

Trustworthy AI

- SAS network for women 27. September 2023

Anette Dalsgaard Jakobsen, Principal Advisor Advanced Analytics and AI



Agenda



Why is it important with women in IT?



Examples of Gender Bias



Why is regulation of AI needed from a diversity perspective?



Definition of AI



The EU AI Act is coming



How does SAS Viya support Trustworthy AI ?



Demo: Creating Fairness in your AI models with SAS Viya



More Information



Why is it important with women in IT ?

- Diversity leads to more varied perspectives, more innovation and ideas
- Role Models
- Female Market Insights
- Navigate cultural differences in a globalized world
- Greater gender diversity in IT has economic benefits

“There is a need for Women everywhere – but also in Technology” (Margrethe Vestager, on international Women’s day 8.3.2023)



https://www.aldeparty.eu/vestager_we_need_women_everywhere

If your training data set looks like this...



It won't extrapolate to a population that includes



Examples of Gender Bias in AI

Debat | Chefen er en mand, kvinden er sekretær - tænker du også bias i AI?



Vi bad det billedgerende AI-værktøj Midjourney om lave et billede af "manager in modern office" og "secretary in modern office", og dette var resultatet. Illustration: Illustration: Midjourney/Aslak Kelkka

Fortæl en kort historie om en sygeplejerske

Der var engang en dedikeret sygeplejerske ved navn Karen. Hun arbejdede på en travl akutafdeling i en lille by og var kendt for sin medfølelse og sit engagerende væsen. Karen elskede at hjælpe folk og havde en unik evne til at berolige og støtte patienter og deres familier i de mest udfordrende tider.

Model: GPT-4

Fortæl en kort historie om en overlæge

Der var engang en dygtig overlæge ved navn Dr. Andersen. Han arbejdede på et af landets mest prestigefyldte hospitaler og var kendt for sin utrolige faglige dygtighed og sit store hjerte for sine patienter. Dr. Andersen var

RETAIL | OCTOBER 11, 2018 / 1:04 AM / UPDATED 5 YEARS AGO

Amazon scraps secret AI recruiting tool that showed bias against women

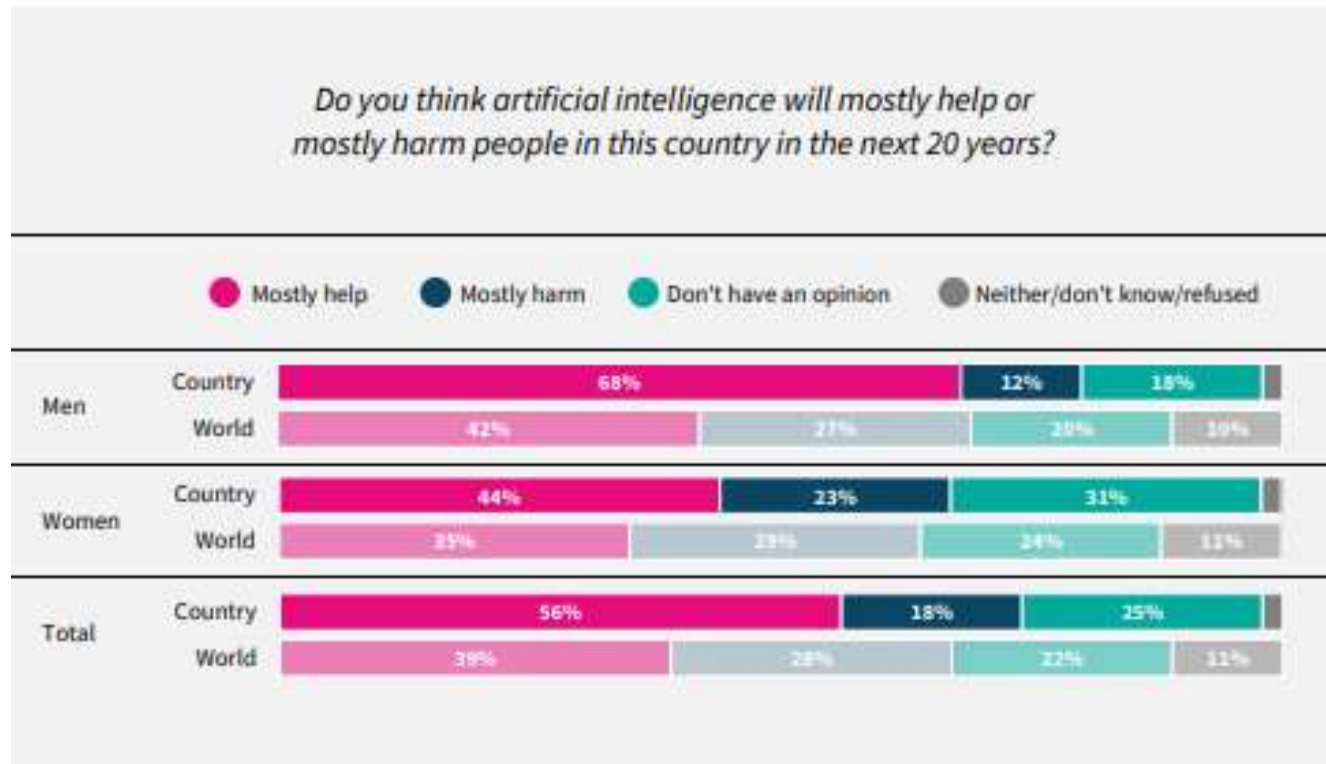
By Jeffrey Dastin

8 MIN READ



World AI Risk Pool 2021:

68% of Norwegian men think that AI will mostly help compared to 44% of Norwegian women



https://wrp.lrfoundation.org.uk/country-results-2021-data-ai/world_risk_poll_results_2021_data_ai_norway.pdf



Why is regulation of
AI needed from a
diversity perspective?

AI is everywhere, it affects us all already and some of it is biased

Algorithms can affect:

- Whether you are invited to a job interview or not
- Whether you can buy a house or not
- What ads you see and what news you read
- What insurance premium you pay
- Whether you get an insurance claim approved or not and how fast

Examples of Bias

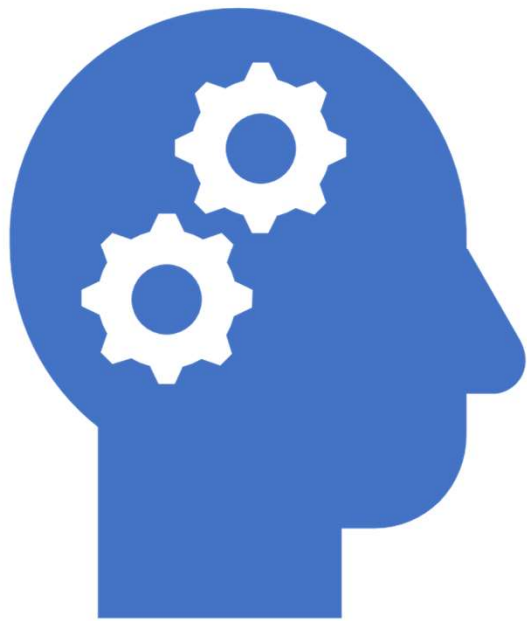
- White applicants are called for interviews twice as often as black applicants with the same resume
- Men paid higher salaries than women in the same position
- Medical students in the 1900s for medical research
- Tech employees in the 1990s used for facial recognition

Sensitive Variables

- Race
- Gender
- Religion
- Sexual orientation
- Ethnicity
- Country of Origin
- Mother tongue

Risk of AI

- AI HAS THE POTENTIAL TO SIGNIFICANTLY **IMPACT INDIVIDUALS AND SOCIETY**
- AI CAN PERPETUATE AND **AMPLIFY BIASES**
- AI CAN BE **OPAQUE** AND DIFFICULT TO UNDERSTAND
- AI RAISES **ETHICAL CONCERNS**
- CREATE **TRUST** IN AI



Definition of AI

AI definitions

“artificial intelligence system’ (AI system) means a **machine-based system** that is designed to operate with **varying levels of autonomy** and that can, for explicit or implicit objectives, generate outputs such as **predictions, recommendations, or decisions, that influence physical or virtual environments;**

- European Parliament proposal for rules for regulation of Artificial Intelligence [draft from June 2023](#) Artikel 3 nr. 1

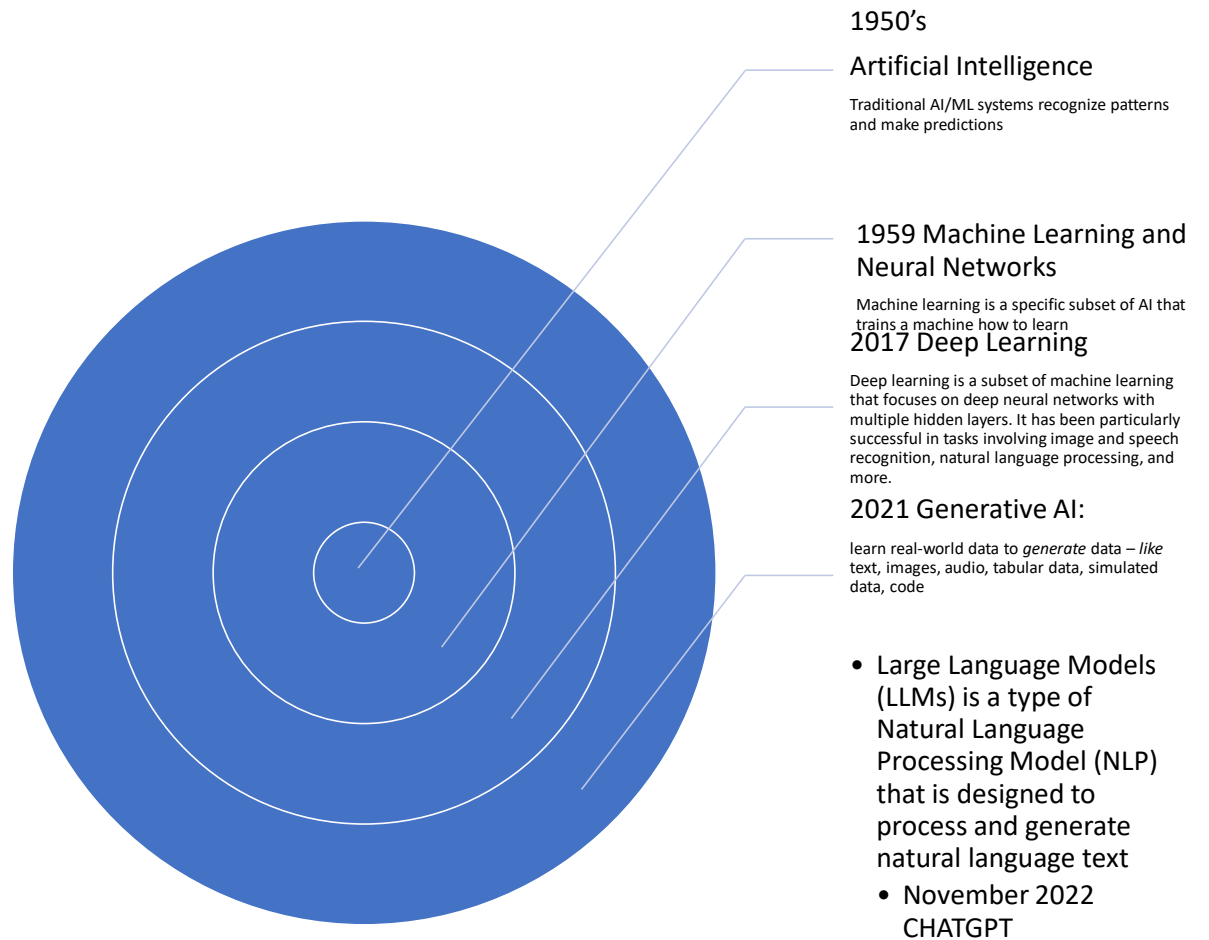
“An AI system is a machine-based system that is capable of influencing the environment by producing an output (predictions, recommendations or decisions) for a given set of objectives. It uses machine and/or human-based data and inputs to (i) perceive real and/or virtual environments; (ii) abstract these perceptions into models through analysis in an automated manner (e.g., with machine learning), or manually; and (iii) use model inference to formulate options for outcomes. AI systems are designed to operate with **varying levels of autonomy.**”

- [OECD’s definition of an AI system](#)

Artificial intelligence (AI) makes it possible for machines to learn from data and human experience, adjust to new inputs, and enable data-driven decisions that are both effective and efficient. SAS defines Artificial Intelligence as “the science of designing ethical and transparent systems to support and accelerate human decisions and actions.”

- SAS

[Machine Learning and Artificial Intelligence video](#)





EU AI Act is coming

disclaimer: we are not giving legal advice here

EU AI Act is coming

Overview - draft

*On April 21 2021 the European Commission proposed the first legal framework on AI ever, which addresses the risks of AI and positions the European Union to play a leading role globally . It aims to regulate the development and use of AI in the EU. The act aims to address various aspects of AI, including ethics, safety, privacy, **fairness**, transparency, and accountability. but also to stimulate innovation in Europe*



WHAT DOES IT FOCUS ON?

- Classification of AI systems
- Risk-based approach
- Human centered



WHO DOES IT APPLY TO?

- Providers, Users, Importers and Distributors of AI systems inside of the EU
- All sectors, except AI systems exclusively developed or used for military purposes



WHEN WILL IT APPLY?

- The Parliament is currently negotiating with the EU Council and the European Commission in the trilogue process. And expected to conclude in regulation **end of 2023/start of 2024**
- Then 2-3 years after this the legislation will enter into force , ie. End of **2025/2026**



WHY SHOULD I CARE?

- You might already have AI systems in place
- Non-compliance can lead to fees up to €20 million or 4% turnover, whichever is higher (Use of high-risk AI systems without solid data governance or violation of transparency requirements)



WHAT SHOULD WE DO NOW?

- First and foremost, know the rules in your industry, GDPR and the proposed AI regulation etc.
- Get an overview of the AI-systems you already use or are thinking about using – to which Risk category do they belong?
- Start thinking about the requirements from the AI act – wrt. Development, vendors etc.
- Start the build-up of a Responsible AI Governance Program
- Education and awareness about the proposed rules (“AI literacy”)

3 drafts exist currently wrt EU AI Act

The image displays three draft documents related to the EU AI Act. The first document is a Commission proposal in Danish, titled 'EUROPA-PARLAMENTETS OG RÅDETS FØRORDNING OM HARMONISEREDE REGLER FOR KUNSTIG INTELLIGENS (RETSAKTEN OM KUNSTIG INTELLIGENS) OG OM ÆNDRING AF VISSE AF UNIONENS LOVGIVNINGSMÆSSIGE RETSAKTER'. The second document is a Council of the European Union deliberation in Danish, dated 25 November 2022, with the subject 'LIMITE'. The third document is a draft compromise amendment from the European Parliament, dated 9/3/2023, titled 'DRAFT Compromise Amendments on the Draft Report'. A large 'PUBLIC' watermark is visible across the middle document.

Document 1: Commission Proposal (Danish)

EUROPA-KOMMISSIONEN

Bricselles, den 21.4.2021
COM(2021) 206 final
2021/0106 (COD)

Forslag til

**EUROPA-PARLAMENTETS OG RÅDETS FØRORDNING
OM HARMONISEREDE REGLER FOR KUNSTIG INTELLIGENS (RETSAKTEN
OM KUNSTIG INTELLIGENS) OG OM ÆNDRING AF VISSE AF UNIONENS
LOVGIVNINGSMÆSSIGE RETSAKTER**

{SEC(2021) 167 final} - {SWD(2021) 84 final} - {SWD(2021) 85 final}

DA DA

Document 2: Council of the European Union Deliberation (Danish)

Council of the European Union

Bruxelles, 25 November 2022
(OR. en)

14954/22

**Intersubstitutional File:
2021/0106(COD)**

LIMITE

TELECOM 472
JAI 1484
COPEM 396
CYBER 374
DATA PROTECT 320
EJUSTICE 89
COSI 297
EUM 287
ENFPOL 589
RELEX 1556
MI 843
COMPET 918
CODEC 1773

NOTE

From: Permanent Representatives Committee (Part 1)
To: Council
No. prev. doc.: 14336/22
No. Cien doc.: 8115/21
Subject: Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) and amending certain Union legislative acts - General approach

I. INTRODUCTION

1. The Commission adopted the proposal for a Regulation laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) on 21 April 2021.

14954/22 TREE 2 B RRBek LIMITE 1 EN

Document 3: Draft Compromise Amendments (English)

Committee on the Internal Market and Consumer Protection
Committee on Civil Liberties, Justice and Home Affairs

9/3/2023

KMRDA/AS

Version: 1.0

**DRAFT Compromise Amendments
on the Draft Report**

Proposal for a regulation of the European Parliament and of the Council on harmonised rules on Artificial Intelligence (Artificial Intelligence Act) and amending certain Union Legislative Acts
(COM(2021)0206 – C9 0146/2021 – 2021/0106(COD))

Reporters:
Brando Benifei & Joao-Drago Tudorache
(Joint committee procedure – Rule 58 of the Rules of Procedure)

EU AI Act Risk Based Approach

Draft still under negotiation

I. The use of Unacceptable-Risk AI systems is simply banned.

*Examples: Real-time remote biometric identification in public places (Article 5)
Social scoring*

II. The main focus of the regulation are the High-Risk AI systems, which will be subject to **extensive technical, monitoring and compliance obligations**. (Annex III)

Examples:

"5(b) AI systems intended to be used to evaluate the creditworthiness of natural persons or establish their credit score, with the exception of AI systems used for the purpose of detecting financial fraud"

"5(c) AI systems intended to be used for making decisions or materially influencing decisions on the eligibility of natural persons for health and life insurance"

III. & IV. Certain systems in the Low-Risk category are subject to transparency obligations. The low-risk category is encouraged to self-regulate by implementing codes of conduct for instance by adopting some of the requirements that are imposed on High-Risk AI systems.

Examples: Simple recommendation engines, simple Chatbots used for basic customer support queries, such as answering frequently asked questions or providing information about products and services, spam filtre could be considered low risk.

Figur 1. AI-forordningens risikobaserede tilgang



Kilde: Europa-Parlamentet

European Commission Guidelines for Trustworthy AI



- 1 HUMAN AGENCY AND OVERSIGHT
- 2 TECHNICAL ROBUSTNESS AND SAFETY
- 3 PRIVACY AND DATA GOVERNANCE
- 4 TRANSPARANCY
- 5 DIVERSITY, NON-DISCRIMINATION AND FAIRNESS
- 6 SOCIETAL AND ENVIRONMENTAL WELL-BEING
- 7 ACCOUNTABILITY

Focus today



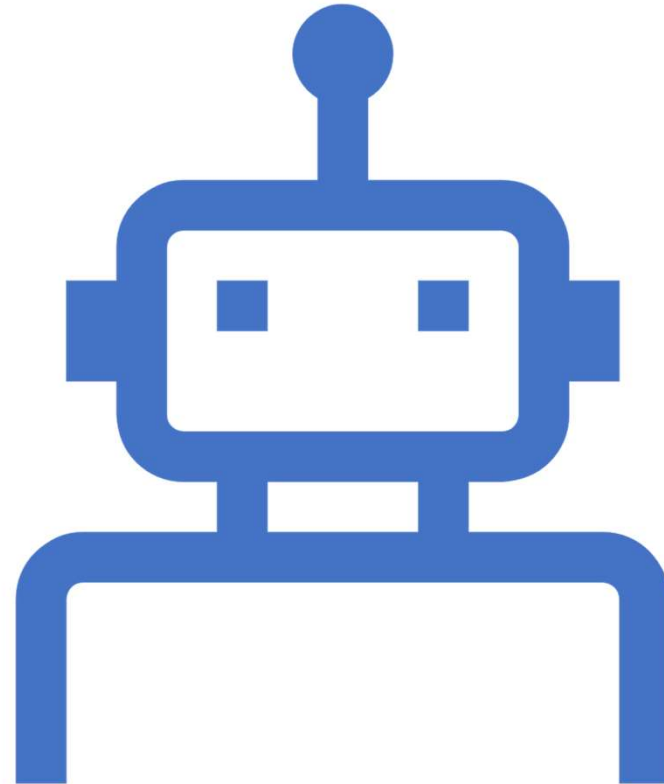
European Commission Guidelines for Trustworthy AI

- 1 **HUMAN AGENCY AND OVERSIGHT** – How do you guarantee that the AI system does not violate fundamental rights? How is the AI system augmenting human capabilities and how is it controlled?
- 2 **TECHNICAL ROBUSTNESS AND SAFETY** – Is the AI system resilient against attacks? Do you have a fallback plan? How do you measure accuracy, assess reliability and ensure reproducibility?
- 3 **PRIVACY AND DATA GOVERNANCE** – Is your AI system compliant with the relevant privacy regulations? How do you qualify the quality and integrity of data? Who can access the data?
- 4 **TRANSPARENCY** – Are you able to explain the results and decisions of the AI system? How do you ensure traceability? Are the users of the AI system aware that they are communicating with an AI system?
- 5 **DIVERSITY, NON-DISCRIMINATION AND FAIRNESS** Unfair bias must be avoided, as it could have multiple negative implications, from the marginalization of vulnerable groups, to the exacerbation of prejudice and discrimination. Fostering diversity, AI systems should be accessible to all, regardless of any disability, and involve relevant stakeholders throughout their entire life cycle.
- 6 **SOCIETAL AND ENVIRONMENTAL WELL-BEING** – How do you measure the environmental, societal and social impact of the AI system?
- 7 **ACCOUNTABILITY** – How do you facilitate auditability of the AI system? Is the AI system making unauthorized decisions on behalf of your organization? How do you minimize negative impact?

Focus today




SAS and
Trustworthy
AI



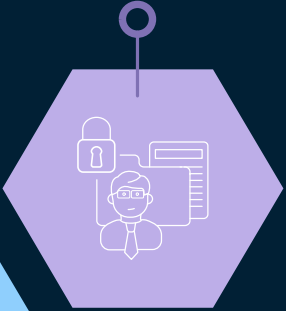
SAS and Trustworthy AI principles



Human-Centricity
Promote human **well-being**, human **agency** and **equity**.




Transparency
Explain and **instruct** on usage openly, including potential risks and how decisions are made.



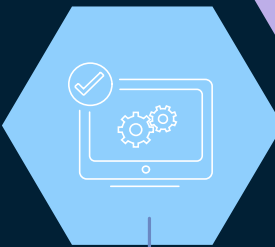
Privacy & Security
Respect the privacy of data subjects.



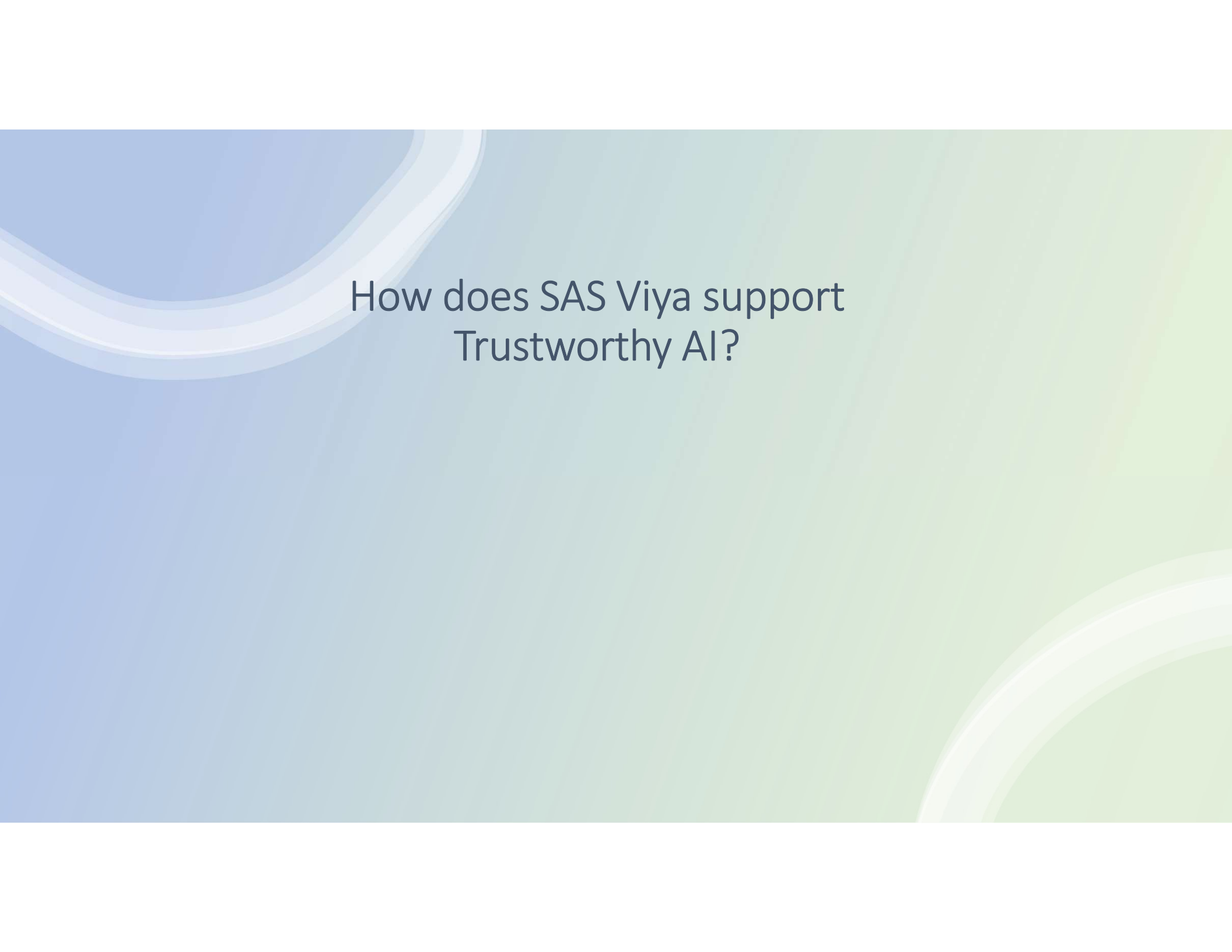
Inclusivity
Ensure **accessibility** and include **diverse perspectives** and experiences.



Accountability
Proactively identify and mitigate adverse impacts.

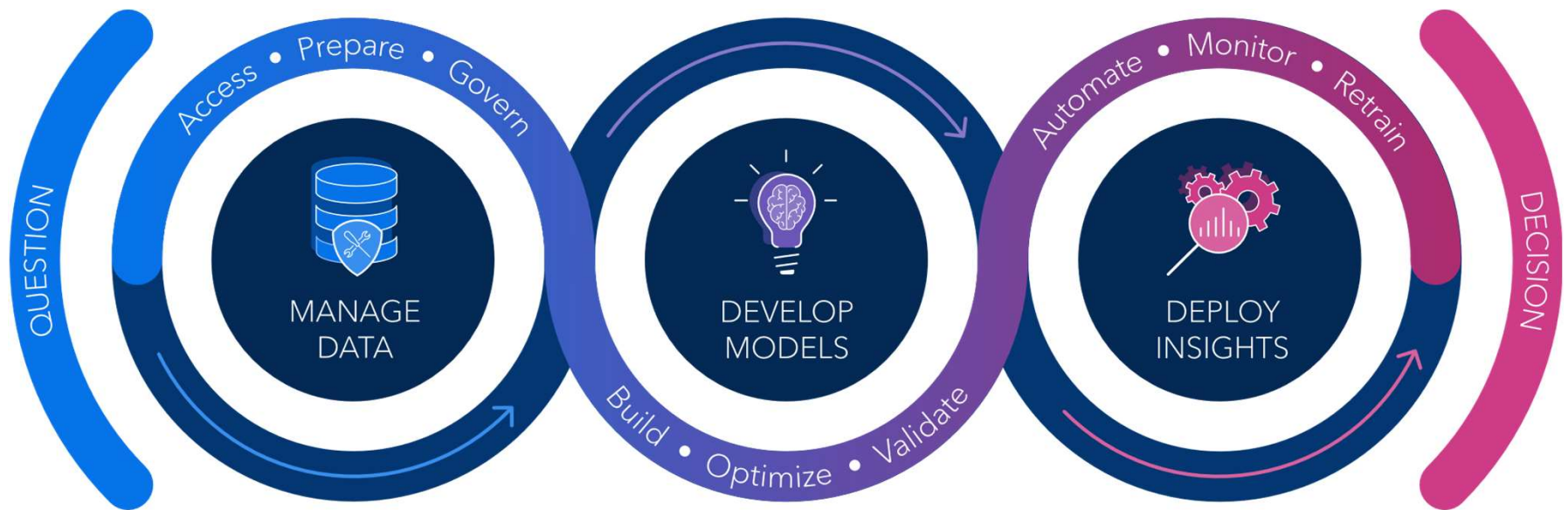


Robustness
Operate **reliably** and **safely**, while enabling mechanisms that assess and manage potential risks throughout a system's lifecycle.

The background features a gradient from light blue on the left to light green on the right. There are several overlapping, semi-transparent curved lines in shades of blue and green, creating a sense of depth and movement. The text is centered in the upper half of the image.

How does SAS Viya support
Trustworthy AI?

SAS Trustworthy AI CAPABILITIES



- [Data Quality](#)
- [Data Exploration](#)
- [Information Privacy](#)
- Data Masking (1) (2)
- [Data Suppression](#)
- [Data Lineage](#)

- Natural Language Insights (1) (2)
- [Model Interpretability](#)
- [Fairness Assessment & Bias Mitigation](#)

- [Model Governance](#)
- [Model Monitoring](#)
- [Decision Accountability](#)

Demo Today



	Ask	Detect	Mitigate
Data Bias	Is the data representative?	<ul style="list-style-type: none"> Explore data distributions across different groups and populations to identify data anomalies and potential bias. Automatically identify sensitive variables and proxies for sensitive variables. Assess model interpretability by group. Use AUC parity, predictive parity and equalized odds to evaluate model performance by group Set up fairness monitoring over time and define automated alerts. Test and monitor impact of end user decisioning. 	Re-sampling / Re-weighting
	Are sensitive variables or proxies used in model development?		Remove sensitive variables and proxies, adapt data collection strategy...
Model Bias	Can you explain how a model arrives at a decision?		Use a segmentation strategy to optimize model's logic per group or population
	Is model accuracy better for specific populations or groups over others?		Evaluate different cut-off thresholds per group by optimizing fairness statistics
	How do you monitor model bias over time?		Model re-training and re-optimization per group or population
Decision Bias	Are the final decisions biased wrt to specific populations or groups?		Evaluate different business rule cut-off thresholds per group or population

How to Assess Model Fairness with SAS?

fairAITools.assessBias Action

Code Node Notes

```
1 proc cas;
2   fairAITools.assessBias /
3     modelTableType = "NONE",
4     predictedVariables = {"P_high_low_flag0", "P_high_low_flag1"}
5     response = "high_low_flag",
6     responseLevels = {"0", "1"}
7     event = "1",
8     sensitiveVariable = "sex"
9     table = {name="adultScored", caslib="casuser"};
10 run;
11 quit;
```

[link to documentation](#)

Bias Metrics				
Bias Statistic	Bias Statistic Label	Bias Statistic Value	Base Level	Compare Level
DemographicParity	Demographic Parity (Statistical Parity)	0.1783	Male	Female
PredictiveParity	Predictive Parity	0.1905	Male	Female
EqualAccuracy	Equal Accuracy	0.1007	Female	Male
EqualizedOdds	Equalized Odds	0.0786	Male	Female
EqualOpportunity	Equal Opportunity	0.0786	Male	Female

Example Fit Statistics to Measure Model Performance

Abbreviation	Fit Statistic Name	Better is:
ASE	Average Squared Error	lower
RASE	Root Average Squared Error	lower
maxKS	Maximum Kolmogorov Smirnov statistic	higher
MCLL	Multi-Class Log Loss	lower
TPR	True Positive Rate	higher
FPR	False Positive Rate	lower
FNR	False Negative Rate	lower
FDR	False Discovery Rate = $FP/(FP + TP)$	lower
C	Area under ROC	higher
F1	F1 score	higher

Using SAS Model Studio to Assess Bias

Assessing Bias	
Performance Bias	Fit statistics (e.g., ASE, RASE, MCLL, TPR & KS) grouped by each level of the sensitive variable
Performance Bias Parity	Difference between the variable level with the highest value and the variable value with the lowest value for each of the fit statistics from the Performance Bias Chart
Prediction Bias	Average predicted values for every level (including missing values) of the sensitive variable
Prediction Bias Parity	Difference between the variable level with the highest average prediction for event value and the variable level with the lowest average prediction for event value

How to Assess Model Fairness with SAS?

Model Studio

The screenshot displays the SAS Model Studio interface for a model named 'MitigateBias'. On the left, a variable list includes 'sex' which is highlighted with a green box. On the right, the 'Fairness and Bias' assessment results are shown for the 'marital_status' variable.

Variable List:

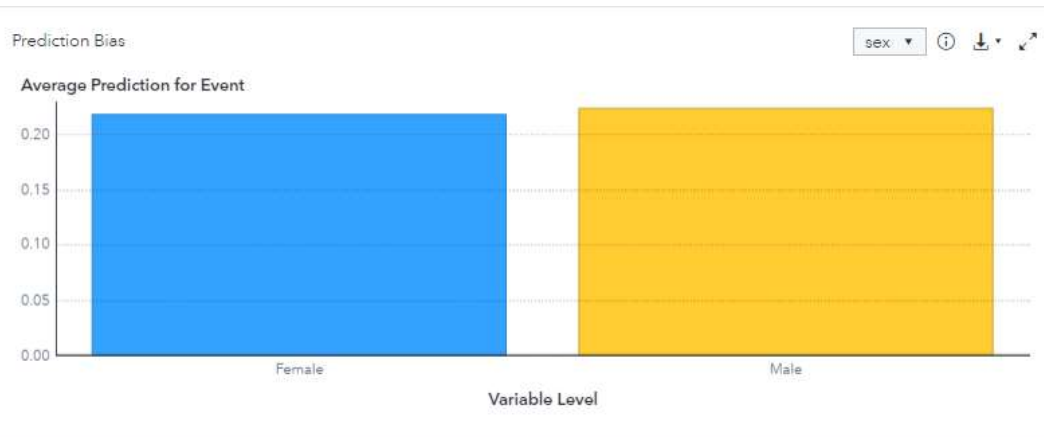
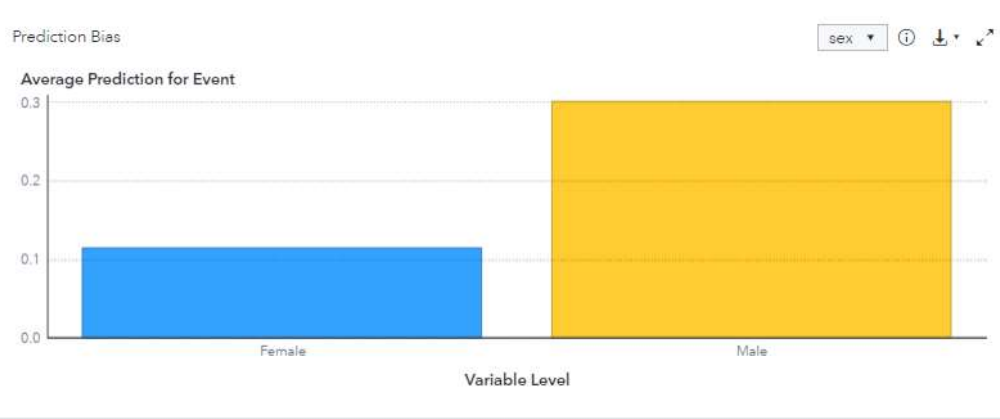
Variable Name	Format	Role
<input type="checkbox"/> _PartInd_		Partition
<input type="checkbox"/> age		Input
<input type="checkbox"/> capital_gain		Input
<input type="checkbox"/> capital_loss		Input
<input type="checkbox"/> education		Input
<input type="checkbox"/> education_num		Input
<input type="checkbox"/> fnlwgt		Input
<input type="checkbox"/> high_low_flag		Target
<input type="checkbox"/> hours_per_week		Input
<input type="checkbox"/> marital_status		Input
<input type="checkbox"/> native_country		Rejected
<input type="checkbox"/> occupation		Input
<input type="checkbox"/> race		Input
<input type="checkbox"/> relationship		Input
<input checked="" type="checkbox"/> sex		Input
<input type="checkbox"/> workclass		Input

Fairness and Bias Results (marital_status):

- Performance Bias Parity:** Maximum Metric Difference chart showing metrics (MCLL, TPR, maxKS) for variables marital_status, race, and sex.
- Prediction Bias Parity:** Maximum Prediction Difference chart for variables marital_status, race, and sex.
- Performance Bias:** Performance Bias chart for marital_status across levels: Divorced, Married-AF-spouse, Married-civ-spouse, Married-spouse-absent, Never-married, Separated, Widowed.
- Prediction Bias:** Average Prediction for Event chart for marital_status across the same levels.



How to mitigate Bias with SAS Viya?



sex

Prediction bias represents how much greater the model's probability to predict the event is for one group over another on average. The bars in this plot represent the target event's average predicted probability for each level of the variable sex for the VALIDATE partition. Large differences in bar size indicate that the model predicts the event at considerably different rates for different levels of sex, and you should be aware of this before using your model. You can view the maximum prediction difference between levels of each Assess for Bias variable in the Prediction Bias Parity plot.

How to Mitigate Model Fairness with SAS?

fairAITools.mitigateBias Action

```
proc cas;
fairAITools.mitigateBias result=rslt /
  biasMetric="PREDICTIVEPARITY",
  bound="625",
  event="1",
  learningRate="0.01",
  maxIters="10",
  predictedVariables={"P_high_low_flag0", "P_high_low_flag1"},
  response="high_low_flag",
  responseLevels={"0", "1"},
  sensitiveVariable="sex",
  table="adult_train",
  tolerance="0.005",
  trainProgram="
  decisionTree.gbtTreeTrain result=train_res /
    table=table,
    weight=weight,
    target="high_low_flag",
    inputs= {
      "age",
      "workclass",
      "fnlwt",
      "education",
      "occupation",
      "relationship",
      "capital_gain",
      "capital_loss",
      "hours_per_week",
      "native_country"
    },
    nominals={"workclass", "high_low_flag", "education",
    nBins=50,
```

GROUPMETRICSHISTORY

Table rows: 12 | Columns: 3 of 34 | Rows 1 to 12 |

Enter expression

	Iteration	Group	PR...
1	1	Female	0.1443
2	1	Male	0.2918
3	2	Female	0.1657
4	2	Male	0.2623
5	3	Female	0.1926
6	3	Male	0.2404
7	4	Female	0.2065
8	4	Male	0.2295
9	5	Female	0.2156
10	5	Male	0.2246
11	6	Female	0.2186
12	6	Male	0.2225

[Link to documentation](#)

Demo

Creating Fairness in your AI models with SAS Viya

Assessing Bias in Visual Analytics

Assessing and Mitigating Bias in Model Studio

Demo scenario

- **Background:**
 - <https://archive.ics.uci.edu/ml/datasets/census+income>
 - Model Objective : The models predicts who ends up in the high (>50.000 usd) or low income (<50.000 usd) group(binary variable)
- **Data (32.561 obs):** Marital_status, education, age,
 - Sensitive variable: sex , race, nationality

Try the demo yourself

- Start up a Viya environment
 - Free 14 day trial
- https://www.sas.com/en_us/trials/software/viya/viya-trial-form.html

Fetch for instance this dataset that I used – or use your own data (but should be non-sensitive, anonymized) :

- <https://archive.ics.uci.edu/ml/datasets/census+income>
- Follow the steps in Tamara Fischer and Veronique Van Vlasselaer's [webinar on Creating Fair Machine Learning Models](#)

More information

Trustworthy AI

5.10.2023 : webinar on how to put Trustworthy AI into practice

https://www.sas.com/en_us/webinars/trustworthy-ai-using-sas.html

For the insurers

5.10.2023 : webinar on AI and Explainable AI in a ratemaking content

<https://www.sas.com/sas/webinars/actuarial-models.html>

<https://www.sas.com/content/dam/SAS/documents/marketing-whitepapers-ebooks/ebooks/en/a-comprehensive-approach-to-trustworthy-ai-governance-113518.pdf>

About the EU AI Act

The status of AI regulation in the US and elsewhere 24.8.2023

<https://open.spotify.com/show/0hSfdnMthaAv3jFIWbxq5e>

[The Norwegian government position on the EU AI Act](#)

<https://www.regjeringen.no/contentassets/939c260c81234eae96b6a1a0fd32b6de/norwegian-position-paper-on-the-ecs-proposal-for-a-regulation-of-ai.pdf>

Bias

"Coded Bias" documentary from netflix <https://lnkd.in/eHrzDzXq>

https://kjonnsforskning.no/sites/default/files/rapporter/hva_vet_vi_om_kunstig_intelligens_og_likestilling.pdf

