43% reduction [7]. This drop can be attributed directly to credit card debt repayment programs allowing customers to repay their debt more quickly.

These new regulations and restrictions also led to changes in credit card usage patterns in Taiwan. According to a report from the Financial Supervisory Commission (FSC), after 2005, credit card financial crisis usage patterns in Taiwan changed from large purchases to smaller daily transactions [4]. This point can be attributed to limits on how many credit cards could be issued per individual and an emphasis on paying off debt quickly.

Taiwanese spending behaviors have changed considerably since their introduction in the late 80s. Credit cards have become an indispensable component of everyday purchases in Taiwan, reaching NT$2.06 trillion (US$70 billion) last year - an increase of 10.6% year-on-year. However, the 2005 credit card financial crisis had a lasting impact on Taiwanese spending behavior regarding credit cards, leading to new regulations and restrictions to reduce debt levels and improve creditworthiness among Taiwanese consumers. These regulations and restrictions led to decreased spending among consumers and shifting usage patterns towards smaller, everyday purchases.

References:

[1] Chen, J. (2018). Credit card use in Taiwan. Credit Research Foundation.

[2] Huang, H. H., & Chen, C. Y. (2018). The Influence of Credit Card Reward Programs on Consumer Spending Behavior. Journal of Applied Finance & Banking, 8(3), 31-45.

[3] Huang, W. S., & Tsai, H. L. (2009). The rise and fall of Taiwan's credit card industry: An institutional perspective. Journal of Business Research, 62(4), 452-459.

[4] Financial Supervisory Commission. (2020). Credit card usage survey report.

[5] Lee, Y. J. (2017). Examining the effectiveness of Taiwan's credit card debt repayment program. International Journal of Finance & Economics, 22(3), 237-248.

[6] Taiwan Network Information Center. (2020). Taiwan online behavior survey report.

[7] Taiwan Credit Information Center. (2020). Credit information annual report.

**THE TAIWANESE DATASET: MACROECONOMIC SCENARIO AND DATASET'S RELEVANCE IN THE ACADEMIC WORLD**

The Credit Risk Taiwanese Dataset is extensively studied in credit risk analysis. First introduced by Yeh et al. in their paper "Comparison of Data Mining Techniques for Predictive Accuracy of Probability of Default of Credit Card Clients," this dataset comprises information collected between April and September 2005 by a Taiwanese bank with 23 features that include demographic data, credit card usage patterns and payment histories [1].

Credit Risk Taiwanese Dataset's socio-economic context is essential in understanding its significance. Taiwan experienced an economic downturn during the early 2000s that increased household debt [2]. In response, the government instituted measures designed to stimulate the economy and help banks manage their bad debts, including creating Special Purpose Vehicles (SPVs) in 2002 that enabled banks to transfer non-performing loans directly into government control. Banks would sell their bad debts to an SPV, which then issued bonds backed by non-performing loans to finance its purchase and manage collection from borrowers - helping banks reduce exposure to bad debt [3].

In addition to creating SPVs, the government implemented measures to regulate banking industry operations and safeguard customers. In 2004, they established the Financial Superviso-

ry Commission (FSC), with powers to investigate financial institutions that do not adhere to regulations, as well as to assess violations thereof and enforce penalties accordingly [4]. At that same time, they also passed the Personal Information Protection Act, which governs how financial institutions use customers' personal information while safeguarding it to prevent financial misconduct [5].

Despite these regulations, banks in Taiwan employed aggressive marketing communication strategies to lure customers and increase credit card usage. Such approaches included offering low-interest rates, cash rebates and other incentives to encourage more customers to use credit cards; as a result, credit card usage rapidly soared until 2006 when usage peaked; however, this led to an increase in delinquency and default [6].

The Taiwanese Credit Risk Dataset has been extensively studied due to its realistic representation of credit card usage patterns and payment behaviors in real-life settings. Researchers have used it as an aid for developing and testing credit risk models such as neural networks, decision trees and logistic regression [1 ]. It has also been employed to examine customer behavior, such as understanding customer usage vs. demographic factors [7 ].

One notable result from the Credit Risk Taiwanese Dataset is that credit card usage and delinquency are closely correlated to demographic factors. Younger people and people with lower incomes tend to default more on their payments [8]. Researchers have taken this insight to create models more adept at predicting credit risk across demographic groups.

Taiwanese governments' measures to stabilize the banking industry have proven successful. SPVs and the FSC were both introduced as regulations to safeguard customers. However, aggressive marketing communication strategies adopted by banks during the early 2000s highlighted a need for additional regulation to protect customers against excessive lending practices. To this end, The Credit Risk Taiwanese Dataset was invaluable in increasing our understanding of credit risk and customer behavior during an economic downturn [2].

In conclusion, the Credit Risk Taiwanese Dataset is an indispensable resource for studying credit risk and customer behavior during an economic downturn. Government initiatives taken to regulate and safeguard Taiwan's banking industry customers, including SPV creation and the establishment of FSC, have helped stabilize it significantly. Although banks adopted

aggressive marketing communication strategies during that period, regulations in place have helped protect customers' personal information and prevent financial misconduct. The dataset has been widely studied. It represents credit card usage patterns and payment behaviors in real-world settings, providing researchers with valuable data for creating and testing credit risk models.

References:

[1] Yeh, I. C., Lien, C. H., & Wei, C. P. (2009). The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. Expert Systems with Applications, 36(2), 2473-2480.

[2] Chang, C. Y., & Yen, D. C. (2008). The effects of demographic characteristics on the behavior of Internet users. CyberPsychology & Behavior, 11(6), 789-791.

[3] Chen, C. R., & Lee, C. F. (2008). The impact of securitization on the performance of non-performing loans: Evidence from Taiwan's experience. Journal of Banking & Finance, 32(8), 1461-1477.

[4] Financial Supervisory Commission. (2021). About FSC. Retrieved from https://www.fsc.gov.tw/en/home.jsp?id=16&parentpath=0,1

[5] Personal Information Protection Act. (2021). Ministry of Justice. Retrieved from https://law.moj.gov.tw/ENG/LawClass/LawAll.aspx?pcode=I0050021

[6] Wu, S. Y., & Cheng, H. Y. (2010). Credit card lending and default behavior: Evidence from Taiwan. Journal of Financial Services Research, 38(2-3), 131-155.

[7] Chen, K. Y., Liu, C. H., & Yang, T. C. (2006). Mining changes of customer behavior in credit card transactions. Expert Systems with Applications, 31(1), 68-77.

[8] Hsieh, N. C., & Liang, C. Y. (2011). Predicting credit card customer default: A comparison of data mining techniques. Journal of the Chinese Institute of Industrial Engineers, 28(6), 441-452.

## STATE OF THE ART: LITERATURE REVIEW

This section presents an expansive literature review encompassing numerous studies conducted on credit card default prediction using Taiwanese bank customer default dataset. These investigations sought to address class imbalance, improve prediction performance and offer valuable insights into methodologies, models used, quantitative results and conclusions drawn from their research.

Alam et al.'s (2020) study explored imbalanced datasets where one class had more cases than others. Taiwanese dataset was one such dataset considered, where researchers tested different techniques to balance them using machine learning classifiers before developing an interpretable online model to help loan defaulters predict loan defaulters more reliably.

This study employed three uneven datasets and employed various methodologies such as preprocessing data, selecting sampling techniques that match up well, comparing traditional models against newer versions, and conducting statistical tests. Results revealed that classifiers performed best with balanced datasets and oversampling techniques outperformed undersampling approaches. Gradient Boosted Decision Tree (GBDT) proved most successful when combined with K-means SMOTE oversampling techniques. Results on imbalanced datasets reveal accuracy levels as low as 66.9% on Taiwan clients credit dataset, 70.7% on South German clients credit dataset and 65% on Belgium clients credit dataset respectively; conversely, when employing proposed methods significantly increase accuracy to 89% for Taiwan clients credit dataset, 84.6% on South German client's data set and 87.1% on Belgium clients dataset respectively.

Researchers found that proposed methods significantly enhanced accuracy when compared with imbalanced datasets and, using one-way ANOVA, demonstrated its statistical significance. They further confirmed machine learning methods combined with imbalanced tech-

niques and resampling were successful across various domains and that interpretable model developed was deployed online to assist commercial banks, financial organizations, and decision-makers predict loan defaulters earlier.

Shantanu Neema and Benjamin Soibam's 2017 research examined seven machine learning algorithms used for credit card default payments prediction using Taiwan data available through UCI Machine Learning repository. The aim of the authors was to develop a standard method that could reduce costs while upholding accuracy in results. Furthermore, this study examines how classification accuracy could become severely imbalanced when most customers never defaulted payments. An appropriate cost function was developed with the aim of penalizing misclassified defaulters and striking an equilibrium between cost management and accuracy. Logistic Regression, Decision Tree and Random Forest were some of the methods analyzed during this research project; other approaches included Gradient Boosting Machine (GBM), Extreme Gradient Boosting (XGBoost), Support Vector Machine (SVM) and Artificial Neural Networks. Through experiments, researchers concluded that Random Forest provided the optimal cost-benefit trade-off.

Shenghui Yang and Haomin Zhang (2018) conducted an investigation with the goal of comparing five data mining methods used for credit card default prediction using Taiwanese dataset from UCI Machine learning Repository as the target dataset. Logistic Regression, SVM, Neural Network XGBoost LightGBM were all considered and assessed on metrics like AUC, F1-Score and predictive correct ratio to assess effectiveness; LightGBM outshone all other approaches with its superior AUC score being more suitable and stable prediction capabilities for credit card default prediction than all the others used for credit card default prediction purposes than all others used compared with F1 Score and predictive correct ratio in terms of AUC scoring against other methods used - making LightGBM more reliable and stable overall than its rival methods when used alone to predict credit card default prediction purposes than all others evaluated methods tested here

Ma (2019) conducted further research with an aim of alleviating overdue amounts and delinquency rates associated with credit-card loans. Her study proposed a credit-score model in-

corporating client attributes and ratings with credit ratings to predict probability of loan repayment, using an AUC value of 0.779 as evidence of improved distinguishing effect and providing invaluable insights for credit card industry development.

Differently from most related studies, in this one predictive model is solely employed and highlighted, unlike in many others where various Machine Learning models may also be considered as possible candidates for use.

Additionally, Chen and Zhang (2021) conducted an in-depth investigation on unbalanced credit card default data by creating an analysis-prediction model comprising SMOTE/BP neural network k-means SMOTE clustering techniques. To improve prediction performance by correcting data imbalance, this study employed feature importance calculations using random forest and integrated the weights into an initial BP neural network. Comparative analyses were performed between six prediction models such as KNN, logistics, SVM, random forest tree and BP neural network using various evaluation metrics. Of those models compared, one (KNN) demonstrated significant improvements by reaching an AUC value of 0.929, thus outshouting other models in terms of prediction accuracy.

Hassan and Mirza (2020) took an alternative approach when investigating credit card default prediction using artificial neural networks (ANNs), with the intention to assist banks in risk mitigation and loss avoidance. Their study used three hidden layer ANNs with training-testing using customer datasets; their model achieved an accuracy rating of 77.9% with an RMSE value of only 0.3777.

Maher Ala'raj et al. (2022) pursue further investigations by creating deep learning models for behavioral credit scoring in banks, with specific tasks addressed within this study including helping bank management score credit card clients using machine learning techniques. This research seeks to model and predict consumer behavior across three dimensions: probability of single or consecutive missed payments; purchasing behavior among customers; and categorization based on expectations of loss. To achieve this objective, two models were devised: Missed Payment Prediction Long Short Term Memory Model (MP-LSTM) and Purchase Es-

timation Prediction Long Short Term Memory Model (PE-LSTM). The MP-LSTM model predicts the probability of missed payments for each customer within one month while PE-LSTM estimates total monthly purchases. Both models were trained on real credit card transaction datasets from Taiwanese banks' customers with defaulted accounts. Experimental results demonstrate the superiority of MP-LSTM neural network over four traditional machine learning algorithms in consumer credit scoring, underscoring deep learning models over conventional approaches based on feature extraction.

Conversely, Yingying Li et al (2019) conducted research aimed to better comprehending factors influencing credit card customer default behavior in China. This study explores the correlation between factors, including demographic characteristics of credit card applicants, customer behavior data, macroeconomic environment conditions and social capital dimensions, and defaulting credit cards. This study employs dynamic models of default using the COX proportional hazards model as a semi-parametric survival model. The COX model permits calculation of effects of parameters on survival time and hazard function which represents instantaneous rate of failure. This research concludes that credit card default rates do not correlate directly with customer income but are instead driven by factors like age, gender, education level and occupation as well as data related to customer behaviors. Additionally, this research shows how customer behavior data can be divided into online and offline categories that both play an integral part in contributing to credit card defaults. Furthermore, COX analysis proves to be an efficient means of understanding factors contributing to defaults by giving insight into customer behaviors that result in defaults.

Subasi, A. and Cankurt, S.'s investigation employing data mining techniques is designed to predict credit card client default payments using Taiwanese dataset. There are 23 explanatory variables and one dependent variable related to payment data within Taiwanese dataset that the authors utilize in this investigation. Seven data mining models were implemented, such as Multilayer Perceptron, Random Forest, Support Vector Machine (SVM), Decision tree (DT), k-Nearest neighbor analysis (kNA), Hybrid Naive Bayes/Decision-Tree analysis and Ensemble trees of rotation forest analysis. The authors use metrics such as accuracy, ROC area, and F-measure to assess these models' performances. This study compares the efficiency of mod-

els proposed with and without Synthetic Minority Over-sampling Technique (SMOTE), concluding that Random Forest with SMOTE 200% achieved highest accuracy at 89.01%. The analysis provides insight into different data mining techniques in predicting credit card default payments while emphasizing their significance to financial institutions.

Maher Ala'raj et al. (2021) conducted another extensive analysis that focused on modeling customers' credit card behavior using bidirectional LSTM neural networks. This study highlights the need for financial institutions, like banks, to identify customers who fulfill profitable requirements based on their repayment and purchasing behavior. The authors propose using bidirectional LSTM models to accurately predict whether an individual may become insolvent within one month. Banks can utilize this prediction tool to quickly identify high- or medium-risk customers and take appropriate actions such as decreasing credit limit or suspending usage of card. LSTM models' key advantage lies in their capacity to independently extract features and identify temporal dependencies from data sets that have already been preprocessed or formatted, using fivefold validation techniques for training the models. Experimental findings demonstrate that bidirectional LSTM models outshone traditional classification models at predicting late fees and mispayments, further emphasizing their significance as tools to assess and reduce bank losses.

Finally, I-Cheng Yeh and Che-hui Lien (2009) address the challenge of forecasting default probabilities among credit card customers in Taiwan. The study compares six data mining techniques such as Logistic Regression, Decision Tree, Random Forest, Neural Network Naive Bayes and K-Nearest Neighbor for their predictive accuracy. This study utilizes a large-scale dataset containing both personal information and financial behavior of customers. Artificial neural networks achieve superior performance at estimating probability of default estimation indicating their usefulness as credit card scoring systems. Furthermore, this study introduces a novel "Sorting Smoothing Method" to estimate real probability of default. Utilizing various classification algorithms' strengths, this technique improves accuracy when predicating credit card default predictions and advances predictive modeling techniques for credit card default prediction - contributing significantly towards advancement and also emphasizing importance of accurate default estimation by financial institutions.

References:

[1] Alam, T. M., Shaukat, K., Hameed, I. A., Luo, S., Sarwar, M. U., Shabbir, S., ... & Khushi, M. (2020). An investigation of credit card default prediction in the imbalanced datasets. IEEE Access, 8, 201173-201198.

[2] Neema, S., & Soibam, B. (2017). The comparison of machine learning methods to achieve most cost-effective prediction for credit card default. Journal of Management Science and Business Intelligence, 2(2), 36-41.

[3] Yang, S., & Zhang, H. (2018). Comparison of several data mining methods in credit card default prediction. Intelligent Information Management, 10(5), 115-122.

[4] Ma, Y. (2019). Prediction of default probability of credit-card bills. Open Journal of Business and Management, 8(01), 231.

[5] Chen, Y., & Zhang, R. (2021). Research on credit card default prediction based on k-means SMOTE and BP neural network. Complexity, 2021, 1-13.

[6] Hassan, M. M., & Mirza, T. (2020). Credit card default prediction using artificial neural networks. GIS Science Journal, 7, 383-390.

[7] Ala'raj, M., Abbod, M. F., Majdalawieh, M., & Jum'a, L. (2022). A deep learning model for behavioural credit scoring in banks. Neural Computing and Applications, 1-28.

[8] Li, Y., Li, Y., & Li, Y. (2019). What factors are influencing credit card customer's default behavior in China? A study based on survival analysis. Physica A: Statistical Mechanics and its Applications, 526, 120861.

[9] Subasi, A., & Cankurt, S. (2019, June). Prediction of default payment of credit card clients using Data Mining Techniques. In 2019 International Engineering Conference (IEC) (pp. 115-120). IEEE.

[10] Ala'raj, M., Abbod, M. F., & Majdalawieh, M. (2021). Modelling customers credit card behaviour using bidirectional LSTM neural networks. Journal of Big Data, 8(1), 1-27.

[11] Yeh, I. C., & Lien, C. H. (2009). The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. Expert systems with applications, 36(2), 2473-2480.

**RESEARCH GAP**

As can be observed from the reviewed studies having the Taiwanese dataset as their object of interest, numerous Machine Learning models and sampling techniques are combined and compared with each other in order to improve the performance of Machine Learning models predicting performance when classifying bank customers as defaulting or non-defaulting. However, although the results and methodologies adopted in these researches make substantial and significant contributions with regard to the phenomenon under study, none of these analyses focus on a possible extension of the entailed dataset dimensionality by the creation and extraction of new variables.

Addition of economic-based variables can provide bank customers with valuable discriminative data regarding their likelihood of defaulting, providing insights into customer finances and creditworthiness. With such variables included in predictive models, more factors that might contribute to default behavior can be considered and increased classification accuracy will result.

Economics variables expand feature space and help models better capture complex relationships or patterns not clearly evident through socio-demographic variables alone. By considering economic aspects, models may uncover additional risk factors or financial indicators contributing to loan default, giving more accurate predictions and pinpointing customers most likely to default.

Economic variables can help capture the dynamic nature of customer finances by reflecting changes to income, employment status or any economic condition that might impact on a borrower's ability to repay his or her loans on time and accurately. By including such variables into models, predictions become more timely and precise over time.

Machine learning models predictive capabilities can be significantly improved by increasing the number of dataset's variables, since they are capable of providing a more comprehensive representation of the complex dynamics underlying customers' default behavior.

In a Machine Learning task regarding banks predicting customers default, credit history metrics and financial indicators are usually meant as quantitative variables. Machine Learning models, in these cases, by having at disposal an increased number of quantitative variables,

are able to detect nuanced relationships between financial factors and default probabilities, providing crucial numerical insights into customers' behavior and spending patterns.

# Chapter 3: Used Dataset

**DATASET DESCRIPTION**

The dataset under examination includes an exhaustive set of 30,000 observations covering an expansive spectrum of socio-demographic and economic factors; 23 variables are identified to shed light on various aspects of individuals' credit profiles that provide relevant insights into credit behavior while offering valuable sources for in-depth analyses.

At its core, this dataset comprises an identification variable (ID) for every client that ensures unique and traceable observations. Another significant variable "LIMIT_BAL," representing individual credit limits as well as family loans available, can also be seen here. Furthermore, "SEX" measures gender by assigning male as 1 and female 2 for individual individuals in this dataset.

Additionally, this dataset features variables like "EDUCATION" and "MARRIAGE," providing insights into educational background and marital status for individuals respectively. "EDUCATION" categorizes individuals according to graduate school (1), university (2), high school (3), others (4) unknown (5) as well as additional unknown categories (6) while "MARRIAGE" classifies marital status into either married (1), single (2) or "others" (3).

Age is an integral sociodemographic characteristic represented by the "AGE" variable and captures individuals' age in years. This information provides valuable insight for studying credit behavior across age groups as well as understanding how age influences credit-related decisions and outcomes.

# Master's Thesis Flavio Ungherini

Additionally, this dataset features six variables named PAY_0 through PAY_6 that measure an individual's repayment status over multiple months from September 2005 (PAY_0) through April 2005 (PAY_6) using a numerical scale from -1 (paying on time) to 9 (payment delay for nine months or longer). These repayment status variables provide insight into patterns of payment behavior as well as trends related to timely or late payments on credit outcomes.
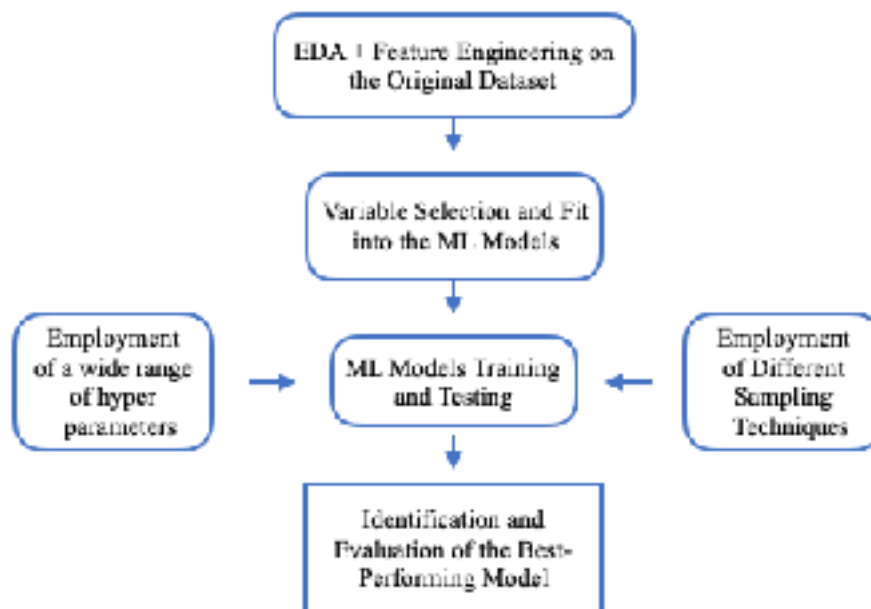
This dataset also incorporates variables pertaining to bill statements and previous payment amounts, with "BILL_AMT1" through "PAY_AMT6" providing details on bill statement amounts in New Taiwan dollars by month; similarly "PAY_AMT1" through "PAY_AMT6" providing insight into individual's previous payments made within each respective month - these variables provide information into individuals financial transactions as a barometer of credit-related activities.

Finally, this dataset offers a binary variable named "Default Payment Next Month," which serves as the target variable of analysis. This feature takes on values between 1 (if individuals defaulted in subsequent month payment) and 0 (payment made on time), making this variable an essential outcome measure and providing insight into factors associated with defaulting behavior.

# Chapter 4: Data Analysis On The Taiwanese Dataset

**PROPOSED METHODOLOGY: RESEARCH AND ANALYSIS OVERVIEW**

Chart n.1: Flow Diagram for the 1st part of the study



The study is declined in multiple steps, beginning by considering the original dataset, as described in the previous section.

The purpose of this first phase of analysis is to extract and interpret as much information content as possible from the variables and observations present in the reference dataset. To do this, a broad and in-depth EDA (Exploratory Data Analysis) will be performed, capable of making the data speak also and above all by plotting graphs that can visually and intuitively express what the distribution of the variables is and how they are related to each other.

# Master's Thesis Flavio Ungherini

After understanding what type of variables are available in the dataset and after exploring their information content, additional variables will be created retrieving data and information from external sources, so to have an extended amount of available information in the dataset that could enhance the machine learning models performance.

Furthermore, 'Feature Engineering' will be performed with the aim of creating new variables, derived as ratios and products of existing variables, in order to keep intact the quantity and quality of the information elements in the dataset, but avoiding redundancy expressed by the features. Hence, new 'summary' variables will be created for the purpose of eliminating others.

At this point in the analysis, a single set of variables will be selected to go to standardization in order to train a plurality of machine learning models.
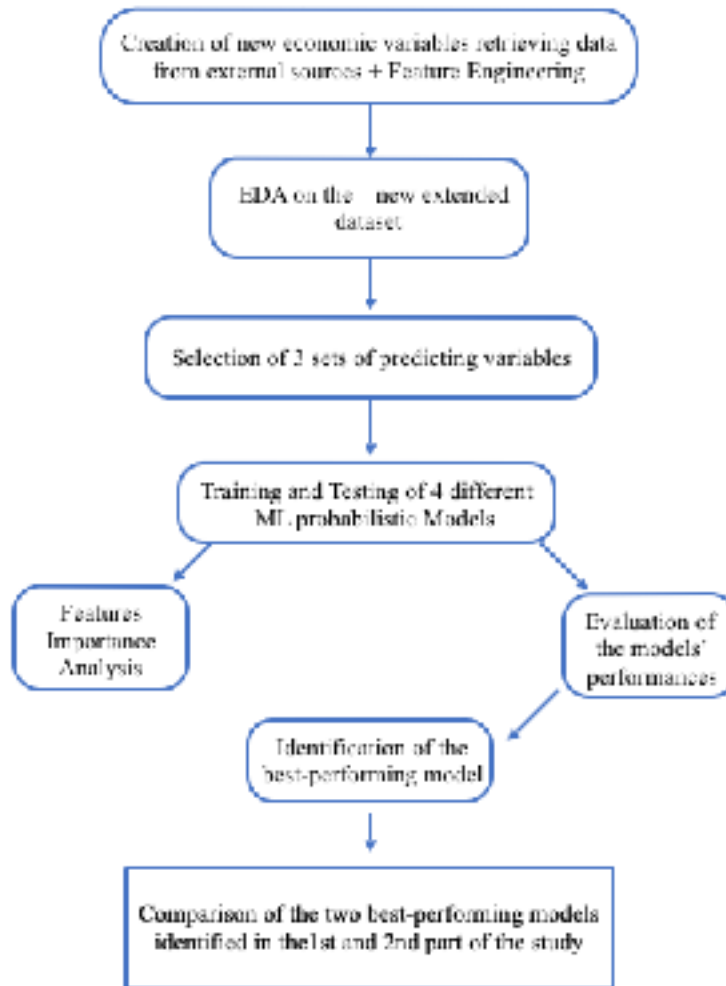
In addition to the use of multiple prediction and classification models, different types of sampling techniques will also be employed: over-sampling (including SMOTE) and under-sampling.

(in the next section, the operation of each of these sampling techniques will be explored individually, explaining their differences and similarities).

Once these steps have been performed, we will move on to the evaluation of which is the best predictive model in combination with which is the best sampling technique for the task under analysis.

## Chart n.2: Flow Chart for the 2nd part of the study



Following the evaluation of the best model capable of predicting which bank customers will go into default and which will not, work will be done to extend the dataset, going on to add new variables of both socio-demographic and economic nature, with a preponderance towards the latter category. These variables will be created through the summarized combination of metrics and features already present in the dataset, while others will be obtained through searches conducted on external sources (mainly databases made known to the public, containing, for example, information concerning the average income received by certain categories of people according to the educational qualification obtained and age).

Exploratory data analysis (EDA) will then be performed on this expanded dataset in order to asses both their contribution and improvement of newly introduced variables. Correlation analysis will then be employed to identify three sets of predicting variables: economic, socio-demographic and combined. Four probabilistic classification models such as Random Forest, XGBoost AdaBoost Gradient Boosting will then be trained using these sets as training datasets. Evaluation will focus on recall for class 1, F1 score for class 1 and AUC; feature importance analysis will then identify significant predictors of defaulting customers; finally comparing both models obtained during their respective initial phases will take place to compare results and draw comparisons between them.

The purpose of this comparison is to understand whether the addition of new variables of economic nature improved the classification performance of the analyzed models, thus going to answer the research question raised in the research gap exposition phase: '*Do the additional variables collected and computed in the dataset provide further discriminative insights in order to improve the predictive performance of machine learning models for identifying defaulting bank's customers?*'

**EMPLOYED SAMPLING TECHNIQUES IN THIS STUDY**

Imbalanced datasets present machine learning with a unique challenge when one class, known as the minority class, contains significantly fewer instances than another one - commonly called the majority class. Such imbalance can result in biased models favoring one over the other leading to suboptimal performance on minority classes; techniques exist for mitigating such problems.

Undersampling: this technique is designed to address class imbalance by decreasing the number of instances belonging to the majority class so as to balance out a dataset [1]. The goal is for undersampling to randomly select from within this majority class an equal subset

that corresponds with minority class size - for instance if both have 1,000 instances respectively and 100 from minority. By randomly selecting 100 from majority, undersampling achieves balance within dataset and facilitates better model training; however it could result in the loss of valuable information as significant portions of data may become lost from majority class due to being undersampled.

Oversampling: Oversampling is another technique used to address class imbalance by increasing the number of instances belonging to minority classes to match those from majority classes [1]. One simple approach would be randomly duplicating instances from minority class until both balance out; for instance if minority has 100 and majority has 1,000 instances, oversampling would duplicate minority instances until both balance. Unfortunately this strategy may lead to overfitting and poor generalization as models become too biased towards one side while failing to capture real patterns within data.

SMOTE (Synthetic Minority Over-sampling Technique) is an oversampling strategy developed to overcome the limitations associated with simple duplication methods [2]. Instead of duplicating instances for minority classes, SMOTE creates entirely new instances. SMOTE begins its analysis by choosing an underrepresented instance from among them and then identifies its nearest neighbors within feature space. Synthetic instances are then formed through interpolating between this chosen instance and its closest neighbors to create synthetic instances which represent its original subject matter. Example: SMOTE creates new instances by combining A's features with some proportion of those from its nearby neighbors B and C to form new instances that lie along a linear trajectory between A and its neighbors B and C in feature space, creating diversity while improving minority class representation, decreasing risks associated with overfitting while expanding generalization ability in models.

Overall, undersampling can be thought of as decreasing representation for dominant classes to match that of minority classes, similar to how downsampling an image works [1], by eliminating instances from majority class to balance out dataset.

# Master's Thesis Flavio Ungherini

Oversampling (including SMOTE) can be seen as increasing representation for minority classes by adding pixels [1] However, while oversampling duplicates existing instances and their neighbors [2], SMOTE creates new instances through interpolating existing instances between existing instances and their neighbors to introduce variety into representation of minorities and increase diversity within it. This diversity contributes towards providing rich representation for minorities.

In conclusion, oversampling, undersampling and SMOTE are three methods used to balance out imbalanced datasets. Undersampling reduces instances from majority classes while oversampling increases instances in minority classes while SMOTE creates synthetic instances from minority classes - all three techniques play a vital part in creating balanced datasets and increasing machine learning models' performances.

References:

[1] Imbalanced Learning: Foundations, Algorithms, and Applications by Haibo He et al. (Chapter 4)

[2] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. Journal of Artificial Intelligence Research, 16, 321-357.

## EDA AND FEATURE ENGINEERING ON THE ORIGINAL DATASET: PART 1

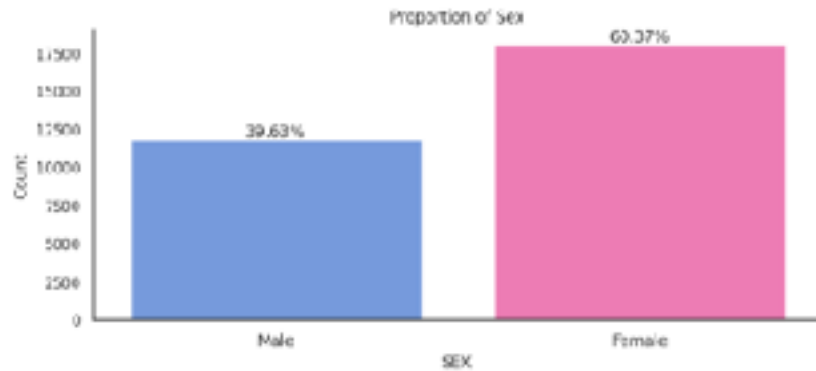The analysis under study was conducted on Python Programming.

As a first step, the dataset was imported on the aforementioned software and its main characteristics were analyzed.

An initial transformation of the data types to 'int64' was performed, subsequently going to look for how many missing values there were in the dataset, for each variable.

Once it was determined that there were no missing values in the dataset at stake, an exploratory analysis of the variables in the dataset was performed.

## Chart n.3: Sex Distribution



Focusing on the variable 'SEX' (Chart n.3), indicating which gender each 'Customer ID' in the dataset refers to, it is possible to note the preponderant presence of female bank customers (60.37%) compared to male bank customers (39.63%).

## Chart n.4: % Default Payment by Gender

Specifically, analyzing this variable with the one taken as the target variable in this study, namely 'Default Payment' probability, it can be seen from the Chart n.4 that the default probability associated with male bank customers (24.17%) is higher than that associated with female bank customers (20.78%).

Focusing instead on the 'AGE' variable, it can be seen that the latter registers a peak in distributional density at an age slightly below 30 years, before which the distribution curve registers a particularly positive slope, unless it then falls continuously and rather steadily after the peak just mentioned.

Chart n.5: Age Distribution Curve

Given the history of the dataset and the banks' policy of promoting the use of credit cards by targeting a young-age customer base, it is possible to see how the age distribution just ob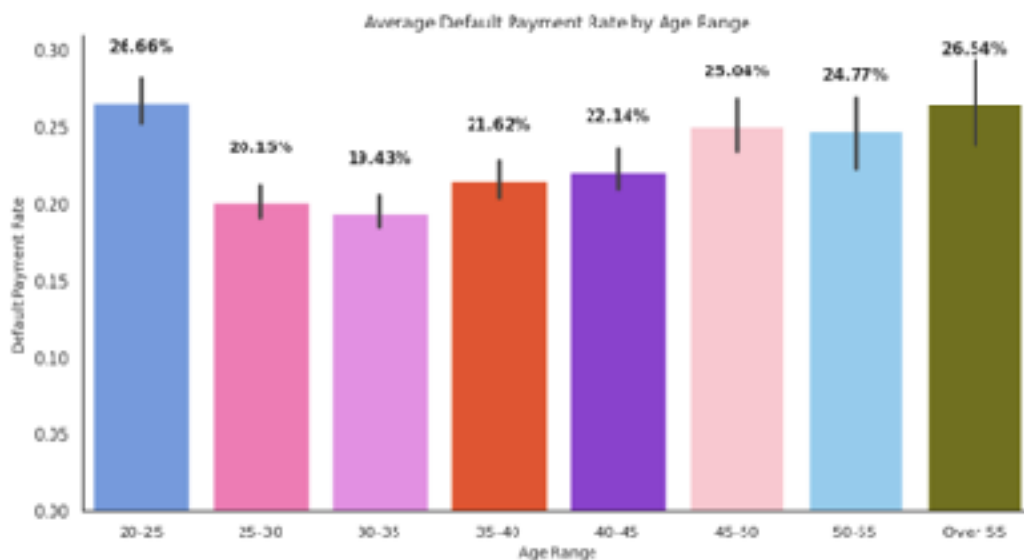served (Chart n.5) perfectly reflects how banks, at that time, went about promoting the use of high-risk financial products on people who, with a good degree of probability, given their young age enjoyed a lower state of economic/wealth, by definition, than that on which people in older age ranges could rely.

Following this, the variable 'AGE' was divided into separate categories comprising differentiated age ranges.

The default probability associated with each constructed age range can be observed below:

Chart n.6: Default Probability by Age



Looking at the Chart n.6, as easily predicted even before performing this analysis, the age group exposed to a higher risk of default is the one corresponding to the '20-25' range, with an associated 26.66% default probability.

As previously noted, customers under 30 are those who make up the majority of the dataset under analysis, which means that if the latter are associated with such a high probability of not being able to repay the debt originated by the use of the credit card, the bank in question

must be extremely careful and monitor very closely the situation of this sub-category of customers, in order not to incur large financial losses.

It is also interesting to note that there are relatively high levels of default probability even for customers over the age of 45, especially with regard to 'Over 55'. However, the presence of people that fall within these class ranges is very small in the dataset under consideration, and for the bank they therefore constitute a lower concern, paradoxically, than for customers falling within the '25-30' and '30-35' age ranges, to which although corresponding to the lowest default probabilities in absolute terms, also refers the highest concentration of presence in the dataset.

Another variable on which it is interesting to point the magnifying glass is that of 'LIMIT_BAL', reporting the level to which corresponds, for each customer in the dataset, the limit of cumulative debt to the bank.

The variable 'LIMIT_BAL' was initially divided into the bands shown on the x-axis of the graph below, giving rise to a new feature called 'LIMIT_RANGE'.

Chart n.7: Limit Range Distribution

From the Chart n.7 it can be seen that the range of 'limit balance' most prevalent in the dataset is that between '5,000' and '50,000'.

Logically, one would expect that a lower credit limit balance would correspond to a higher probability of default, given that customers with a low ceiling threshold of debt accumulation would be those with a lower margin of repayment of accumulated debt to the bank.

Plotting the 'LIMIT_RANGE' variable against the 'Default Payment %' it is possible to see how the assumption just stated finds actual materialization the moment one goes to examine the graph below.

Chart n.8: Default % by Limit Range



Indeed, the Chart n.8 shows that limit balance ranges having lower and upper bounds gradually increasing correspond to a lower probability of default.

So, in addition to being the most prevalent 'Limit Range' in the dataset, the corresponding '5,000 - 50,000' range is the one with which the highest absolute probability of default is as-

sociated, amounting to 31.79%.

Following the observation of these data, the distribution of the variable 'Marriage' was considered, in order to understand its characteristics and its links with the other variables in the dataset.

To a 'Marital Status' = 1 corresponds the status of 'Single,' while to an equal 'Marital Status' = 2 that of 'Married.' On the other hand, as far as class 0 is concerned, this is referred to by a status equal to 'N/A', while that related to class 3 by a status corresponding to 'Others'.

## Chart n.9: Marriage Variable Distribution



Given the sparse presence of the last two classes mentioned within the dataset, these were aggregated within a single class, equal to 3, into which customers with a marital status equal, with a good degree of probability, to divorce are assumed to converge.

Observing which are the probabilities of default corresponding to each of these marital status categories, it can be seen from Chart n.10 that married people, having generally good financial stability due to the financial help that can be taken away from the partner in case of need, have an associated '% of Default' lower than that of 'Single' and 'Divorced' people.

**73**

Chart n.10: % of Default by Marital Status



Examining instead what the distribution of the target variable is on the dataset under analysis and exploration, Chart n.11 shows that it is characterized by a marked class imbalance.

Chart n.11: Overall Default %

The 22.1 percent of the customers in the dataset are classified as 'defaulting customers,' while the 77.9 percent of them are classified as 'non-defaulting customers.'

Subsequently, in order to have a broader and all-encompassing view of what are the correlations among the variables in the dataset, a correlation matrix was taken into analysis from which it emerges, quite clearly, how many of the variables in the dataset are strongly correlated with each other: from this observation, therefore, comes the need to reduce the dimensionality of the dataset by going to create variables that group the information content of multiple features within a more limited number of them.

Chart n.12: Original Dataset's Correlation Matrix

Thus, three new variables are created: 'REM_DEBT' (i.e., the remaining debt associated with each customer), 'pay_delay' (expressing the total accumulated delay in repaying the debt each customer owes the bank), and 'USAGE_IDX' (reporting to which extent each customer has taken advantage of the line of credit taken out with the bank).

Specifically, 'REM_DEBT' = 'TOT_BILL_AMT' / 'TOT_PAY_AMT', where the numerator is obtained through the summation of all 'BILL_AMT' present for each record, just as the denominator is obtained through the summation of all 'PAY_AMT' present in the dataset for each customer.

Instead, 'USAGE_IDX' = 'AVG_BILL_AMT' / 'LIMIT_BAL', where 'AVG_BILL_AMT' = 'TOT_BILL_AMT' / 6, as cash flows for each customer over six separate and consecutive months are considered in the dataset.

Lastly, 'pay_delay' = 'PAY_1' + 'PAY_2' + 'PAY_3' + 'PAY_4' + 'PAY_5' + 'PAY_6', in which 'PAY_i' expresses how many months are in arrears in payments.

The correlation matrix including the newly created variables appears as shown in the Chart n.13:

Chart n.13: Correlation Matrix after Feature Engineering



**BEST-PERFORMING MODEL: PART 1**

As a result of these observations, the variables selected as 'predicting variables,' and thus standardized in order to be taken into consideration as such, are as follows:'AGE,'EDUCA-TION,'SEX,'LIMIT_BAL,'MARRIAGE_dummy_1,'MARRIAGE_dummy_2,'REM_DEBT' , 'pay_delay' , 'USAGE_IDX.

The two variables expressing marital status, resulting from a variable dummying operation, exclusively consider 'Single' and 'Married' customers.

**77**

# Master's Thesis Flavio Ungherini

Given the critical issues that having an unbalanced dataset may entail in the context of a classification analysis performed by machine learning models, the class imbalance was adjusted by going to adopt, in combination with the models trained to classify banks' defaulting customers, three different sampling techniques: under-sampling, over-sampling and SMOTE.

The models used in this first phase of the study are: KNN, Gaussian Naive Bayes, Support Vector Machine, Random Forest Classifier , Logistic Regression and the Decision Tree Classifier.

The model that performs better is the Random Forest (Table n.1), adopting a SMOTE-type sampling technique and internal cross validation with a split number of 5.

The set of hyper parameters used to achieve this performance is as follows: ("n_estimators": [2, 3, 4], "max_depth": [3, 4], "max_features": [3,4].

These hyperparameters were optimized using GridSearchCV, an iterative algorithm which tests each possible combination and evaluates classifier's performance via cross-validation. Then, the splitting procedure was executed using scikit-learn's train_test_split function for testing split samples.

This function randomly splits data with an initial test size of 20% and random seed of 10. After initial split, SMOTE algorithm is applied to address class imbalance in training data; SMOTE generates synthetic samples of the minority class as needed in order to balance distribution evenly across classes.

Table n.1: Random Forest Classification Report

|  | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| 0 | 0.85 | 0.85 | 0.85 | 4683 |
| 1 | 0.46 | 0.46 | 0.46 | 1317 |
|  |  |  |  |  |
| Accuracy |  |  | 0.76 | 6000 |
| Macro Avg | 0.65 | 0.65 | 0.65 | 6000 |
| Weighted Avg | 0.76 | 0.76 | 0.76 | 6000 |

This model's performance can interpreted also looking at the Chart n.14:

Chart n.14: Random Forest Confusion Matrix



Given the nature of this task, the main objective consists in minimizing the number of 'false negatives', currently equal to 723 out of 6.000. The number of considered bank customers in this case is equal to 6.000 instead of 30.000, since the test set is equal to the 20% of the whole dataset.

**EDA AND FEATURE ENGINEERING ON THE EXTENDED DATASET: PART 2**

# Master's Thesis Flavio Ungherini

The goal of this second part of the analysis is to go on to create additional variables that enrich the information content of the dataset examined, so that the performance of the classification algorithms used for the detection of defaulting customers by banks can be improved. Building on what was observed in the concluding part of the first part of the analysis, the 'Random Forest' is the best-performing model so far. Since it is a probabilistic model, the latter returns as output the probability that each customer in the dataset under study has of going into default, as well as the complementary probability of not going into default.

When the probability of default for any observation is greater than 50 percent, then the customer in question is classified as 'defaulting,' and vice versa.

An example of what is the output returned by the Random Forest, as well as by any probabilistic machine learning model, is represented in the Table n.2:

Table n.2: Random Forest Prediction Output

|   | No Defaut | Default |
|---|---|---|
| 0 | 0.96 | 0.04 |
| 1 | 0.45 | 0.55 |
| 2 | 0.94 | 0.06 |
| 3 | 1.00 | 0.00 |
| 4 | 0.93 | 0.07 |

For example, in the probability result associated with the first of the rows in the output provided, it is possible to see that the customer in question has an associated probability of 'non-default' of 96 percent, and a complementary probability of 'default' of 4 percent.

So, it is intuitive to understand how this person will not be considered at risk of default by the bank of which he or she is a customer.

A good way to exploit the probabilities returned in output from the random forest is to create risk bands within the dataset.

For this purpose, four different risk classes were created: 'Low risk', for customers with a default risk of less than 25 percent, 'Medium-Low risk' for people with whom a default risk between 25 percent and 50 percent is associated, 'Medium-High risk' for bank clients with a default risk between 50 percent and 75 percent, and finally 'High risk' for people whose probability of going into default is between 75 percent and 100 percent.

This thresholds setting within the afore-mentioned variable represents a further analysis about the results obtained with the Random Forest that allows to add a new variable to the dataset at stake.

The risk bands thus defined are distributed as shown by the Chart n.15:

## Chart n.15: Risk Categories Classes Overview



The 'Low' and 'Medium-Low' risk bands together amount to 77.24 percent of the total. Those with a particularly low risk of default are even more than half of the total number of clients in

the dataset. However, 8 percent of the observations in the latter are associated with a very high probability of default, above 75 percent.

After dividing the variable 'AGE' into 4 sub-groups, categorizing these sub-groups based on the quartiles of the distribution. Relating this variable to that of risk categories, the Chart n.16 suggests the following insights:

Chart n.16: Risk Categories by Age Categories



'Young' people seem to be the ones that are mostly exposed to a level of risk equal to 'High' or 'Medium-High'. This is in line with the dataset history and the social context at stake. Specifically, out of the four age categories defined, the 29,55% of people with an 'High' default risk category are 'Young'. On the other hand, 'Mature' people, who are the ones with an age ranging from 34 to 46, are the ones who are less likely to default: in fact, the 27,9% of

people with a 'Low' default risk are 'Mature'. Furthermore, only the 17,3% of customers falling within this age category have an high risk default probability.

Instead, creating four different sub-groups for the variable 'USAGE IDX', according to a subdivision made by the feature's quartiles and comparing it with the risk categories defined earlier, the following can be seen:

Table n.3: Risk Categories by Usage Index

| Risk Categories | High | Low | Medium-High | Medium-Low |
|---|---|---|---|---|
| USAGE IDX Category | | | | |
| Low Usage | 58 | 679 | 199 | 564 |
| Medium-Low Usage | 78 | 1027 | 134 | 261 |
| Medium-High Usage | 143 | 809 | 225 | 323 |
| High Usage | 197 | 618 | 332 | 353 |

Table n.3 shows as the people that make an High and Medium-High use of their available credit are exposed to a major default risk than the ones that make a less marked use of the credit granted by the bank. Indeed, 70,66% percent of people belonging to a category risk equal to 'High' has a usage index category equal to 'High' or 'Medium-High'. On the other hand, only the 45% of people associated with a 'Low' category risk make an 'High' or 'Medium-High' usage of the available credit.

One variable on which an in-depth focus has not yet been performed is 'EDUCATION,' which is an extremely important indicator when trying to predict a customer's default. In fact, a higher level of education may suggest, or at least hint at, greater knowledge of the financial instruments used and greater budgeting skills developed through studies, as well as a more controlled propensity to consume commensurate with one's economic possibilities.

This feature contains four sub-groups within it: 1 = 'High school diploma'; 2 = 'Bachelor's degree'; 3 = 'Mater's degree'; 0 = 'Others'.

Table n.4: Risk Categories by Education Levels

| Risk Categories | High | Low | Medium-High | Medium-Low |
|---|---|---|---|---|
| EDUCATION | | | | |
| 0 | 2 | 67 | 3 | 23 |
| 1 | 134 | 408 | 171 | 293 |
| 2 | 268 | 1325 | 482 | 728 |
| 3 | 72 | 234 | 234 | 457 |

As observable from Table n.4, the 63,6% of people with a Master's degree has a low default probability, while the 47,2% of people with a 'Bachelor degree' tends to have a low risk default probability. Moreover, it can noticed that the 40,5% of the customers with a 'High school' diploma tend to have a low risk default probability. These findings prove as the higher the education level, the higher the chances of having a 'Low' default risk probability. Likewise, the probabilities of falling within the 'High default risk' category are distributed as follows: 3,4% for customers with a Master's degree, 9,5% for Bachelor degree graduates, and 13,3% for customers with an 'High school' diploma. Overall, the higher the eduction level, the lower the chances of meeting an high default risk probability.

Up to this point in the EDA, new variables have been derived from existing features, identifying subgroups within them and comparing them to risk categories.
However, the objective of this second part of the analysis is to create new variables by retrieving data from external sources.
A variable of absolute importance for the purpose of classifying a customer as defaulting or not is its perceived income. However, this data is not contained in the dataset under analysis,

effectively limiting what may be the performance capability of the machine learning models employed for classification.

However, thanks to information found within the Taiwan Statistics Bureau, pertaining to wages received in Taiwan in the early 2000s' by age group and level of education, it was possible to reconstruct sixteen distinct levels of estimated monthly income in Taiwan Dollars (TWD).

The sixteen levels of 'Monthly Income Estimate (TWD)' are the result of the combinations of the four levels of education and the four age groups defined in this second part of the analysis. Exploring all the possible combinations resulting from the joint analysis of the four sub-groups of both of these variables, it is possible to identify sixteen distinct micro-groups, to each of which corresponds a different level of estimated monthly income in Taiwanese dollars.

It is certainly interesting to visualize the income levels that have been defined in the new variable 'Monthly Income Estimate (TWD)' with the risk categories previously defined:

Chart n.17: Monthly Income and Risk
Categories' Aggregated View

# Master's Thesis Flavio Ungherini

In general, from Chart n.17 it is possible to observe how the majority of wages are above 39,231 TWD monthly; focusing exclusively on incomes above the figure just indicated, given their massive presence within the dataset under study, it is possible to observe how increasing levels of income correspond to a lower level of 'High' risk of default.

After creating this new variable, the goal is to create new ones that can further assist the machine learning models employed for this task in performing the classification analysis.

One feature that could play an important role in this regard, much used by both banks and companies operating in all types of industries in order to understand what the economic value of their customers is, is CLV (Customer Lifetime Value).

In order to create this new field within the dataset, two new objects were defined in the work environment, as well as a new variable within the dataset itself.

Specifically, first, the new variable 'Annual Revenue' was created, obtained through the simple multiplication of the 'Monthly Income Estimate (TWD)' by the twelve months that constitute a year.

Next, the 'Retention Rate' was defined for each customer in the dataset. In order to derive this new feature, the variable 'REM_DEBT' was taken into analysis, which represents the measure of the debt owed by each customer to the bank.

After grouping the variable by index level = 0, and after calculating the average 'remaining debt' for each group of customers, the retention rate is defined and used together with the 'Monthly Income Estimate (TWD)' in the Customer Lifetime Value cascade for each customer.

The reason why the idea of remaining debt is associated with the retention rate is that the amount of remaining debt owed by each customer to the bank certainly has a significant influence on what the quality and durability of the relationship between the two parties may be over time. Intuitively, it is logical to expect that the lower the amount of debt owed by the customer to the bank, the higher the likelihood that the two parties will continue to have a lasting relationship over time.

In order to calculate CLV, however, it is also necessary to define a 'Discount Rate,' which is assumed to be 10 percent, according to a rule very frequently repeated in the calculation of

**86**

customer lifetime value.

Hence, at the conclusion of all these steps, the 'Customer Lifetime Value' is calculated as follows:

$$CLV = \text{'Annual Revenue (TWD)'} * (1 - \text{'Retention Rate'}) / (\text{'Discount Rate'})$$

Once defined, 5,700 unique values are counted for this variable.

Next, a new feature, named 'CLV_groups,' is created from the variable just defined, with the objective of identifying four subgroups of Customer Lifetime Value on the basis of a division made by quartiles.

Thus, the following four categories are identified: 'Negative CLV', 'Low CLV', 'Medium CLV' and 'High CLV'.

Once these subgroups are defined and related to the risk categories, it can be seen that 85.5 percent of total clients have 'Negative CLV'. Given the history and socioeconomic background of the dataset, it is not hard to believe that such a high percentage is associated with their level of CLV. In fact, most customers are likely to generate negative economic cash flows for the bank in the future because the bank has invested primarily in acquiring an economically and socially fragile customer base.

In addition, 45% of customers with "high CLV" have a "low probability of credit risk"; this probability is 48% for customers with "medium CLV" and 56% for customers with "low CLV." Similarly, 3.7 percent of customers with a "high CLV" have a "probability of high credit risk"; this percentage is 7 percent for customers with a "low CLV."

Focusing exclusively on customers with a positive CLV, Chart n.18 is rendered:

Chart n.18: Credit Risk Distribution Within
CLV's Levels



Among the customers with a positive CLV, those with a higher presence in the dataset are those with an associated 'Low CLV' level. In addition, it is interesting to note that those with a 'High CLV' are exposed to a particularly low risk of default.

Next, by comparing the levels these categories of CLV levels with the four classes related to 'EDUCATION' levels, it is noted how the 95.5 percent of customers with a 'High CLV' has at least a 'Bachelor Degree' study title, while out of the customers with a 'Low CLV' only the 69.8 percent has at least a 'Bachelor Degree'. It's interesting to notice as the there are no people with a study title inferior than the High School diploma among the customers with a 'High CLV'. At the same time, it's relevant to observe that the maximum number of people with a study title inferior than a High School diploma belongs to the 'Low CLV' customer category. Taking advantage of the variables currently in the dataset, it is possible to create new variables that will add to the information content and dimensionality of the dataset.

# Master's Thesis Flavio Ungherini

It is possible to use the two features 'Annual Revenue (TWD)' and 'REM_DEBT' in order to create two new variables expressing the 'Customer Profitability' and the 'Debt-to-Income Ratio' of each customer.

Specifically, 'Customer Profitability' is obtained as the difference between the 'Annual Revenue' in Taiwanese dollars and the 'Remaining Debt' already available in the dataset.

In addition to this, 'Debt-to-Income Ratio' can be defined as the ratio of 'Remaining Debt' to 'Annual Revenue'.

Finally, it is possible to create a new variable, 'Payment Ratio,' obtained as follows: ('Limit Balance' - 'Remaining Debt') / ('Limit Balance').

This variable expresses how much of the debt from credit card use has been returned to the bank by each customer.

In addition to the metrics just defined in the dataset, a KPI that not only banks but every type of company looks at with particular attention is ROI (return on investment).

ROI allows, in this house for the bank whose customers are being analyzed, to understand whether the investment made in acquisition and retention on each individual customer is matched by a positive or negative ROI. The higher the ROI on a customer, the more worthwhile the investment on that person by the bank.

Circumscribing this consideration to the history of the dataset examined, the main investment made by the bank in question refers to advertising expenses incurred to incentivize credit card consumption on each customer acquired.

However, for a creation of the 'ROI' feature as precise, accurate and close to reality as possible, it is deemed necessary to go through a very complex and structured data retrieval process.

'Return on Investment' was calculated as follows within this dataset:

$$ROI = (\text{'Customer Profitability'}) / 6.02$$

6.02 TWD is the estimated advertising budget amount spent by the bank object of study in 2005. According to the data provided by this Statista report: https://www.statista.com/sta-

**89**

tistics/308842/ad-spend-credit-card-issuers-usa/, it is possible to consider as a benchmark the eight major American credit card issuers advertising expenses in the US in 2021 (in US Dollars). Those expenses were subsequently averaged, divided by the number corresponding to the American population above the age of 18 in 2021 (77,8% of the total US population: the American population as a whole was equal to roughly 332 mln inhabitants, 258.3 mln of which above the age of 18). Thus, dividing the average advertising spend of those eight credit card issuers in the US in 2021, equal to 76 mln US dollars, by the aforementioned population, the resulting average amount of economic resources spent on every population's credit card eligible person is equal to 0.2942 US dollars. Given an exchange rate, at the moment of writing, equal to 1$ = 30.74 TWD, then 0.2942 US dollars = 9.04 TWD. Since the dataset taken at stake relates to a phenomenon occurred in 2005, the 9.04 TWD amount was proportionated to 6.02 TWD, on the base of the effective inflation rate in Taiwan in 2005 (equal to 2.31%, data taken from Statista at this link: https://www.statista.com/statistics/727598/inflation-rate-in-taiwan/).

The formula used to obtain the estimated 2005 advertising spend amount is the following one: (9.04) / (1+0.0231) ^ (2021-2005) = 6.02 TWD.

Lastly, it is necessary to specify that these big eight credit card issuers were considered as a benchmark for this analysis since the bank object of study is assumed to be a big player in the financial services industry in Taiwan.

In the graph represented below, the distribution of customers with positive and negative ROIs within the dataset can be observed in Chart n.19:

## Chart n.19: ROI's Overview



Positive and Negative Values Distribution

Positives

16.9%

83.1%

Negatives

Given the aggressive and economically demanding marketing campaigns carried out by banks in Taiwan in 2005 and the low credit-worthiness of a great part of the credit cards holders in that Country, the Return on Investment (where for investment is intended the estimated amount of money spent on advertising by credit card issuers) mainly takes negative values: specifically, to the 83,1% of customers in the dataset corresponds a negative ROI. After extending the dimensionality of the dataset by adding variables in the dataset that did not originally exist, we can take a look at the correlation matrix of the new dataset.

Chart n.20: New Dataset's Correlation
Matrix

testX_copy
correlation plot (Pearson)

Of all the variables in the matrix shown in Chart n.20, a portion of them will be employed among predicting variables used to train the proposed machine learning models to the classification of defaulting customers.

The goal at this point in the analysis is to see if the newly created variables can improve the classification performance observed in the first part of the study.

To do this, three separate sets of predicting variables will be constructed: one consisting of economic variables only, one consisting of socio-demographic variables only, and a third set consisting of an intersection between these two sets, that is, containing both economic and socio-demographic variables.

The set of economic-only predicting variables counts within it the following features, none of which were initially present in the original dataset: 'CLV', 'Monthly Income Estimate (TWD)', 'Debt-to-Income Ratio', 'Payment Ratio', 'ROI (TWD)'.

The set of variables targeting only socio-demographic variables includes the following features: 'Age', 'Sex', 'Education', 'Marriage_dummy_1','Marriage_dummy_2'

Finally, the more extended set of predicting variables, obtained as a fusion between the two groups just exposed, consists of: 'Age', 'Sex', 'Education', 'Marriage_dummy_1', 'Marriage_dummy_2', 'Monthly Income Estimate (TWD)', 'Debt-to-Income Ratio', 'Payment Ratio', 'ROI (TWD)'.

The Machine Learning models that will be trained from the sets of variables just outlined are all probabilistic ones, since in the first phase of the study the best-performing model turned out to be the Random Forest, which is precisely probabilistic in nature.

The models used in this phase of the study are: Random Forest, XGBoost, AdaBoost and Gradient Boosting.

# Chapter 5: Results

**COMPARATIVE ANALYSIS OF THE USED MODELS: FUNCTIONALITIES, SIMILARITIES, DIFFERENCES AND ADVANTAGES**

This section presents an overview of four machine learning models - Random Forest, XGBoost, AdaBoost and Gradient Boosting - that were utilized to predict banks' defaulting customers. Their functions and characteristics are discussed extensively along with any algorithmic approaches or similarities or differences they possess, along with any advantages these models possess over alternative ways of anticipating bank customer defaulting behaviors.

# Master's Thesis Flavio Ungherini

1. Random Forest: Random Forest is an ensemble learning method that utilizes multiple decision trees in combination to make predictions[1]. It works through bootstrap sampling and feature randomization. In Random Forest, each tree is trained independently on randomly chosen subsets from training data before collectively aggregating all predictions through majority voting or average prediction for classification or regression tasks respectively[1]. By randomly choosing subsets and features among trees within Random Forest, its diversity ensures more efficient handling of high dimensional data while decreasing overfitting[1].

Random Forest can bring several advantages to predictive modeling. First, its automated feature selection identifies only relevant features at each split split point to quickly select those most pertinent to prediction [1]. Second, by aggregating multiple trees' predictions together for improved generalization performance, Random Forest also helps analysts gain a clear picture of relative contribution of different variables during prediction [1]. Lastly, Random Forest allows an analyst to accurately gauge importance estimates between features to better comprehend relative contributions across variables[1].

2. XGBoost (eXtreme Gradient Boosting): XGBoost is an optimized implementation of gradient boosting that addresses limitations found with traditional gradient boosting methods by employing a more robust algorithmic framework. In essence, this framework sequentially builds an ensemble of weak learners (usually decision trees) by minimizing loss functions[2]. XGBoost stands out by being capable of capturing second-order gradients for more precise estimation of model parameters[2].

The XGBoost algorithm begins by initializing its model with a constant prediction value and then calculates its gradient and second-order derivative with respect to this prediction, providing valuable information regarding direction and curvature of loss surfaces[2]. On each iteration, weak learners such as decision trees are fitted onto negative gradients of loss functions in an effort to reduce errors produced by previous models [2] with each learning rate controlling contribution by each learner[2]. As new weak learners are added into its model, its overall prediction value remains constant[2].

XGBoost stands out among gradient boosting methods by offering several enhancements over traditional gradient boosting approaches. First, it incorporates regularization techniques such as shrinkage and column subsampling that prevent overfitting while improving model generalization[2]. Second, parallel processing speeds model training considerably faster while providing greater scalability of large datasets[2]. Thirdly, missing value mechanisms automatically learn to impute missing data using available information[2]. These features have made XGBoost popular both within machine learning competitions as well as industry applications[2]. These features contribute greatly towards its immense popularity within various domains[2].

3. AdaBoost (Adaptive Boosting): AdaBoost is an ensemble method which combines several weak learners into one strong learner[3]. AdaBoost works by iteratively adjusting weights based on classification performance of training samples to enable subsequent weak learners to focus their efforts on misclassified ones [3]. AdaBoost's central principle resides in its emphasis of difficult samples thereby improving performance on challenging cases[3].

AdaBoost starts by assigning all training samples equal weights, then uses weak learners such as decision stumps (decision trees with one split) to train on these weighted samples in each iteration[3]. Misclassified samples then have their weight increased for subsequent iterations so the model focuses its efforts on challenging instances while increasing its ability to handle complex datasets[3]. This adaptive weighting scheme gives AdaBoost its edge for complex datasets[3].

AdaBoost generates its final prediction by aggregating all weighted votes of weak learners[3]. Each weak learner's weights reflect its classification performance; those that achieve greater accuracy receive additional weight[3]. Combining multiple weak learners into a stronger learner helps AdaBoost strengthen its generalization capability more than an individual weak learner might[3].

4. Gradient Boosting: Gradient Boosting is an ensemble method which employs iterative addition of weak learners - typically decision trees - until its performance meets an established

loss function [4]. Gradient Boosting optimizes model performance by targeting residual errors made by previous weak learners[4]. Gradient Boosting utilizes sequential training of weak learners to correct past miscalculations made by earlier weak learners to gradually lower overall prediction error [4].

Gradient Boosting algorithms start their models off with constant prediction values [4]. They then compute a negative gradient of loss function with respect to current model predictions, providing direction on where predictions should be adjusted in order to minimize loss[4]. A weak learner, often decision trees, is fitted on to this negative gradient capturing patterns left by prior models[4]. Once again the model updates by adding learners using learning rates which determine their contribution towards creating final predictions[4].

Gradient Boosting offers greater freedom when choosing its loss function, allowing the model to be tailored specifically to each task requirement[4]. It can be utilized for both regression and classification tasks with various data types and prediction objectives in mind[4], thanks to its iterative nature and emphasis on residuals enabling it to capture complex relationships within data, leading to highly accurate predictions when tuned and optimized accordingly[4].

Analogies and Differences: All four models--Random Forest, XGBoost, AdaBoost and Gradient Boosting--can be classified as ensemble learning techniques that use multiple weak learners to form one strong learner. All four use decision trees as weak learners[1][2][3][4], yet there remain notable distinctions among them.

Random Forest constructs each tree independently through bootstrap sampling and feature randomization, in an effort to minimize overfitting and ensure robust predictions. By contrast, XGBoost, AdaBoost, and Gradient Boosting sequentially build trees; subsequent trees aim at correcting errors made by predecessor trees. These sequential models learn from errors they made earlier while improving predictive performance as time goes by.

Gradient Boosting (XGBoost) and Gradient Boosting share an inherent similarity: updating models by minimizing specified loss functions using negative gradients, updating models us-

ing regularization techniques, parallel processing capabilities and handling missing values as described here[2][4]. However, XGBoost adds further advantages, including regularization techniques, parallel processing capacities and handling missing values[2].

AdaBoost takes an alternative approach by adapting training sample weights, placing more importance on misclassified instances to boost its model's performance[3]. With its adaptive weighting scheme, AdaBoost excels at handling large datasets or challenging instances[3].

One key difference lies in the capability of XGBoost and Gradient Boosting to directly address missing values while Random Forest and AdaBoost require preprocessing or imputation methods in order to effectively handle missing data[2].

In conclusion, while Random Forest, XGBoost, AdaBoost and Gradient Boosting all employ similar decision tree approaches as ensemble learning algorithms, each algorithm employs distinct techniques when building its ensemble models. Random Forest models boast superiority over alternative approaches for predicting bank customer default, including robustness, automatic feature selection and prediction aggregation. XGBoost stands out as an innovative improvement over conventional gradient boosting methods with regularization, parallel processing and the handling of missing values. AdaBoost features challenging instances and adaptive weighting which is great for complex datasets while Gradient Boosting takes an iterative approach and emphasizes residuals to capture complex relationships while producing highly accurate predictions.

References:

[1] Breiman, L. (2001). Random forests. Machine learning, 45(1), 5-32.

[2] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794.

[3] Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. Journal of computer and system sciences, 55(1),

119-139.

[4] Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. The Annals of Statistics, 29(5), 1189-1232.

**EMPLOYED METRICS FOR THE MODELS PERFORMANCE EVALUATION**

As part of an evaluation of a machine learning model's ability to classify bank customers who default, one important metric to monitor is its "Recall for Class 1". This measures whether all instances of defaulted customers - also referred to as the positive class - were accurately identified by its model.

Recall (also referred to as Sensitivity or True Positive Rate) measures how accurately models identify actual positive instances (defaulting customers). Recall for Class 1 is designed specifically to measure how successfully models identify defaulting customers; often an important focus for banks.

As soon as a bank begins dealing with defaulting customers, having an accurate model with high recall for class 1 is paramount; misclassifying defaulting customers as non-defaulters (false negatives) could have severe repercussions for them and their bank. A false negative occurs when an inaccurate model predicts someone as non-defaulting when in reality they have actually defaulted and can lead to financial losses due to late risk assessment or debt collection procedures being performed on these individuals.

Prioritizing recall for class 1 customers can help banks identify more defaulting customers quickly, thus minimizing risks related to missed potentially problematic cases. Recall should be balanced against other evaluation metrics like precision which measures proportion of predicted positives that actually turn out true positives - this depends on each bank's specific goals and constraints for operations.

Recall for class 1 is one measure to evaluate how accurately machine learning models can detect defaulting customers; "F1-score for class 1" provides another balanced evaluation metric which takes precision and recall into consideration.

"F1-score" is a measure of model performance which incorporates precision and recall into one value; specifically, "F1-score for class 1" measures the model's success in classifying defaulting customers.

Utilizing both F1-score for class 1 and Recall for class 1 helps address the tradeoff between precision and recall. Precision indicates the proportion of predicted positive instances that turn out to be truly positive, while Recall represents how accurately real positive instances were identified as being such.

Classifying defaulting customers requires high recall to prevent false negatives and capture as many defaulters as possible; however, such high recall could also increase false positives - non-defaulters being misclassified as defaulting - leading to unnecessary actions, like blocking accounts or denying credit to people who may never actually default.

By taking into account the F1-score, a model's overall performance can be evaluated using precision and recall measures, in combination. This measure serves as a balanced check on models which achieve both precision and recall; rewarding models which meet this objective. By finding this balance point between capturing all defaulting customers while minimizing false positives.

Overall, while recall for class 1 places an emphasis on capturing defaulting customers, F1-score for class 1 provides a balanced evaluation that simultaneously measures precision and recall; offering banks more nuanced insight into how successful their models have been at classifying banks' defaulting customers.

# Master's Thesis Flavio Ungherini

As part of any bank classification task, making use of AUC (Area Under the ROC Curve) metrics can prove very helpful in accurately classifying defaulting customers across classification thresholds.

A Receiver Operating Characteristic (ROC) curve can be created by plotting true positive rate against false positive rate at various threshold settings and measuring it against AUC (area under this curve), where AUC ranges between 0-1; higher AUC indicates superior model performance in distinguishing defaulting from non-defaulting customers.
 Here are a few reasons why the AUC metric can help with this task:

-Evaluation across Multiple Thresholds: The Area Under Curve (AUC) measures how well a model performs at various classification thresholds, providing a more complete picture of its discriminative ability. It's particularly relevant when the cost of false positives and false negatives cannot easily be quantified, or when decision thresholds may vary based on individual contexts.

-Robustness to Class Imbalance: The AUC metric is generally resistant to class imbalance issues commonly encountered when performing default prediction tasks, like banking datasets where defaulters make up a minority class and thus affect evaluation metrics negatively. AUC serves as a reliable measure of whether models can rank defaulting customers higher than non-defaulting customers across various class distributions.

-Model Comparison: Area Under Curve (AUC) allows for simple and direct comparison between various models or variants of one model, providing one scalar value that summarizes overall discriminatory power to facilitate straightforward evaluations between algorithms or techniques.

-Threshold Selection Guidance: ROC curve and AUC provide insights into the tradeoff between true positive and false positive rates, providing decision makers with guidance for selecting an optimum threshold based on risk tolerance and operational needs; selecting one can balance accurate identification of defaulting customers while minimizing false positives.

Overall, AUC provides an effective measure for gauging a model's classification performance across thresholds; provides robustness to class imbalance; allows model comparison; and assists threshold selection. In doing so it complements other evaluation metrics like recall and F1-score in improving understanding and decision-making when classifying bank customers who have defaulted.

As part of their assessment of machine learning models used to classify bank customers in default, banks take an approach where multiple metrics and weighted combination thereof are taken into consideration and assigned different weights of importance - this results in the creation of the Weighted Importance Metric which takes into account different weightings for recall for class 1, F1 score for class 1 and AUC metrics - then ranks model performances according to this new 'Weighted Importance' metric in descending order.

Assigning weights of importance to evaluation metrics reflects the priorities and preferences when assessing model performance. By assigning these weights, certain metrics can be highlighted based on their significance to particular requirements; here, the highest weight (0.65 is given to "Recall for class 1," signifying its significance for correctly identifying defaulting customers while at the same time minimizing false negatives to capture as many defaulting customers as possible.

A weight of 0.30 assigned to F1-score for Class 1 indicates a balanced measure which takes into account precision and recall equally, suggesting an attempt at finding an equilibrium between correctly identifying defaulting customers while also minimizing false positives.

At last, AUC received an estimated weight of 0.05 to reflect that while distinguishing defaulting and non-defaulting customers is essential to any evaluation effort, this metric also serves to evaluate overall model performance; yet given less weight.

Ranking model performances by their Weighted Importance in descending order allows for identification of models which excel across all metrics according to weights specified, pro-

viding an integrated view that considers different aspects of model performance while taking preferences into consideration.

By employing the "Weighted Importance" metric and ranking, models can be prioritized and compared based on their ability to detect defaulting customers while striking an equitable balance between precision and recall as well as distinguishing defaulters from non-defaulters.

**EVALUATION OF ML MODELS FOR DEFAULT PREDICTION: IDENTIFYING SUPERIOR PERFORMERS WITH THE 'WEIGHTED PERFORMANCE METRIC'**

Table n.5: Classification Models

Performance Evaluation

| Index | Model Name | Weighted Performance | Recall for class 1 | F1 score for class 1 | AUC score |
|-------|-----------|---------------------|-------------------|---------------------|-----------|
| 1 | **ADA (Eco.2)** | **0.55** | **0.62** | **0.39** | **0.6** |
| 2 | RF (Demo) | 0.5355 | 0.56 | 0.47 | 0.61 |
| 3 | ADA (Demo.1) | 0.528 | 0.58 | 0.4 | 0.62 |
| 4 | ADA (Eco.1) | 0.512 | 0.54 | 0.43 | 0.64 |
| 5 | ADA (Eco-Demo.1) | 0.497 | 0.55 | 0.37 | 0.57 |
| 6 | XGB (Demo) | 0.4945 | 0.5 | 0.46 | 0.63 |

| 7 | ADA (Eco-Demo.2) | 0.491 | 0.54 | 0.37 | 0.58 |
|---|---|---|---|---|---|
| 8 | GB (Eco) | 0.488 | 0.54 | 0.36 | 0.58 |
| 9 | XGB (Eco-Demo) | 0.3515 | 0.34 | 0.34 | 0.57 |
| 10 | RF (Eco) | 0.323 | 0.31 | 0.31 | 0.57 |
| 11 | GB (Demo) | 0.3225 | 0.28 | 0.37 | 0.59 |
| 12 | RF (Eco-Demo) | 0.319 | 0.29 | 0.34 | 0.57 |
| 13 | GB (Eco-Demo) | 0.3185 | 0.28 | 0.36 | 0.57 |
| 14 | XGB (Eco) | 0.29 | 0.26 | 0.31 | 0.56 |

Table n.5 presents an evaluation of four machine learning models employed in this study, each categorized according to its "Weighted Performance" metric. The purpose is to identify models which demonstrate superior predictive capabilities while accounting for the complex nature of the problem.

For all the models at stake, a SMOTE sampling technique was used, as well as an inner cross validation with external hold-out.

In this table, 'ADA' = AdaBoost, 'RF' = Random Forest, 'XGB' = XGBoost, 'GB' = Gradient Boosting.

To every model is associated a set of predicting variables, contained in parentheses.

Specifically, 'Eco' indicates the set of economic predicting variables, 'Demo' the set of socio-demographic predicting variables and, lastly, 'Eco-Demo' refers to the group of predicting features including both economic and socio-demographic variables.

In the case of 'ADA' (AdaBoost), it can be noticed that the numbers '1' and '2' are also included in the parentheses together with the predicting variables set label.

When '2' appears in the parentheses, it means that the AdaBoost was trained using a wider range of hyperparameters than the one in which '1' appears together with the predicting variables set label.

ADA(DEMO.2) appeared to be model with an higher 'Weighted Performance', as well as with an higher 'Recall for class 1' score. Although, having this model a feature importance equal to 100% for the 'EDUCATION' feature, this model was considered as anomalous and then discarded from the present analysis.

As a consequence, ADA(ECO.2) stands out among the models by scoring 0.55 on its "Weighted Performance" scale; this indicates its exceptional predictive prowess compared to its counterparts. ADA(ECO.2)'s performance indicates its adept incorporation of multiple economic predicting variables enabling accurate default risk predictions; by capitalizing on such features ADA(ECO.2) has shown superior ability in understanding nuanced relationships between economic factors and default probabilities.

Examination of evaluation metrics provides additional insight into the strengths of tested models. Notably, ADA(ECO.2)'s recall score for class 1 stands out; at 0.62 this reflects its ability to correctly recognize instances belonging to positive class (customers at risk of defaulting payments), thus successfully minimizing false negatives which is essential in helping identify as many at-risk customers as possible.

F1-score for class 1, which provides an objective balance measure that considers both precision and recall, highlights RF(DEMO)'s effective performance with its 0.47 score; this indicates its success at both accurately identifying positive instances while minimizing false positives occurrence. Furthermore, its F1 score shows its reliability for making reliable predictions while simultaneously maintaining low misclassification rates; making this balanced performance particularly relevant when precision and recall carry equal weight in certain circumstances.

Additionally, when looking at AUC scores (area under receiver operating characteristic curve), ADA(ECO.1) stands out with an outstanding AUC_score of 0.64 - the area under receiver operating characteristic curve serves as an aggregate measurement of discriminatory power and ability to distinguish positive from negative instances - its superior AUC score suggests robust discriminative capabilities when looking across economic predictability vari-

ables; its high AUC score also highlights how effectively differentiate default risks to enable more precise risk analyses and informed decision-making processes.

These findings demonstrate the numerous strengths displayed by all evaluated models: ADA(ECO.2) excels in terms of weighted performance and recall; RF(DEMO) stands out with its favorable F1 score; while ADA(ECO.1) displays superior discriminative capabilities. Such findings underscore the value of using multiple metrics when assessing models as it provides more nuanced understandings into each models' predictive ability and highlights any individual models' respective strengths or weaknesses.

Based on an evaluation of the table provided, it is evident that the model with the highest "Weighted Performance" value outperforms all others for identifying customers at risk of defaulting. ADA(ECO.2) excelled with a "Weighted Performance" score of 0.55; its predictive capabilities outshone those of all other models evaluated. By prioritizing models with higher weighted Performance values banks can make more accurate predictions and identify customers likely to default payments more accurately and identify those likely to do so sooner.

# Chapter 6: Discussions

### COMPARATIVE ANALYSIS OF FEATURE IMPORTANCES IN ML MODELS: EXPLORING ANALOGIES AND DISCORDANCIES BETWEEN ECONOMIC AND SOCIO-DEMOGRAPHIC PREDICTORS

This section investigates any similarities or discrepancies observed among multiple Machine Learning models concerning feature importances. Analysis centers around three sets of predicting variables - Economic and Demographic, Demographic and Economic. Individual feature importances were then assessed within each set using random Forest (RF), XGBoost

(XGB), AdaBoost (ADA) and Gradient Boosting (GB) models - to reveal either consistent patterns or divergent trends among them in regards to predictor importance in relation to target variables ("Default Binary"). Findings highlighted various predictor importances by these models while emphasizing particulars.

Once the ensemble classifiers used (Random Forest, XGBoost, AdaBoost or Gradient Boosting) have been trained, feature importances can be computed using GridSearchCV's best estimator's "feature_importances_ attribute". These measurements of feature importance produce measures based on how each feature contributes to improving ensemble models' performance or loss reduction. Feature importances are computed as the mean of accumulation of the impurity decrease within each tree. The importance values are calculated by combining individual feature importances across all models in an ensemble. For instance, Random Forest determines importance by taking the reduction in 'Gini impurity' obtained during split process across all decision trees as measured by individual features during splitting process across decision trees; similarly XGBoost measures importance according to frequency used and improvement seen in loss function improvement from using each feature for splitting; AdaBoost uses weighted sum of classification errors attributable to each feature and Gradient Boosting evaluates significance by cumulative reduction in the loss function achieved by each feature across all boosting iterations.

Economic and Socio-Demographic Predictors:

These tables were built combining the features importances results obtained when evaluating every each of the four machine learning models used for classifying defaulting customers employing this set of predicting variables.

# Master's Thesis Flavio Ungherini

Economic and Socio-Demographic predictors:

Table n.6: Random Forest (Eco-Demo)

| Feature Name | Importance |
|---|---|
| Payment Ratio | 0.18367 |
| ROI (TWD) | 0.146542 |
| Debt-to-Income Ratio | 0.134697 |
| Monthly Income Estimate (TWD) | 0.126327 |
| AGE | 0.100684 |
| EDUCATION | 0.097644 |
| MARRIAGE_dummy _2 | 0.0954562 |
| MARRIAGE_dummy _1 | 0.0742288 |
| SEX | 0.0407513 |

Table n.7: XGBoost (Eco-Demo)

| Feature Name | Importance |
|---|---|
| Marriage_dummy_2 | 0.323491 |
| EDUCATION | 0.241986 |
| Marriage_dummy_1 | 0.18265 |
| Monthly Income Estimate (TWD) | 0.090508 |
| SEX | 0.079663 |
| ROI (TWD) | 0.025388 |
| Payment Ratio | 0.022138 |
| Debt-to-Income Ratio | 0.019218 |
| AGE | 0.014958 |

Table n.8: AdaBoost (Eco-Demo)

| Feature Name | Importance |
|---|---|
| AGE | 0.57500 |
| Monthly Income Estimate (TWD) | 0.10700 |
| Payment Ratio | 0.09600 |
| ROI (TWD) | 0.0615 |
| Marriage_dummy_2 | 0.04900 |
| Marriage_dummy_1 | 0.045500 |
| Debt-to-Income Ratio | 0.044500 |
| EDUCATION | 0.0145 |
| SEX | 0.00700 |

Table n.9: Gradient Boosting (Eco-Demo)

| Feature Name | Importance |
|---|---|
| Payment Ratio | 0.168448 |
| Monthly Income Estimate (TWD) | 0.157131 |
| Marriage_dummy_1 | 0.141565 |
| ROI (TWD) | 0.129077 |
| Debt-to-Income Ratio | 0.117491 |
| Marriage_dummy_2 | 0.089161 |
| AGE | 0.073761 |
| EDUCATION | 0.073207 |
| SEX | 0.05016 |

Analysis of Economic and Demographic predictor variables shows their different degrees of importance across models. The RF model (Table n.6) pinpoints "Payment Ratio" as one of its key predictors; this measure compares monthly payments against outstanding balance. This suggests that payment habits play a pivotal role in whether or not individuals default on loans. Return on Investment in Taiwanese Dollars indicates the profitability of investments as another indicator for default prediction. "Debt-to-Income Ratio" is considered an important indicator, reflecting borrowers' ability to balance debt payments against income. These findings align with economic intuition by emphasizing financial aspects associated with default prediction.

Contrarily, the XGB model (Table n.7) gives high priority to variables related to marital status and education; specifically "MARRIAGE_dummy_2", an indicator for married individuals, as an influential predictor suggesting they might be less prone to defaulting than not married borrowers. "EDUCATION" ranks second as an influential predictor, suggesting that higher educational levels may decrease risk of default and therefore highlight socio-demographic considerations as potential predictors beyond pure financial considerations in predicting defaults. As indicated above, marital status does have some bearing in predicting default risk, although at a lesser scale than before.

ADA (Table n.8) assigns great significance to both "AGE" and Monthly Income Estimate (TWD), with particular attention given to "AGE." Borrowers' age has long been recognized as a reliable indicator of default risk, providing insight into factors like stability, experience and financial responsibilities that might contribute to it. "Monthly Income Estimate (TWD)" measures the borrower's ability to meet loan payments and is consistently one of the top essential features across all models. Monthly Income Estimate (TWD) ranks second as an influential variable for both ADA and GB (Table n.8, Table n.9); these findings underscore the key role demographic and economic variables play in default prediction; variations across these models highlight their unique predictive strengths.

Socio-Demographic Predictors:

Table n.10: Random Forest (Demo)

| Feature Name | Importance |
|---|---|
| AGE | 0.579094 |
| EDUCATION | 0.230482 |
| Marriage_dummy_1 | 0.0849234 |
| Marriage_dummy_2 | 0.0527912 |
| SEX | 0.052709 |

Table n.11: XGBoost (Demo)

| Feature Name | Importance |
|---|---|
| Marriage_dummy_2 | 0.370964 |
| Marriage_dummy_1 | 0.245708 |
| EDUCATION | 0.202043 |
| SEX | 0.120036 |
| AGE | 0.051768 |

Table n.12: AdaBoost (Demo)

| Feature Name | Importance |
|---|---|
| AGE | 0.670000 |
| EDUCATION | 0.150000 |
| Marriage_dummy_1 | 0.150000 |
| Marriage_dummy_2 | 0.025000 |
| SEX | 0.005000 |

Table n.13: Gradient Boosting (Demo)

| Feature Name | Importance |
|---|---|
| Marriage_dummy_2 | 0.548043 |
| Marriage_dummy_1 | 0..238428 |
| EDUCATION | 0.097552 |
| SEX | 0.06421 |
| AGE | 0.051768 |

Analysis of Demographic Predictors has revealed both consistent trends and variations across models. The RF model (Table n.10) highlighted "AGE" as its most influential predictor, emphasizing its significance for default prediction; this matches up well with older borrowers having greater financial security and lower risks of default than younger borrowers. Next to

"AGE", "EDUCATION" became another influential predictor; further suggesting higher education may correlate to reduced likelihoods of default and further still, "MARRIAGE_dummy_1", representing single individuals was given significant weight due to affective variables influencing default risk.

In the XGB model (Table n.11), feature importances are more evenly distributed among demographic predictors. "MARRIAGE_dummy_2", representing married individuals, holds the greatest significance; suggesting they are likely to engage in responsible financial behaviors that lower default risks than unmarried borrowers. Furthermore, "MARRIAGE_dummy_1", representing unmarried borrowers, also reinforces this idea that marital status influences default prediction. Furthermore "EDUCATION" receives great emphasis indicating higher educational attainment is linked with decreased likelihood of default while "SEX", representing gender as having less significance as far as default risk considerations go.

Similar trends emerge across both ADA (Table n.12) and GB (Table n.13) models, with "AGE" frequently identified as an indicator. Meanwhile, emphasis may vary between marital status (married_dummy_1) and marital_status_2 for default risk prediction purposes, suggesting marital status plays some part but to different degrees across models. Furthermore, "EDUCATION" plays an integral part of default prediction with differing degrees of significance reaffirming its impact.

Among the four predictions made with respect to economic, demographic and socio-demographic predictor variables respectively 'SEX' was shown as having the lowest importance.

Demographic predictors continue to demonstrate their usefulness for default prediction and provide valuable insights into borrowers' characteristics that indicate default risk.

Economic Predictors:

### Table n.14: Random Forest (Eco)

| Feature Name | Importance |
|---|---|
| Monthly Income Estimate (TWD) | 0.274384 |
| Payment Ratio | 0.253461 |
| CLV | 0.173810 |
| ROI (TWD) | 0.158394 |
| Debt-to-Income Ratio | 0.139951 |

### Table n.15: XGBoost (Eco)

| Feature Name | Importance |
|---|---|
| Monthly Income Estimate (TWD) | 0.509423 |
| ROI (TWD) | 0.174011 |
| Payment Ratio | 0.118588 |
| Debt-to-Income Ratio | 0.108223 |
| CLV | 0.089754 |

### Table n.16: AdaBoost (Eco)

| Feature Name | Importance |
|---|---|
| Monthly Income Estimate (TWD) | 0.507500 |
| Payment Ratio | 0.187500 |
| CLV | 0.132500 |
| Debt-to-Income Ratio | 0.092500 |
| ROI (TWD) | 0.080000 |

### Table n.17: Gradient Boosting (Eco)

| Feature Name | Importance |
|---|---|
| Monthly Income Estimate (TWD) | 0.452716 |
| Payment Ratio | 0.239845 |
| CLV | 0.119987 |
| ROI (TWD) | 0.104456 |
| Debt-to-Income Ratio | 0.082997 |

Economic predictor feature importances reflect both similarities and variance among models, providing insights into their relative importance in foretelling default. "Monthly Income Estimate (TWD)" consistently stands out among models as the primary economic indicator, underscoring how income levels play such an integral part of meeting loan obligations and avoiding default. Furthermore, The RF model (Table n.14) places great significance on "Payment Ratio", or reflecting proportion of monthly payments to outstanding balance; higher ratios indicate greater creditworthiness of borrowers.

However, in contrast to these models, the XGB model (Table n.15) attaches greater significance to "ROI (TWD)," suggesting that potential profits of investments contributes directly to default prediction. Furthermore, "ROI (TWD)" stands out as an innovative aspect of this model that considers financial gains from investments as one factor influencing default behavior of borrowers; its inclusion highlights this unique characteristic and sets itself apart from RF and ADA (Table n.16) models which give investment returns comparatively lesser importance than this variable suggests.

Further, the RF model gives relatively higher weight to "CLV", or Customer Lifetime Value; an economic indicator used to capture long-term profitability associated with each borrower and predict defaults. Notably, however, CLV was not included as an important predictor in XGB model due to differing modeling approaches and weight given for this economic variable.

Finally, it is noteworthy to observe how Random Forest, AdaBoost and Gradient Boosting models all assign equal ranking importance for economic predicting variables despite having differing feature importance scores.

Overall, Monthly Income Estimate (TWD) consistently comes out as key; however, different modeling strategies and the predictive capacity attributed to various economic factors vary considerably, placing emphasis on "Payment Ratio", ROI (TWD), or CLV depending on which economic variables come under focus.

# Master's Thesis Flavio Ungherini

Analysis of feature importances across several Machine Learning models has highlighted both similarities and variances in how different variables are assigned predictive importance, with selection and interpretation of predictors differing across models indicating their individual contribution towards predicting target variables. Such findings highlight the need to carefully select features, interpret models accurately, compare and contrast results of various models to gain an in-depth knowledge of predictive dynamics for customer intelligence analysis in the financial services industry.

## COMPARISON BETWEEN THE BEST-PERFORMING MODELS OF THE TWO EXECUTED STUDIES

The best predictive model identified in the first part of the analysis performed is the Random Forest, with a set of predictor variables consisting of the following features: 'USAGE_IDX', 'pay_delay', REM_DEBT', 'AGE', 'LIMIT_BAL', 'EDUCATION', 'SEX', 'MARRIAGE_-dummy_1', 'MARRIAGE_dummy_2'.

In contrast, the best classification model identified in the second part of the study is, on a general scale, AdaBoost. In particular, this model performs better classification of defaulting customers when variables of an exclusively economic nature are used and when an extended range of hyperparameters is used.

In both cases, the sampling technique used is SMOTE; another point of contact between the two models is that an inner cross validation with external hold out is used for both.

Importantly, the Random Forest identified as the best performing model in the first part of the study was trained by a number of observations equal to 80% of 30,000, or 24,000 records.

In contrast, the AdaBoost that provides the best performance in the second part of the study carried out is trained on the 80% of 6,000 observations, or 4,800.

Since SMOTE is employed in both models, there is actually an increase in the number of training data in both cases. However, a substantial difference persists between the two about the number of data used to train the models.

**114**

So, dwelling on this mere observation, the likelihood that the Random Forest above performs better than AdaBoost is very high, the latter being trained by considering considerably fewer data.

Looking, however, at the evaluation metrics considered to determine which was the better classification model, we come to a different conclusion from the one just assumed.

The two models will be compared by taking only 'Recall for class 1' as the reference metric. In this case, it is not necessary to use the 'Weighted Performance' metric used previously, as the latter was developed with the intention of comparing the performance of a large number of classification models.

In this case, however, since the number of machine learning models being compared is two, then the latter can be compared on the basis of what is expressed by the respective 'Recall for class 1,' which is by the way the metric that weighs the most in the calculation of 'Weighted Performance.'

The Random Forest classification report in question is as follows:

Table n.18: Random Forest Classification Report

|  | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| 0 | 0.87 | 0.82 | 0.85 | 4683 |
| 1 | 0.48 | 0.57 | 0.52 | 1317 |
|  |  |  |  |  |
| Accuracy |  |  | 0.77 | 6000 |
| Macro Avg | 0.67 | 0.70 | 0.68 | 6000 |
| Weighted Avg | 0.79 | 0.77 | 0.78 | 6000 |

In this case, Table n.18 indicates that the 'Recall for class 1' equals 0.57.

**115**

In contrast, as noted earlier, the 'Recall for class 1' related to the AdaBoost with exclusively economic prediction variables is equal to 0.62.

So, looking at this metric as the one of greatest interest and relevance in classifying defaulting customers, it can be established that the AdaBoost with exclusively economic prediction variables performs better than the Random Forest identified as the best performing model in the first part of the study.

Interestingly, the economic variables used as predicting features by the AdaBoost under consideration are all variables created on top at the original dataset, that is, developed through the combination of variables that originally existed in the dataset and through the retrieval of data from external sources, as happened for 'Annual Revenue (TWD)'.

In conclusion, it is possible to establish how the variables 'CLV', 'Monthly income Estimate (TWD)', 'Debt-to-Income Ratio', 'Payment Ratio' and 'ROI (TWD)' are the discriminative variables that guarantee additional predictive value on the analysis conducted.

Specifically, the variables in question were created on-top of the original dataset, thus suggesting how the inclusion and creation of new 'economic' variables guarantee an improvement in the performance that the machine learning models employed in this study have in identifying bank's customers who will go into default.

# Chapter 7: Conclusions, Limits and Future Research

**MANAGERIAL IMPLICATIONS AND CONCLUSIONS**

This study's results have significant managerial ramifications for banking industries, specifically regarding credit risk evaluation and default prediction. By utilizing an expansive dataset combining socio-demographic and economic variables, banks can gain valuable insights into customers' behavior while improving their ability to identify individuals likely to default on loans. It highlights both economic factors in addition to traditional socio-demographic variables as well as potential advantages feature engineering can bring for improving machine learning performance models.

Economic indicators, like income levels, employment stability and debt-to-income ratios provide banks with a comprehensive view of customers' finances. Armed with this data, banks are better informed to make decisions regarding loan approvals, interest rates and credit limits that reduce risks more efficiently. By monitoring economic variables closely banks can detect customers experiencing difficulties financially or changes that impact them economically quickly so proactive measures may be taken against potential default risks that arise as a result of such challenges or changes.

AdaBoost was named as the highest performing model using exclusively newly created economic variables, demonstrating its success as feature engineering contributes to credit risk modeling. By adding more variables that capture economic information that may improve accuracy of classification of defaulting customers. Therefore, banks should allocate sufficient resources toward developing and refining variables tailored to their customer base to enhance predictive capabilities and effectively forecast default risks.

Furthermore, this study emphasizes the significance of ongoing monitoring and updating of credit risk models. As economic conditions shift over time, banks must regularly reassess and

adapt their models. By regularly including new economic variables into models or training them with new variables to stay abreast of emerging risks while strengthening risk management strategies altogether - banks will remain well equipped to deal with emerging market dynamics or changing customer profiles.

Furthermore, this study can serve as a valuable resource for other financial institutions looking to strengthen their credit risk assessment practices. By adopting similar methodologies and including economic variables tailored specifically for their customer base, banks can improve predictive capabilities and make more informed decisions regarding loan approvals and risk mitigation. Likewise, its findings serve as a roadmap for financial institutions seeking to enhance credit risk modeling techniques as part of an industry wide improvement in risk management practices.

Overall, this study underscores the significance of including newly generated economic variables through feature engineering into credit risk modeling. Leveraging these variables and applying advanced machine learning techniques such as AdaBoost allows banks to improve their ability to identify customers at risk of default and develop more effective risk management strategies leading to improved loan portfolio quality, reduced credit losses and ultimately an improvement in financial performance. Its application extends far beyond its focus bank alone - its implications could apply broadly among financial institutions seeking to enhance their credit risk analysis abilities.

**LIMITATIONS AND FUTURE RESEARCH**

Although this study made substantial advancements to credit risk assessment and default prediction, several limitations and areas for future research should be recognized. Neural networks have proven highly useful classification techniques across numerous domains making their absence notable. Therefore, future studies could explore using neural network models to

**118**

enhance identification of customers likely to default further, leading to greater predictive accuracy as well as deeper insight into complex patterns within a dataset.

Despite attempts were made to approximate missing data such as Annual Revenue via related variables, approximating them with relative variables created from scratch, estimates may fail to capture the true variables' values accurately, resulting in potential bias and lack of precision while adding new variables to the original dataset.

Moreover, this study examined data collected from one Taiwanese bank only; thus the findings may not apply directly to other banks or regions due to significant variations in customer socio-demographics and economic characteristics between banking systems and cultural contexts. Therefore, future research should involve replicating this study with datasets from multiple banks/countries to validate findings as well as explore any variations in predictive performance of models.

As with any quantitative study, there may be limitations regarding data quality and availability. The dataset used may have some inaccuracies that skew results and conclusions negatively; future research should address these gaps by using comprehensive, high-quality datasets from multiple sources to gain a fuller picture of customer creditworthiness.

As this research has made significant advances to credit risk assessment and default prediction, several limitations and areas for future study have become clear. Investigating neural network models further, validating findings across diverse datasets and banking systems and addressing data quality issues all offer promising avenues for investigation - by addressing these constraints researchers can further advance credit risk modeling towards more precise risk management practices in banking industries worldwide.

# Master's Thesis Flavio Ungherini

# Extended Abstract

Over many decades, banks have used technological innovation to reinvent customer interactions. ATMs were introduced during the 1960s; electronic card payments came later; in 2000s online banking gained widespread use before mobile-based "banking on the go" became widespread by 2010s.

AI-fueled digital age has flourished due to a host of factors: falling data storage costs, better access and connectivity for all, rapid advances in AI technologies enabling greater automation and efficiency gains across various tasks when implemented effectively.

Forbes conducted research in 2021 which concluded that daily internet user activity produced 2.5 quintillion bytes of data daily; each person creating approximately 1.7 megabytes per second [1].

Covid-19's pandemic outbreak resulted in widespread demand for financial services to increase significantly online. Banks made massive investments into Artificial Intelligence adoption as they strive to meet customers' evolving requirements and satisfy customers.

KPMG conducted a survey in 2021 which demonstrated that 83% of players operating within the financial industry used artificial intelligence of some form, placing financial services as second only behind manufacturing on this KPI measure.

Financial services players' increasing adoption of AI is further evidenced when looking at industry growth rates: an estimated 31.5% annual average rate acceleration in AI investments by 2027 could reach up to USD 54.7 billion [2].

Artificial Intelligence technologies can significantly increase revenues and lower fixed and variable costs through greater customer customization, higher automation rates, reduced error rates and better resource exploitation; new unrealized opportunities may emerge through improved data processing capacities that yield meaningful insight from vast quantities of information [3].

# Master's Thesis Flavio Ungherini

As banks seek to minimize financial losses, increase revenues, deliver superior services to their customers and boost satisfaction levels, they have increasingly turned to Artificial Intelligence and Machine Learning techniques for performing several strategic, operational and daily functions such as: credit scoring, fraud detection and prevention, identity verification, Customer Lifetime Value analysis and anti-churn policies creation, risk management and AML (Anti-Money Laundering).

Banks increasingly depend on AI technologies for customer service enhancements, automation of routine tasks, fraud detection faster and risk management more effectively. According to an Accenture research, banks' investments in AI should increase by an estimated 30-35% annually over three years up until 2022 [4]. This amount approximately equals $5.6 billion.

Artificial Intelligence adoption in the banking industry became particularly relevant and fundamental after the 2008 financial crisis, with the Basel 2 Accord update required that banks are mandated to develop and adopt credit scoring and profiling systems which enable them to predict in advance and track customer creditworthiness and repayment behavior once economically bound through contractual relationships.

Before granting credit to an applicant requesting a loan or a mortgage, banks must ensure that customers can repay any sum borrowed in line with contractually stipulated terms without creating additional difficulties for themselves as lenders: Artificial Intelligence and Machine Learning play a key role in helping banks to carry out those kind of assessments and considerations, providing the decision makers with valuable insights about the banks' existing and potential customers.

AI and ML applications in the banking industry are not limited to this case: they encompass numerous implementations and include its exploitation to assess Customer Lifetime Value (CLV), an integral metric for banks that measures customer profitability over their lifespan. AI allows them to quickly process large volumes of customer data while making accurate predictions regarding future customer behavior - helping identify which customers might be more profitable, what products or services to offer them and tailor marketing campaigns accordingly.

Predictive modeling is one way banks use AI to assess customer lifetime value (CLV). Predictive models use statistical algorithms to examine customer data and predict future behavior. Predictive modeling also allows banks to predict which customers may leave and provide proactive measures for keeping those customers .

Moreover, banks can leverage customer lifetime value (CLV) data to tailor offers specifically to each of their customers based on transaction history, purchase behavior and other variables : this allows banks to offer exclusive travel rewards or discounts when planning future journeys.

Banks can utilize customer lifetime value (CLV) scores to identify their most valued customers and deliver personalized customer service. Customers with high CLV may receive priority service or even be assigned a dedicated account manager.

Despite Machine Learning adoption in banking industry represents a unique opportunity, concerns and legal constraints about its usage have been raised over the years.

Specifically, AI and ML in banking have raised concerns over whether these technologies threaten human jobs, since they could lead to job displacements for human employees. Although, these technologies are proven to enhance productivity and decision-making by helping employees focus more on tasks requiring human judgment, creativity, or problem-solving skills .

Additionally, AI and ML use in banking is predicted to open up job opportunities requiring new skills such as data analysis, coding and digital literacy. According to research from McKinsey Global Institute, these new jobs could create 20-50 million global job openings by 2030 [5]. For banking industry employees to meet these demands of an evolving job market, they will require investment in training and up-skilling programs designed to prepare employees for these jobs.

Artificial Intelligence adoption in the banking industry also poses challenges regarding compliance with data protection regulations, such as GDPR in Europe. AI and ML can assist banks in GDPR compliance by automating data classification, discovery, retention, and subject access requests. At the same time, this regulation also sets constraints regarding the nature of the data that can be processed by banks (e.g. sensitive data can't be processed) [6].

# Master's Thesis Flavio Ungherini

The General Data Protection Regulation can be certainly considered the most advanced, complete and structured existing regulation worldwide, but it has applicability exclusively in Europe.

As a consequence, other countries adopt different data protection regulations: for example, in Asia AI is used for data validation and security measures in banking industries, adhering to local regulations such as PDPA in Singapore, PDPO in Hong Kong, and PDPA in Malaysia [7].

Overall, Machine learning algorithms can comb through large volumes of customer data to use machine learning to develop predictive models that estimate the likelihood of borrower default on loans. By taking into account various data points such as credit history, income, employment status and other financial considerations,, AI algorithms can provide more accurate credit risk assessments than traditional credit scoring methods and thus better manage loan portfolios while decreasing exposure to NPLs (Non-Performing Loans) [8]. This point helps banks effectively manage loan portfolios while decreasing exposure to non-performing loans.

Artificial Intelligence can assist banks in more efficiently managing their collections and recovery processes by automating the processes for identifying delinquent borrowers, predicting their willingness and ability to pay, and suggesting recovery solutions . By streamlining collections and recovery processes using AI solutions, banks can reduce the time and costs associated with managing NPLs.

One of the primary functions of banks is to provide loans to individuals and corporations to fulfill their credit needs. However, sometimes borrowers default on these loans, leading to substantial financial losses for banks worldwide. Many large and prominent banks worldwide have experienced losses due to customer defaults, which has seriously affected their financial stability and reputation.

Relevant historical cases in which financial losses suffered from banks could have been avoided are the following ones: The 2008 Global Financial Crisis [9], The Asian Financial Crisis of 1997 [10] and The European Financial Debt Crisis of 2009-2012 [11].

# Master's Thesis Flavio Ungherini

In the case of the 2008 global crisis mentioned above, AI could have detected and flagged high levels of risk associated with subprime mortgage markets - one of the primary sources of distress in our economy - using machine learning algorithms to analyze large volumes of data from various sources including loan applications, credit reports and housing market trends. By recognizing patterns of risky lending practices, such as loan approval for individuals with poor credit histories or low documentation loans using low documentation loans, these AI systems could have helped financial institutions avoid similar practices.

In order to address these challenges and to prevent customers defaults negative consequences, the Taiwanese Credit Risk Dataset has been extensively studied in credit risk analysis, due to its realistic representation of credit card usage patterns and payment behaviors in real-life settings.

Before discussing the dataset structure and the proposed study methodology, since it contains data about bank's customers credit card usage, it is essential to introduce how the credit card usage behavior in Taiwan has evolved since their introduction in the country.

Credit cards first made their debut in Taiwan during the late 80s, initially used only for major purchases like appliances. Since then however, their usage has increased to include everyday expenses like groceries and dining out due to credit card convenience and online shopping: this shift can be attributed to rewards programs providing quick payments with incentives like cash back or airline miles as part of convenient payment schemes that lend itself well for everyday spending needs. Credit card spending reached NT$2.06 trillion in 2019 (a 10.6% year-on-year growth) [12].

After Taiwan's 2005 financial crisis, new regulations and restrictions were placed upon credit card issuers in order to reduce debt and boost consumer creditworthiness. Measures included debt repayment programs, transparent disclosure of interest rates and fees and limits on the number of cards one could possess, as well as limits on personal use of multiple cards at one time. As a result of these regulations there was significant shift in spending patterns after crisis with increased use for smaller day-to-day purchases rather than large transactions such as those conducted online or via mobile phones.

Taiwanese debt repayment programs were highly successful, as average per-borrower debt dropped 43% between 2005 and 2019 [13]. This reduction can be attributed to programs which allowed customers to repay debt more rapidly.

After having clarified which has been the evolution of the credit card behavior adopted by banks' customers in Taiwan, the structure of dataset object of study can be introduced: it comprises information collected between April and September 2005 by a Taiwanese bank with 23 features that include demographic data, credit card usage patterns and payment histories.

In that period, banks in Taiwan employed aggressive marketing communication strategies to lure customers and increase credit card usage. Such approaches included offering low-interest rates, cash rebates and other incentives to encourage more customers to use credit cards; as a result, credit card usage rapidly soared until 2006 when their usage peaked; however, this led to an increase in banks' customers delinquency and default rates, since the customer bases attracted by banks was particularly fragile on an economic stability standpoint.

Despite this dataset has been extensively studied in the academic world with the aim to develop and compare several machine learning models capable of predicting customers' default, none of the conducted researches has proposed a method based on the creation of a wide amount of new variables through information retrieving from external sources and feature engineering.

This study has the objective of integrating new discriminative features in the dataset at stake, proposing a methodology that has the purpose to answer to the following research question:

'*Do the additional variables collected and computed in the dataset provide further discriminative insights in order to improve the predictive performance of machine learning models for identifying defaulting bank's customers?*'

# Master's Thesis Flavio Ungherini

This study involves multiple steps. First, the original dataset is reviewed: an exploratory data analysis (EDA) is carried out in order to extract and interpret information from variables and observations; graphs are utilized as visual displays of distribution of variables and their relationships between observations.

Next, additional variables are created through collecting external source data in order to enhance machine learning models' performance. Feature engineering techniques are then utilized in order to combine existing variables together into new ones while still preserving information quality while preventing redundancy; summary variables are introduced as necessary to avoid redundancy.

Subsequently, variables for standardization to train multiple machine learning models are selected and standardized before sampling techniques such as oversampling ( including SMOTE ) and under-sampling are implemented and explained further in subsequent sections. Finally, predictive models with optimal predictive power and sampling strategies for an analysis task are assessed for evaluation purposes.

After evaluation, the dataset is expanded with socio-demographic and economic variables created from existing metrics and features as well as external sources (mainly public databases) with which information has been gleaned from. EDA is performed on this expanded dataset in order to evaluate their contribution and improvement, while correlation analysis identifies three sets of predicting variables having 'economic', 'socio-demographic' and both 'economic and socio-demographic ' nature.

Four probabilistic classification models (Random Forest, XGBoost, AdaBoost and Gradient Boosting) are trained on identified sets of predicting variables as training datasets for evaluation; Recall for class 1, F1-score for class 1 and AUC are focused upon while feature importance analysis helps uncover significant predictors for defaulting customers. Finally, initial models obtained are compared with each other to compare and interpret results.

# Master's Thesis Flavio Ungherini

Overall, AdaBoost emerges as the best-performing machine learning model among the ones considered in this study: specifically, the group of predicting variables that yields the best classification performance is the one exclusively composed by features of 'Economic' nature.

Comparing the performance of this model with the best classification model using only predicting features from the original dataset, it can be seen that the integration of these new variables within the dataset object of study was decisive in improving the performance of classification and identification of defaulting customers for the bank.

Specifically, the newly created 'economic' features identified as discriminative are the following: 'Monthly Income Estimate (TWD)', 'Payment Ratio', 'CLV', 'Debt-to-Income Ratio' and 'ROI (TWD)', in which CLV stands for Customer Lifetime Value and ROI (TWD) stands for the Return on Investment in Taiwanese dollars.

Next, an analysis is performed to determine whether there are similarities and discrepancies between the feature importances established by the four machine learning models used in this study.

This comparison shows that, among the predicting variables of both economic and socio-demographic nature combined, the economic features generally tend to assume a more pronounced importance than the socio-demographic features.

Among the economic variables, however, the variables that are most often identified as most important for the four probabilistic models in question are, respectively: 'Monthly income estimate (TWD),' 'Payment Ratio,' and 'CLV,' with the first of the three mentioned exhibiting a much greater degree of importance than all the other features considered.

The study's results provide managerial insights for banks, especially with regard to credit risk evaluation and default prediction. Banks can utilize comprehensive datasets encompassing socio-demographic and economic variables to gain greater insights into customer behaviors, as well as improve their ability to recognize individuals likely to default on loans. Economic indicators, like income levels, employment stability and debt-to-income ratios provide banks with a comprehensive view of customers' finances and enable better decision-making regarding loan approvals, interest rates, credit limits and loan renewals. Banks also benefit by clo-

**128**

sely tracking these indicators to detect financial challenges or changes quickly affecting customers and taking precautionary steps against default risks quickly and proactively.

Hence, banks should devote resources toward developing variables tailored specifically for their customer base to increase predictive capabilities and efficiently forecast default risks.

The default risk prediction helps banks to have a more thorough understanding and control on their cash flows, having the possibility to comprehend in advance how to stem the negative financial effects arising from the missed payments of their customers and, most importantly, to implement strategies to dismiss deteriorated credits.

Despite the numerous managerial implications and insights suggested by this study's results, it acknowledges numerous limitations and areas for future investigation. Neural network models could provide significant improvements in default prediction accuracy while increasing predictive accuracy overall. Furthermore, approximating missing data with relative variables may introduce bias or imprecision highlighting a need for more precise methods of approximating missing data.

Due to different customer characteristics and cultural contexts, its findings may not apply directly across banks or regions. Therefore, future research should utilize datasets from multiple banks or countries so as to validate findings and assess model performance across diverse contexts.

Lastly, data quality and availability could pose limitations that impede accuracy in results, so future research must utilize comprehensive datasets from multiple sources in order to gain a fuller picture of customer creditworthiness.

References:

[1] https://www.forbes.com/sites/forbestechcouncil/2021/08/02/understanding-generation-data/?sh=3a3bb7b136b7

[2] https://info.kpmg.us/news-perspectives/technology-innovation/thriving-in-an-ai-world/ai-adoption-accelerated-during-pandemic.html

[3]https://www.mckinsey.de/~/media/McKinsey/Industries/Financial%20Services/Our%20Insights/
AI%20bank%20of%20the%20future%20Can%20banks%20meet%20the%20AI%20challenge/AI-bank-of-the-future-Can-banks-meet-the-AI-challenge.pdf

[4] https://www.accenture.com/_acnmedia/PDF-112/Accenture-Banking-on-AI.pdf#zoom=50

[5] McKinsey Global Institute, "Skill Shift: Automation and the Future of the Workforce"
https://www.mckinsey.com/featured-insights/future-of-work/skill-shift-automation-and-the-future-of-the-workforce.

[6] Kshetri, N. (2020). Using Artificial Intelligence to Meet GDPR Requirements. Communications of the ACM, 63(4), 56-63.

[7] Yap, W. Y., Tan, C. H., & Ong, T. S. (2021). Artificial intelligence in data protection: A review of GDPR, PDPA and CCPA. Journal of Information Privacy and Security, 17(2), 61-75.

[8] Blanchard, O. (2019). Macroeconomics (7th ed.). Pearson.

[9] Amato, J. D., & Fantacci, L. (2015). The financial crisis: lessons and challenges. Springer

[10] Blommestein, H., & Spencer, P. (2001). Corporate Governance and Banking Stability in East Asia. OECD Journal: Financial Market Trends, 2001(1), 89-104.

[11] Gros, D. (2012). The European Financial Stability Facility: What it is and what it does. CEPS Policy Brief, (270).

[12] Financial Supervisory Commission. (2020). Credit card usage survey report.

[13] Taiwan Credit Information Center. (2020). Credit information annual report.