









Nydia Remolina Assistant Professor of Law, Singapore Management University







510

213 RB (1999) - 1990

12/101

TIME

THE GOVERNMENT OF THE GRAND-DUCHY OF LUXEMBOURG Ministry of Finance



BUS UNI COMOS



# CONTENT

1. Conceptual Foundations

2. The AI Governance Ecosystem: Frameworks, policies and regulations

3. Al Governance in Singapore and Beyond

4. Gaps in the Al Governance Ecosystem





# **CONCEPTUAL FOUNDATIONS**

1. The "Al" Ecosystem







# **CONCEPTUAL FOUNDATIONS**

## 2. The Regulatory Definition of AI



An AI system is a machine-based system that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments. Different AI systems in their levels of autonomy and adaptiveness after deployment



#### 2. The Regulatory Definition of AI



• Article 3, Al Act:

'Al system' means a machine-based system that is designed to operate with **varying levels of autonomy** and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, **infers**, **from the input it receives**, **how to generate outputs** such as **predictions**, **content**, **recommendations**, **or decisions** that can influence physical or virtual environments.

# **CONCEPTUAL FOUNDATIONS**

#### 2. The AI Value Chain

Stage	Key Players	Responsibilities
1. Al Development	Developers, Al Providers	<ul><li>Develop AI models</li><li>Ensure compliance</li></ul>
2. Al Deployment	Companies, Financial Institutions	<ul> <li>integrate Al into products/services and offer them to users</li> <li>Ensure regulatory alignment</li> </ul>
3. AI Use	End Users (Businesses, Consumers)	<ul> <li>Interact with AI</li> <li>Use AI systems but do not develop or modify them</li> </ul>
4. Al Oversight & Regulation	Regulators, Authorities	<ul> <li>Monitor compliance</li> <li>Enforce penalties - Classify Al risks</li> </ul>



# THE AI GOVERNANCE ECOSYSTEM

**AI governance:** frameworks, policies, and regulations that guide the responsible development, deployment, and use of artificial intelligence (AI). It ensures AI systems are ethical, transparent, accountable, and aligned with societal values while minimizing risks. (Gasser and Almeida, 2017; Wachter, Mittelstadt and Floridi, 2017).

# Market- DrivenState-DrivenRights-Driven• US• China• EU

Source: Bradford, 2023



# **REGULATORY AND POLICY RESPONSES**



Year	Initiative
2019	National Al Strategy
2020	Model Al Governance Framework
2023	AI Verify Governance and Testing Framework
2023	Guidelines on Data Privacy in Al
2023	Discussion Paper on GenAl
2023	Veritas Toolkit
2023	Project Mindforge
2024	MAS paper on Generative Al Risks
2024	Project Moonshot
2024	Model Al Governance Framework for Generative Al







## **PDPC Model AI governance (Singapore)**

- Model AI Governance Framework (2020) 2<sup>nd</sup> Ed)
  - Guidance re key ethical and governance issues for AI solutions
  - Centred on human-centricity, good data accountability practices, and creating open and transparent communication

#### **Guiding Principles**



Decisions made by AI should be **EXPLAINABLE, TRANSPARENT & FAIR** 



#### From Principles to Practice





making

Appropriate

involvement

of harm to

individuals

Minimise the risk

degree of human





your

Measures

- organisation SOPs to monitor
- and manage risks
- Staff training

Determining the Level of Human Involvement in Al-





Risk-based

Minimise bias in

approach to

explainability,

robustness and

regular tuning

data and model

measures such as



Interaction and Communication

- Make Al policies known to users
- Allow users to provide feedback, if possible Make
- communications easy to understand



## **Generative AI (Singapore)**

#### Discussion paper on Gen AI (7 Jun'23)

 6 key risks – (i) Mistakes & Hallucination, (ii) Privacy & Confidentiality, (iii) Disinformation, Toxicity & Cyber threats at scale, (iv) copyright erosion, (v) embedded biases, and (vi) value misalignment

#### Proposed Model AI Governance Framework for Gen AI (16 Jan'24)

- Complements traditional model Al governance framework.
- Aims to:
  - provide systematic and balanced approach to foster trusted AI ecosystem; and
  - Address Gen AI concerns while facilitating innovation.





## AI and Data Analytics in Finance: Overview of MAS' role

- MAS plays active role in artificial intelligence ("AI") and Data Analytics use in financial Industry
  - Evangelist (SFF, Grants, Accelerator)
  - Governor/Supervisor (FEAT, Veritas, Digital Advisory Services)
  - Adopter (COSMIC, SGFinDex, Project Ellipse)
- Selected Key developments
  - Prototype of "Project Ellipse" launched on BIS Open Tech platform in Mar'22
  - Publication of Veritas whitepapers on assessment methodologies for FEAT principles in Feb'22
  - S\$180m National Artificial Intelligence (AI) Programme in Finance in Nov'21
  - Inclusion of Investment Holding Data in SGFinDex in Nov'21 following Dec'20 launch
  - Announcement of COSMIC data sharing platform for ML,TF, and PF data in Oct'21



#### **Guiding Principles in Finance**

• FEAT

AS Monetary Authority of Singapore

- AIDA use cases should consider FEAT (Fairness, Ethics, Accountability and Transparency) principles
  - Fairness: Justifiability, Accuracy and Bias
  - Ethics: At least equivalent and aligned to human decisions
  - Accountability: Internal, and External
  - Transparency: Communications to data subjects, Explainability re mechanics and/or Consequences
- Calibration to be driven by materiality, e.g. role AIDA play in decision-making, complexity of AIDA model, etc
- Impacts of AIDA also to be considered, e.g. stakeholder, monetary, financial, regulatory impacts, etc

Fairness	<ul> <li>P1: Individuals or groups of individuals are not systematically disadvantaged through AIDA-driven decisions, unless these decisions can be justified.</li> <li>P2: Use of personal attributes as input factors for AIDA-driven decisions is justified.</li> <li>P3: Data and models used for AIDA-driven decisions are regularly reviewed and validated for accuracy and relevance, and to minimise unintentional bias.</li> <li>P4: AIDA-driven decisions are regularly reviewed so that models behave as designed and intended.</li> </ul>
Ethics	<ul> <li>P5: Use of AIDA is aligned with the firm's ethical standards, values and codes of conduct.</li> <li>P6: AIDA-driven decisions are held to at least the same ethical standards as human-driven decisions.</li> </ul>
Accountability	<ul> <li>P7: Use of AIDA in AIDA-driven decision-making is approved by an appropriate internal authority.</li> <li>P8: Firms using AIDA are accountable for both internally developed and externally sourced AIDA models.</li> <li>P9: Firms using AIDA proactively raise management and Board awareness of their use of AIDA.</li> <li>P10: Data subjects are provided with channels to enquire about, submit appeals for and request reviews of AIDA-driven decisions that affect them.</li> <li>P11: Verified and relevant supplementary data provided by data subjects are taken into account when performing a review of AIDA-driven decisions.</li> </ul>
Transparency	<ul> <li>P12: To increase public confidence, use of AIDA is proactively disclosed to data subjects as part of general communication.</li> <li>P13: Data subjects are provided, upon request, clear explanations on what data is used to make AIDA-driven decisions about the data subject and how the data affects the decision.</li> <li>P14: Data subjects are provided, upon request, clear explanations on the consequences that AIDA-driven decisions may have on them.</li> </ul>



#### **Generative Al**

- Project Mindforge Gen AI Risk Framework for Financial Sector (15 Nov'23)
  - Whitepaper detailing the risk framework
  - GenAI risk framework, with 7 risk dimensions: (a) Accountability and Governance, (b) Monitoring and Stability, (c) Transparency and Explainability, (d) Fairness and Bias, (e) Legal and Regulatory, (f) Ethics and Impact, and (g) Cyber and Data Security.
- Phase 1: The report found that the existing FEAT principles remained broadly relevant, though enhancements to existing principles and additional considerations to be taken into account were further recommended. Steps and additional guardrails for financial sector players to mitigate generative AI risks were also recommended.



List of MindForge Consortium Members



## **Generative AI – Key element in the Whitepaper**

#### Project Mindforge – Gen AI Risk Framework for Financial Sector (15 Nov'23)

- Guardrails to mitigate Gen-AI related risks: human-in-the-loop, due diligence on third-party generated AI systems, value alignment ex-post control, multidisciplinary approach, use the AI Verify Foundation.
- The Whitepaper recognizes new complexities into vendor-FI relationships.
- The Risks are assessed in different lifecycle stages:
  - -System Context and Desing
  - -Data Acquisition
  - -Model Onboarding and Build
  - -Deployment and Monitoring
  - -Model Use and Output.







#### Sandboxes

- Data protection:
- Information Commissioner's Office (UK)
- Norwegian Data Protection Agency
- Superintendency of Industry and Commerce (Colombia)
- CNIL Sandbox Initiative for Health Data and Privacy-by-Design (France)
- IMDA-PDPC Data Regulatory Sandbox (Singapore)
- Financial regulatory sandbox:
  - FCA: 4.5%
    - Cohort 1: 0/18
    - Cohort 2: 2/24
    - Cohort 3: 0/18
    - Cohort 4: 1/29
    - Cohort 5: 1/29
    - Cohort 6: 3/22
    - Cohort 7:0/13
  - MAS: 0%
    - 0 before 2018

#### Gen AI Evaluation Sandbox for Trusted AI (31 Oct'23)

Anchored by the catalogue, a compilation of technical testing tools

Provides baseline of evaluation tools of Gen AI products.



#### The problem of discrimination in AI credit scoring

#### Model 1: OLS, using observations 1-10745 Dependent variable: TARGET

	Coefficient	Std. Error	t-ratio	p-value			
const	0.0375632	0.0857195	0.4382	0.6612			
GENDER	0.0123190	0.00476760	2.584	0.0098	***		
CREDIT_AMOUNT	-2.18188e-08	7.52012e-09	-2.901	0.0037	***		
AGE OF CAR	0.000512744	0.000188460	2.721	0.0065	***		
AGE OF CLIENT	-2.33543e-06	7.57000e-07	-3.085	0.0020	***		
EDUCATION_TYPE	-0.0223241	0.00486574	-4.588	< 0.0001	***		
REGION_RATING_SCORE	0.0237853	0.00434910	5.469	< 0.0001	***		
NUMBER_OF_CHILDREN	-0.00766924	0.00300480	-2.552	0.0107	**		
LOG_ANNUITY_AMOUNT	0.0345452	0.00657663	5.253	< 0.0001	***		
LOG_DAYS_EMPLOYED	-0.0131367	0.00235096	-5.588	< 0.0001	***		
LOG_INCOME	-0.0188018	0.00596179	-3.154	0.0016	***		
Mean dependent var	0.061145	S.D. dependent	var	0.239607			
Sum squared resid	606.5705	S.E. of regressi	on	0.237717			
R-squared	0.016629	0.016629 Adjusted R-squared					
F(10, 10734)	18.15168	P-value(F)		2.27e-33			
Log-likelihood	196.0840	Akaike criterion	n	-370.1680			
Schwarz criterion	-290.0638	Hannan-Quinn		-343.1498			

#### Model 2: OLS, using observations 1-10745 Dependent variable: TARGET

	Coefficient	Std. Error	t-ratio	p-value	
const	-0.00475901	0.0848536	-0.05608	0.9553	
CREDIT_AMOUNT	-2.56982e-08	7.46237e-09	-3.444	0.0006	***
AGE_OF_CAR	0.000518673	0.000188582	2.750	0.0060	***
EDUCATION_TYPE	-0.0237196	0.00478362	-4.959	< 0.0001	***
REGION_RATING_SCORE	0.0243988	0.00434767	5.612	< 0.0001	***
NUMBER_OF_CHILDREN	-0.00617184	0.00296081	-2.085	0.0371	**
LOG_ANNUITY_AMOUNT	0.0345107	0.00658071	5.244	< 0.0001	***
LOG_DAYS_EMPLOYED	-0.0152294	0.00228051	-6.678	< 0.0001	***
LOG_INCOME	-0.0161516	0.00580110	-2.784	0.0054	***
Mean dependent var	0.061145	S.D. dependent v	ar	0.239607	
Sum squared resid	607.5180	S.E. of regression	n	0.237880	
R-squared	0.015093	0.015093 Adjusted R-squared			
F(8, 10736)	20.56542	P-value(F)		3.23e-31	
Log-likelihood	187.6984	Akaike criterion		-357.3968	
Schwarz criterion	-291.8570	Hannan-Quinn		-335.2910	

- Holding all other variables constant, being male, on average, increases the likelihood of default by 1.23%.
- Effects of removing 'GENDER': the negative correlation between 'GENDER' and 'EDUCATION\_TYPE' causes the coefficient of 'EDUCATION\_TYPE' to decrease when 'GENDER' is excluded.
- The "proxy" variables have significant correlations with other variables (eg. Education and Income according to the correlation coefficient matrix), making it virtually impossible to eradicate the gender discrimination through excluding variables.
- The exclusion of all "protected" variables and their proxies in the model is ineffective in eliminating the discrimination.



#### **Metrics - Veritas (Singapore)**

- Collaboration between MAS and FIs to develop framework for responsible use of AI
- 5 Whitepapers re assessment methodology for FEAT principles
- Open-source toolkit for automation of the fairness metrics assessment, with visualisation and plug-ins to integrate with FI's IT systems
  - In credit scoring: Veritas provides a qualitative and a mathematical approach to fairness.

Installation	
The easiest way to install veritastool is to download it from <b>PyPI</b> . It's going to install the library itself and its prerequisites as well. It is suggested to create virtual environment with requirements.txt file first.	
pip install veritastool	C
Then, you will be able to import the library and use its functionalities. Before we do that, we can run a test fu on our sample datasets to see if our codes are performing as expected.	nction
<pre>from veritastool.util.utility import test_function_cs test_function_cs()</pre>	C
Output:	
Evaluate: 100%	
Evaluation of credit scoring performed normally	
Initialization	
You can now import the custom library that you would to use for diagnosis. In this example we will use the Co Scoring custom library.	redit
<pre>from veritastool.model.modelwrapper import ModelWrapper from veritastool.model.model_container import ModelContainer from veritastool.usecases.credit_scoring import CreditScoring</pre>	Q

ш	PILLO	11 11.4		$\mathbf{D}\mathbf{D}$	$\mathbf{D}$	<u>a</u>		$\mathbf{O}$	<b>d</b>		22												
Yo	ou can a	lso to	ggle th	e wic	lget	to vie	ew ye	our re	esult	s in	a int	eract	tive v	/isua	alizat	ion	forma	ıt.					
<pre>cre_sco_obj.evaluate(visualize = True)</pre>													C										
0	utput:																						
	Model Type	: Classific	ation				Sam	ole Weig	ght	Rejec	tion Inf	erence							Model	Name: Cred	it_Sci	oring	
	Protected	Feature:	SEX (pri	vileged	group :	= [[1]])					¥ F	Priority: I	Benefit		Impact	Nor	nal	Concern: E	ligible	Type: Diffe	rence	e	
	Fairness															Pe	rforman	се					
	Metric							Assess	sment								Assessme	ent					
		Equa	І Орро	pportunity Fair Accu									curac	racy									
	Value		057		Threshold								Value										
		-0	.057 ±	0.020							0.10	00				<b>0.764</b> ± 0.009							
					Fair	ness	Metri	c Ass	essm	ient						Performance Metrics Fairness Metrics							
	0.2															r						-	
																Actual positives $-P - TP + FN$ Actual positives $N = TN + FP$							
	0.1	0.1												_	Base rate = $BR = \frac{P}{r_{rr}}$								
		and the second														Positive rate	$= PR = \frac{TP + FP}{P + M}$			1			
	0.0 gine															Negative rat	$r = NR = \frac{TN + FN}{P + N}$			1			
	>															True Positiv	$Rate = TPR = \frac{1}{1}$	7P 1P+FN					
	-0.1	-0.1													True Negative	Rate $-TNR = \frac{T}{TN}$	N +FP		1				
															False Positive	Rate = FPR = $\frac{51}{FP_1}$	7N 7N		1				
	-0.2													Paise Negative	Rate = $FNR = \frac{1}{FN}$ tive Value = $PPV$	+79							
		ا م		5	5	~	5	5	5	- 2	- si	۱ ط	5	1			Negative Pred	ictive Value = NPV	$TF + FF = \frac{TN}{TN + FN}$				
		Parit	Opp FPR Parit	TNR	FNR Parit	PPV Parit	Parit	FDR Parit	FOR	Equa	eg I Odc	Grou	AUC	ig-los Parity			False Omission	$Rate = FOR = \frac{1}{78}$	FN I+FN			•	
											equa	By		2		•	Palen Dierenar	e rate = FND = -	<i>n</i>		Þ		
				-	Drin	nan / H	stric		brock	dd Dam													
					Prin	ndry Me	suric		nresho	nd Rang	16												



# A.I. Verify (Singapore)

- Launched in May 2022
- AI Governance Testing Framework and Toolkit
- Validates against a set of principles through standardised tests vis. a set of open-source testing solutions
- set of process checklist for self-assessment.
- Open source and supported by AI Verify Foundation launched in June 2023.
- The foundation is further venturing into Generative AI assurance and testing with Project Moonshot, announced on 31 May 2024.
- The foundation's work aligns with the "Operations Management" aspect of the Traditional Al Framework and the "Trusted Development and Deployment" and "Testing and Assurance" aspects of the Gen Al Framework.



## A.I. Verify (Singapore)



#### **Enabling Manual and Automated Red-Teaming**

Project Moonshot facilitates manual and automated redteaming, incorporating automated attack modules based on research-backed techniques to test multiple LLM applications simultaneously.

Red-Teaming allows the adversarial prompting of LLMs to induce them to behave in a manner incongruent with their design.

As Red-Teaming conventionally relies on humans, it is hard to scale. Project Moonshot has developed some attack modules that enable automated prompt generation, which allows automated red teaming.





## The European Union Act – A Risk Based Approach



- Risk management model based on the following classification of AI
  - systems resulting in an unbearable risk: social scoring
  - high-risk systems: credit scoring
  - low or minimum-risk systems (not regulated)
  - general-purpose AI models and general-purpose AI models with systemic risk



#### **Data Protection Regulation**

• EU: GDPR contains transparency rights but with limitations applicable to algorithmic credit scoring. E.g. Art. 13-14: if the controller draws inferences, notification duties may be avoided.

• SG: PDPA Fifth Schedule allows organizations to decline access to opinion data kept solely for evaluative purposes (e.g., determining whether any contract should be continued)



- PDPC Decision 2021:
- HSBC: Redacted Data was opinion data auto-generated by HSBC's proprietary algorithm that determined an individual's suitability for a credit card by analysing data from various sources



## **Topics for Discussion**

- Growing Role of Third-Party Technology Vendors: Shared responsibility between vendors and financial institutions in technology risk management?
- New era of online scams (APP fraud)
- Revisiting Sector-specific Rules in the era of AI (Digital advisory services)