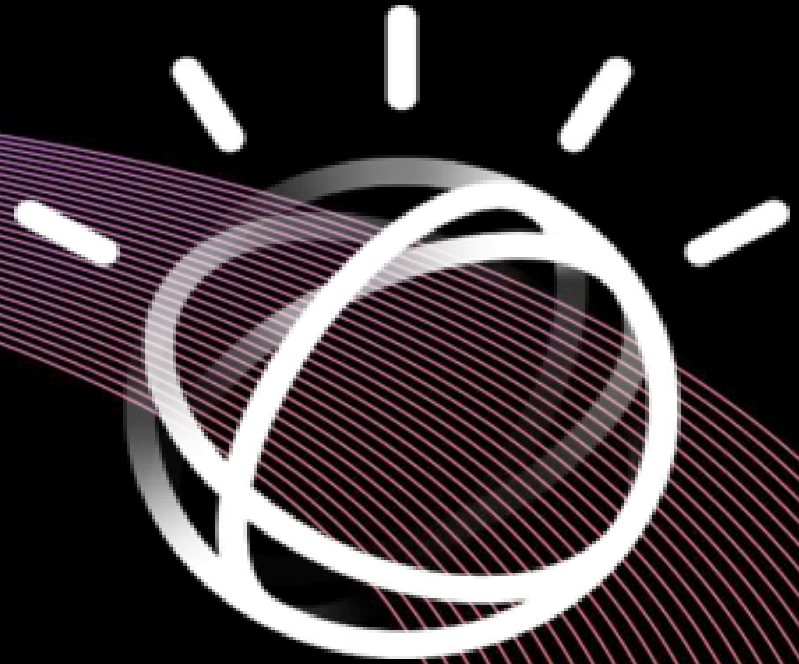# Trustworthy AI with IBM

Brian Snyder
Sr Data Science Technical Specialist
bsnyder@us.ibm.com

IBM

# Rising concerns on trust in AI decisions

**YouTube sued for using AI to racially profile content creators**

They claim YouTube's algorithms discriminate against black users

**BlackRock shelves unexplainable AI liquidity models**

Risk USA: Neural nets beat other models in t...

**Data science during COVID-19: Some reassembly required**

Most likely, the assumptions behind your data science model or the patterns in your data did not survive the coronavirus pandemic. Here's how to address the challenges of model drift

**Can AI models respond to black swan events like COVID-19?**

**Over-Segmenting In Financial Services Is So Over - Bye, Bye**

Sections ☰                    The Washington Post                    Get 1 year for $29
Democracy Dies in Darkness

Business

**Apple Card algorithm sparks gender bias allegations against Goldman Sachs**

**Amazon scraps secret AI recruiting tool that showed bias against women**

**EFF to HUD: Algorithms Are No Excuse for Discrimination**

BY JAMIE WILLIAMS, SAIRA HUSSAIN, AND JEREMY GILLULA | SEPTEMBER 26, 2019

# ...and regulators are catching up

## USA
SR 11–7 requires model risk management for all models in financial services

2019—Proposal for Algorithmic Accountability Act

2021 – US Govt. National AI Initiative Act

## Canada
2017—National AI Strategy launched. Impact Analysis

2020—All public agencies must do an impact analysis for AI models

## European Union
2021 – Draft regulation for trust in AI development

2019—Guidelines for AI development

## Mexico
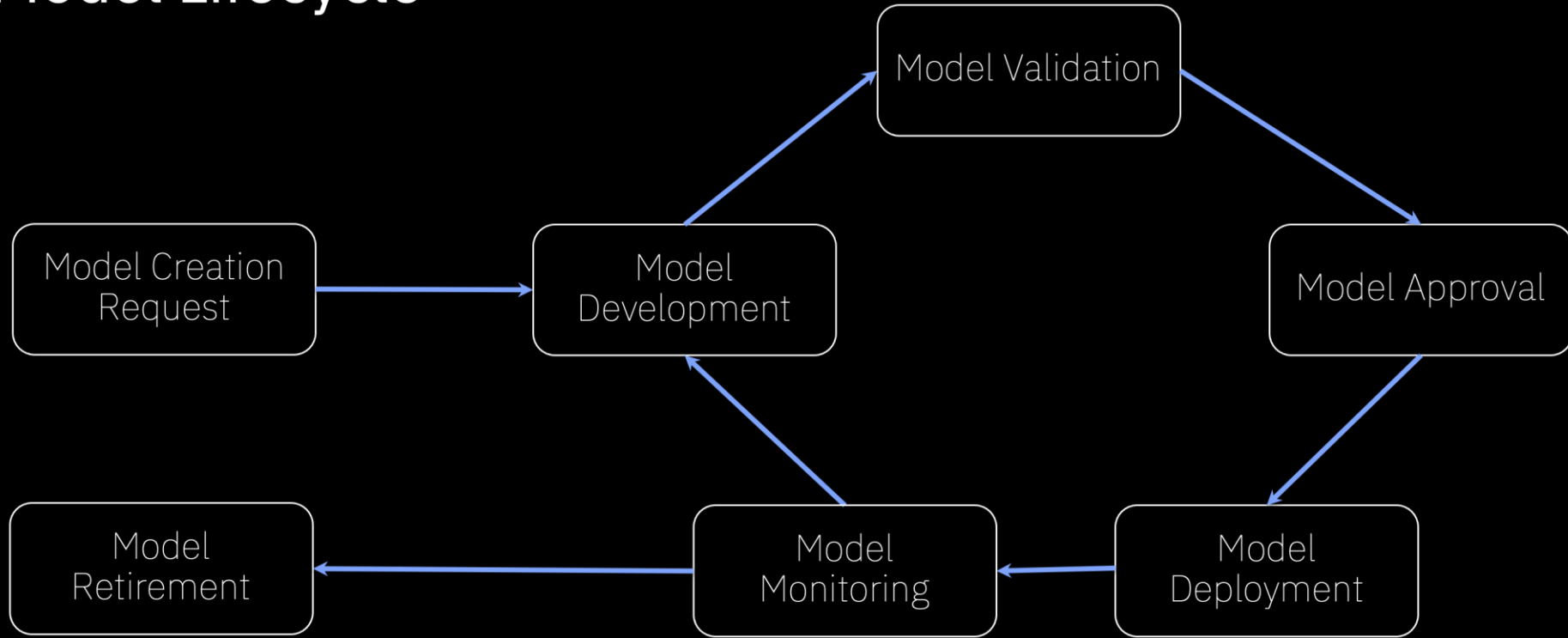2018—General principles for AI development in the government

## Partnerships on AI
Partnership between tech companies to study best practices and impact of AI

## AI Now Institute
NYU research center focused on social implications of AI
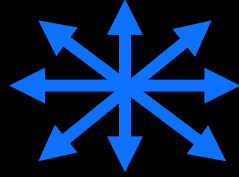
# Model Lifecycle

# Aspects of Trustworthy AI
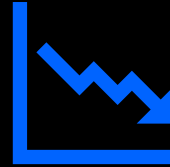
**Fairness**

Impartial and addressing bias

*Are **privileged groups** at a systematic advantage compared to other groups?*

**Robustness**

Models need to perform well across the lifecycle, handle exceptions effectively, enable confidence in systems outcomes

*Are relevant performance **metrics** monitored over time?*

**Drift**

Changes in input data cause model to make inaccurate decisions.

*Do anomalies exist between training data and data ranges or combinations seen in real life?*

**Explainability**

Easy to understand outcomes/decisions

*Why did the AI arrive at an outcome? At what point would the outcome have been different?*

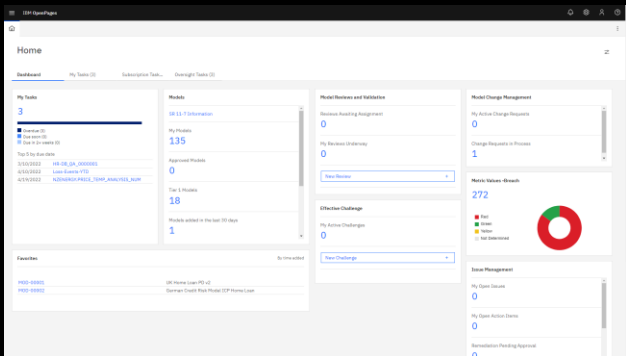**Transparency**

Open to inspecting facts and details

*Can we increase understanding of why and how AI was created?*
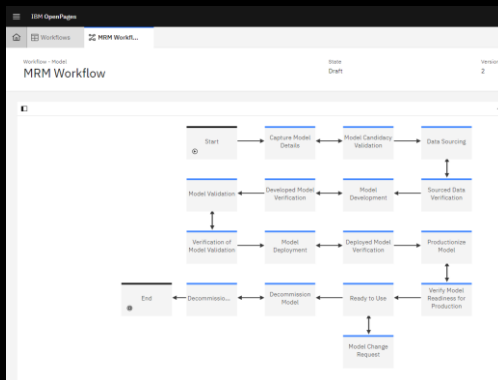
**Common AI Use Cases**

- Lending: Loan Origination, Loan Default
- Collections
- Claims Processing
- Underwriting

- Targeted Marketing Campaigns
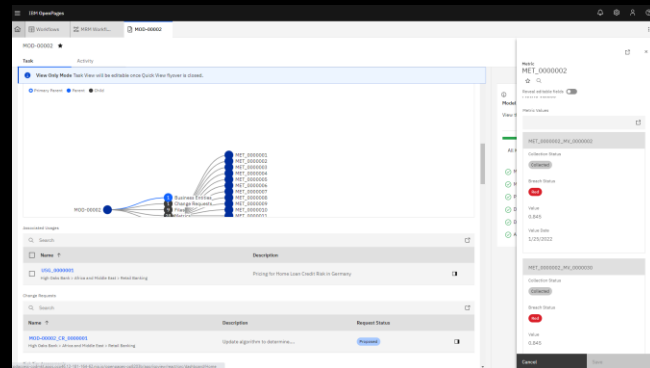- Segmentation
- Customer Management
- HR

# Trustworthy AI Demo
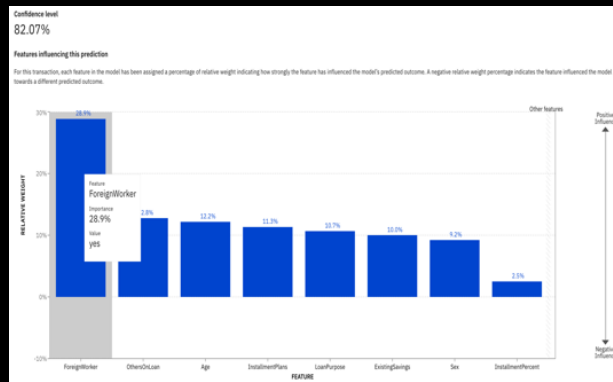


Enterprise Inventory Dashboard



Enterprise Workflow
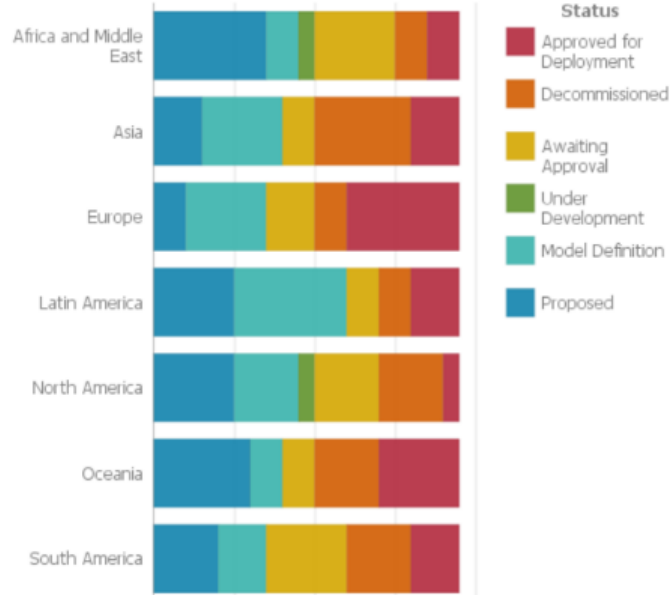


Variety of Risk Metrics Captured For Reporting



Drift in data consistency, Drift in accuracy, Bias Detection



Local and contrastive explanations

# Model Status by Region



**Model Status by Region**

| # of Models | Proposed | Model Definition | Under Development | Awaiting Approval | Decommissioned | Approved for Deployment | Total |
|---|---|---|---|---|---|---|---|
| Africa and Middle East | 7 | 2 | 1 | 5 | 2 | 2 | |
| Asia | 3 | 5 | | 2 | 6 | 3 | |
| Europe | 2 | 5 | | 3 | 2 | 7 | 19 |
| Latin America | 5 | 7 | | 2 | 2 | 3 | 19 |
| North America | 5 | 4 | 1 | 4 | 4 | 1 | 19 |
| Oceania | 6 | 2 | | 2 | 4 | 5 | 19 |
| South America | 4 | 3 | | 5 | 4 | 3 | 19 |
| **Total** | **32** | **28** | **2** | **23** | **24** | **24** | **133** |

**Status**
- Approved for Deployment
- Decommissioned
- Awaiting Approval
- Under Development
- Model Definition
- Proposed

# Change Request by Region



**Change Requests by Region**

| # of Requests | Proposed | Awaiting Review | Approved - Awaiting Implementation | Implemented | Approved - Awaiting Assessment | Total |
|---|---|---|---|---|---|---|
| Africa and Middle East | | | | | 1 | 1 |
| Europe | | | | 17 | | 17 |
| Latin America | 1 | | | 8 | | 9 |
| North America | 5 | 3 | 7 | 14 | | 29 |
| South America | | 7 | 2 | 2 | 1 | 12 |
| **Total** | 6 | 10 | 9 | 41 | 2 | 68 |

# ModelOps Data Science and Trustworthy AI as a Team - DEMO

# US-based Multinational Bank

Upon movement of all proof-of-concept projects into production, the **bank will have the ability to govern all AI projects** using their existing technology and skill investment so that existing business units do not need to change their current systems and reskill their employees.

**\*IBM Cloud Pak for Data** provides organizations a transversal and centralized view of the AI lifecycle through **an integrated platform that covers the three key functions of model build, model deploy, and model management.**

**Business problem**

This bank uses many tools and systems for model development and deployment, making model governance challenging. They were unable to **ensure models comply with enterprise policies, identify inefficiencies, provide standardized regulatory reporting, learn and scale best practices**

**Solution**

IBM used its leading enterprise insights platform (**Cloud Pak for Data**) with its open architecture enabling a **smooth integration of models developed and deployed on other platforms** from a governance perspective
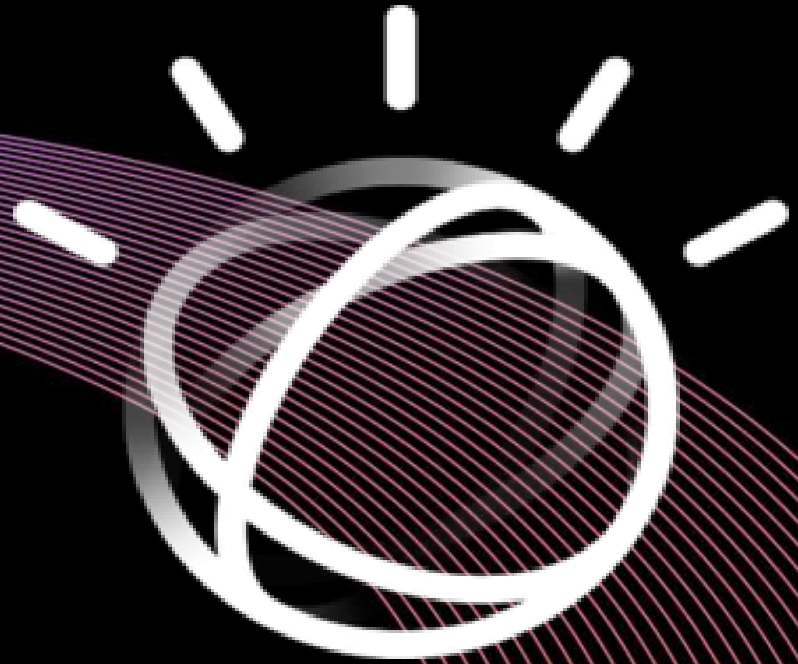
29 October 2021

# Trustworthy AI
# with IBM

Brian Snyder
Sr Data Science Technical Specialist
bsnyder@us.ibm.com

**THANK YOU!**

IBM

# Demo Snapshots

# Gartner Magic Quadrant for Data Science and Machine Learning Platforms



2021 Magic Quadrant

CHALLENGERS | LEADERS

Alteryx
IBM 2020
SAS
IBM 2021
IBM
Dataiku
MathWorks
TIBCO Software
Databricks
Microsoft
DataRobot
Google
Amazon Web Services
KNIME
H2O.ai
RapidMiner
IBM 2019
Domino
Cloudera
Alibaba Cloud
Samsung SDS
Altair
Anaconda

NICHE PLAYERS | VISIONARIES

ABILITY TO EXECUTE

COMPLETENESS OF VISION

As of Jan 2021        © Gartner, Inc

Supports multiple tasks across the data science life cycle, including:
- Problem and business context understanding
- Data ingestion
- Data preparation
- Data exploration
- Feature engineering
- Model creation and training
- Model testing
- Deployment
- Monitoring
- Maintenance
- Data and model governance
- Explainable artificial intelligence (XAI)
- Business value tracking
- Collaboration

https://www.gartner.com/doc/reprints?id=1-25DOZD29&ct=210304&st=sb

IBM **OpenPages**

# Home

**Dashboard**  My Tasks (4)  Subscription Task...  Oversight Tasks (3)

## My Tasks

**4**

- ■ Overdue (4)
- ■ Due soon (0)
- ■ Due in 2+ weeks (0)

Top 5 by due date

| | |
|---|---|
| 1/14/2022 | MOD_0000032 |
| 3/10/2022 | HR-DB_QA_0000001 |
| 4/10/2022 | Loss-Events-YTD |
| 4/19/2022 | NZENERGY.PRICE_TEMP_ANALYSIS_NUM |

## Models

SR 11-7 Information

My Models
**135**

Approved Models
**0**

Tier 1 Models
**18**

Models added in the last 30 days
**0**

## Model Reviews and Validation

Reviews Awaiting Assignment
**0**

My Reviews Underway
**0**

New Review +

## Effective Challenge

My Active Challenges
**0**

New Challenge +

## Model Change Management

My Active Change Requests
**0**

Change Requests in Process
**0**

## Metric Values -Breach

**196**

- ■ Red
- ■ Green
- ■ Yellow
- ■ Not Determined

## Issue Management

My Open Issues
**0**

My Open Action Items
**0**

## Favorites

By time added

| | |
|---|---|
| MOD-00001 | UK Home Loan PD v2 |
| MOD-00002 | German Credit Risk Model ICP Home Loan |

14

# Metric Values (161)   Breach Status : Red ✕

Search   □ Active Only   ⬆   New +

| | Name ↑ | Description | Metric Owner | Breach Status | |
|---|---|---|---|---|---|
| □ | MET_0000002_MV_0000002<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000030<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000104<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000118<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000142<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000156<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000170<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000184<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000198<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000002_MV_0000212<br>High Oaks Bank > Africa and Middle East > Retail Banking | Fairness score | | Red | ▢ |
| □ | MET_0000003_MV_0000003<br>High Oaks Bank > Africa and Middle East > Retail Banking | Quality score | | Red | ▢ |
| □ | MET_0000003_MV_0000031<br>High Oaks Bank > Africa and Middle East > Retail Banking | Quality score | | Red | ▢ |

Choose Favorite Model

# Home

## Dashboard | My Tasks (4) | Subscription Task... | Oversight Tasks (3)

### My Tasks

**4**

- Overdue (4)
- Due soon (0)
- Due in 2+ weeks (0)

Top 5 by due date

| | | |
|---|---|---|
| 1/14/2022 | MOD_0000032 | |
| 3/10/2022 | HR-DB_QA_0000001 | |
| 4/10/2022 | Loss-Events-YTD | |
| 4/10/2022 | | |
| 4/19/2022 | NZENERGY.PRICE_TEMP_ANALYSIS_NUM | |

### Models

SR 11-7 Information

My Models
**135**

Approved Models
**0**

Tier 1 Models
**18**

Models added in the last 30 days
**0**

### Model Reviews and Validation

Reviews Awaiting Assignment
**0**

My Reviews Underway
**0**

[ New Review + ]

### Effective Challenge

My Active Challenges
**0**

[ New Challenge + ]

### Model Change Management

My Active Change Requests
**0**

Change Requests in Process
**0**

### Metric Values -Breach

**196**

- Red
- Green
- Yellow
- Not Determined

### Favorites

By time added

| MOD-00001 | UK Home Loan PD v2 |
|---|---|
| MOD-00002 | German Credit Risk Model ICP Home Loan |

### Issue Management

My Open Issues
**0**

My Open Action Items
**0**

16

## Model Details

**Model**
# MOD-00002 ★

**Model Status**
Draft

**Final Tier**
Tier 2

Task | Activity

🔍 | Reveal editable fields ⚪ | *Required *Modified

## General ⓘ

**Name** *
MOD-00002

**Description**
German Credit Risk Model ICP Home Loan

**Model Status**
Draft

**Version**
1

**Model or Non-Model** *
Model

**Candidate Status**
Confirmed

**Machine Learning Model**
Yes

**Monitored with Watson OpenScale**
Yes

**Additional Description**
A UK Home Loan PD model based on client history analysis with Experian Credit Data. Version 2 revised the data set from Experian and enhanced model accuracy. Covers Owner-Occupied and Buy-to-Let.

**Model Category**
Retail Credit Risk

**Basel Model**
No

**Measurement Type**
Probability of Default

## Ownership

**Model Owner** *
👤 jblanco

**Model Delegate**
👤 System Administrator

**Model Type**
Logistic Regression

**Third Party Model**
No

**Model Provider**
IBM

---

ⓘ

**Model General View**

View the critical information about this model.

All Key Items (6) ⌄

✓ Model or Non-Model *
✓ Model Owner *
✓ Proposal Original Date *
✓ Definition Original Date *
✓ Development Original Date *
✓ Approval Original Date *

Use preconfigured workflow or in this case create your own.

# Fairness and Artificial Intelligence in Banking

Use of AI/ML creates many risks – bias, lack of explainability, operational, … – with implications for both compliance and risk management

Regulators issued guidance regarding these risks, but the regulatory landscape continues to evolve as new evidence of consumer impact emerges

**Fairness and bias**

**Complexity**

**Explainability**

**Alternative Data**

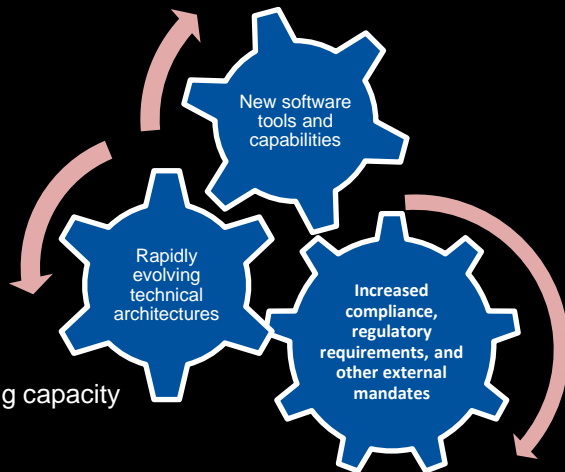# Very specific challenges relate to the nature of a data science model.

> The technology and regulations, and how they are applied, are constantly evolving.

Traditional (SAS, Matlab, IBM, etc.)

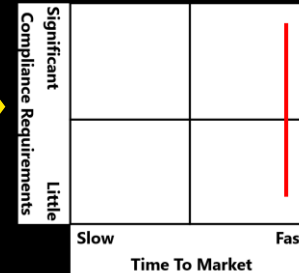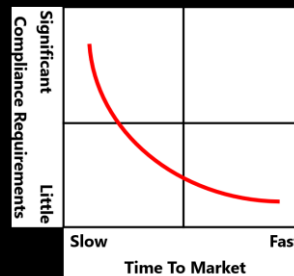ML-centric (IBM Watson, Datarobot, RapidMiner, H2O, etc.)

New entrants (Algorithmia, ModelOp, Modzy, etc.)

How to ensure the necessary *velocity* and *quality* in a changing, multi-system and multi-source environment?

New software tools and capabilities

Rapidly evolving technical architectures

Increased compliance, regulatory requirements, and other external mandates

Cloud

Unlimited computing capacity

Containerization

Security
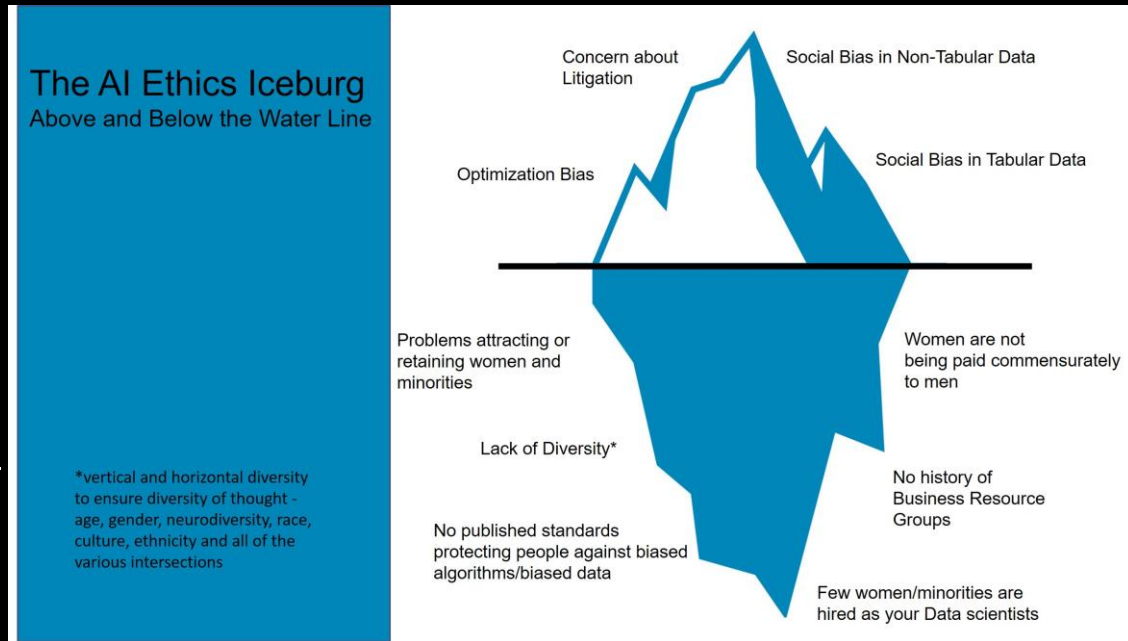
DataOps, DevOps, MLOps

FRB SR 11-7

OSFI E-23

ECB TRIM Guideline

World Bank Credit Scoring Guidelines

OECD AI Principles

IOSCO Consultation report on AI
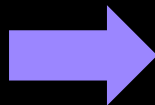
CA Consumer Privacy Act

# Operationalizing AI requires organizations to demonstrate responsible, explainable use of AI

# AI

Driven by data ➔ validity <u>not guaranteed</u>
Models ➔ not code, often <u>unexplainable</u>
Probabilistic ➔ non-deterministic, <u>uncertain</u>
Not just focus on version control, not an
<u>assembly line</u> ➔ Model development MUST
be part of lifecycle management



The AI Ethics Iceburg
Above and Below the Water Line

Concern about Litigation

Social Bias in Non-Tabular Data

Optimization Bias

Social Bias in Tabular Data

Problems attracting or retaining women and minorities

Women are not being paid commensurately to men

Lack of Diversity*

No published standards protecting people against biased algorithms/biased data

No history of Business Resource Groups

Few women/minorities are hired as your Data scientists

*vertical and horizontal diversity to ensure diversity of thought - age, gender, neurodiversity, race, culture, ethnicity and all of the various intersections

Model drift, bias and risk can pose significant liabilities and damage

➔ AI fairness is a corporate social responsibility