

Data Science Challenges

COMMON PITFALLS & ETHICS

About me



Short Bio

10+ years working in Telecom in Marketing & Customer Value Management. Now leading Data Science team enabling NOS to make better decisions using the latest analytical artillery. Nuno holds one MSC in Electrotechnical Engineering (University of Coimbra), a Msc of Data Analytics (University of Porto) and is a PHD candidate of the MAP-i Doctoral Program in Computer Science **studying Fairness in ML**.

“Trustworthy AI” - PhD Motivation

New cycle of learning on a broad AI/ML scope;
Tech skills & tools refresh;
Business relevancy;
Social mission;

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

Why are we here today?

At the end of session you should have new knowledge to start applying on your Data Science Projects

1. Improve Data Science Projects execution

Support your development team in planning DS activities making sure that critical path includes tasks to mitigate common pitfall risks;

2. Move Ethics out of the lecture theatre and put it into practice

Apply ethical framework in your Data Science projects learning the concepts and identifying case studies patterns;

3. Promote Fairness discussions at Model development

Ask the right questions together with all the stakeholders involved in a Data Science Product development / maintenance by measuring and discussing the reason of biases and fairness on your models.

Managing expectations

1. Online interaction is challenging

A few work group sessions are included on the agenda and I'm counting on active participation through the session:)

2. Communication

To be inclusive concepts intuition will be provided preferably in English (and less in math notation);

3. Not Comprehensive lecture on tech DS

More focus given around predictive and prescriptive analytics – there are many more technical AI disciplines that can be used to create value in a DS Project;

4. Ethics in a Data Science Context

Ethics is a multi-disciplinary topic however in this presentation is discussed mostly in an applied algorithmic policy context.

Mini-breakout sessions

1	2	3	4
XXXX	XXXX	XXXX	XXXX
XXXX	XXXX	XXXX	XXXX
XXXX	XXXX	XXXX	XXXX

- We'll have a mini breakout session for each module to experiment with lectured concepts;
- Groups are supposedly mixing different teams to maximize diversity! (any improvements please let me know)
- We'll use Zoom breakout feature and I'll set a counter – there'll be a briefing before to make sure we're on same page;
- We'll do the discussions together using jamboard – there'll be a link provided for each challenge.
- Presentation (and related materials for challenges) in a share for your convenience:

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

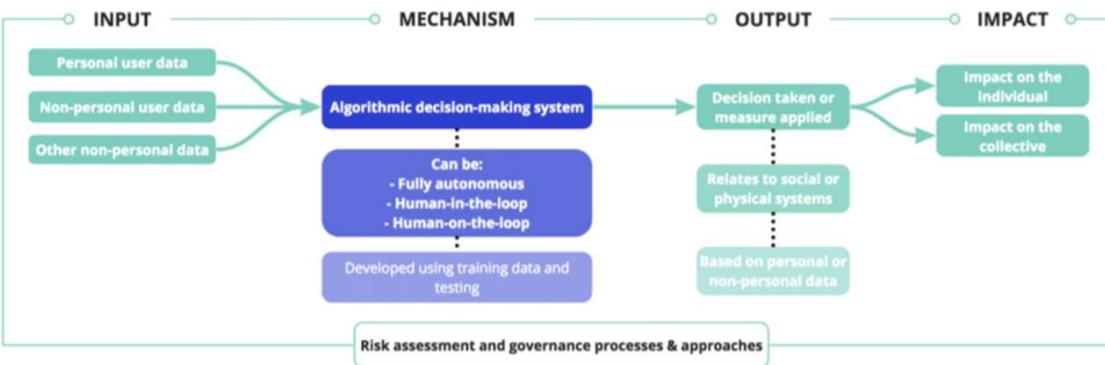
Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		



Algorithmic Decision Making

Also known as predictive analytics, can be seen as statistical risk assessment, i.e., we aim to predict the probability of an outcome (a given event occurs in the future).

Working Definition



A software system – including its testing, training and input data, as well as associated governance processes – that, autonomously or with human involvement, takes decisions or applies measures relating to social or physical systems on the basis of personal or non-personal data, with impacts either at the individual or collective level.[1]

Applications

Scoring applications - e.g. credit, fraud, insurance, hiring, college admission, bail determination, preventive healthcare, etc

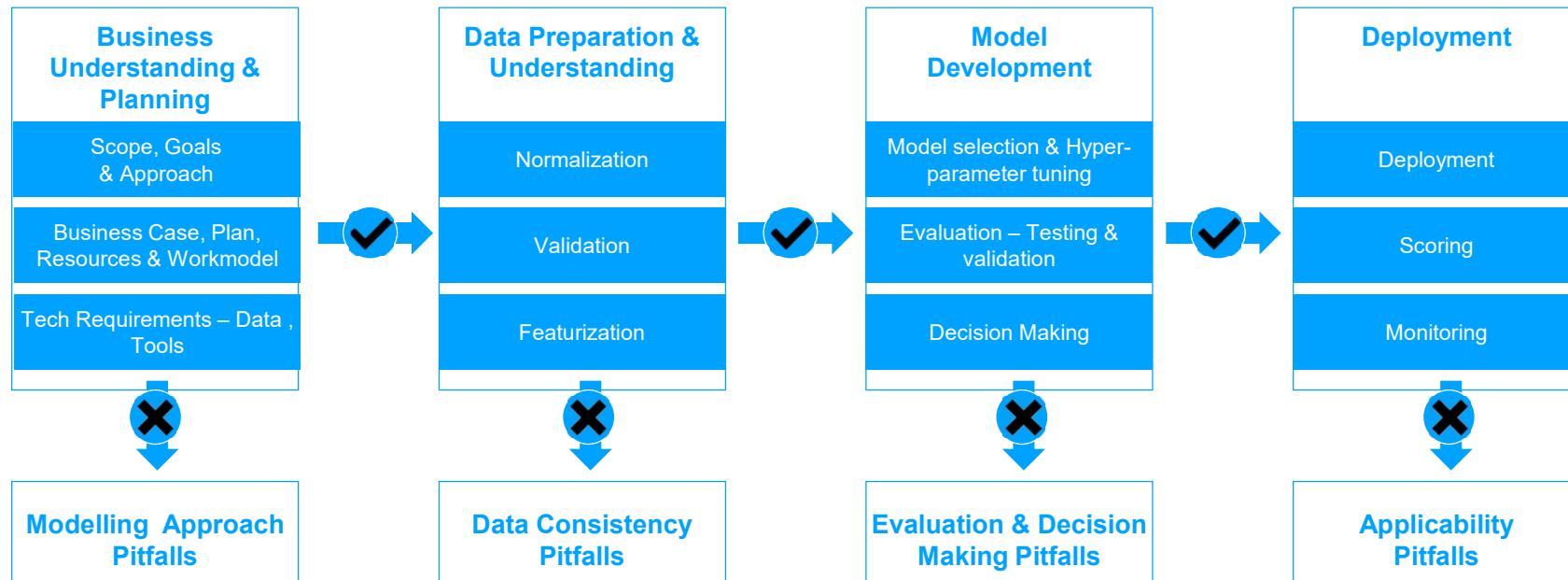
- Different types of **search engines**, including general, semantic, and meta search engines;
- **Aggregation applications**, such as news aggregators, which collect, categorise and regroup information from multiple sources into one single point of access;
- **Forecasting, profiling and recommendation applications**, including targeted advertisements, selection of recommended products or content, personalised pricing and predictive policing;
- **Content production applications** (e.g. algorithmic journalism)
- **Filtering and observation applications**, such as spam filters, malware filters, and filters for detecting illegal content in online environments and platforms.
- Other 'sense-making' applications, crunching data and drawing insights.



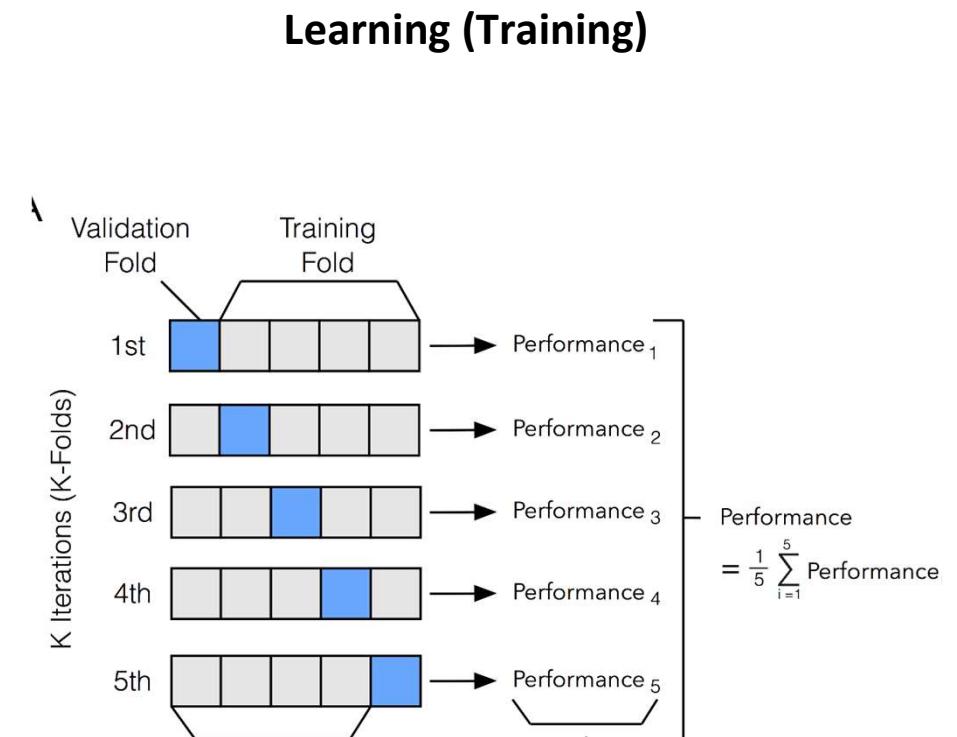
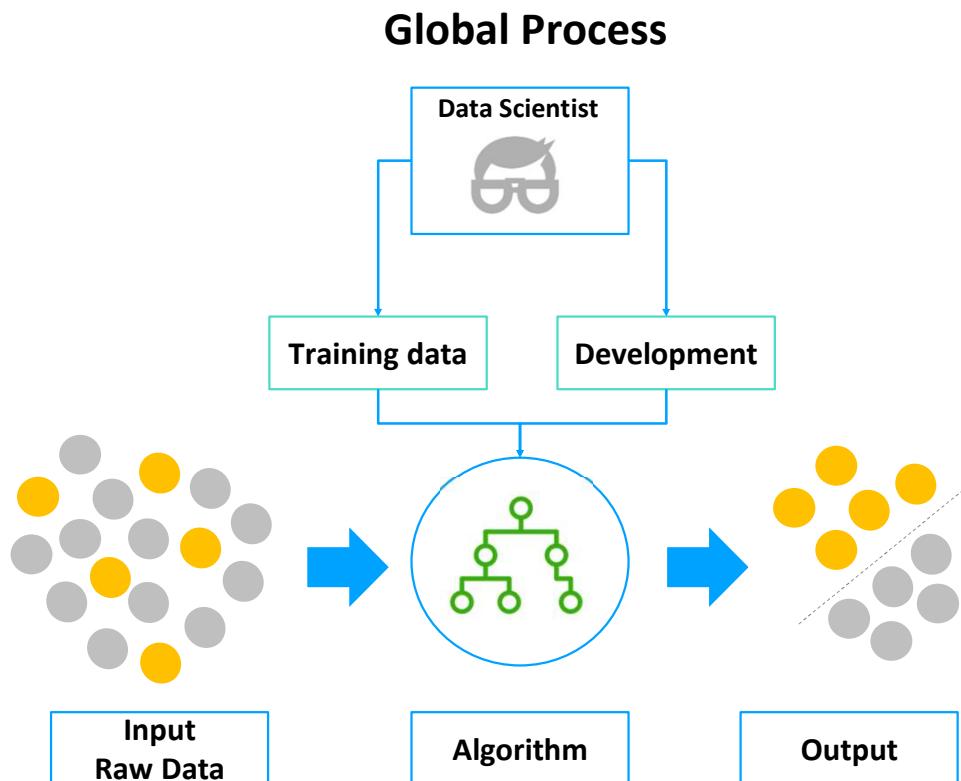
Common Pitfalls in Predictive Analytics

Data Science is a multidisciplinary field that may fail in unexpected ways during development or operational phase

Illustrative

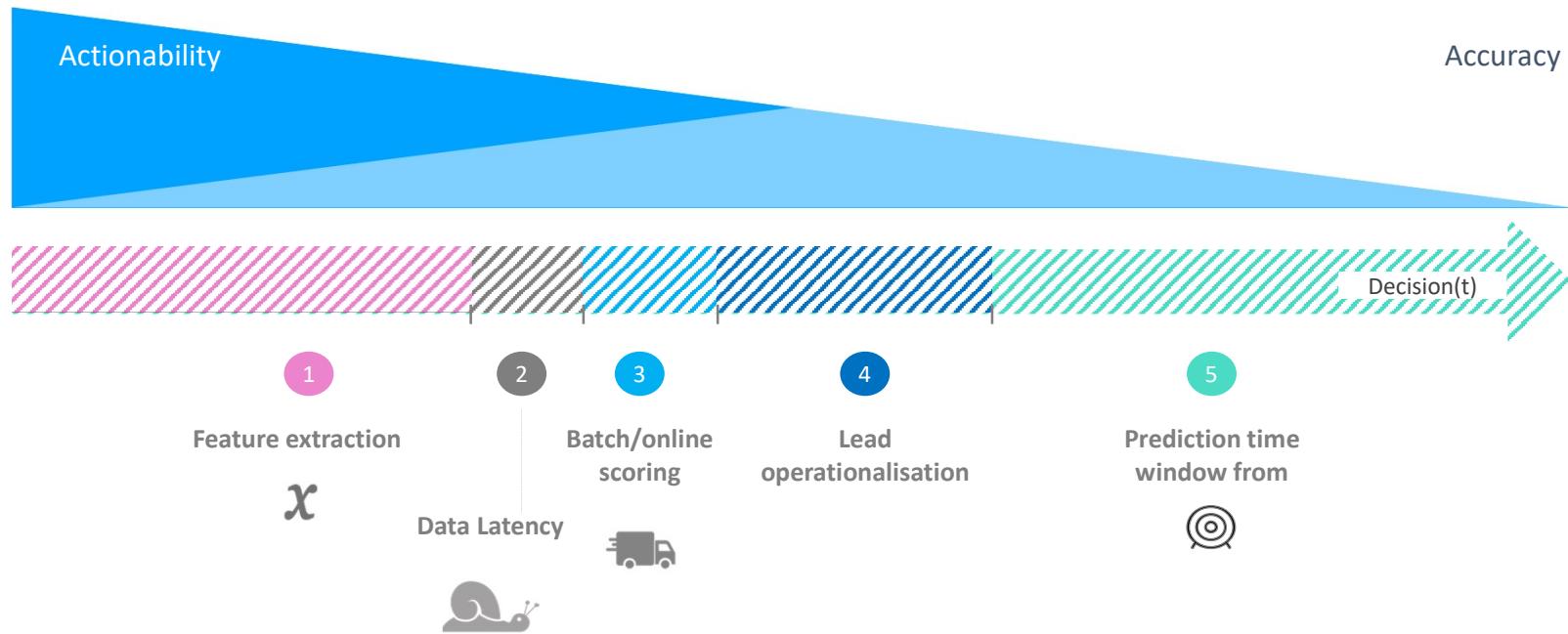


Classification Refresher



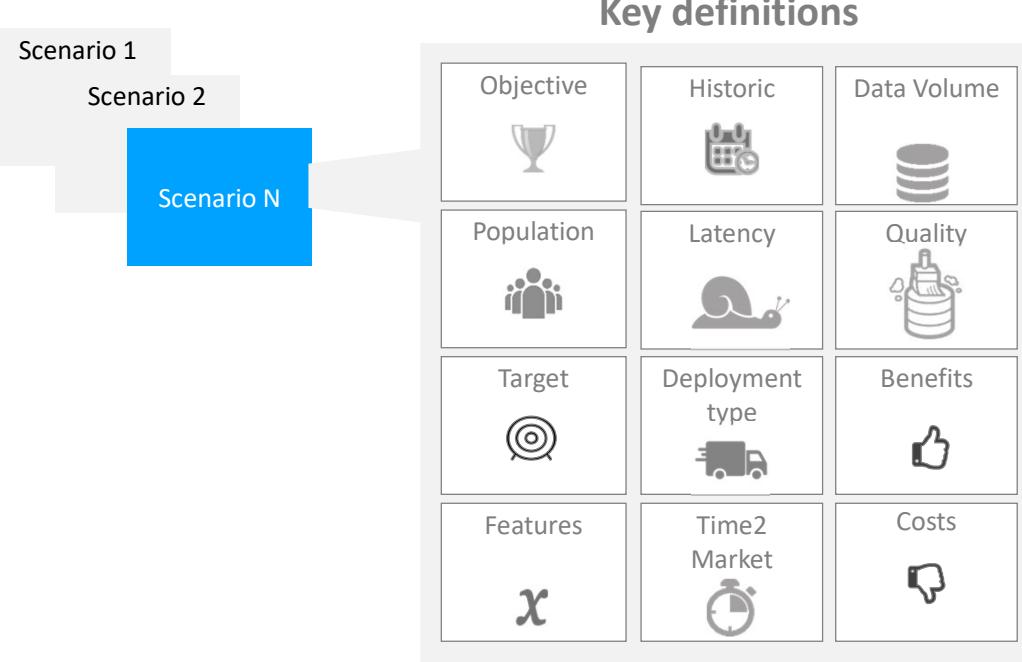
Classification Refresher

There are a few degrees of freedom that need to be discussed to set-up the modelling approach



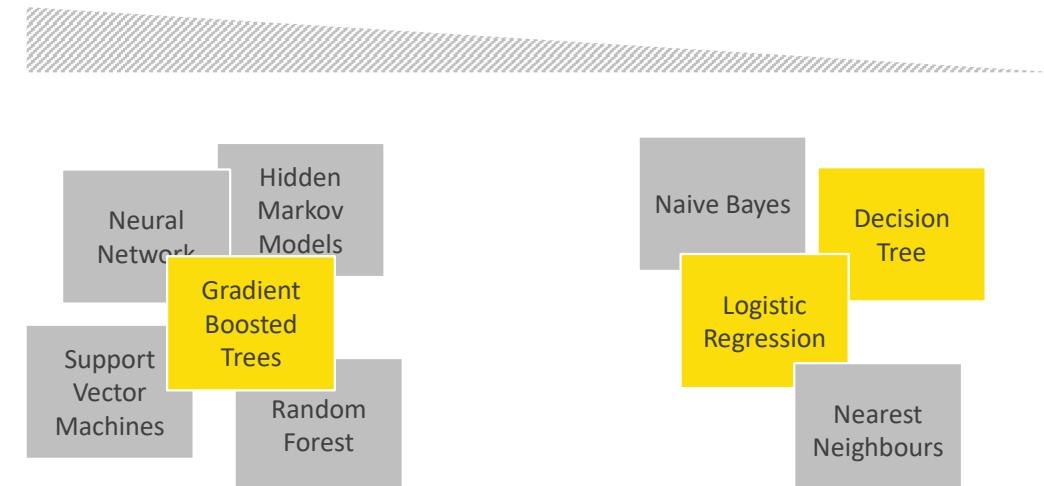
Classification Refresher

Scenario Evaluation



Algorithm Selection

+Prediction Accuracy

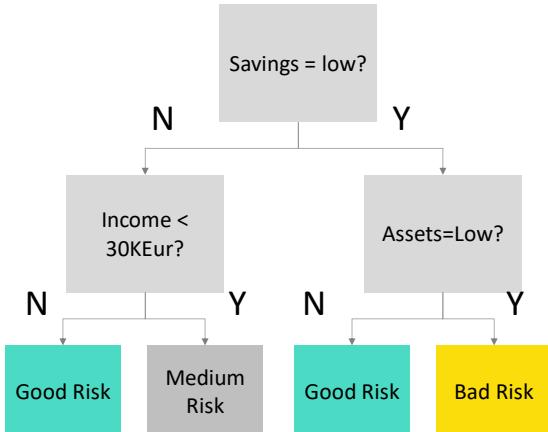


+ more than 200 Models available in programmatic frameworks

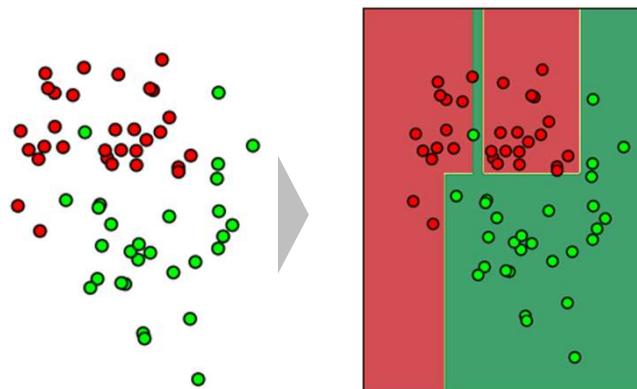
To be
detailed

Classification Refresher: Decision Tree

A. Intuition



B. Decision boundary



- Decision trees are built creating decision nodes that select the best variable at each step that best separates each class in the training dataset.

- Space is divided in “rectangles” and each one is labelled with the more dominant k-class.
- Non-Linear relationships can be found on data.

C. Configuration & Strengths/Weaknesses

Parametrization

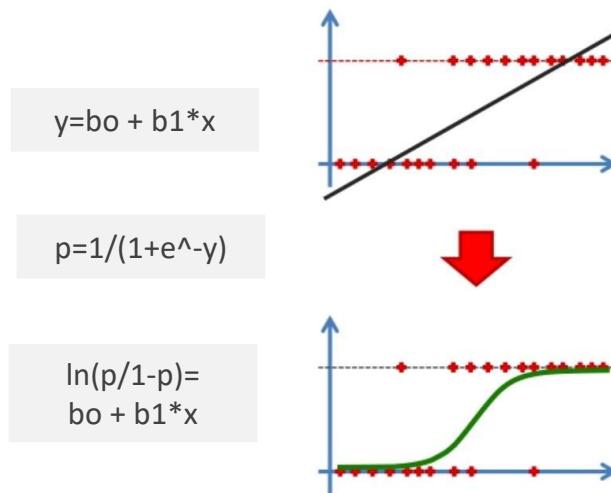
- Depth of tree
- Minimum samples Split
- Minimum samples leaf
- Maximum no features

Strengths / Weaknesses

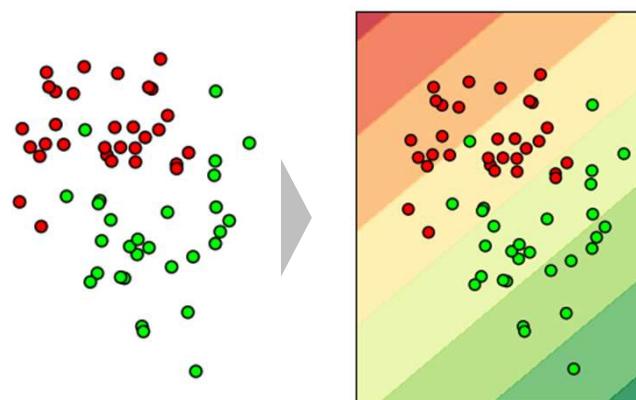
- | |
|---|
| (+) Highly Interpretable |
| (+) Transparent; |
| (+) little effort in data preparation (eg: can deal both continuous and discrete data); |
| (-) less accurate while compared to other available predictors; |

Classification Refresher: Logistic Regression

A. Intuition



B. Decision boundary



- We can think of Logistic Regression as an extension of the linear regression model but where the output is fitted through a sigmoid into a value (prob) ranging from 0 to 1.

C. Configuration & Strengths/Weaknesses

Parametrization

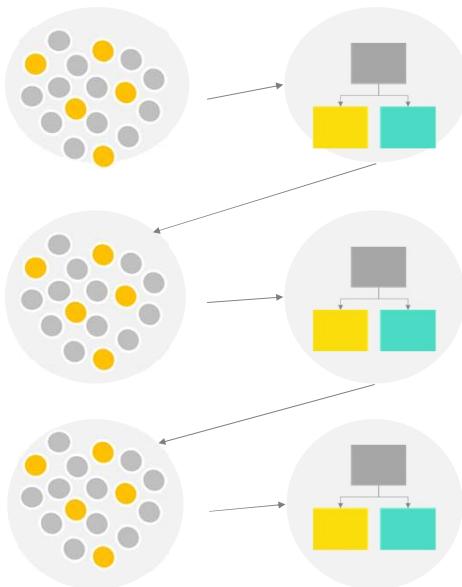
- Regularization parameters to control trade-off between overfit and simplicity (underfit)

Strengths / Weaknesses

- | | |
|---|--|
| <ul style="list-style-type: none"> (+) simple and efficient; (+) low variance; (+) it provides probability score for observations; | <ul style="list-style-type: none"> (-) doesn't handle large number of features; (-) requires feature processing (eg: non-linear features, categorical, multicollinearity, etc) |
|---|--|

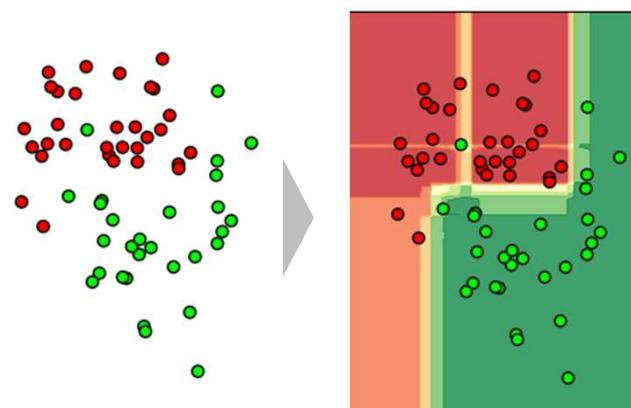
Classification Refresher: Gradient Boosted Trees

A. Intuition



- Data is fitted first in a simple tree model, then subsequent models focus on resolving classification errors;
- In the end all the predictors are combined through weights in a single model.

B. Decision boundary



- Decision boundaries are more complex since multiple trees are fitted on areas with prediction error is greater;

C. Configuration & Strengths/Weaknesses

Parametrization

- Number of trees;
- Learning Rate
- Maximum depth of tree;
- Min number samples to split tree node;
- Min number samples for a node;
- Max number features to consider for a split;

Strengths / Weaknesses

- | |
|---|
| <p>(+) accuracy – extensions of the boosting idea have a track record of winning Kaggle competitions;</p> <p>(+) similar with DT - little effort in data preparation</p> <p>(-) interpretability</p> <p>(-) takes more time to train since it's sequential</p> <p>(-) more prone to overfit</p> |
|---|



Classification Refresher: Evaluation Metrics

Evaluation metric is a number that measures the performance that your machine learning model when it comes to assigning observations to certain classes.

		True condition			
Total population	Condition positive	Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$	
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
	True positive rate (TPR), Recall, Sensitivity, probability of detection, Power = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$	F_1 score = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
	False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$		

Most Common Derived Ratios - https://en.wikipedia.org/wiki/Evaluation_of_binary_classifiers



Classification Refresher: Evaluation Metrics

Rank	Score	Label	Predict
1	0.997	1	1
2	0.993	1	1
3	0.986	1	1
4	0.982	1	1
5	0.971	0	0
6	0.965	1	0
7	0.964	0	0
8	0.961	0	0
9	0.953	0	0
10	0.932	1	0
11	0.918	0	0
12	0.873	0	0
13	0.854	0	0
14	0.839	0	0
15	0.777	0	0
16	0.723	0	0
17	0.634	0	0
18	0.512	0	0
19	0.487	0	0
20	0.473	0	0

Predicted Positive: 4
Predicted Negative: 16

Population

- Total Label Positives: 6
- Total Label Negatives: 14
- Prevalence $6/20=0.3$

For threshold > 0.980 or top k=4:

Predicted Condition	True Condition		
	Total Population	Positive	Negative
Positive	TP 4	FP 0	
Negative	FN 2	TN 14	

True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$
False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$

- True Positive Rate (Recall) = $4/6 = 0.66$
- False Positive Rate: $0/14=0$
- False Negative Rate = $2/6=0.33$
- True Negative Rate = $14/14=1.0$
- Precision = $4/4 =1.0$

Adapted from Pedro Saleiros's "Tutorial: Fairness in decision-making with AI: a practical guide & hands-on tutorial using Aequitas
<https://www.youtube.com/watch?v=yOR71zBm3Uc>



Classification Refresher: Evaluation Metrics

Rank	Score	Label	Predict
1	0.997	1	1
2	0.993	1	1
3	0.986	1	1
4	0.982	1	1
5	0.971	0	1
6	0.965	1	1
7	0.964	0	1
8	0.961	0	1
9	0.953	0	1
10	0.932	1	1
11	0.918	0	0
12	0.873	0	0
13	0.854	0	0
14	0.839	0	0
15	0.777	0	0
16	0.723	0	0
17	0.634	0	0
18	0.512	0	0
19	0.487	0	0
20	0.473	0	0

Predicted Positive: 10
Predicted Negative: 10

Population

- Total Label Positives: 6
- Total Label Negatives: 14
- Prevalence $6/20=0.3$

For threshold > 0.920 or top k=10:

Predicted Condition	True Condition		
	Total Population	Positive	Negative
Positive	TP 6	FP 4	
Negative	FN 0	TN 10	

True positive rate (TPR),
Recall, Sensitivity,
probability of detection,
Power
 $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$

False positive rate (FPR),
Fall-out,
probability of false alarm
 $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$

False negative rate (FNR),
Miss rate
 $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$

Specificity (SPC),
Selectivity, True negative
rate (TNR)
 $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$

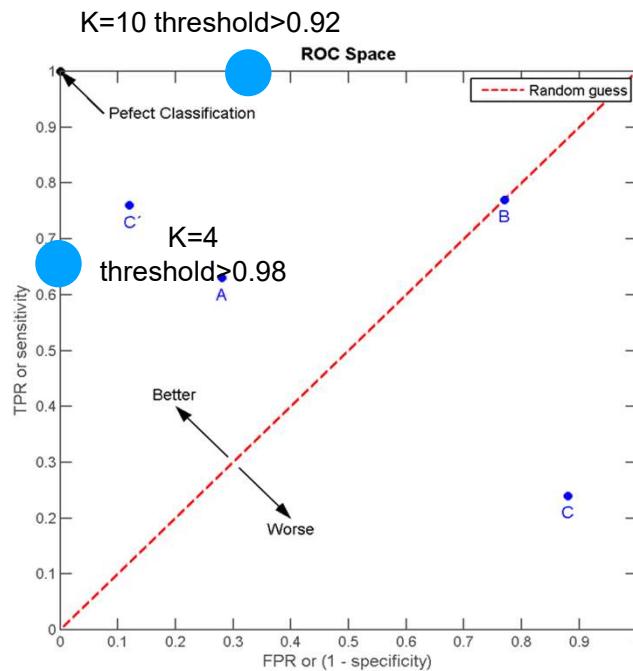
- True Positive Rate (Recall) = $6/6 = 1.0$
- False Positive Rate: $4/14=0.29$
- False Negative Rate = $0/6=0$
- True Negative Rate = $10/14=0.71$
- Precision = $6/10 =0.6$

Adapted from Pedro Saleiros's "Tutorial: Fairness in decision-making with AI: a practical guide & hands-on tutorial using Aequitas
<https://www.youtube.com/watch?v=yOR71zBm3Uc>



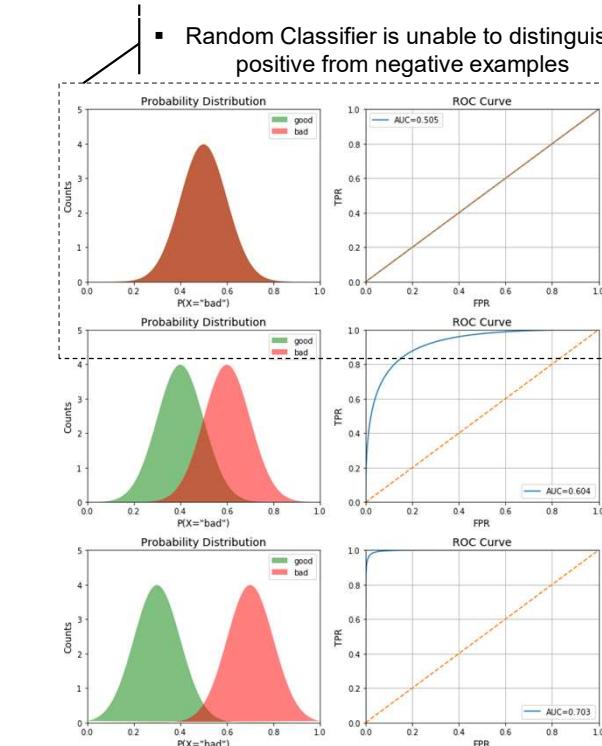
Classification Refresher: Evaluation Metrics

A. ROC COORDINATES



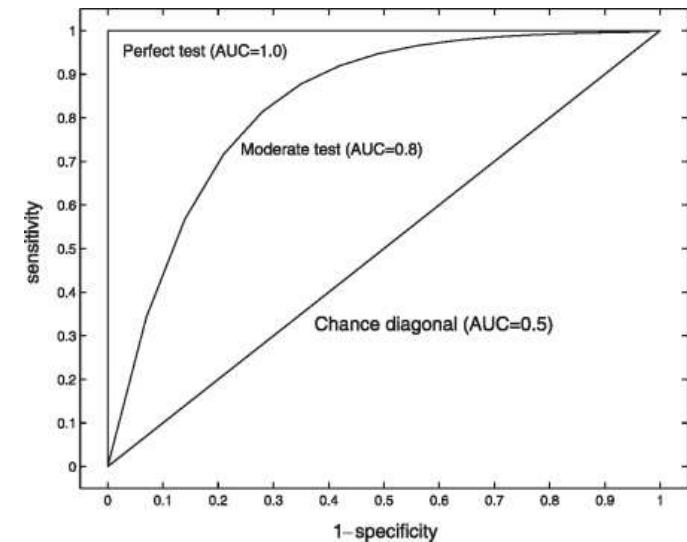
- To get ROC coordinates we define a threshold for 1/0 predictions and then we calculate the False Positive rate and Recall values for that threshold.

B. ROC PROFILES



- ROC curve is an analytical tool to assess model performance (if your dataset is not highly unbalanced)

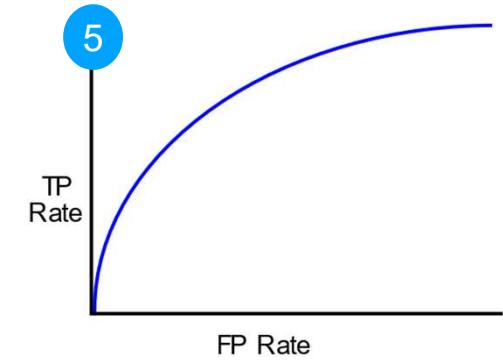
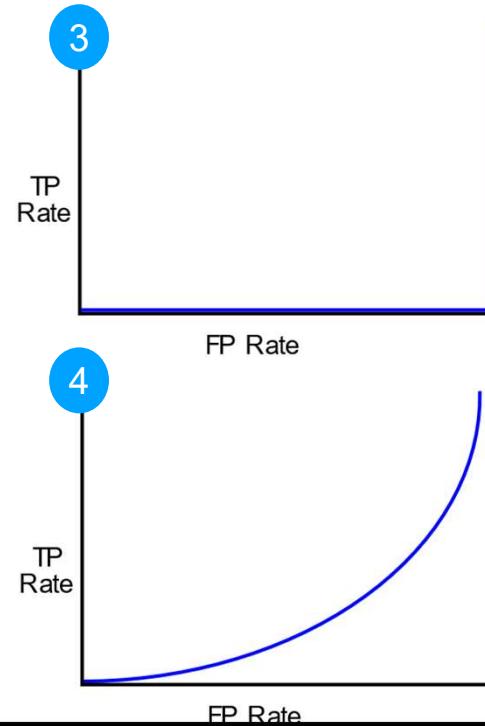
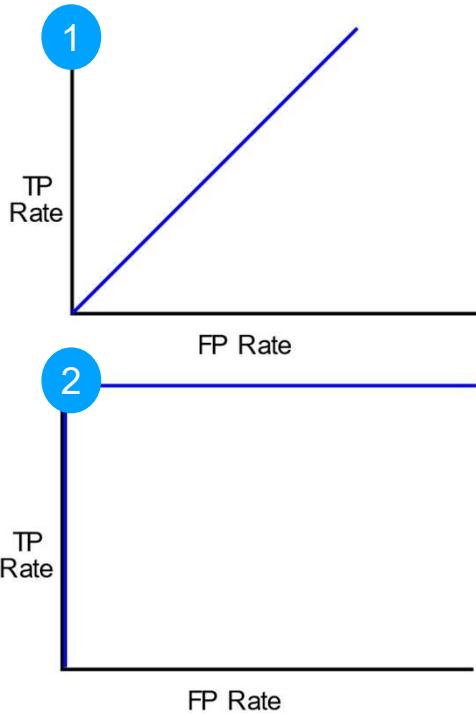
C. AUC



- AUC is the area under the ROC expresses how much a model can distinguish two classes along the score distribution in a given sample.

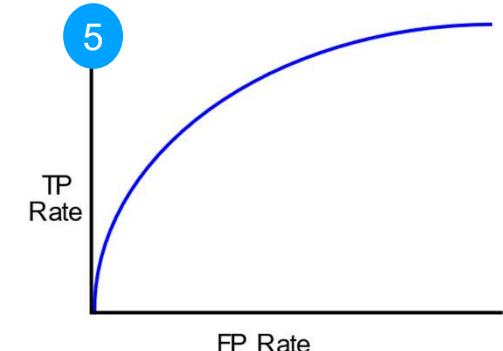
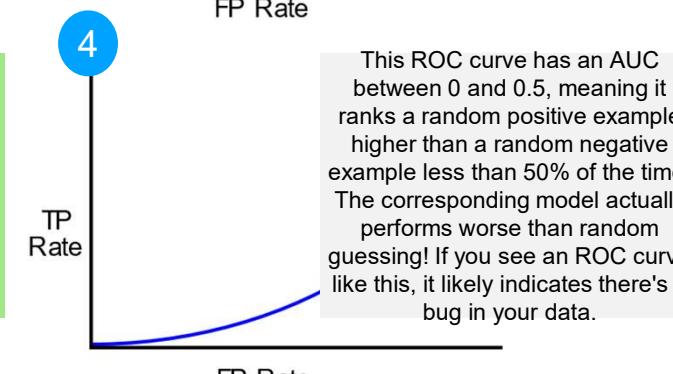
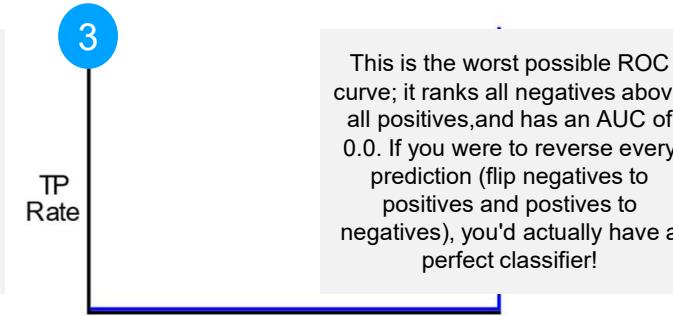
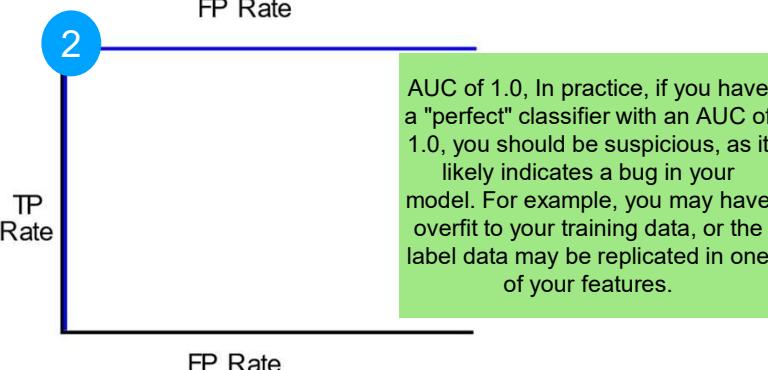
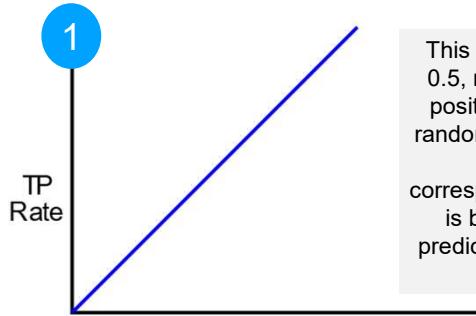
Classification Refresher: Evaluation Metrics

Which of the following ROC curves produce AUC values greater than 0.5?



Classification Refresher: Evaluation Metrics

Which of the following ROC curves produce AUC values greater than 0.5?





Classification Refresher: Evaluation Metrics

Cheat Sheet

Binary classification performances measure cheat sheet
 Damien François - v1.0 - 2009 (damien.francois@uclouvain.be)

Confusion matrix for two possible outcomes p (positive) and n (negative) <table border="1"> <thead> <tr> <th colspan="2"></th><th colspan="2">Actual</th><th rowspan="2">Total</th></tr> <tr> <th colspan="2"></th><th>p</th><th>n</th></tr> </thead> <tbody> <tr> <th rowspan="2">Predicted</th><th>p'</th><td>true positive</td><td>false positive</td><th>p</th></tr> <tr> <th>n'</th><td>false negative</td><td>true negative</td><th>n</th></tr> <tr> <th>total</th><td>P'</td><td>N'</td><td></td><td></td></tr> </tbody> </table> <p>Classification accuracy $(TP + TN) / (TP + TN + FP + FN)$ Error rate $(FP + FN) / (TP + TN + FP + FN)$</p>			Actual		Total			p	n	Predicted	p'	true positive	false positive	p	n'	false negative	true negative	n	total	P'	N'			True positive rate: proportion of actual positives which are predicted positive $TP / (TP + FN)$ True negative rate: proportion of actual negative which are predicted negative $TN / (TN + FP)$	Youden's index: arithmetic mean between sensitivity and specificity $sensitivity - (1 - specificity)$ Matthews correlation: correlation between the actual and predicted $(TP \cdot TN - FP \cdot FN) / ((TP+FP)(TP+FN)(TN+FP)(TN+FN))^{1/2}$ comprised between -1 and 1	(Cumulative) Lift chart: plot of the true positive rate as a function of the proportion of the population being predicted positive, controlled by some classifier parameter (e.g. a threshold)	
		Actual		Total																							
		p	n																								
Predicted	p'	true positive	false positive	p																							
	n'	false negative	true negative	n																							
total	P'	N'																									
Paired criteria	Precision: (or Positive predictive value) proportion of predicted positives which are actual positive $TP / (TP + FP)$ Recall: proportion of actual positives which are predicted positive $TP / (TP + FN)$	Positive likelihood: likelihood that a predicted positive is an actual positive $sensitivity / (1 - specificity)$ Negative likelihood: likelihood that a predicted negative is an actual negative $specificity / (1 - sensitivity)$	Discriminant power: normalised likelihood index $\sqrt{3} / \pi \cdot (\log(sensitivity / (1 - specificity)) + \log(specificity / (1 - sensitivity)))$ <1 = poor, >3 = good, fair otherwise	Relationships $sensitivity = recall = true positive rate$ $specificity = true negative rate$ $BCR = \frac{1}{2} \cdot (sensitivity + specificity)$ $BCR = 2 \cdot Youden's\ index - 1$ $F\text{-measure} = F\text{measure}$ $Accuracy = 1 - error\ rate$																							
Combined criteria	BCR: Balanced Classification Rate $\frac{1}{2} (TP / (TP + FN) + TN / (TN + FP))$ BER: Balanced Error Rate, or HTER: Half Total Error Rate: $1 - BCR$ F-measure: harmonic mean between precision and recall $2 \cdot (precision \cdot recall) / (precision + recall)$ F_β-measure: weighted harmonic mean between precision and recall $(1+\beta^2) TP / ((1+\beta^2) TP + \beta^2 FN + FP)$	AUC: The area under the ROC is between 0 and 1	Graphical tools ROC curve: receiver operating characteristic curve : 2-D curve parametrized by one parameter of the classification algorithm, e.g. some threshold in the « true positive rate » space AUC: The area under the ROC is between 0 and 1	References Sokolova, M. and Lapalme, G. 2009. A systematic analysis of performance measures for classification tasks. <i>Inf. Process. Manage.</i> 45, 4 (Jul. 2009), 427-437. Demsar, J.: Statistical comparisons of classifiers over multiple data sets. <i>Journal of Machine Learning Research</i> 7 (2006) 1-30																							
Sensitivity: proportion of actual positives which are predicted positive $TP / (TP + FN)$ Specificity: proportion of actual negative which are predicted negative $TN / (TN + FP)$	 	The harmonic mean between specificity and sensitivity is also often used and sometimes referred to as F-measure.																									



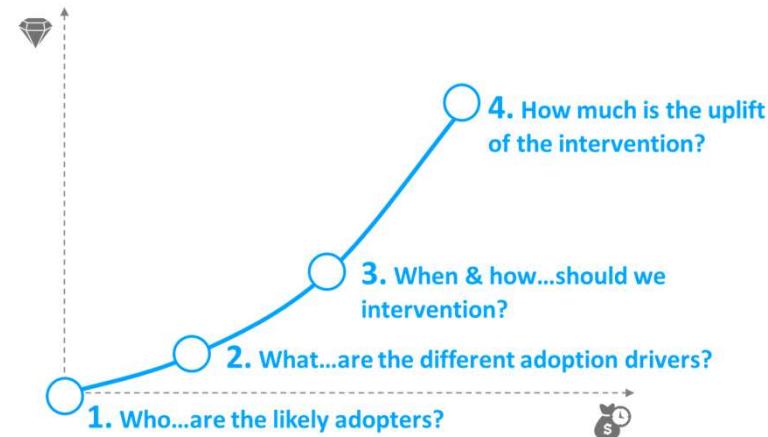
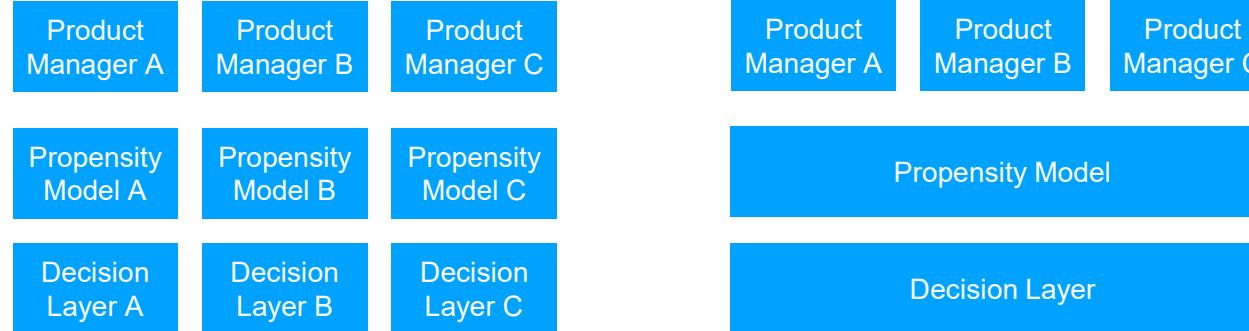
Pitfall #1:Lack of breadth and depth in modelling approach

Can result in proliferation of similar pattern use cases that will take resources for new opportunities (or the inclusion of more analytical features)

Pitfall: Data Science products can mirror your organization functions spending development/maintenance resources on a similar problem

Example: Next-Best-Offer related projects can have various stakeholders but the pattern is always the same: Customer has product X and we want to offer product Z,Y

Solution: (usually) go for a global product that addresses a fundamental pattern and invest in layers of analytical sophistication/experimentation as a better strategy to deliver more value.





Pitfall #2: Not looking at data before modelling

Each model can have different set of assumptions that have to be met, looking at data helps us to understand its characteristics and form an hypothesis of which model is best to capture those characteristics.

Pitfall: Not checking model assumptions can lead to poor results – “garbage in, garbage out”

Example: depending on the model the a not met assumption can affect training behavior

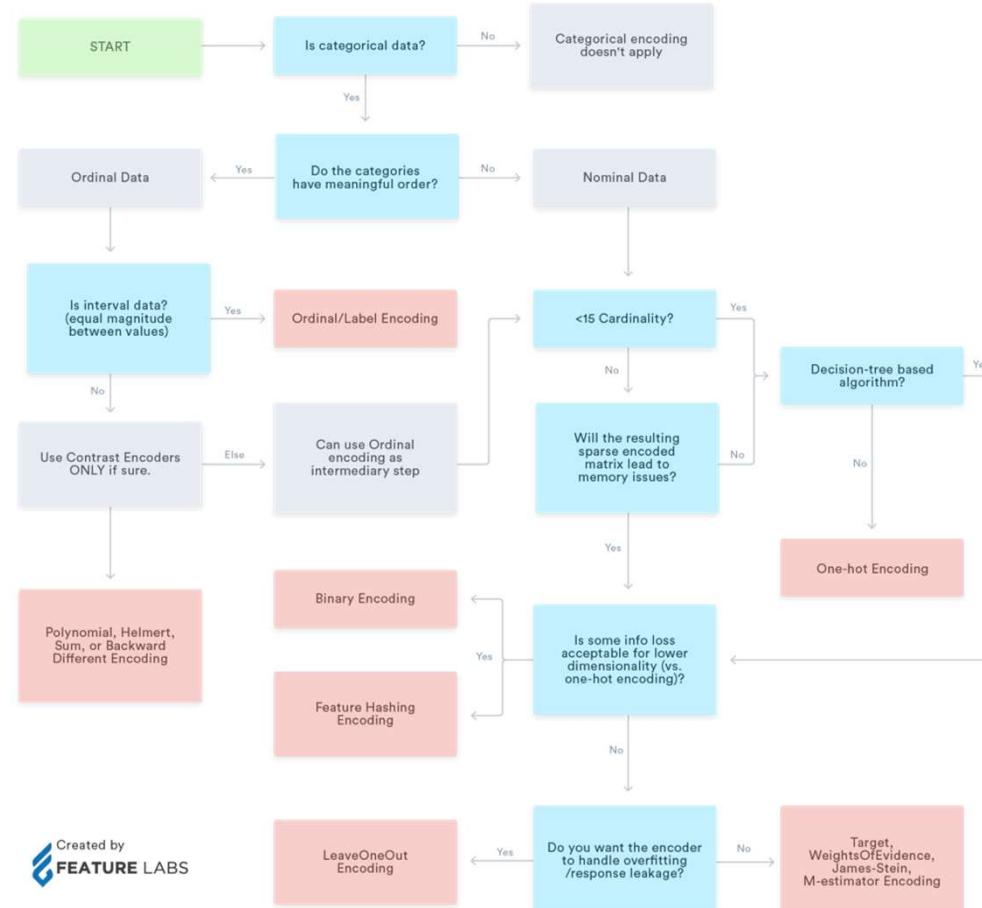
Solution: conscious development safeguards to address these matters (usually a baseline model more robust)

Model(s)	Assumption	Impact	Solution
Regression	Independence between features	Model will not calculate correctly coefficients and will not be able to identify which x variables have most statistical influence on the y var.	Plotting correlations between variables to diagnose and apply feature selection if issue is detected.
Distance (or Gradient Descent) based models	Features are in same scale	<ul style="list-style-type: none"> In distance based models features with higher scale will dominate solution; In models with gradient descent optimization will may end up taking long time to learn. 	Scale features using a normalization process such as Z (better for outliers!) or Min-Max normalization.
SVMs, NNets, GBTs	Input features are numerical.	<ul style="list-style-type: none"> Training process will left out categorical variables if they're not encoded properly. 	Experiment with different encoding methods depending on categorical/ordinal data, cardinality and overfitting handling.



Pitfall #2: Not looking at data before modelling

Categorical Encoding Methods Cheat-Sheet

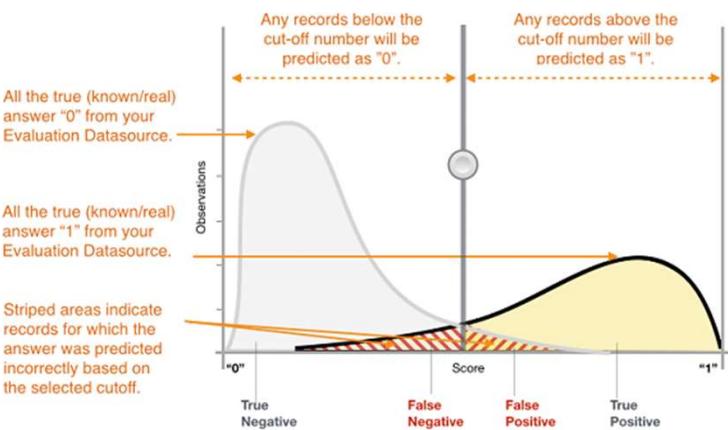




Pitfall #3: Out-of-the-Box Binary Classification

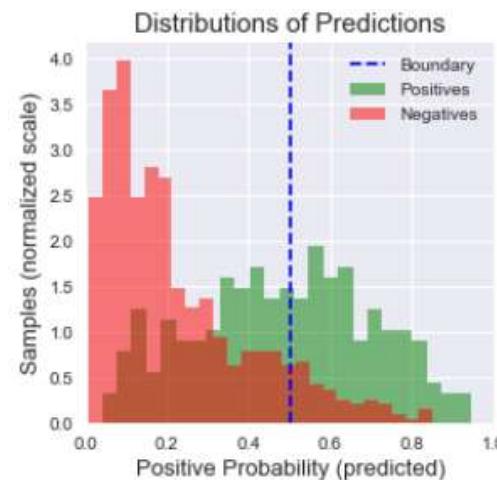
Decision-Making is different than just classification (e.g. images or text). Consequential decision-making implies an action informed by the predictions and this actions are constrained by a budget or capacity.

Pitfall: Using out-of-the-box binary predictions assume that if the model score is > 0.5 then it is 1, else it is 0.



AWS Dev Guide/Binary Classification –
<https://docs.aws.amazon.com/machine-learning/latest/dg/binary-classification.html>

Example: Predicting Bank Loan Default – Positive Class = Payment Default



German Credit Risk Dataset – Jupyter Notebook

Solution: score all instances and select those above a defined threshold either for batch or online predictions.

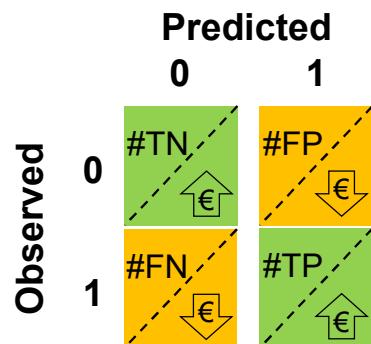
Rank	Score	Label	Predict				
1	0.997	1	1				
2	0.993	1	1				
3	0.986	1	1				
4	0.982	1	1				
5	0.971	0	0				
6	0.965	1	0				
7	0.964	0	0				
8	0.961	0	0				
9	0.953	0	0				
10	0.932	1	0				
11	0.918	0	0				
12	0.873	0	0				
13	0.854	0	0				
14	0.839	0	0				
15	0.777	0	0				



Pitfall #4: Missing link betw. model threshold and value

xxxxxxxxx

Pitfall: Value created by the model is not optimized accordingly with the model score threshold

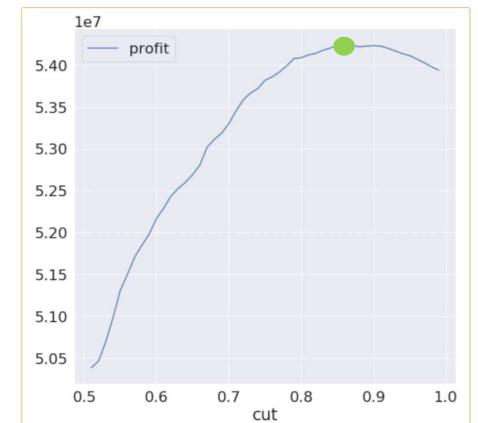


Example: Loan application illustrates which costs are critical in threshold (or model) selection – desirably with a low number of False negatives.

TN Revenue	FP Cost
<ul style="list-style-type: none"> Average annual income for a payer 12.000Eur	<ul style="list-style-type: none"> Opportunity cost of denying a loan to a payer -12.000Eur
FN Cost	TP Cost
<ul style="list-style-type: none"> Amount lost due to sanction a defaulter's loan -10.00.000	<ul style="list-style-type: none"> Forfeiting application processing fee -3.000Eur

$$\text{Net Revenue} = \#TN * \text{Revenue} - \#FP * \text{Cost} - \#FN * \text{Cost} - \#TP * \text{Cost}$$

Net Revenue vs Probability Predictions



Solution: Plotting the Net Revenue vs Threshold reveals the best cut-off to maximize profit.

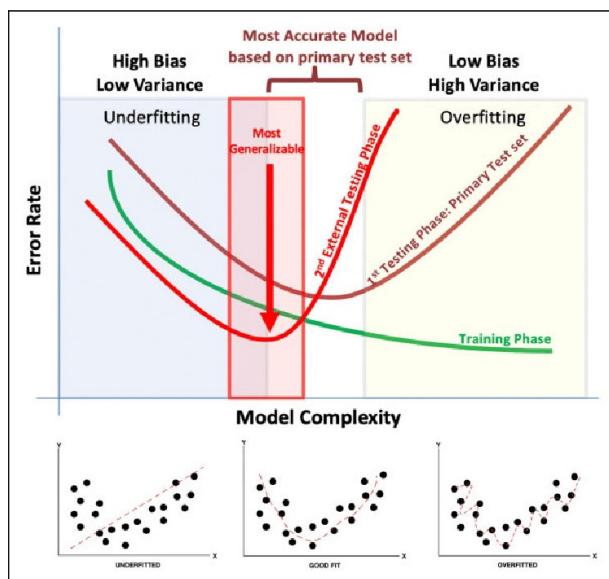




Pitfall #6: Diagnose bias and variance trade-off in learning

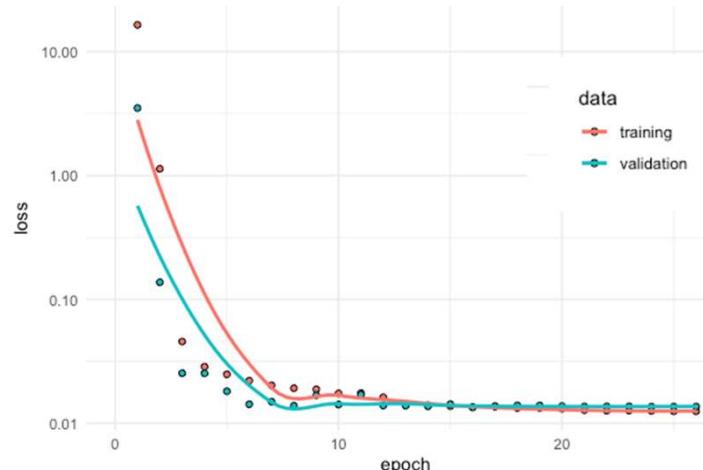
The bias-variance dilemma is a widely known problem in the field of machine learning, looking at the profile of learning curves a proper diagnostic can be found and solutions to fix the issue.

Pitfall: Not assessing properly model behavior to improve the learning task



Example: "Optimal Fit" is the goal of the learning algorithm

Example of learning curve showing near optimality assuming we have adequately minimized the loss score.



Solution: In this scenario there's no need to continue training (or it'll lead to overfitting).

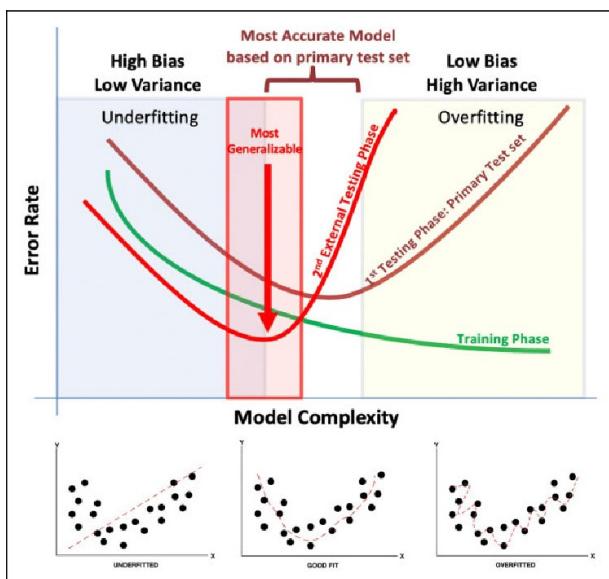
- The plot of training loss decreases to a point of stability.
- The plot of validation loss decreases to a point of stability.
- The generalization gap is minimal (nearly zero in an ideal situation).



Pitfall #6: Diagnose bias and variance trade-off in learning

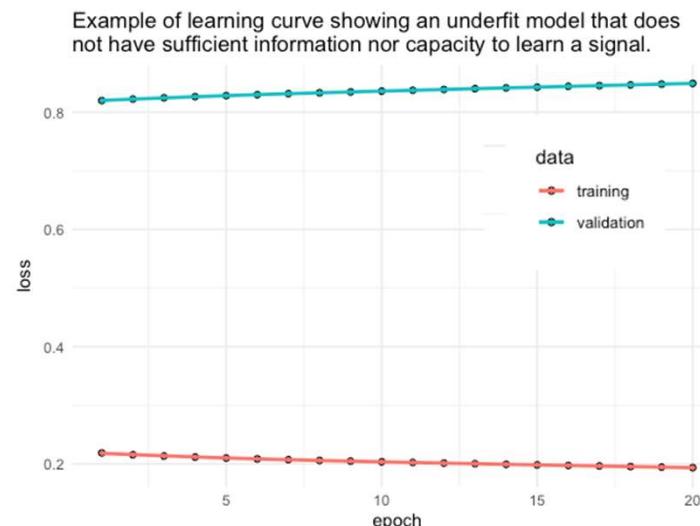
The bias-variance dilemma is a widely known problem in the field of machine learning, looking at the profile of learning curves a proper diagnostic can be made and solutions to fix the issue.

Pitfall: Not assessing properly model behavior to improve the learning task



Example: "Underfit"

is a model that hasn't learned the training dataset to get a low error value – flat curve.



Solution: In this scenario model isn't capable of learning.

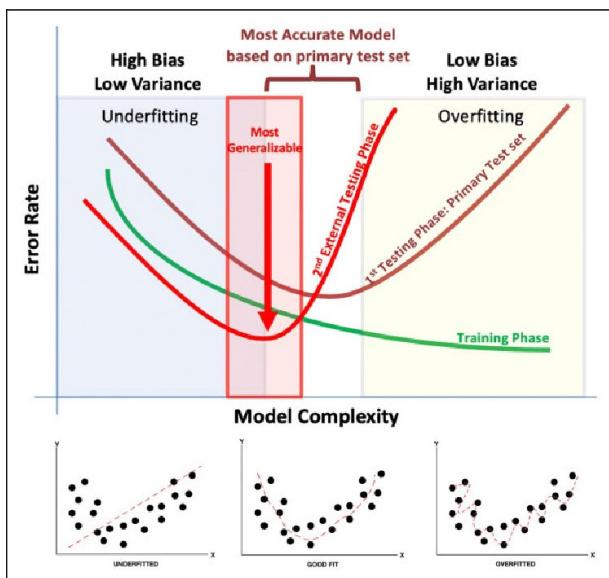
- Add more observations. You may not have enough data for the existing patterns to become strong signals.
- Add more features
- If you have explicit regularization parameters specified (i.e. dropout, weight regularization), remove or reduce these parameters;
- Model capacity may not be large enough to capture and learn existing signals – hyperparameter optimization



Pitfall #6: Diagnose bias and variance trade-off in learning

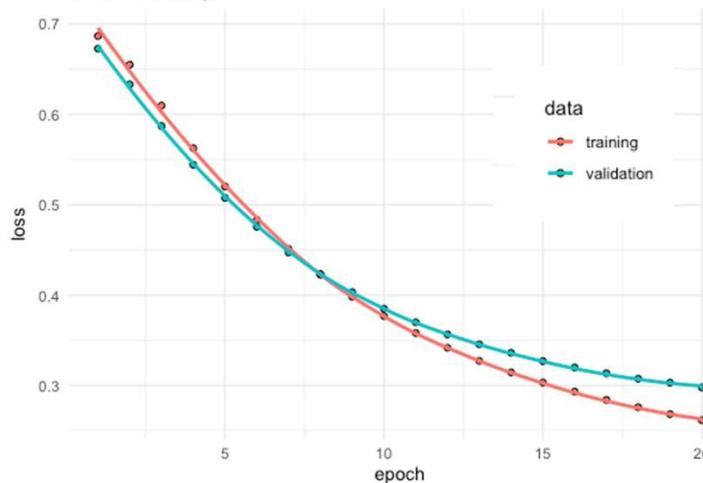
The bias-variance dilemma is a widely known problem in the field of machine learning, looking at the profile of learning curves a proper diagnostic can be made and solutions to fix the issue.

Pitfall: Not assessing properly model behavior to improve the learning task



Example: "Underfit" learning curves still on decrease trend.

Example of learning curve showing an underfit model that requires further training.



Solution: In this scenario learning process was stopped prematurely

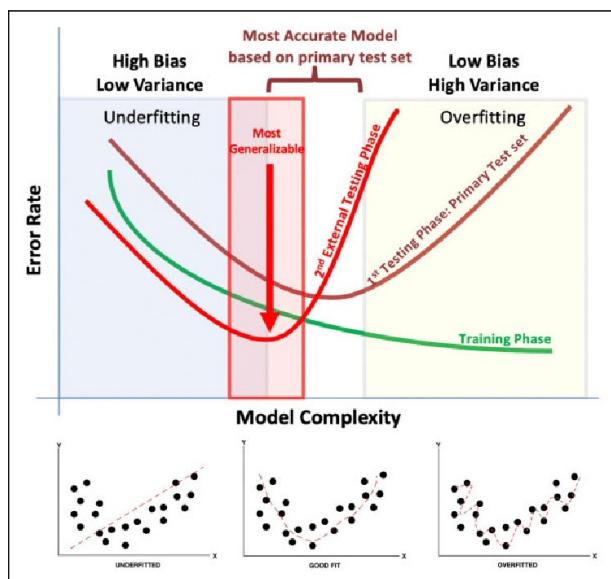
- Increase the number of epochs until the validation curve has stopped improving.
- If it is taking a long time to reach a minimum for the validation curve, increase the learning rate;



Pitfall #6: Diagnose bias and variance trade-off in learning

The bias-variance dilemma is a widely known problem in the field of machine learning, looking at the profile of learning curves a proper diagnostic can be made and solutions to fix the issue.

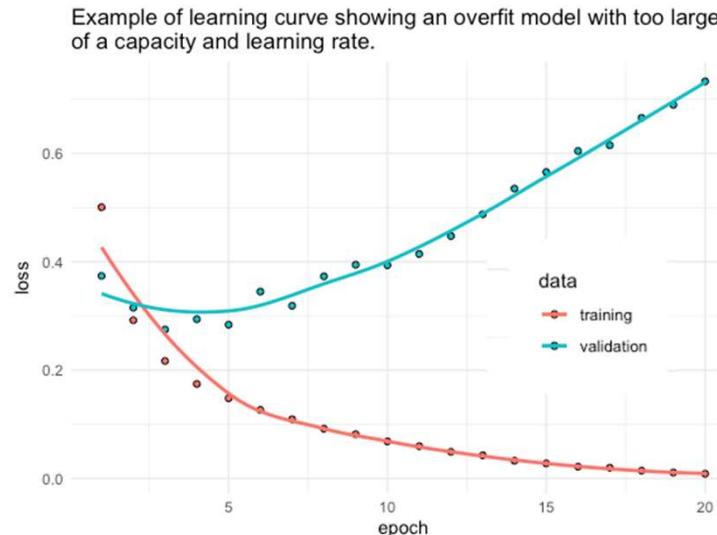
Pitfall: Not assessing properly model behavior to improve the learning task



Example: "Overfit"

Training loss continues to decrease while validation loss has a strong U-shape curve

Solution: overfitting is not necessarily a bad thing signals model extracted all what had to be learned up to a certain point



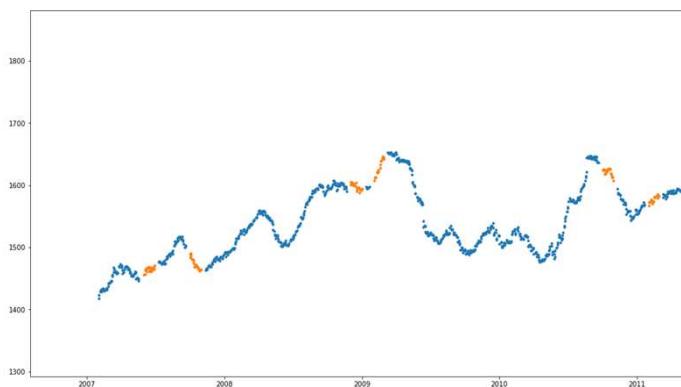
- Apply regularization techniques depending on the algorithm used – eg: if it is NN reducing the number and/or size of hidden layers;
 - Early stopping the model once the validation curve has stopped improving.- using the model parameters for the best training metric;



Pitfall #7: Time series validation

Standard cross-validation methods do not reflect sequential discovery of time series, creating a risk of “future leakage”

Pitfall: data with a sequential structure must be accounted for in model validation otherwise results in production can be quite different from development

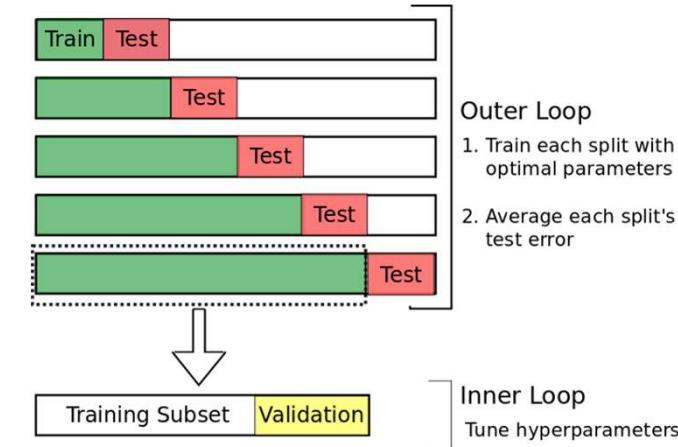


Example: estimating Customer usage/churn or revenue/lifetime value using time-series data is a typical use case

Customer	Day	Usage
1	1	321
1	2	313
1	3	233
1	4	244
1	5	140
1	6	80
1	7	30
1	8	2
1	9	0

Solution: ensure that cross-validation folds are created respecting the time sequence to get model validation results as close as production setting.

Nested Cross-Validation



Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		



Group Challenge – Practice Metric Evaluation Selection

Goal: Discuss what evaluation metric(s) should be used in a classification problem translating the technical metric in the context of each challenge.

Challenge Briefing

- Get Together with your Group in the **dedicated Zoom link** to discuss the questions related with your Use Case **[10min]**;
- **Come back to the main zoom meeting** and:
 - 1) present your solution for your use case in <insert jamboard link here>**[5min]**
 - 2) when you're not presenting, **you're invited to comment the other's group use cases solutions.**



#1 Breathalyzer Tests

A breathalyzer registers someone's blood alcohol content to tell if they are "over the limit" or "under the influence" of alcohol. It is typically used at roadside police stops to determine if someone is legally able to drive.

1. What is the positive class?
2. What does 85% recall mean?
3. What does 75% precision mean?
4. What is the most important metric?

#2 Should we unlock a phone?

We are building a facial recognition algorithm to allow people to unlock their phone. If the phone recognizes the person as the authorized user, it will unlock the phone. If it doesn't recognize the user, it will prompt them to try again or try an alternative method (such as a passphrase).

#3 Detect malicious programs

When running a program for the first time, we are running some information about the program (such as where it was downloaded, size of the executable, etc) through a classifier. If the program is deemed safe, it will run. If it is deemed unsafe, the user will be prompted to confirm that the program is safe before running.

#4 De-duplicate records

You are writing a deduplication algorithm. Its goal is to flag entries in your database that are duplicates of existing records. It is more complicated than checking if two records are identical (Which is an easy problem), but instead tries to assess if differences are meaningful. Records that our algorithm detects as duplicates will be reviewed, and if we are sure they are duplicates, they are removed.

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

“the most important thing is not life, but the good life.”

Socrates, 399 BC

Technology shapes how we're seeking “good life” and it's not ethically neutral.

Aren't we all Ethical already?

Strong culture organisations already have Ethics as a value however they're not specifically catered to social and ethical processes needed to develop AI systems

**<Insert your company
mission & values ☺ >**

...**Data Ethics** is an extension that **evaluates moral problems related with Data, algorithms and development practices** in order to provide **morally good solutions aligned with the organisation values**.

It should go from high-level abstraction concepts down to to **practical steps that can support creation & maintenance of AI systems that help us make decisions leading to fair and equitable outcomes**.



Why Data Ethics?

There are different potential harms caused by AI systems that can be managed by a solid Data Ethics practice



Bias and Discrimination

Systems that gain insights from existing structures and societal dynamics can reproduce and reinforce patterns of inequality.



Denial of Individual Autonomy, Recourse, and Rights

AI systems can potentially automate decision making processes and this can complicate accountability since it can be distributed through the system,



Non-transparent, Unexplainable, or Unjustifiable Outcomes

Depending on the ML model there may be challenging to understand the resulting algorithmic rules which can hide traces of discrimination, bias, inequity or fairness.



Invasions of privacy

AI systems are developed frequently using personal data which can be captured, extracted or handled in a way without knowledge or proper consent of the data subject.



Isolation and Disintegration of Social Connection

Societies are built on relationships of trust, empathy and mutual understanding, as AI gets more pervasive human to human interaction is reduced and everything that comes along.



Unreliable, Unsafe, or Poor-Quality Outcomes

Negligent design, development or deployment of AI systems can lead to outcomes that hurt individuals and undermine public trust in the responsible use of societally beneficial AI tech.



[1] Goodman, B. F. (2016). European Union regulations on algorithmic decision-making and a “right to explanation.” ICML Workshop on Human Interpretability in Machine Learning
[2] In <https://www.accenture.com/us-en/insights/technology/testing-ai>
[3] https://www.bcg.com/publications/2020/europe-can-catch-up-in-ai-but-must-act-today?utm_medium=Email&utm_source=esp&utm_campaign=most_innovative_companies&utm_description=early&utm_topic=covid&utm_geo=global&utm_content=202007&utm_usertoken=CRM_76ee82224b32d3ff8ede5bdf349924b2f40a5cd2

Why does Data Ethics matter to organizations?

Data Ethics Practice can be a new competitive advantage

- "...the data controllers should assure implementation of measures that prevents, inter alia, discriminatory effects...also the “right-to-explanation” is deeply embedded on GDPR in Articles 13-15, specifying that individuals (or “data subjects”) have the right to know which data is being collected and how it is being used and what for." [1]
- There's a nation race for AI adoption (where Portugal is lagging) taking an ethical approach will be a key enabler for that adoption in our economy as a growth driver [3]
-  **GDPR Compliance**
-  **Consumers Trust**
-  **PT Economic Growth**
-  **Savvy Modelling**
- "Business executives today want to win customers' trust as growth comes from customers trusting the business. And, aligned to that customers have to trust you. You need to ensure that your systems are supporting that trust imperative. And that's exactly why businesses need to care... customers trust businesses which have verifiable, explainable, trustworthy systems" [2]
- Going through the process of auditing a Model for an explanation to better understand “Fairness” dimensions leads us to an insight journey, building internal trust on the DS practice. Also, DS practice will be challenged to another level of sophistication since explanations and auditing/mitigation bias are cutting edge technical matters that can be leveraged to create competitive value.

Are there any Data Ethics frameworks?

Organizations are introducing policy frameworks to enhance and regulate AI based services or products

The relationship between trust in AI and trustworthy machine learning technologies

FAT* '20, January 27–30, 2020, Barcelona, Spain

Table 1: Trustworthy technology classes related to FAT* frameworks. X - no mention, ✓ - mentioned, ✓✓ - emphasised

Framework	Year	Document Owner	Entities	Country	Fairness	Explainability	Safety	Auditability
Top 10 principles of ethical AI	2017	UNI Global Union	Ind	Switzerland	✓	✓	✓	✓
Toronto Declaration	2018	Amnesty International	Gov, Ind	Canada	✓✓	✓	X	✓
Future of work and Education for the Digital Age	2018	T20: Think 20	Gov	Argentina	✓✓	✓	✓	✓
Universal Guidelines for AI	2018	The public voice coalition	Ind	Belgium	✓✓	✓	✓	✓
Human Rights in the Age of AI	2018	Access Now	Gov, Ind	United States	✓✓	✓	✓✓	✓
Preparing for the Future of AI	2016	US national Science, and Technology Council	Gov, Ind, Acad	United States	✓	✓	✓✓	✓
Draft AI R&D Guidelines	2017	Japan Government	Gov	Japan	X	✓	✓✓	✓
White Paper on AI Standardization	2018	Standards Administration of China	Gov, Ind	China	✓	X	✓✓	✓✓
Statements on AI, Robotics and "Autonomous" Systems	2018	European Group on Ethics in Science and New Technologies	Gov, Ind, Acad	Belgium	✓	✓	✓	✓✓
For a Meaningful Artificial Intelligence	2018	Mission assigned by the French Prime Minister	Gov, Ind	France	✓	✓	X	✓✓
AI at the Service of Citizens	2018	Agency for Digital Italy	Gov, Ind	Italy	✓	✓	✓	✓
AI for Europe	2018	European Commission	Gov, Ind	Belgium	✓	✓	✓	✓
AI in the UK	2018	UK House of Lords	Gov, Ind	United Kingdom	✓✓	✓	✓	✓
AI in Mexico	2018	British Embassy in Mexico City	Gov	Mexico	✓	X	✓	✓
Artificial Intelligence Strategy	2018	German Federal Ministry of Education, Economic Affairs, and Labour and Social Affairs	Gov, Ind	Germany	✓	✓	✓	✓✓
Draft Ethics Guidelines for Trustworthy AI	2018	European High Level Expert Group on AI	Gov, Ind, Civ		✓✓	✓	✓✓	✓
AI Principles and Ethics	2019	Smart Dubai	Ind	UAE	✓✓	✓	✓✓	✓
Principles to Promote FAT* AI in the Financial Sector	2019	Monetary Authority of Singapore	Gov, Ind	Singapore	✓	✓	X	✓
Tenets	2016	Partnership on AI	Gov, Ind, Acad	United States	✓	✓	✓	✓
Aisomar AI Principles	2017	Future of Life Institute	Ind	United States	✓	✓	✓	✓
The GNI Principles	2017	Global Network Initiative	Gov, Ind	United States	X	✓	✓	✓
Montreal Declaration	2018	University of Montreal	Gov, Ind, Civ	Canada	✓✓	✓✓	✓✓	✓
Ethically Aligned Design	2019	IEEE	Ind	United States	✓	✓	✓	✓✓
Seeking Ground Rules for AI	2019	New York Times	Ind, GeP	United States	✓	✓	✓	✓
European Ethical Charter on the Use of AI in Judicial Systems	2018	Council of Europe (CEPEJ)	Gov	France	✓	✓	✓	✓
AI Policy Principles	2017	ITI	Gov, Ind	United States	✓	✓	✓✓	✓
The Ethics of Code	2019	Sage	Ind	United States	✓	X	X	✓
Microsoft AI Principles	2018	Microsoft	Ind	United States	✓	✓	✓✓	✓
AI at Google: Our Principles	2018	Google	Ind	United States	✓	✓	✓✓	✓
AI Principles of Telefónica	2018	Telefónica	Ind	Spain	✓	✓	✓	X
Guiding Principles on Trusted AI Ethics	2019	Telia Company	Ind	Sweden	✓	✓	✓✓	✓
Declaration of the Ethical Principle for AI	2019	IA Latam	Ind	Chile	✓	✓	✓✓	✓

Toreini, E., Aitken, M., Coopamootoo, K., Elliott, K., Zelaya, C. G., & van Moorsel, A. (2019). The relationship between trust in AI and trustworthy machine learning technologies.

Framework	Year	Document Owner	Entities	Country	Fairness	Explainability	Safety	Audibility
Microsoft AI Principles	2018	Microsoft	Ind	United States	✓	✓	✓✓	✓
AI at Google: Our Principles	2018	Google	Ind	United States	✓	✓	✓✓	✓
AI Principles of Telefónica	2018	Telefónica	Ind	Spain	✓	✓	✓	X
Guiding Principles on Trusted AI Ethics	2019	Telia Company	Ind	Sweden	✓	✓	✓✓	✓
Declaration of the Ethical Principles for AI	2019	IA Latam	Ind	Chile	✓	✓	✓✓	✓

The principles within each theme are:

Privacy:

- Explainability
- Transparency
- Open Source Data and Algorithms
- Notification when Interacting with an AI
- Notification when AI Makes a Decision about an Individual
- Regular Reporting Requirement
- Right to Information
- Open Procurement (for Government)

Fairness and Non-discrimination:

- Non-discrimination and the Prevention of Bias
- Fairness
- Inclusiveness in Design
- Inclusiveness in Impact
- Representative and High Quality Data
- Equality

Accountability:

- Accountability
- Recommendation for New Regulations
- Impact Assessment
- Evaluation and Auditing Requirement
- Verifiability and Replicability
- Liability and Legal Responsibility
- Ability to Appeal
- Environmental Responsibility
- Creation of a Monitoring Body
- Remedy for Automated Decision

Human Control of Technology:

- Human Control of Technology
- Human Review of Automated Decision
- Ability to Opt out of Automated Decision

Professional Responsibility:

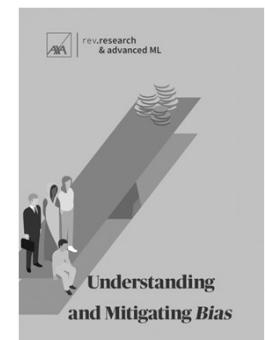
- Multistakeholder Collaboration
- Responsible Design
- Consideration of Long Term Effects
- Accuracy
- Scientific Integrity

Promotion of Human Values:

- Leveraged to Benefit Society
- Human Values and Human Flourishing
- Access to Technology

Further inform
methodology
Artificial Intel
in Ethical and
(Berkman Kle
cyber.harvar

Also other organisations are starting more focused and tactical:



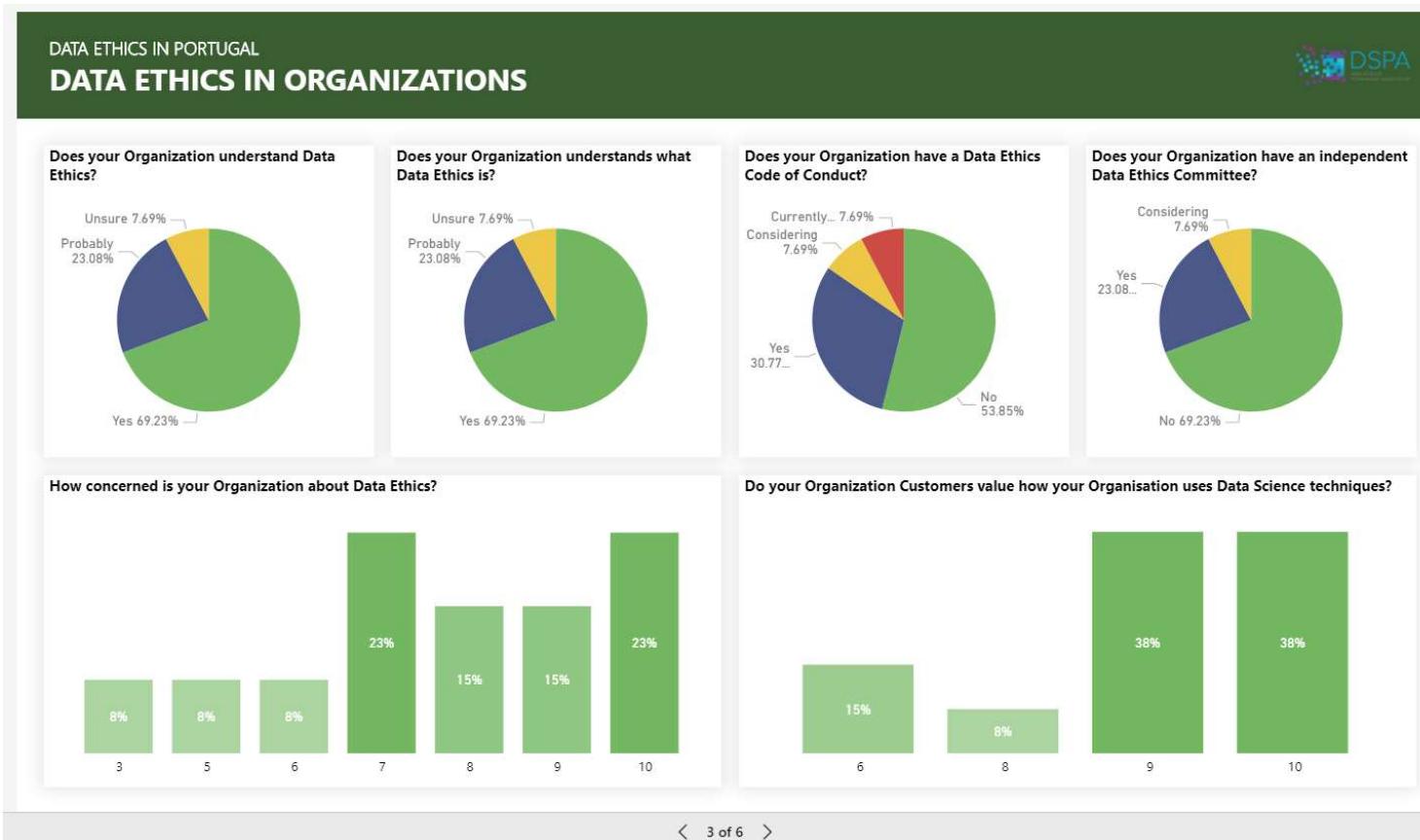
https://axa-rev-research.github.io/static/AXA_Booklet_Bias.pdf

Fjeld, J., Achten, N., Hilligoss, H., Nagy, A., Srikumar, M., Dushkin, R., Tolmatov, A., & Onatsik, K. (2020). Principled Artificial Intelligence: Mapping Consensus in Ethical and Rightsbased Approaches to Principles for AI.



And in practice? Are we doing it?

The abundance of Data Ethics frameworks doesn't translate to many PT organizations in practice



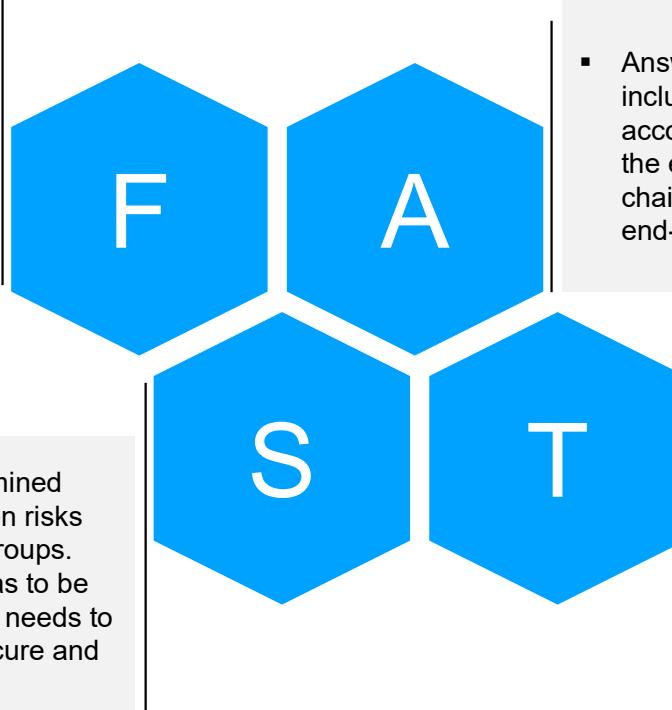


So, where to start from?

A structured framework will organize the main principles. The Alan Turing Institute has a detailed guide with the FAST principles:

Fairness

- AI systems that process personal data must be designed to meet a minimum threshold of discriminatory non-harm. This means that all the tech end-to-end pipeline is designed and implemented from data up to the outcomes that causes to individuals in an unbiased way.



Accountability

- Answerability and auditability has to be included by design, requiring accountable humans in the loop across the entire design and implementation chain, with defined processes that allow end-to-end oversight and review.

Sustainability

- Social impact has to be determined through bringing to light unseen risks that can affect individuals or groups.
- Also technical sustainability has to be ensured for the AI system that needs to be safe, accurate, reliable, secure and robust.

Transparency

- AI system needs to be transparent in terms of being able to explain the design/implementation process but also the in terms of product, ie: the contente and justificationof the outcome.

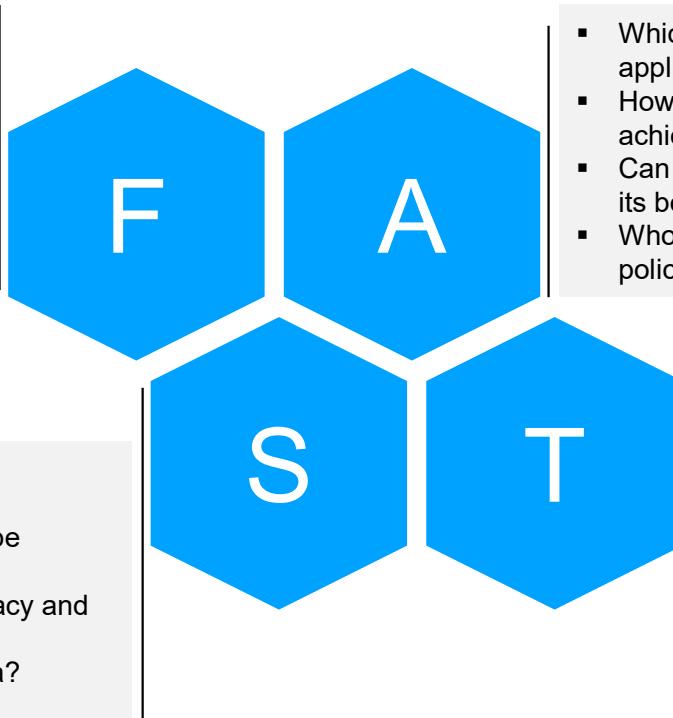


So, where to start from?

A structured framework will organize the main principles, the Alan Turing Institute has a detailed guide with the FAST principles:

Fairness

- How have you identified and minimized any bias in the data or model?
- How was any potential bias coming from the development team was identified and then mitigated?



Accountability

- Which laws and regulations might be applicable to this project?
- How is ethical accountability being achieved?
- Can we trace back the causal origins of its behavior?
- Who's accountable for actions and policies encoded on the system?

Sustainability

- How might the legal rights of organizations and individuals be impinged by our use of data?
- How might an individuals' privacy and anonymity be impinged by via aggregation and linking of data?

Transparency

- How transparent does the model needs to be and how is that achieved?
- What are likely misinterpretations of results and what can be done to prevent those misinterpretations?
- How transparency helps to debug/improve system?
- How can help tailor interventions?



What happens when something goes wrong?

A few Data Ethics cautionary tales that made “success” on social media...

Hiring [1]

Amazon scraps secret AI recruiting tool that showed bias against women

Pricing [2]

Staples website displays different prices to people depending on distance from a rival brick-and-mortar store

Justice [3]

An algorithm used by courts to assist in bail decisions was biased against non-white defendants



[1] <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

[2] <https://www.wsj.com/articles/SB1000142412788732377204578189391813881534>

[3] <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>



What happens when something goes wrong?

A few Data Ethics cautionary tales that made “success” on social media...

Self Driving Cars [1]

The uncertainty around assigning liability in AV highlights the desirability of a coordinated effort by policymakers and manufacturers to devise organisational structures and regulatory measures which support a fair criteria for culpability attribution

Social Welfare Tax [2]

The robodebt automated welfare recovery scheme matched annualised pay information from the Australian Taxation Office and income data to claw back overpaid welfare payments... But one-in-five debt letters sent were based on false information. The Robodebt class action bought by Gordon Legal has been settled at a cost to the government of around \$1.2 billion.

Health Care [3]

AI system as only advisory and expect that accountability is thereby secure, because human clinicians still make the final decision. dilemma with two equally undesirable choices. Either clinicians spend the time to develop their own opinions as to the best course of action, meaning that the artificial intelligence system adds little value; or clinicians must accept the advice blindly, further weakening both the control and epistemic conditions of moral accountability.



[1] Legal issues in automated vehicles: critically considering the potential role of consent and interactive digital interfaces
<https://www.nature.com/articles/s41599-020-00644-2>

[2] <https://7news.com.au/business/finance/robodebt-class-action-everything-you-need-to-know-about-the-12-billion-lawsuit-brought-by-gordon-legal-c-1595625>

[3] <https://www.who.int/bulletin/volumes/98/4/19-237487/en/>



What happens when something goes wrong?

A few Data Ethics cautionary tales that made “success” on social media...

Data Brokers ^[1]

The type of location data in question is collected from millions of phones, with most people unaware that their movements are being tracked this way and unable to find out who has access to that information. There are few laws regulating location data companies, and government agencies have used this to their advantage, spending millions to gain access to this information.

Facebook ^[2]

Facebook breached data protection laws by failing to keep users' personal information secure, allowing Cambridge Analytica to harvest the data of up to 87 million people without their consent worldwide. The now-defunct firm worked for the Trump presidential campaign and used the data to influence several elections around the world.

Google ^[3]

Google has agreed to pay a record \$170 million penalty to settle accusations that YouTube broke the law when it knowingly tracked and sold ads targeted to children ...The settlement involves the largest-ever penalty under the Children's Online Privacy Protection Act, which YouTube allegedly violated by collecting user information from kids to fuel its behavioral advertising business.



[1] <https://www.vox.com/recode/22038383/dhs-cbp-investigation-cellphone-data-brokers-venn/>

[2] <https://www.npr.org/2019/10/30/774749376/facebook-pays-643-000-fine-for-role-in-cambridge-analytica-scandal?t=1607184810121>

[3] <https://edition.cnn.com/2019/09/04/tech/google-youtube-ftc-settlement/index.html>



What happens when something goes wrong?

A few Data Ethics cautionary tales that made “success” on social media...

Uber & Lyft [1]

*Uber and Lyft blamed for San Francisco's congested streets...
Lyft and Uber caused a 51 per cent increase in vehicle hours on the road – far outweighing the impact of other causes like changes in population and employment.*

Predatory Ads [2]

...company's artificial intelligence technology “enables us to efficiently serve prospective students as well as customize our outreach and engagement with current students to help them stay and succeed in school... settlements worth over \$500 million last year to resolve major actions pursued by the Federal Trade Commission and 49 state attorneys general alleging deceptive practices.

Apple [3]



The [@AppleCard](#) is such a fucking sexist program. My wife and I filed joint tax returns, live in a community-property state, and have been married for a long time. Yet Apple's black box algorithm thinks I deserve 20x the credit limit she does. No appeals work.

8:34 PM · Nov 7, 2019

28.1K 10.4K people are Tweeting about this



[1] <https://www.dezeen.com/2018/10/18/uber-lyft-transportation-network-companies-blamed-san-francisco-traffic-congestion/>

[2] <https://www.republicreport.org/2020/biden-must-cut-off-taxpayer-billions-to-predatory-colleges-that-ruin-students-lives/>

[3] <https://www.nytimes.com/2019/11/10/business/Apple-credit-card-investigation.html>



Trust as a competitive differentiator in an AI economy

Transparency reinforces trust and Customer willingness to share data, that by its turn allows AI to flourish new use cases, as opposing to hinder innovation and growth

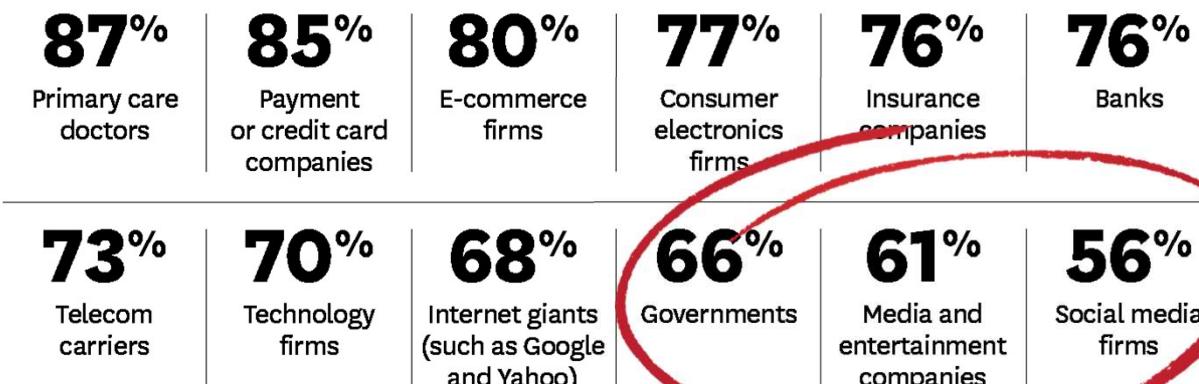
Survey Methodology

46 companies surveyed representing seven categories of business around the world. asked Consumers to rate the firms on the following scale:

- **completely trustworthy** - would freely share sensitive personal data with a firm because they trust the firm not to misuse it;
- **trustworthy** - would “not mind” exchanging sensitive data for a desired service;
- **untrustworthy** - would provide sensitive data only if required to do so in exchange for an essential service;
- **completely untrustworthy** - would never share sensitive data with the firm.

Do They Trust You with Their Data?

Percentages of consumers who said that each category of organization was “trustworthy” or “completely trustworthy” when it came to making sure that personal data was never misused.



SOURCE TIMOTHY MOREY, THEODORE “THEO” FORBATH, AND ALLISON SCHOOP
FROM “CUSTOMER DATA: DESIGNING FOR TRANSPARENCY AND TRUST,” MAY 2015

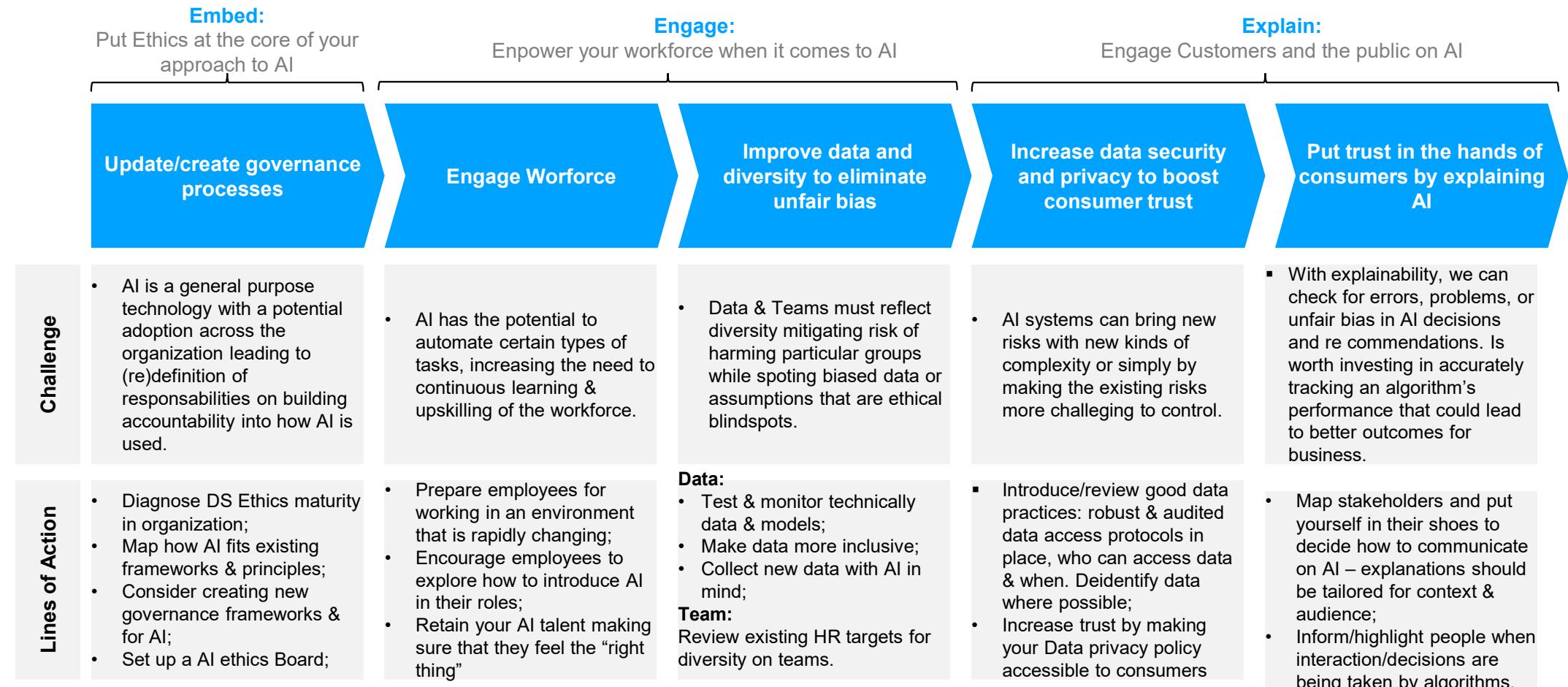
© HBR.ORG



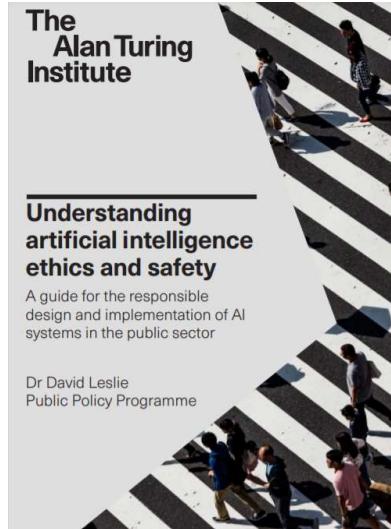
How to go from principles to ground practice?

Depending on your DS Ethics maturity you may complete current practice or focus in specific “pain” areas

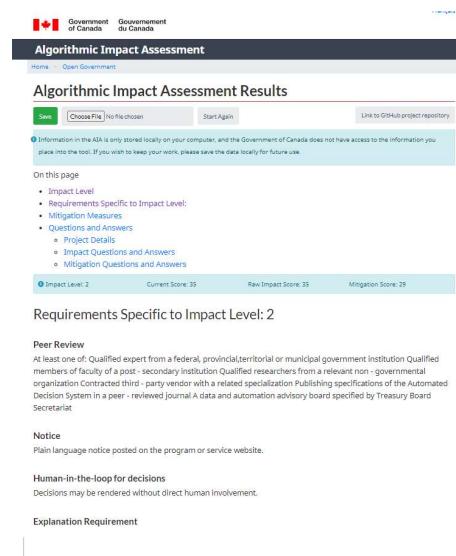
Not Exhaustive



DATA SCIENCE CHALLENGES – COMMON PITFALLS & ETHICS



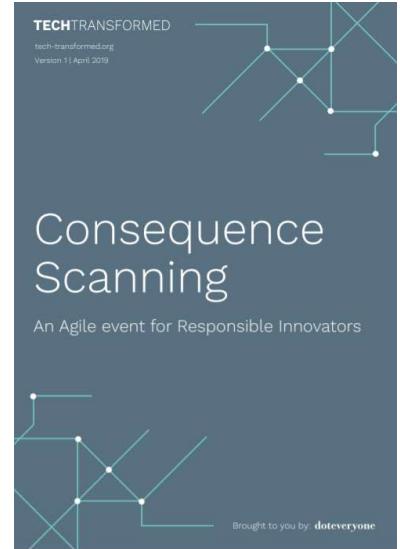
- Thorough guide on values, principles and Process Based Framework with human centered implementation processes. [1]



- Self assessment tool with impact score plus a series of recommendations of mitigation requirements and measures.



- Other self-assessment tool to Help you cook the project prepared by the UK government [3]



- There's also a free kit prepared by TechTransformed to upgrade your agile practice to consider impact of products/services on people & society.

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

Group Challenge – Practice Data Ethical Issues Identification

Goal: Identify ethical issues on an AI system and discuss potential harms for each involved stakeholder.

Challenge Briefing

- Spend some time to read the article;
- Get Together with your Group in the **dedicated Zoom link** to discuss the questions **[20min]**;
- Come back to the main zoom meeting and:
 - 1) present your view to the audience – write it on the jamboard
<insert jamboard link> [5min/group]
 - 2) when you're not presenting, **you're invited to comment the other's group ethical issues.**



<https://www.theguardian.com/technology/2017/sep/07/new-artificial-intelligence-can-tell-whether-youre-gay-or-straight-from-a-photograph>

New AI can guess whether you're gay or straight from a photograph

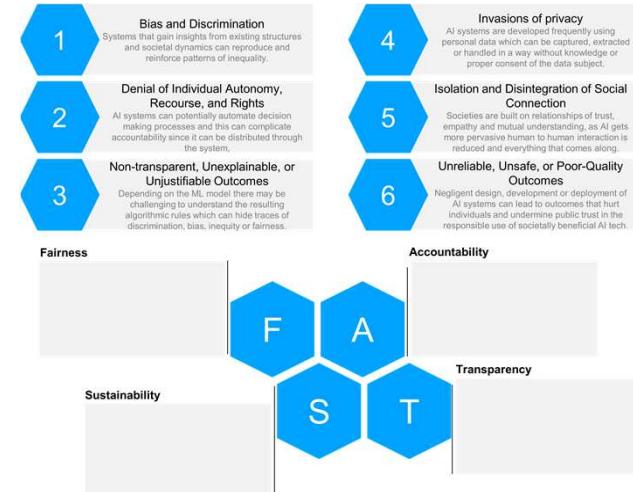
An algorithm deduced the sexuality of people on a dating site with up to 91% accuracy, raising tricky ethical questions



An illustrated depiction of facial analysis technology similar to that used in the experiment. Illustration: Alamy
Artificial intelligence can accurately guess whether people are gay or straight based on photos of their faces, according to new research that suggests machines can have significantly better "gaydar" than humans.

The study from Stanford University - which found that a computer algorithm could correctly distinguish between gay and straight men 81% of the time, and 74% for women - has raised questions about the biological origins of sexual orientation, the ethics of facial-detection technology, and the potential for this kind of software to violate people's privacy or be abused for anti-LGBT purposes.

The machine intelligence tested in the research, which was published in the Journal of Personality and Social Psychology and first reported in the Economist, was based on a sample of more than 35,000 facial images that men and women publicly posted on a US dating website. The researchers, Michal Kosinski and Yilun Wang, extracted features from the images using "deep neural networks", meaning a sophisticated mathematical system that learns to analyze visuals based on a large dataset.



- 1) Which **potential harms** and **related ethical issues** do you identify on the AI system presented on the article?
- 2) Can you think of an **ethical application** of this system?

Group Challenge – Practice Data Ethical Issues Identification

Helper – it could help you to structure your brainstorm output in this format

Article Quote	Potential Harm/Implication	Article Quote	Potential Harm/Implication
F ‘AI can accurately guess whether people are gay or straight based on photos of their faces	Bias & Discrimination Acknowledging sexual orientation as a binary classification in scientific research can reinforce discrimination on groups which are not included.	S	
F		S	
T		A	

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

Agenda

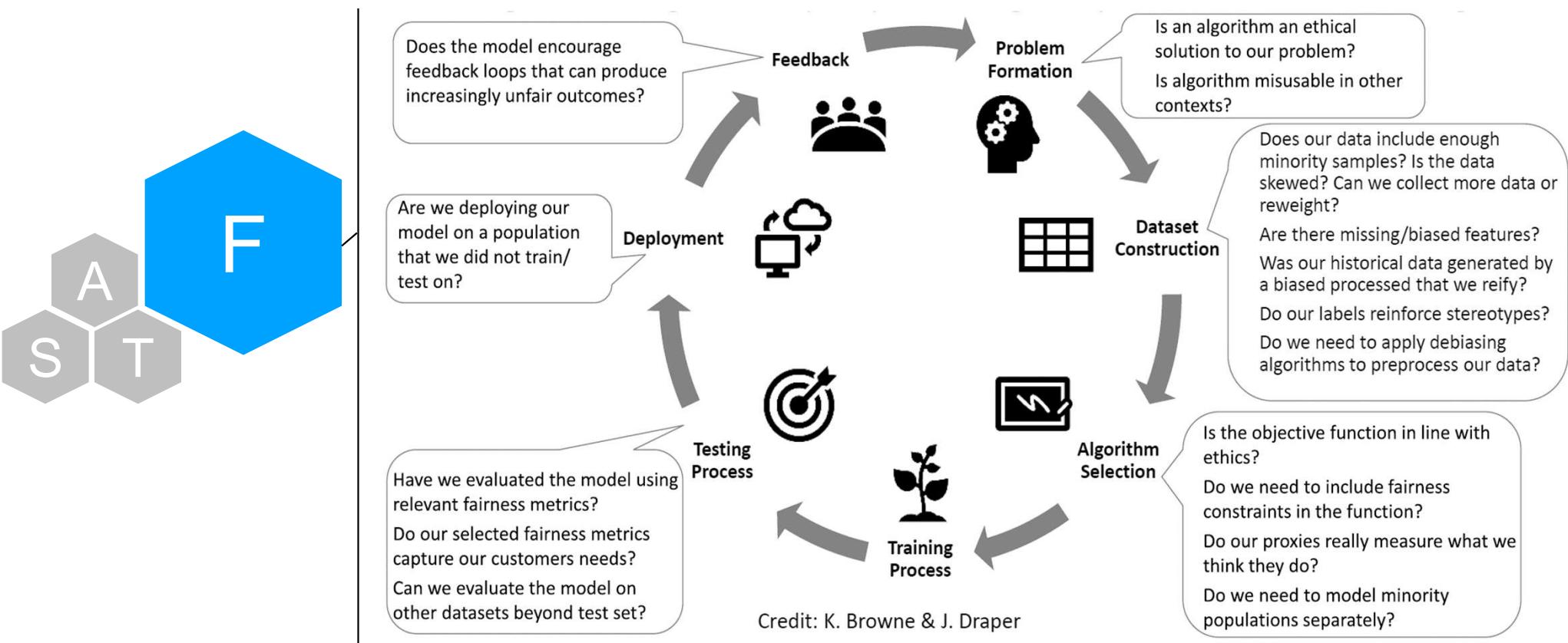
Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		



How do we develop AI systems that help make decisions leading to fair and equitable systems?

Each step on the pipeline may introduce potential bias to the downstream tasks





Data is a social mirror of an unfair world?

Type of Bias in Data		Other Sources of Bias							
Bias Type	Example	Evaluation	Aggregation	Simpson's Paradox	Longitudinal Data Fallacy				
Historical Bias <i>Historical bias is the already existing bias and socio-technical issues in the world and can seep into from the data generation process even given a perfect sampling and feature selection</i>	An example of this type of bias can be found in a 2018 image search result where searching for women CEOs ultimately resulted in fewer female CEO images due to the fact that only 5% of Fortune 500 CEOs were woman —which would cause the search results to be biased towards male CEOs								
Representation Bias <i>Representation bias happens from the way we define and sample from a population</i>	Lacking geographical diversity in datasets like ImageNet is an example for this type of bias, where Western countries such as US and Great Britain are over-represented								
Measurement Bias <i>Measurement bias happens from the way we choose, utilize, and measure a particular feature</i>	In the recidivism risk prediction tool COMPAS, where prior arrests and friend/family arrests were used as proxy variables to measure level of “riskiness” or “crime”—which on its own can be viewed as mismeasured proxies . This is due to the fact that minority communities are controlled and policed more frequently, so they have higher arrest rates.								
Population Bias <i>Population bias arises when statistics, demographics, representatives, and user characteristics are different in the user population represented in the dataset or platform from the original target population</i>	Different user demographics on different social platforms , such as women being more likely to use Pinterest, Facebook, Instagram, while men being more active in online forums like Reddit or Twitter.								
		Evaluation	Aggregation	Simpson's Paradox	Longitudinal Data Fallacy				
		Sampling	Behavioral	Content Production	Linking				
		Temporal	Popularity	Algorithmic	User Interaction				
		Social	Emergent	Self Selection	Omitted Variable				
		Cause Effect Bias	Observer Bias	Funding Bias	Sample Size				

There's no unique mathematical definition of fairness!

Fairness isn't a general concept, it can be concerned “with important opportunities that affect people's life chances” (Solon Barocas · Moritz Hardt, 2017) - **domain specific** - or concerned with “social salient qualities have served as the basis for unjustified and systematically adverse treatment in the past” - **feature specific**.



Domain Specific

4.4. Access to supply of goods and services, including housing

Under EU law, Protection from discrimination in the field of access to the supply of goods and services, including housing, applies to the ground of race under the Racial Equality Directive, and to the ground of sex under the Gender Goods and Services Directive. Paragraph 13 of the Preamble to the Gender Goods and Services Directive gives more precision to prohibition of discrimination, stating that it relates to all goods and services “which are available to the public irrespective of the person concerned as regards both the public and private sectors, including public bodies, and which are offered outside the area of private and family life and the transactions carried out in this context”. It expressly excludes application to ‘the content of media or advertising’ and ‘public or private education’, though this latter exclusion does not narrow the scope of the Racial Equality Directive, which expressly covers education. The Gender Goods and Services Directive also refers to Article 57 of the Treaty on the Functioning of the European Union:

from Handbook on European non- discrimination law

- “Under EU law, Protection from discrimination in the field of access to the supply of **goods and services**, including housing, applies to the ground of race under the Racial Equality Directive, and to the ground of sex under the Gender Goods”

Feature Specific

4. PROTECTED GROUNDS

- 4.1. Introduction
 - 4.2. Sex
 - 4.3. Sexual orientation
 - 4.4. Disability
 - 4.5. Age
 - 4.6. Race, ethnicity, colour and membership of a national minority
 - 4.7. Nationality or national origin
 - 4.8. Religion or belief
 - 4.9. Language
 - 4.10. Social origin, birth and property
 - 4.11. Political or other opinion
 - 4.12. ‘Other status’
- Key points
- Further reading

from Handbook on European non- discrimination law

- “The European non-discrimination directives prohibit differential treatment that is based on **certain ‘protected grounds**', containing a fixed and limited list of protected grounds”

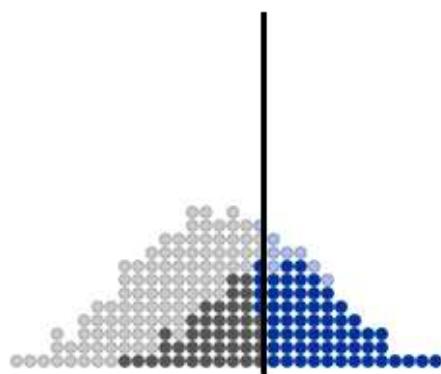


To understand Fairness implications, we can simulate a loan application use case for two distinct populations with different default probabilities...

Blue Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 61

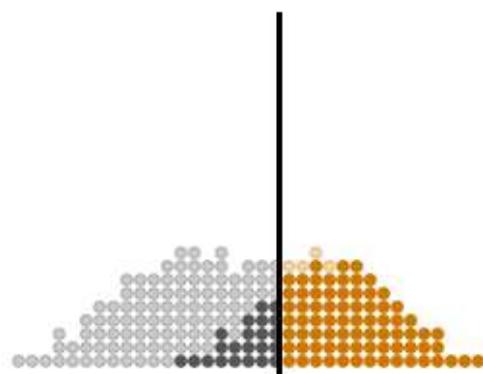


denied loan / would default
denied loan / would pay back

Orange Population

0 10 20 30 40 50 60 70 80 90 100

loan threshold: 50



denied loan / would default
denied loan / would pay back

Different threshold criteria can lead to different Fairness definition and outcome impact.

Loan Strategy Maximize profit with:	Loan Threshold cut-off for loan granting		True Positive Rate percentage of paying applications getting loans		Positive Rate percentage of all applications getting loans		Incorrect loans denied to paying applicants and granted to defaulters		Profit sum of the costs and profits of granted paying and default loans		
	Blue	Orange	Blue	Orange	Blue	Orange	Blue	Orange	Blue	Orange	Total
Max Profit - The most profitable, since there are no constraints. But the two groups have different thresholds, meaning they are held to different standards.	61	50	60%	78%	34%	41%	24%	13%	\$ 12 100	\$ 20 300	\$ 32 400
Group unaware - Both groups have the same threshold, but the orange group has been given fewer loans overall. Among people who would pay back a loan, the orange group is also at a disadvantage.	55	55	81%	60%	52%	30%	21%	21%	\$ 8 600	\$ 17 000	\$ 25 600
Demographic parity - The number of loans given to each group is the same, but among people who would pay back a loan, the blue group is at a disadvantage	60	52	64%	71%	37%	37%	23%	16%	\$ 11 900	\$ 18 900	\$ 30 800
Equal Opportunity - Among people who would pay back a loan, blue and orange groups do equally well. This choice is almost as profitable as demographic parity, and about as many people get loans overall.	59	53	68%	68%	35%	40%	22%	17%	\$ 11 700	\$ 18 700	\$ 30 400

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

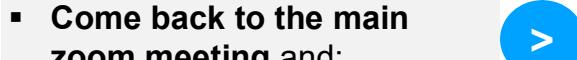


Group Challenge - Which could be the sources of bias a loan use case?

Goal: Go through the process of brainstorming sources of bias in a specific use case

Challenge

- Get Together with your Group in the to read article + discuss the questions related with Use Case [15min];
- Come back to the main zoom meeting and:
 - 1) add & present your solution for your use case in the <insert jamboard link> [5min]
 - 2) when you're not presenting, **you're invited to comment the other's group use cases solutions.**



Case Study - Bank Loan

Description: Acme Bank is developing a system to decide which loan applications to deny based on predicted risk of lender not paying back their loan in time.

- **Goal:** Increase repayment rates for bank loans
- **Data:** Historical loans and payments, credit reporting data, background checks, <other that you may find relevant!>
- **Analysis:** Build model to predict risk of not repaying on time
- **Actions:** Deny loan or increase interest rate/penalties

Question: What are some potential sources of bias in the underlying data?

From KDD workshop

Suggested reading

AI Can Make Bank Loans More Fair

by Sian Townsend

November 06, 2020



Summary. Many financial institutions are turning to AI reverse past discrimination in lending, and to foster a more inclusive economy. But many lenders find that artificial-intelligence-based engines exhibit many of the same biases as humans. How can they address the issue to ensure that biases of the... [more](#)

As banks increasingly deploy artificial intelligence tools to make credit decisions, they are having to revisit an unwelcome fact about the practice of lending: Historically, it has been riddled with biases

<https://hbr.org/2020/11/ai-can-make-bank-loans-more-fair?ab=hero-main-text>



If blue = women and orange = man, how would you decide fairly?

Expert #1. Profit Max

Expert #1 defends that what matters is maximizing profit without restrictions and thus optimize the cut-offs having in account the distributions minimizing loan granting to the women's population.

Expert #2: Group unaware:

What's fair, says Expert #2, is to totally disregard the gender mix of the applicants who are given loans. If no women at all make it into the pool, that's fair so long as people were chosen purely on their merits—what's sometimes called “group unaware” fairness. Otherwise, you'd have to kick out a qualified man to make room for a less qualified woman, because there are a limited number of approvals possible. company should go to great lengths to exclude gender and gender-proxy information from the data set, and then go for the most accurate predictions of who will repay a loan.

#3. Demographic Parity

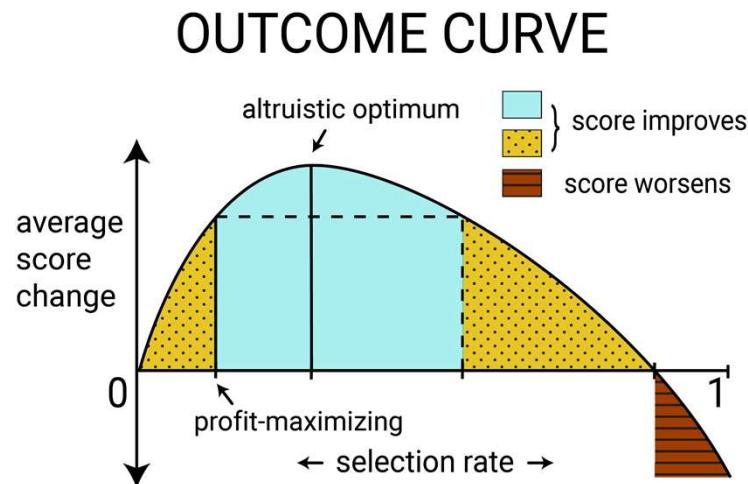
No, says Expert #3. The composition of the set of people who are granted loans should reflect the percentage of applicants: if 30 percent of the applicants are women, then 30 percent of the pool of approved applicants ought to be women.

#4. Equal Opportunity

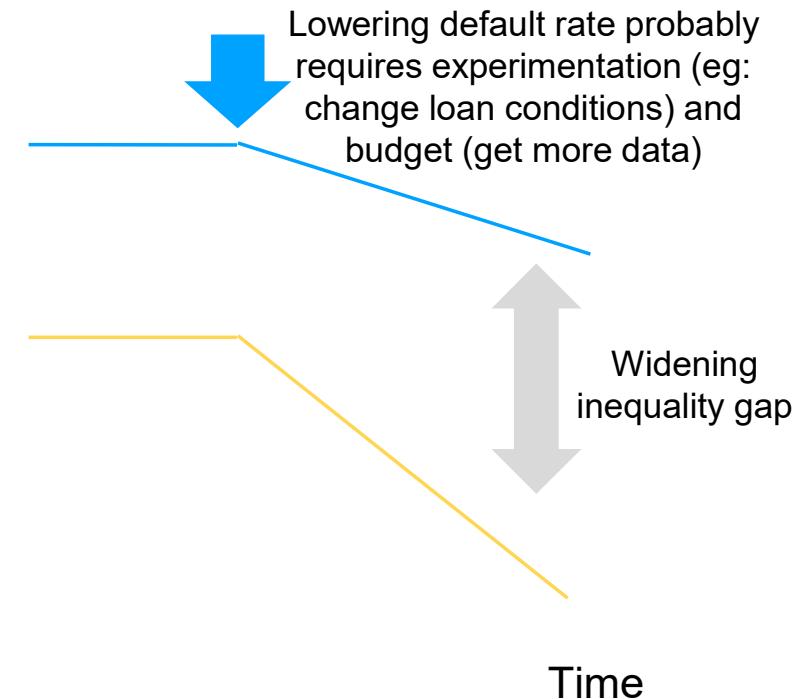
Expert #4 disagrees because if you just want equal proportions of men and women to get loans, you could randomly offer loans to women whether or not the individuals are good risks. Instead, it will be “Equal Opportunity” fairness if the same percentage of men and women who are likely to succeed at loans are given loans. What you don't want is for 90 percent of the male good risks to make it into the acceptance pile, but only 40 percent of the female good risks to do so.



Did you consider the impact of your fairness policy over time?



Loan Default Prevalence



- When enough individuals in a group are granted loans and successfully repay them, the **average credit score in that group is likely to increase**;
- Increasing selection rate up to a point where the average score change is lower than under **unconstrained profit-maximization but still positive**;
- Selection rates in dotted **yellow regions** are causing **relative harm** and if too many individuals are **unable to repay their loans, credit score for the group will decrease**, as is the case in the red striped region.

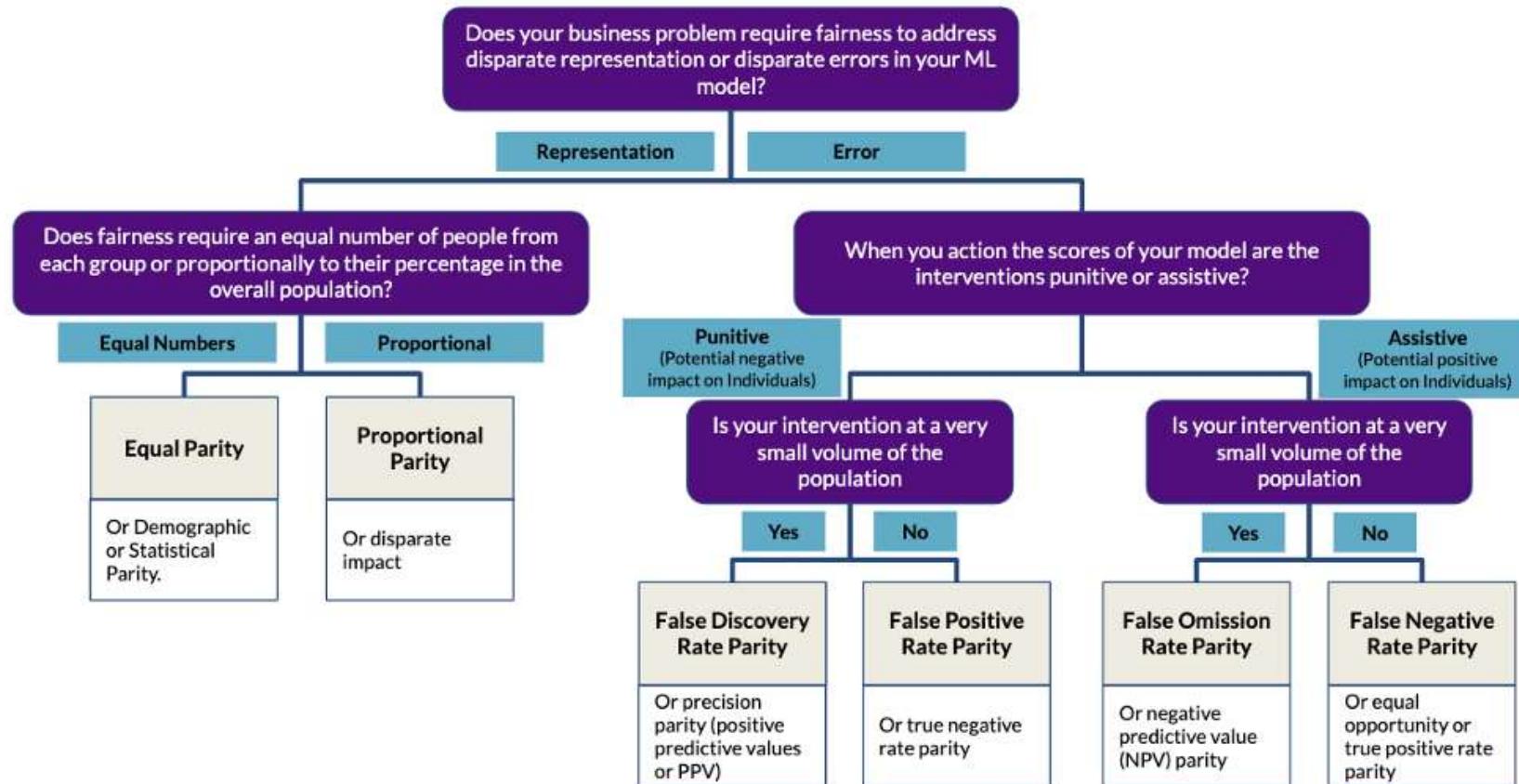
Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		



Which metrics do we care about for our problem?



The Fairness Tree, Adapted from <http://www.datasciencepublicpolicy.org/projects/aequitas/>



Use the fairness tree as a tool to support discussions and evaluate scenarios

Punitive Scenario – Jail Determination

Metric	Probabilistic Term	Metric Translation to Use Case Domain	Caveat																		
Equal Count of False Positives	$P(\text{wrongly jailed} \mid \text{group } i) = C \quad \forall i$	<i>Number of people wrongly jailed from each group is the same,</i>	If groups are not the same size it can be arguable as a fair criteria.																		
Group Size-Adjusted False Positives	$P(\text{wrongly jailed} \mid \text{group } i) = C \quad \forall i$	<i>Just by virtue of the fact that an individual is a member of a given group, what are the chances they'll be wrongly convicted?</i>	<p>% innocents is lower on group B</p> <table border="1"><thead><tr><th>Group</th><th>Pop. Size</th><th>No Bail</th><th>FP</th><th>FP/Popul.</th><th>%Innocent</th></tr></thead><tbody><tr><td>A</td><td>10 000</td><td>100</td><td>10</td><td>0,1%</td><td>90%</td></tr><tr><td>B</td><td>30 000</td><td>100</td><td>30</td><td>0,1%</td><td>70%</td></tr></tbody></table>	Group	Pop. Size	No Bail	FP	FP/Popul.	%Innocent	A	10 000	100	10	0,1%	90%	B	30 000	100	30	0,1%	70%
Group	Pop. Size	No Bail	FP	FP/Popul.	%Innocent																
A	10 000	100	10	0,1%	90%																
B	30 000	100	30	0,1%	70%																
False Discovery Rate (FDR)	$P(\text{wrongly jailed} \mid \text{jailed}, \text{group } i) = C \quad \forall i$	<i>Focused on the group of people which are interventioned - 100/each group jailed on example above – “discovers” the the ratio of FP on the interventioned sub-groups.</i>	Considering two distinct groups A & B with a criminal process with equal FDR and group sized-adjusted false positives: <table border="1"><thead><tr><th>Group</th><th>Pop. Size</th><th>Guilty</th><th>Innocent</th><th>Decision = Jail</th></tr></thead><tbody><tr><td>A</td><td>1 000</td><td>900</td><td>100</td><td>100</td></tr><tr><td>B</td><td>3 000</td><td>300</td><td>2700</td><td>100</td></tr></tbody></table> $\frac{FP_B}{n_B} = \frac{30}{3000} = 1.0\% \quad \frac{FP_A}{n_A} = \frac{10}{1000} = 1.0\%$ $FDR_B = \frac{30}{300} = 10.0\% \quad FDR_A = \frac{10}{100} = 10.0\%$ $FPR_B = \frac{30}{2730} = 1.1\% \quad FPR_A = \frac{10}{10} = 100.0\%$	Group	Pop. Size	Guilty	Innocent	Decision = Jail	A	1 000	900	100	100	B	3 000	300	2700	100			
Group	Pop. Size	Guilty	Innocent	Decision = Jail																	
A	1 000	900	100	100																	
B	3 000	300	2700	100																	
False Positive Rate (FPR)	$P(\text{wrongly jailed} \mid \text{innocent}, \text{group } i) = C \quad \forall i$	<i>For an innocent person, what are the chances they will be wrongly convicted by virtue of the fact that they're a member of a given group?</i>	<ul style="list-style-type: none">There'll be cases where it's needed to balance trade-offs while selecting evaluation metrics																		



Classification Refresher: Evaluation Metrics

Rank	Score	Label	Predict	Skin Color
1	0.997	1	1	non-white
2	0.993	1	1	white
3	0.986	1	1	white
4	0.982	1	1	non-white
5	0.971	0	1	white
6	0.965	1	1	white
7	0.964	0	1	white
8	0.961	0	1	non-white
9	0.953	0	1	non-white
10	0.932	1	1	non-white
11	0.918	0	0	non-white
12	0.873	0	0	white
13	0.854	0	0	white
14	0.839	0	0	white
15	0.777	0	0	white
16	0.723	0	0	non-white
17	0.634	0	0	white
18	0.512	0	0	non-white
19	0.487	0	0	white
20	0.473	0	0	non-white

Population

- Total Label Positives: 6
- Total Label Negatives: 14
- Prevalence $6/20=0.3$

For threshold > 0.920 or top k=10:

Predicted Condition	True Condition		
	Total Population	Positive	Negative
Positive	TP 6	FP 4	
Negative	FN 0	TN 10	

True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$
False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$

- True Positive Rate (Recall) = $6/6 = 1.0$
- False Positive Rate: $4/14=0.29$
- False Negative Rate = $0/6=0$
- True Negative Rate = $10/14=0.71$
- Precision = $6/10 = 0.6$

Adapted from Pedro Saleiros's "Tutorial: Fairness in decision-making with AI: a practical guide & hands-on tutorial using Aequitas <https://www.youtube.com/watch?v=yOR71zBm3Uc>



Classification Refresher: Evaluation Metrics

Rank	Score	Label	Predict	Skin Color
1	0.997	1	1	non-white
2	0.993	1	1	white
3	0.986	1	1	white
4	0.982	1	1	non-white
5	0.971	0	1	white
6	0.965	1	1	white
7	0.964	0	1	white
8	0.961	0	1	non-white
9	0.953	0	1	non-white
10	0.932	1	1	non-white
				Predicted Positive: 10
				Predicted Negative: 10
11	0.918	0	0	non-white
12	0.873	0	0	white
13	0.854	0	0	white
14	0.839	0	0	white
15	0.777	0	0	white
16	0.723	0	0	non-white
17	0.634	0	0	white
18	0.512	0	0	non-white
19	0.487	0	0	white
20	0.473	0	0	non-white

Population

- Total Label Positives: 6
- Total Label Negatives: 14
- Prevalence $6/20=0.3$

For threshold > 0.920 or top k=10:

		Non-White		White			
		True Condition		True Condition			
	Population	Positive	Negative		Population	Positive	Negative
Predicted	Positive	3	2	Predicted	Positive	3	0
Condition	Negative	0	4	Condition	Negative	2	6

Metric	Non-White	White
FPR	$2/6=0.33$	$2/6=0.25$
Recall	$3/3=1$	$3/3=1$
Precision	$3/5=0.6$	$3/5=0.6$
FNR	0	0

Adapted from Pedro Saleiros's "Tutorial: Fairness in decision-making with AI: a practical guide & hands-on tutorial using Aequitas
<https://www.youtube.com/watch?v=yOR71zBm3Uc>

True positive rate (TPR),
 Recall, Sensitivity,
 probability of detection,
 Power
 $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$

False positive rate (FPR),
 Fall-out,
 probability of false alarm
 $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$

False negative rate (FNR),
 Miss rate
 $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$

Specificity (SPC),
 Selectivity, True negative
 rate (TNR)
 $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$



How to assess level of fairness disparity between groups?

1. Compare a given metric with a reference group (e.g. historically favoured group)

$$FPR_g \text{ disp} = \frac{FPR_{a_i}}{FPR_{a_r}} = \frac{\Pr(\widehat{Y}=1|Y=0,A=a_i)}{\Pr(\widehat{Y}=1|Y=0,A=a_r)}$$

2. Define a fairness range for the disparity measure across groups

$$\tau \leq DisparityMeasure_{group_i} \leq \frac{1}{\tau}$$

3. Calculate & evaluate disparate impact

Models	Accuracy	Disparate impact
TF without mitigation	0.73	0.16
TF with adversarial debiasing	0.71	0.42

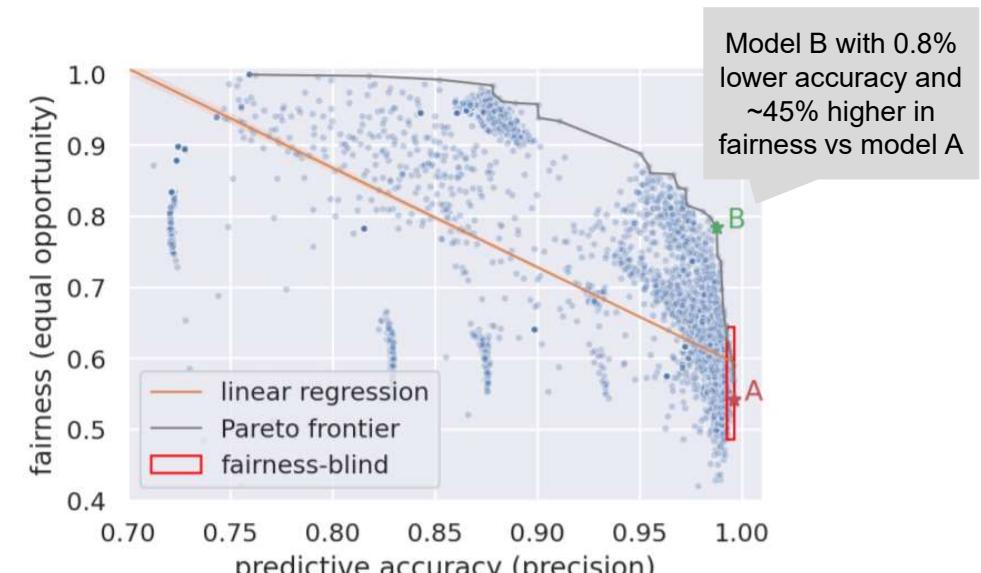
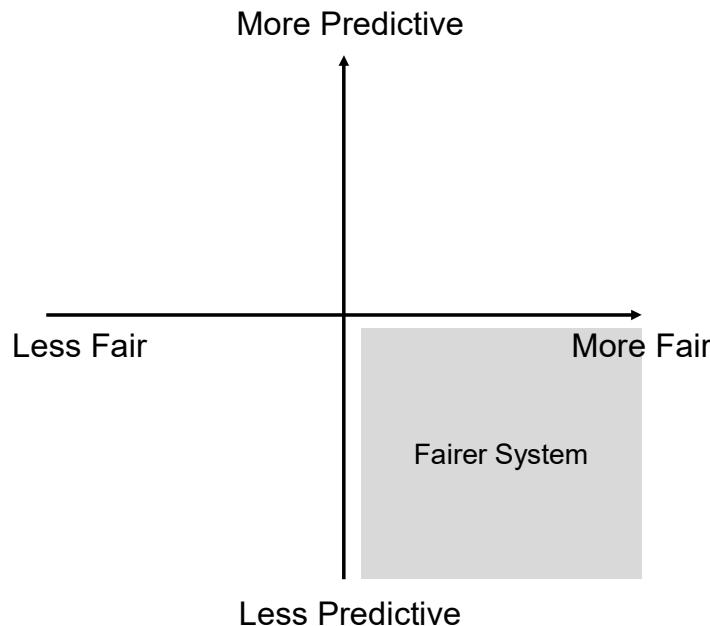


This can be regulated depending on domain/geography

"The four-fifths or 80% rule is described by the guidelines as "a selection rate for any race, sex, or ethnic group which is less than four-fifths (or 80%) of the rate for the group with the highest rate will generally be regarded by the Federal enforcement agencies as evidence of adverse impact..." [1]



Once defined one (or a set) of fairness definition(s) we can do a better judgment of which model best suits us





Fairness is a cross-sectional topic in AI with different domains already benefiting of growing research activity

In a recent survey different domains & subdomains of ML/AI are examined to structure research around state-of-art methods to mitigate bias

Domain	Sub-Domain	Domain	Sub-Domain	Domain	Sub-Domain
ML	Binary Classification	NLP	Word Embedding	DNN	Variational Auto-Encoders
	Regression		Coreference Resolution		Adversarial Learning
	PCA		Language Model		
	Community Detection		Sentence embedding		
	Clustering		Machine translation		
	Graph embedding		Semantic Role Labeling		
	Causal Inference		Classification		



There is a proliferation of open source frameworks to support the development of Fair models

Not Exhaustive

Tool	Latest Commit	Nbr of Commits	Nbr Branches	Nbr Releases	Nbr of Contributors	Github Reference
Tensorflow – what if tool *	30/08/2019	2.776	21	25	176	tensorflow/tensorboard/tree/master/tensorboard/plugins/interactive_inference
SHAP	31/08/2019	964	3	33	64	slundberg/shap
Aequitas	11/07/2019	832	11	35	8	dssg/Aequitas
Fairness Comparison	05/06/2019	503	10	2	10	algofairness/fairness-comparison
Fairest	29/05/2017	474	5	0	7	columbia/fairest
LIME	23/07/2019	433	5	17	34	marcotcr/lime
Fairness Measures	29/12/2017	260	3	0	3	megantosh/fairness_measures_code
AIF360	30/08/2019	149	26	4	16	IBM/AIF360
Fairml	23/03/2017	104	3	0	2	adebayoj/fairml
ThemisML	31/07/2018	65	1	3	1	cosmicBboy/themis-ml
FAT-Forensics	03/05/2019	59	1	0	4	liamdalgl/FAT-Forensics
Microsoft Fair Learn	09/08/2019	37	10	2	6	github.com/microsoft/fairlearn
Audit-AI	25/06/2019	35	2	2	6	pymetrics/audit-ai



There is a proliferation of open source frameworks to support the development of Fair models

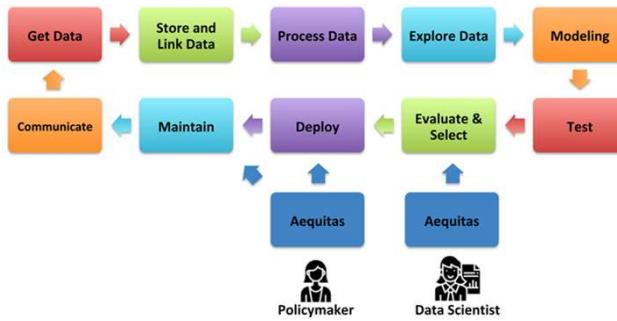
Not Exhaustive

Tool	Latest Commit	Nbr of Commits	Nbr Branches	Nbr Releases	Nbr of Contributors	Github Reference
Tensorflow – what if tool *	30/08/2019	2.776	21	25	176	tensorflow/tensorboard/tree/master/tensorboard/plugins/interactive_inference
SHAP	31/08/2019	964	3	33	64	slundberg/shap
Aequitas	11/07/2019	832	11	35	8	dssg
Fairness Comparison	05/06/2019	503	10	2	10	algofairness/fairness-comparison
Fairest	29/05/2017	474	5	0	7	columbia/fairest
LIME	23/07/2019	433	5	17	34	marcotcr/lime
Fairness Measures	29/12/2017	260	3	0	3	megantosh/fairness_measures_code
AIF360	30/08/2019	149	26	4	16	IBM/AIF360
Fairml	23/03/2017	104	3	0	2	adebayoj/fairml
ThemisML	31/07/2018	65	1	3	1	cosmicBboy/themis-ml
FAT-Forensics	03/05/2019	59	1	0	4	liamdalgl/FAT-Forensics
Microsoft Fair Learn	09/08/2019	37	10	2	6	github.com/microsoft/fairlearn
Audit-AI	25/06/2019	35	2	2	6	pymetrics/audit-ai

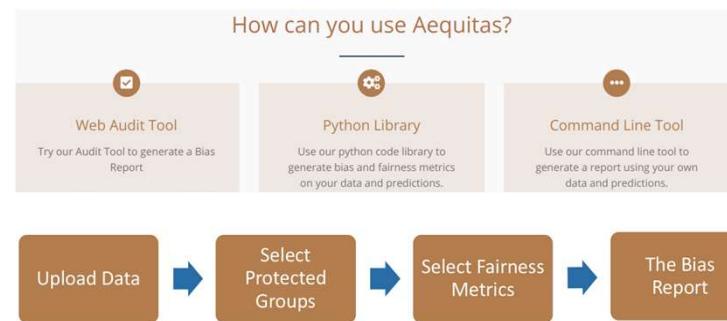
Aequitas – Bias & Fairness Audit

Aequitas is an open-source bias audit toolkit for data scientists, machine learning researchers, and policymakers to audit machine learning models for discrimination and bias, and to make informed and equitable decisions around developing and deploying predictive tools.

What is the Target Audience?



How to use it?



What do we get?

Aequitas		Back to Dashboard	Home	About
The Bias Report				
7214 rows were used to audit bias and fairness.				
Fairness Criteria	Desired Outcome	Reference Group Selected	Unfairly Affected Groups	
Equal Party	Each group is represented equally.	race: Caucasian sex: Male age_cat: 25 - 45	race: Asian Hispanic Other African-American Native American sex: Female age_cat: Less than 25 Greater than 45	
Proportional Party	Each group is represented proportional to their representation in the overall population.	race: Caucasian age_cat: 25 - 45	race: Asian-American Native American Other Asian age_cat: Greater than 45 Less than 25	

- Aequitas was designed having in mind technical Data profiles but also policymakers creating a common ground for discussion.
 - For non-tech profiles a bias report is quite straightforward to get through web app – also there are available python api and batch loading of data.
 - Aequitas provides a detailed audit to your models calculating detailed fairness and bias statistics, mitigation can be done through an informed model selection.



Case Study – Compas Recidivism Risk

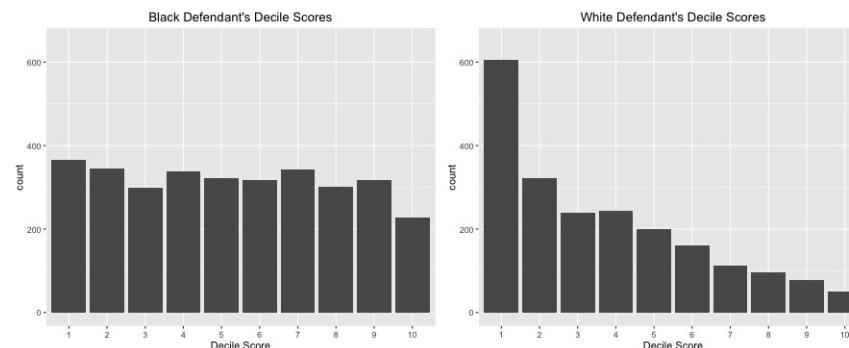
Risk models are used in criminal justice to inform high-stakes decisions in people's life such as bail, sentencing and early release. In May 2016 journalists at Propublica explored the bias in such models, choosing a more widely solution in use – Correctional Offender Management Profiling for Alternative Sanctions (Compas)

COMPAS Scores & Profiles



- Scores of individuals with inconsistent criminal record were found to be unreliable by Propublica journalists.

COMPAS scores for “Risk of Recidivism”



- Propublica showed that scores for white defendants were skewed toward lower-risk categories, while black defendants were evenly distributed across scores.

COMPAS Fairness Metrics

Fairness Metric	Caucasian	African American
False Positive Rate (FDR)	23%	45%
False Negative Rate (FNR)	48%	28%
False Discovery Rate (FDR)	41%	37%

- Propublica and Northpoint (Compas vendor) were arguing over different Fairness Metrics



http://aequitas.dssg.io/audit/27py3o00/compas_for_aequitas/report-1.html#false-positive-rate-parity-span-red-initfailedspan-red-end

Case Study – Compas Recidivism Risk

Checking audit results on the Aequitas Bias Report example for Compas use case we can find same results as reported by Propublica

Audit Results: Group Metrics Values

race

Attribute Value	Group Size Ratio	Predicted Positive Rate	Predicted Positive Group Rate	False Discovery Rate	False Positive Rate	False Omission Rate	False Negative Rate
African-American	0.51	0.66	0.59	0.37	0.45	0.35	0.28
Asian	0	0.0	0.25	0.25	0.09	0.12	0.33
Caucasian	0.34	0.26	0.35	0.41	0.23	0.29	0.48
Hispanic	0.09	0.06	0.3	0.46	0.21	0.29	0.56
Native American	0	0.0	0.67	0.25	0.38	0.17	0.1
Other	0.05	0.02	0.21	0.46	0.15	0.3	0.68

[Go to Previous](#)

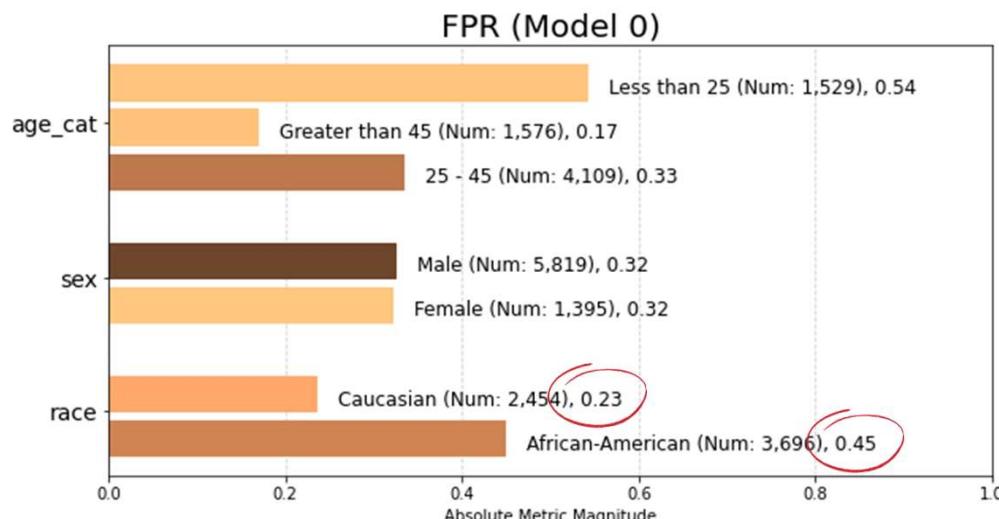
[Go to Top](#)



Case Study – Compas Recidivism Risk

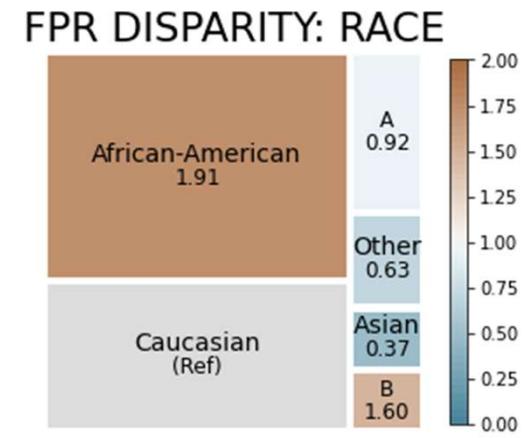
If we need more flexibility, we can easily use aequitas python api to audit our Compas use case

Fairness Metrics for Group above volume threshold



```
fnr = aqp.plot_group_metric(xtab, 'fpr', min_group_size=0.05)
```

Visualizing disparities between groups using Treemap



Not labeled above:
 A: Hispanic, 0.92
 B: Native American, 1.60

```
aqp.plot_disparity(bdf, group_metric='fpr_disparity',  

attribute_name='race', significance_alpha=0.05)
```



There is a proliferation of open source frameworks to support the development of Fair models

Not Exhaustive

Tool	Latest Commit	Nbr of Commits	Nbr Branches	Nbr Releases	Nbr of Contributors	Github Reference
Tensorflow – what if tool *	30/08/2019	2.776	21	25	176	tensorflow/tensorboard/tree/master/tensorboard/plugins/interactive_inference
SHAP	31/08/2019	964	3	33	64	slundberg/shap
Aequitas	11/07/2019	832	11	35	8	dssg/Aequitas
Fairness Comparison	05/06/2019	503	10	2	10	algofairness/fairness-comparison
Fairest	29/05/2017	474	5	0	7	columbia/fairest
LIME	23/07/2019	433	5	17	34	marcotcr/lime
Fairness Measures	29/12/2017	260	3	0	3	megantosh/fairness_measures_code
AIF360	30/08/2019	149	26	4	16	IBM/AIF360
Fairml	23/03/2017	104	3	0	2	adebayoj/fairml
ThemisML	31/07/2018	65	1	3	1	cosmicBboy/themis-ml
FAT-Forensics	03/05/2019	59	1	0	4	liamdalgl/FAT-Forensics
Microsoft Fair Learn	09/08/2019	37	10	2	6	github.com/microsoft/fairlearn
Audit-AI	25/06/2019	35	2	2	6	pymetrics/audit-ai



AIF360 – Bias & Fairness Audit

AI Fairness 360

<https://github.com/IBM/AIF360>

AIF360 toolkit is an open-source library to help detect and remove bias in machine learning models.

The AI Fairness 360 Python package includes a comprehensive set of metrics for datasets and models to test for biases, explanations for these metrics, and algorithms to mitigate bias in datasets and models.

Toolbox

- Fairness metrics (70+)
- Fairness metric explanations
- Bias mitigation algorithms (10)

IBM **CODE**

AIF360

Supported bias mitigation algorithms

Optimized Preprocessing (Calmon et al., 2017)
Disparate Impact Remover (Feldman et al., 2015)
Equalized Odds Postprocessing (Hardt et al., 2016)
Reweighting (Kamiran and Calders, 2012)
Reject Option Classification (Kamiran et al., 2012)
Prejudice Remover Regularizer (Kamishima et al., 2012)
Calibrated Equalized Odds Postprocessing (Pleiss et al., 2017)
Learning Fair Representations (Zemel et al., 2013)
Adversarial Debiasing (Zhang et al., 2018)

Clip slide

Supported fairness metrics

Comprehensive set of group fairness metrics derived from selection rates and error rates
Comprehensive set of sample distortion metrics
Generalized Entropy Index (Speicher et al., 2018)

```

graph TD
    RD[Raw Data] --> SP[Standard Pre-processing]
    SP --> TD[Training Data]
    SP --> TD[Testing Data]
    TD --> ST[Standard Training<br/>in-processing algorithm]
    TD --> CUT[Classifier Unit Tests<br/>Accuracy/Discrimination]
    CUT --> C[Classifier]
    C --> PP[post-processing algorithm]
    C --> CME[classifier metric explainer]
    CME --> CM[classifier metric]
    CM --> D{Deploy}
    CM --> DME[dataset metric explainer]
    CM --> P[pre-processing algorithm]
    P --> DM[dataset metric]
    DM --> DME
    DM -- Yes --> P
    DM -- No --> DME
    D --> RI{External Interventions}
    D --> RRR{Reprocess and/or Retrain}
  
```

We can group the existing unfairness mitigation methods in 4 different classes

Not Exhaustive

1

Data Collection

- Better investment in data collection processes as an effective way to mitigate bias and/or better data stewardship processes.

2

Pre-Processing

- Mitigate bias on the dataset typically by learning a data transformation or oversampling minority classes;

3

In-Processing

- Modifications of existing learning algorithms to minimize discrimination during the training stage can be grouped as in-processing techniques

4

Post-Processing

- Performed after model training, it relies typically on the prediction, protected attribute and ground truth label.

- Datasheets for Datasets (Gebru, et al. 2018)
- Bias Source Analysis (I. Chen, Johansson, & Sontag, 2018)

- Optimized Preprocessing (Calmon et al., 2017)
- Manipulating labels (Kamiran & Calders, 2012);
- Reweighting (Kamiran & Calders, 2012);

- Prejudice removal by regularization (T Kamishima, Akaho, & Sakuma, 2011);
- Adversarial de-biasing: (Beutel, Chen, Zhao, & Chi, 2017) 

- Equalized Odds Postprocessing (Hardt et al., 2016)

We can group the existing unfairness mitigation methods in 4 different classes

Not Exhaustive

1

Data Collection

- Better investment in data collection processes as an effective way to mitigate bias and/or better data stewardship processes.

2

Pre-Processing

- Mitigate bias on the dataset typically by learning a data transformation or oversampling minority classes;

3

In-Processing

- Modifications of existing learning algorithms to minimize discrimination during the training stage can be grouped as in-processing techniques

4

Post-Processing

- Performed after model training, it relies typically on the prediction, protected attribute and ground truth label.

- Datasheets for Datasets (Gebru, et al. 2018)
- Bias Source Analysis (I. Chen, Johansson, & Sontag, 2018)

- Optimized Preprocessing (Calmon et al., 2017)
- Manipulating labels (Kamiran & Calders, 2012);
- Reweighting (Kamiran & Calders, 2012);

- Prejudice removal by regularization (T Kamishima, Akaho, & Sakuma, 2011);
- Adversarial de-biasing: (Beutel, Chen, Zhao, & Chi, 2017) 

- Equalized Odds Postprocessing (Hardt et al., 2016)



Adversarial debiasing is one of the most popular in-processing techniques



Intuition

- Train a model that **predicts Y from X, but encodes information from the protected attribute Z in a way that is not recoverable by an adversarial model.**



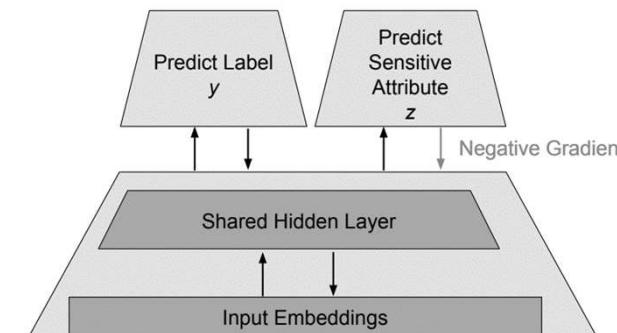
Goal

- Learn a representation $g(X)$ that makes it **easy for f to predict Y, but difficult for a to predict Z**. In optimization terms, we want to **minimize our prediction loss** $L_y(f(g(X)), Y)$ and **maximize the adversarial loss** $L_z(a(g(X)), Z)$.



We let f be our prediction function where $Y=f(g(X))$

Let a be our adversary where $Z=a(g(X))$



$g(X)$ to be the shared learned embedding of our input.



Results

Models	Accuracy	Disparate impact
TF without mitigation	0.73	0.16
TF with adversarial debiasing	0.71	0.42

For a loss of 2pp accuracy we have improved 2.7x the DI fairness metric

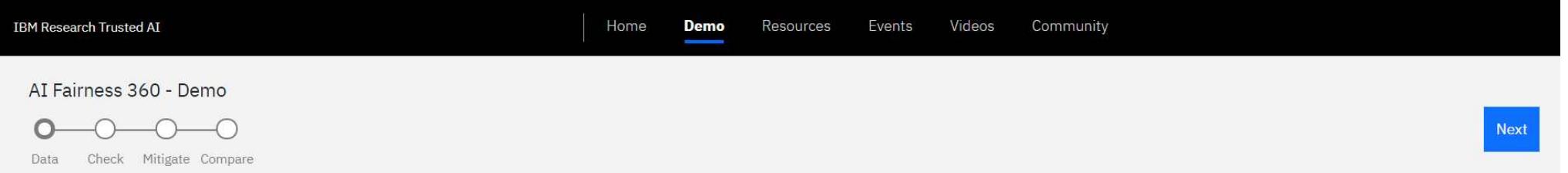
Method is also:

- Flexible, it can be tuned for the desired trade-off between “fairness” and accuracy (or other metric);
- Applicable to different fairness definitions;
- Can support regression and classification.

$$DI = \frac{P_{\text{positive unprivileged}}}{P_{\text{positive privileged}}}$$

AIF360 – Bias & Fairness Audit

We can test the same dataset but using a mitigation method



IBM Research Trusted AI

AI Fairness 360 - Demo

Home **Demo** Resources Events Videos Community

Data Check Mitigate Compare Next

1. Choose sample data set

Bias occurs in data used to train a model. We have provided three sample datasets that you can use to explore bias checking and mitigation. Each dataset contains attributes that should be protected to avoid bias.

Compas (ProPublica recidivism)

Predict a criminal defendant's likelihood of reoffending.

Protected Attributes:

- **Sex**, privileged: *Female*, unprivileged: *Male*
- **Race**, privileged: *Caucasian*, unprivileged: *Not Caucasian*

[Learn more](#)

German credit scoring

Predict an individual's credit risk.

Protected Attributes:

- **Sex**, privileged: *Male*, unprivileged: *Female*
- **Age**, privileged: *Old*, unprivileged: *Young*

[Learn more](#)

Adult census income

Predict whether income exceeds \$50K/yr based on census data.

Protected Attributes:

- **Race**, privileged: *White*, unprivileged: *Non-white*
- **Sex**, privileged: *Male*, unprivileged: *Female*



<https://github.com/Trusted-AI/AIF360/>
<https://www.slideshare.net/AnimeshSingh/aif360-trusted-and-fair-ai>
<http://aif360.mybluemix.net/>

AIF360 – Bias & Fairness Audit

We can test the same dataset but using a mitigation method

The screenshot shows the AI Fairness 360 - Demo interface. At the top, there is a navigation bar with links: Home, Demo (which is underlined), Resources, Events, Videos, and Community. Below the navigation bar, the main content area has a title "AI Fairness 360 - Demo". Underneath the title, there is a horizontal progress bar consisting of four circles: blue (Data), grey (Check), white (Mitigate), and white (Compare). To the right of the progress bar are two buttons: "Back" (grey) and "Next" (blue). Below the progress bar, the text "Dataset: Compas (ProPublica recidivism)" is displayed in a blue box, followed by "Mitigation: none" in a smaller blue box.

2. Check bias metrics

Dataset: Compas (ProPublica recidivism)

Mitigation: none

Protected Attribute: Sex

Privileged Group: **Female**, Unprivileged Group: **Male**

Accuracy with no mitigation applied is 66%

With default thresholds, bias against unprivileged group detected in 4 out of 5 metrics





<https://github.com/Trusted-AI/AIF360/>
<https://www.slideshare.net/AnimeshSingh/aif360-trusted-and-fair-ai>
<http://aif360.mybluemix.net/>

DATA SCIENCE CHALLENGES – COMMON PITFALLS & ETHICS

AIF360 – Bias & Fairness Audit

We can test the same dataset but using a mitigation method

IBM Research Trusted AI

Home **Demo** Resources Events Videos Community

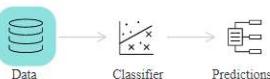
AI Fairness 360 - Demo

Data Check Mitigate Compare Back Next

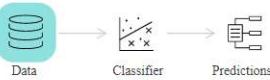
3. Choose bias mitigation algorithm

A variety of algorithms can be used to mitigate bias. The choice of which to use depends on whether you want to fix the data (pre-process), the classifier (in-process), or the predictions (post-process). [Learn more about how to choose.](#)

Reweighting
Weights the examples in each (group, label) combination differently to ensure fairness before classification.



Optimized Pre-Processing
Learns a probabilistic transformation that can modify the features and the labels in the training data.



Adversarial Debiasing
Learns a classifier that maximizes prediction accuracy and simultaneously reduces an adversary's ability to determine the protected attribute from the predictions. This approach leads to a fair classifier as the predictions cannot carry any group discrimination information that the adversary can exploit.



Cookie Preferences



<https://github.com/Trusted-AI/AIF360/>
<https://www.slideshare.net/AnimeshSingh/aif360-trusted-and-fair-ai>
<http://aif360.mybluemix.net/>

AIF360 – Bias & Fairness Audit

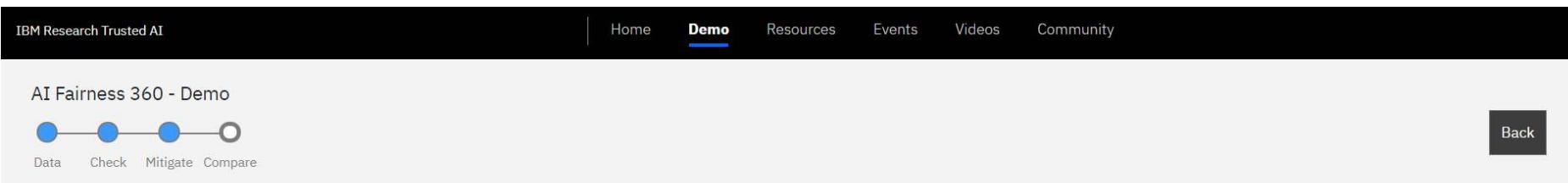
We can test the same dataset but using a mitigation method

IBM Research Trusted AI

Home Demo Resources Events Videos Community

AI Fairness 360 - Demo

Data Check Mitigate Compare Back



The screenshot shows the AIF360 Demo interface. At the top, there's a navigation bar with links for Home, Demo (which is highlighted in blue), Resources, Events, Videos, and Community. Below the navigation bar is a title "AI Fairness 360 - Demo". Underneath the title is a horizontal navigation bar with four items: "Data" (blue circle), "Check" (blue circle), "Mitigate" (blue circle), and "Compare" (white circle). On the right side of this bar is a "Back" button. The main content area is titled "4. Compare original vs. mitigated results". It includes a note about the dataset being Compas (ProPublica recidivism) and the mitigation method used being "Adversarial Debiasing algorithm applied".

4. Compare original vs. mitigated results

Dataset: Compas (ProPublica recidivism)

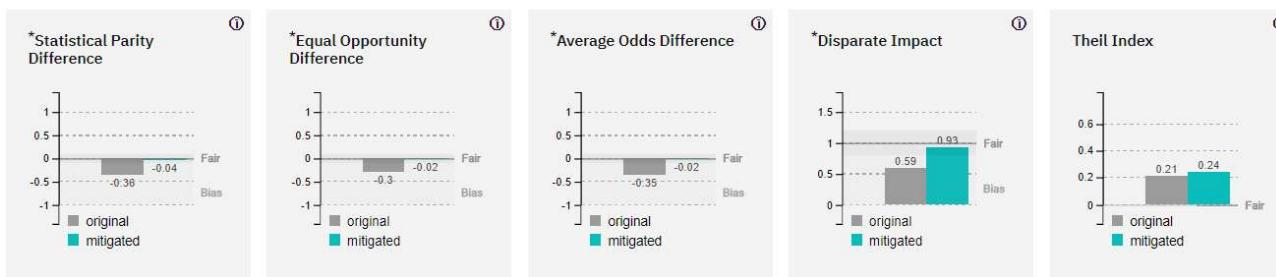
Mitigation: [Adversarial Debiasing algorithm applied](#)

Protected Attribute: Sex

Privileged Group: **Female**, Unprivileged Group: **Male**

Accuracy after mitigation unchanged

Bias against unprivileged group was reduced to acceptable levels* for 4 of 4 previously biased metrics (0 of 5 metrics still indicate bias for unprivileged group)



Protected Attribute: Race

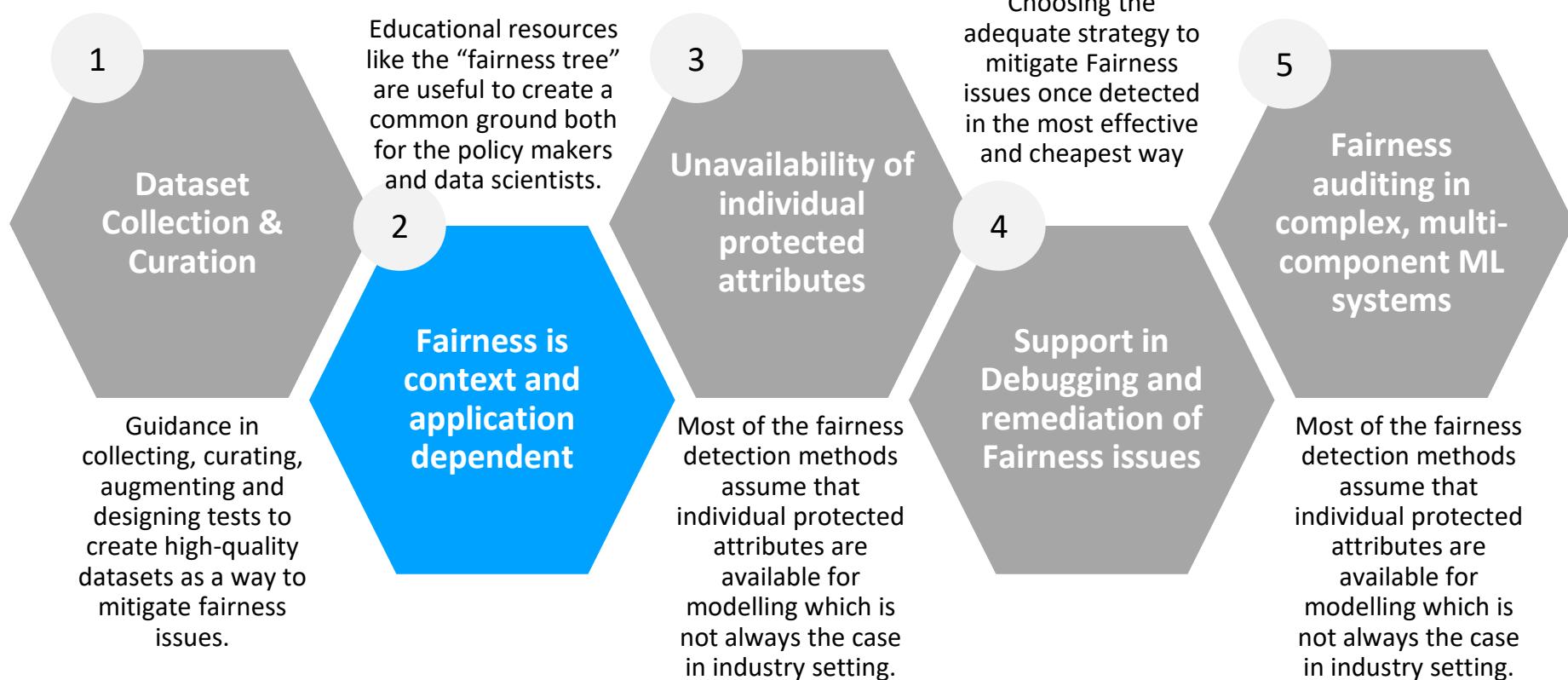
Privileged Group: **Caucasian**, Unprivileged Group: **Not Caucasian**

[Cookie Preferences](#)



Even for the existing (global) practice there are open challenges

A recent study[1] identifies challenges that are not so much covered by existing research topics on Fairness



We need to address and share more industry use-cases patterns

Use Case	Fairness Discussion	Illustrative Existing Research
Pricing / Next-Best-Offer	Fairness Discussion	Illustrative Existing Research
Content Recommendation Systems	<ul style="list-style-type: none"> ▪ In acquisition are we covering all the socio-economical segments? ▪ What pricing/offer is being provided to different groups of Customers and how fair is the current practice? ▪ Which are the possible definitions of a Fair Price? 	<ul style="list-style-type: none"> ▪ Seele, P., Dierksmeier, C., Hofstetter, R., & Schultz, M. D. (2019). Mapping the Ethicality of Algorithmic Pricing: A Review of Dynamic and Personalized Pricing. <i>Journal of Business Ethics</i>. https://doi.org/10.1007/s10551-019-04371-w ▪ Maestre, R., Duque, J., Rubio, A., & Arévalo, J. (2018). Reinforcement Learning for Fair Dynamic Pricing. <i>ArXiv E-Prints</i>, arXiv:1803.09967.
Customer Service & Operations	<ul style="list-style-type: none"> ▪ Is cultural or popularity bias reinforcing recommendations? ▪ Which metrics such diversity could be a more “fair” metric? Is it a trade-off, or can we find items dissimilar but relevant for user’s preferences? ▪ If a product is provided for a household and users are obfuscated how can we define fairness as well? ▪ Should the suppliers have a fair opportunity? 	<ul style="list-style-type: none"> ▪ Porcaro, L., Castillo, C., & Gómez, E. (2019). Music Recommendation Diversity: A Tentative Framework and Preliminary Results. ▪ Yao, S., & Huang, B. (2017). Beyond Parity: Fairness Objectives for Collaborative Filtering. <i>CoRR</i>, abs/1705.08804. http://arxiv.org/abs/1705.08804 ▪ Mehrotra, R., McInerney, J., Bouchard, H., Lalmas, M., & Diaz, F. (2018). Towards a Fair Marketplace: Counterfactual Evaluation of the trade-off between Relevance, Fairness & Satisfaction in Recommendation Systems. 2243–2251. https://doi.org/10.1145/3269206.3272027

Agenda

Today: encounter new ideas, engage through discussions with your peers and reflect on your learning.

Time	Topic	Time	Topic
9:00 am	Welcome, logistics & agenda	11:50 am	Break
9:15 am	Common Pitfalls in Data Science	12:00 am	Fairness in Machine Learning
10:00 am	Group Challenge #1	12:10 am	Group Challenge #3
10:30 am	Break: Grab your coffee!	12:30 am	Audit & Mitigation
10:40 am	DS Ethics	12:55 am	Wrap up & Farewell
11:10 am	Group Challenge #2		

Key takeaways

- Include fairness and equity as part of AI Systems development process;
- Make sure that fairness metrics are translated to the ones that make sense to your context;
- Audit bias & fairness;
- Explore bias reduction strategies;

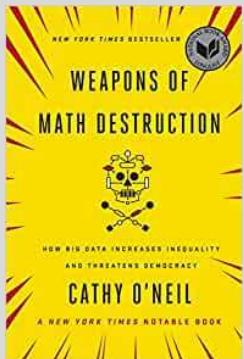
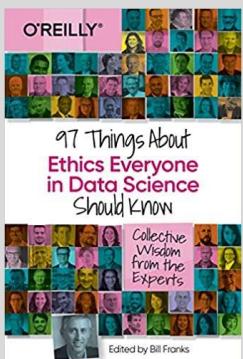
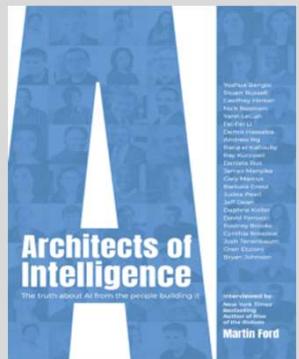
Call2action

What could be ideas to kick-start or accelerate these practices starting tomorrow?

- High-level list your internal DS initiatives – current & about to start;
- Prepare a strutured workshop to go through bias & fairness potential harm impacts;
- Schedule a workshop, discuss and fill list potential harms;
- Structure a brief presentation to management;
- Discuss next steps;

Resources

Inspire



Acquire

- Berkeley CS 294: Fairness in ML
- Cornell INFO 4270: Ethics and policy DS
- Princeton COS 597E: Fairness in ML
- CMU 10718/94889: ML for Public Policy Lab
- KDD 2020 Hands-on Tutorial

Engage

