

Course 1: How Google does ML

Module 4: Inclusive ML

Lesson Title: **Introduction**

Format: Talking head

Inclusive Machine Learning

How Google does ML

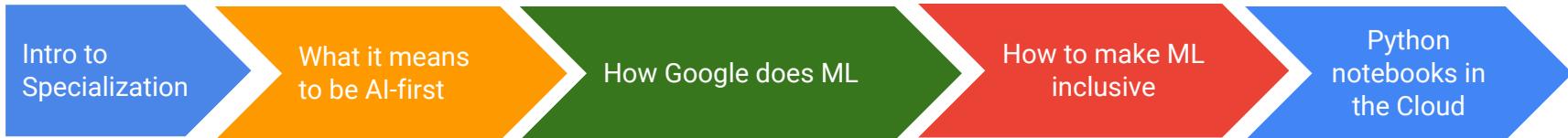
Machine Learning on Google Cloud Platform



Google Cloud

You will learn how to:

- Identify the origins of bias in ML
- Evaluate ML models with biases
- Make models inclusive



Agenda

Machine learning and human bias

Evaluating metrics with inclusion for your ML system

Equality of opportunity

How to find errors in your dataset using Facets

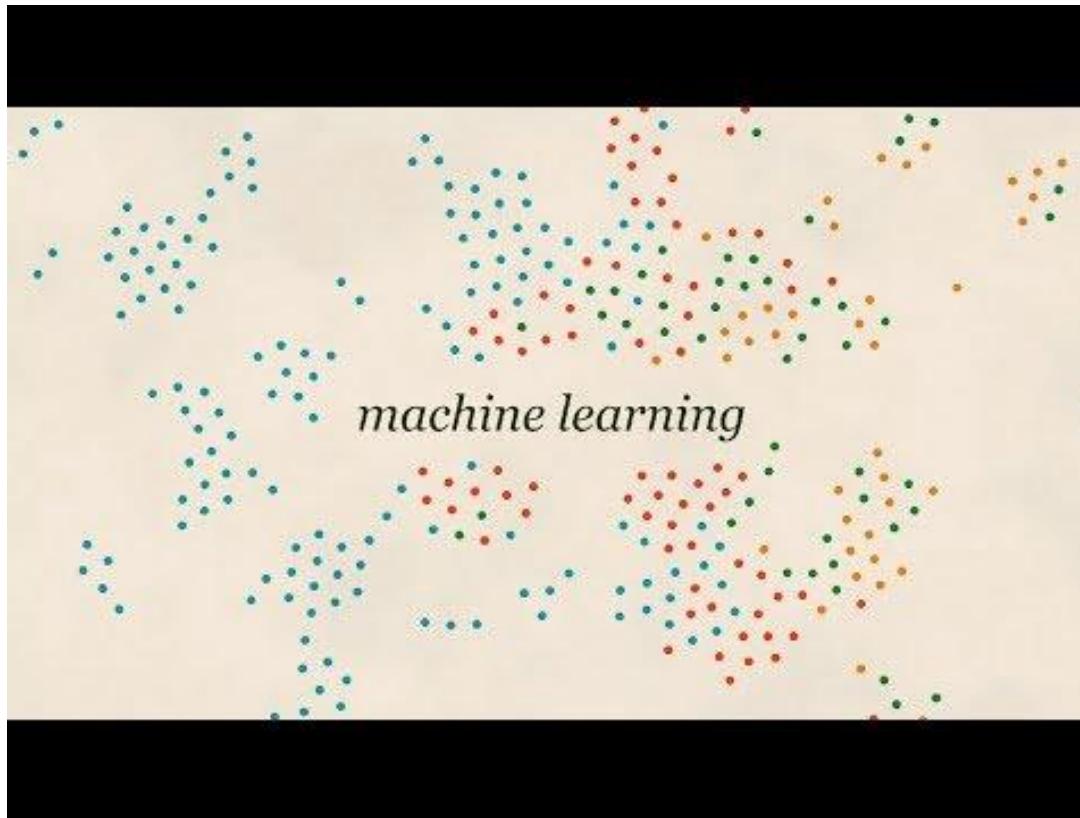
Course 1: How Google does ML

Module 4: How to make ML Inclusive

Lesson Title: **Machine Learning and Human Bias**

Format: Screencast

Human biases lead to biases in machine learning models



Course 1: How Google does ML

Module 4: How to Make ML Inclusive

Lesson Title: **Evaluation Metrics Across Subgroups**

Format: Screencast

Agenda

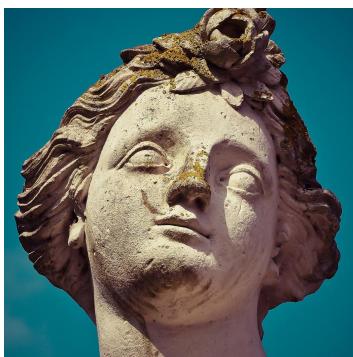
Machine learning and human bias

[Evaluating metrics with inclusion for your ML system](#)

Equality of opportunity

How to find errors in your dataset using Facets

Evaluate your model over subgroups also



The confusion matrix leads to evaluation metric insights

| | | Model Predictions | |
|--------|----------|--|----------|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Label says something exists. Model predicts it. | |
| | Negative | | |

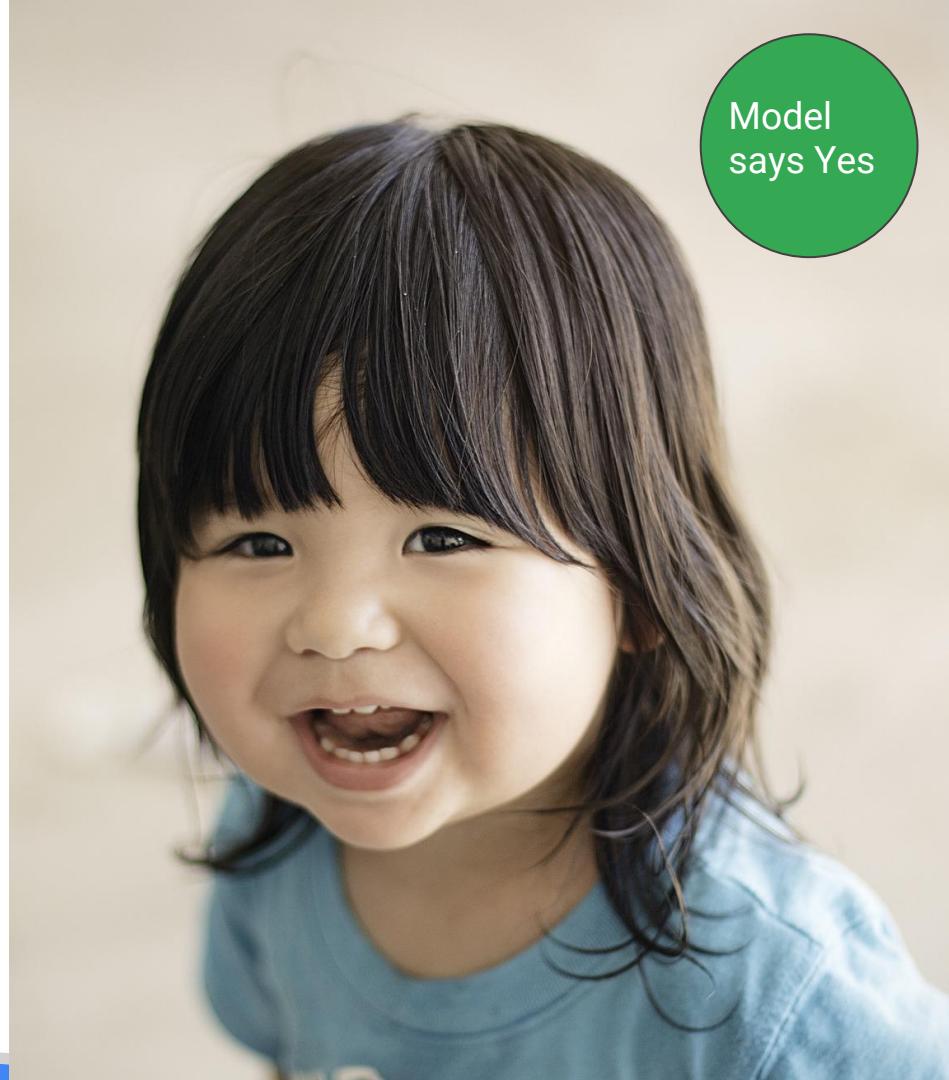


The confusion matrix leads to evaluation metric insights

True Positives (TP)

Label says something exists

Model predicts it exists



Model
says Yes

The confusion matrix leads to evaluation metric insights

| | | Model Predictions | |
|--------|----------|---|---|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Something exists. Model predicts it. | False Negatives (FN) <i>Type II Error</i> Something exists Model doesn't predict it |
| | Negative | |  A photograph of a young child with dark hair, smiling broadly. A small green checkmark icon is positioned to the right of the child's head. |

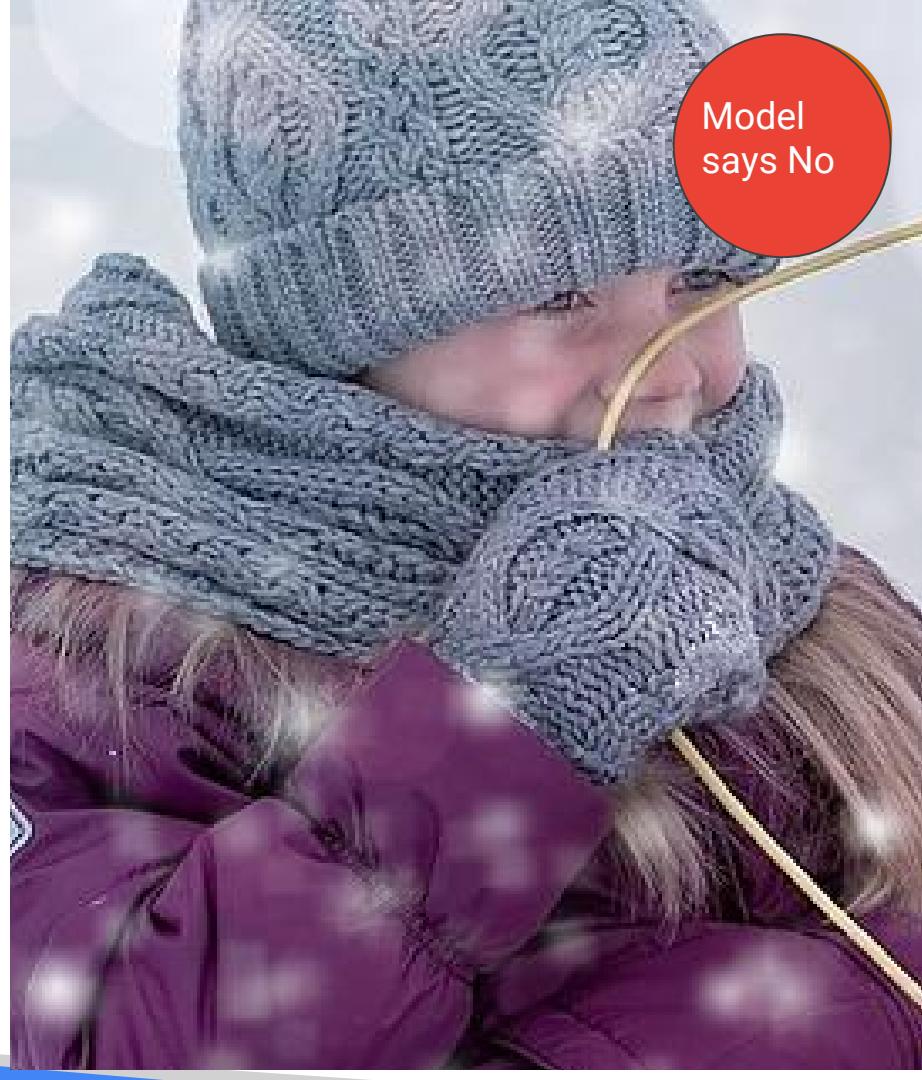
The confusion matrix leads to evaluation metric insights

False Negatives (FN)

Type II Error

Label says something exists

Model predicts it doesn't exist



The confusion matrix leads to evaluation metric insights

| | | Model Predictions | |
|--------|----------|---|---|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Label = something exists. Model predicts it. | False Negatives (FN) <i>Type II Error</i> Label = something exists Model doesn't predict it |
| | Negative | | True Negatives (TN) Something doesn't exist Model doesn't predict it |



The confusion matrix leads to evaluation metric insights

True Negative (TP)

Label says something doesn't exist

Model predicts it doesn't exist



The confusion matrix leads to evaluation metric insights

| | | Model Predictions | |
|--------|----------|--|--|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Something exists. Model predicts it. | False Negatives (FN) <i>Type II Error</i> Something exists. Model doesn't predict it |
| | Negative | False Positives (FP) <i>Type I Error</i> Something doesn't exist Model predicts it | True Negatives (TN) Something doesn't exist Model doesn't predict it |






The confusion matrix leads to evaluation metric insights

False Positives (FN)

Type I Error

Label says something doesn't exist

Model predicts it exists



Model
says Yes

False positives and false negatives errors occur when predictions and labels disagree

| | | Model Predictions | |
|--------|----------|---|--|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) | False Negatives (FN) <i>Type II Error</i> |
| | Negative | False Positives (FP) <i>Type I Error</i> | True Negatives (TN) |

True Positives (TP): Model says Yes (Green Circle). Image: Smiling child.

False Negatives (FN) / Type II Error: Model says No (Red Circle). Image: Child knitting.

False Positives (FP) / Type I Error: Model says Yes (Green Circle). Image: Statue head.

True Negatives (TN): Model says No (Red Circle). Image: Teddy bear.

Course 1: How Google does ML

Module 4: Inclusive ML

Lesson Title: **Statistical Measurements**

Format: Screencast

Evaluation metrics can help highlight areas where machine learning could be more inclusive

| | | Model Predictions | |
|--------|----------|---|--|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Label says something exists. The model predicts it. | False Negatives (FN) <i>Type II Error</i> Label says something exists Model doesn't predict it |
| | Negative | False Positives (FP) <i>Type I Error</i> Label says something doesn't exist Model predicts it | True Negatives (TN) Label says something doesn't exist Model doesn't predict it |

Model says Yes (Green circle, bottom-left)

Model says No (Red circle, top-right)

False negative rate is the fraction of true faces that are not detected by the ML system

| | | Model Predictions | |
|---------------------|----------|---|--|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Label says something exists. The model predicts it. | False Negatives (FN) Type II Error Label says something exists Model doesn't predict it |
| | Negative | False Positives (FP) Label says something doesn't exist. The model predicts it. | True Negatives (TN) Label says something doesn't exist. The model doesn't predict it. |
| False Negative Rate | | $\frac{\text{False Negatives}}{\text{False Negatives} + \text{True Positives}}$ | |

False negative rate is the fraction of true faces that are not detected by the ML system



True Positives (TP)

Label says something exists.
The model predicts it.



False Negatives (FN)

Type II Error
Label says something exists
Model doesn't predict it

$$\text{False Negative Rate} = \frac{\text{False Negatives}}{\text{False Negatives} + \text{True Positives}}$$

False positive rate is the fraction of the faces that the ML model detects that are not really faces

| | | Model Predictions | |
|--------|----------|---|--|
| | | Positive | Negative |
| Labels | Positive | True Positives (TP) Label says something exists. The model predicts it. | False Positive Rate $\frac{\text{False Positives}}{\text{False Positives} + \text{True Negatives}}$ |
| | Negative | False Positives (FP) Type I Error Label says something doesn't exist Model predicts it | |

False positive rate is the fraction of the faces that the ML model detects that are not really faces



True Positives (TP)

Label says something exists.
The model predicts it.



False Positives (FP)

Type I Error
Label says something doesn't exist
Model predicts it

False Positive Rate =

$$\frac{\text{False Positives}}{\text{False Positives} + \text{True Negatives}}$$

Sometimes false positives are better than false negatives

Privacy in images



False positive



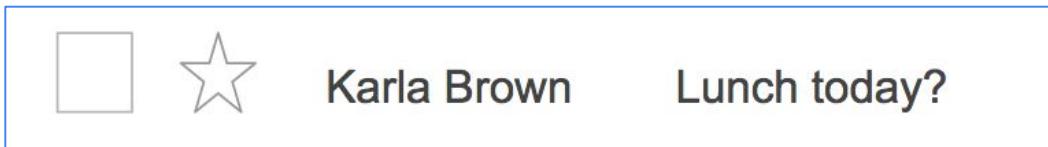
False negative

Sometimes false negatives are better than false positives

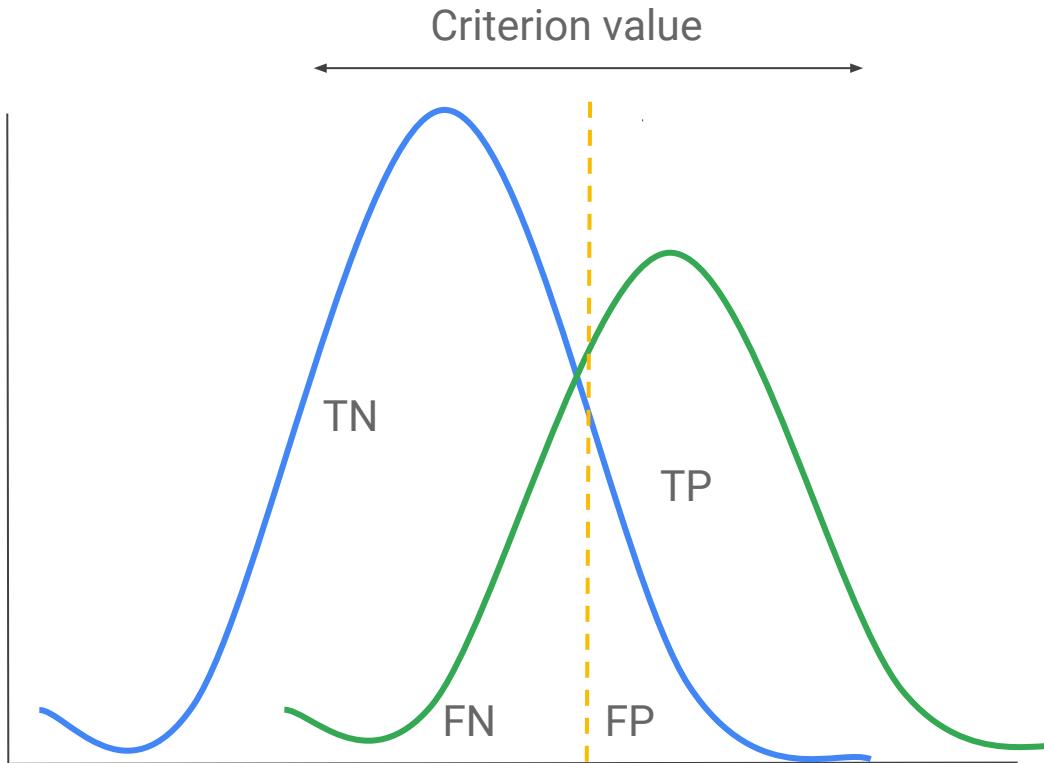
False Negative: E-mail that is SPAM is not caught, so you see it in your inbox.



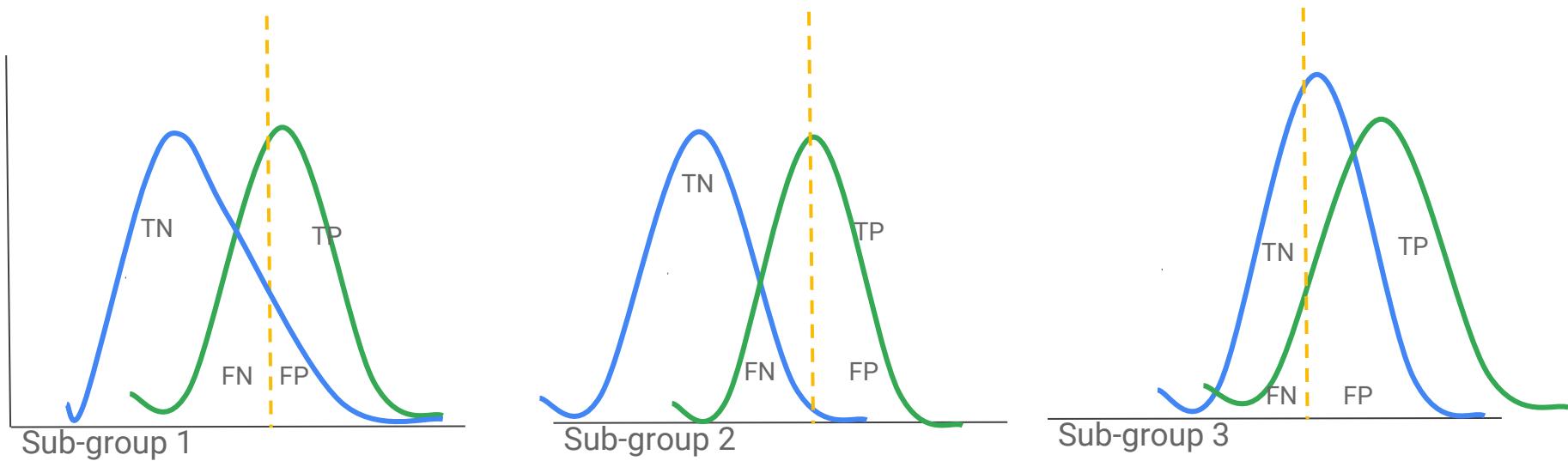
False Positive: E-mail flagged as SPAM is removed from your inbox.

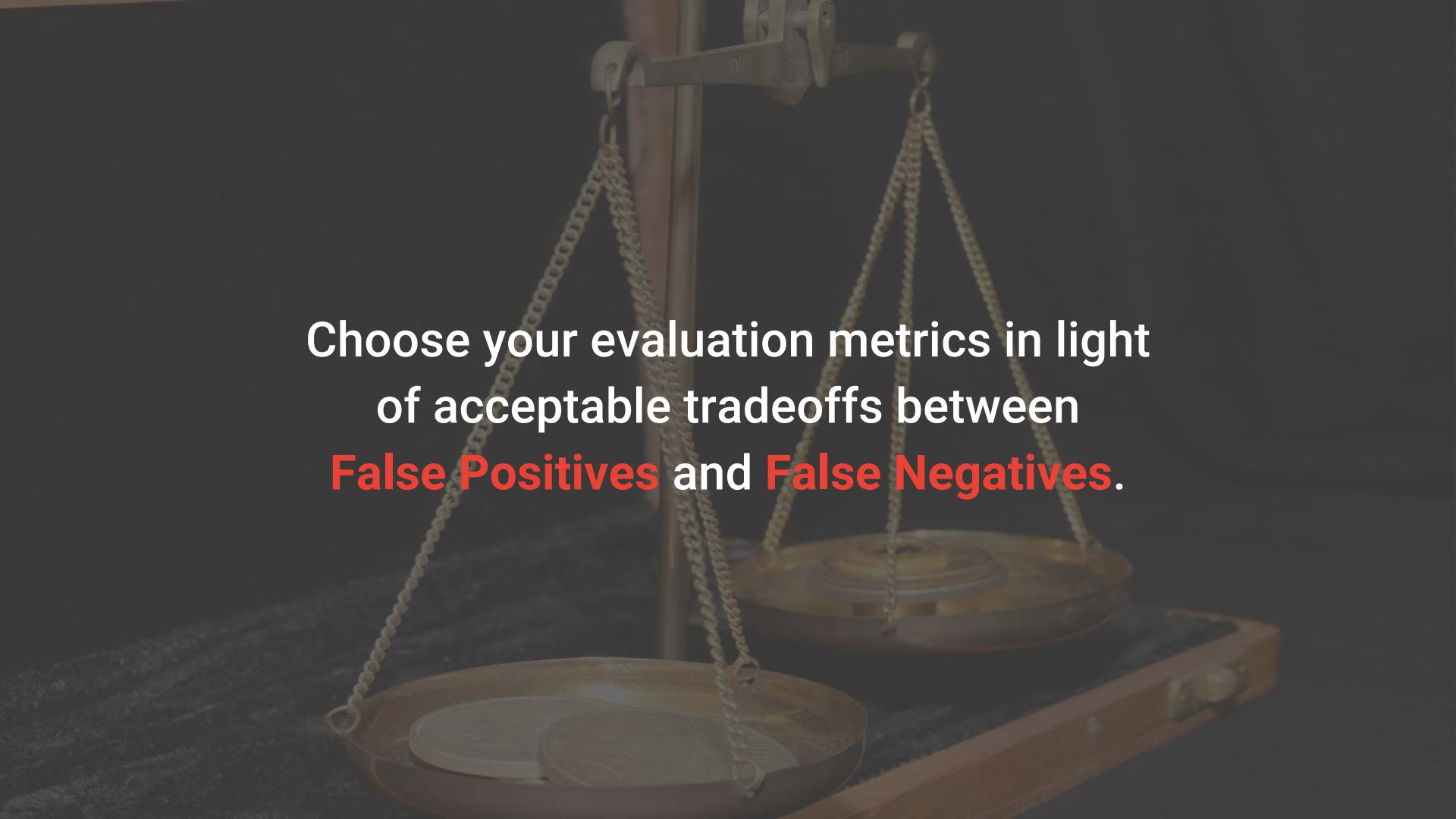


Find the threshold that brings the precision or recall to acceptable values



Check the precision/recall you obtain with that threshold in each of your subgroups





Choose your evaluation metrics in light
of acceptable tradeoffs between
False Positives and **False Negatives**.

Course 1: How Google does ML

Module 4: Human-centered ML

Lesson Title: **Equality of Opportunity**

Format: Screencast

Agenda

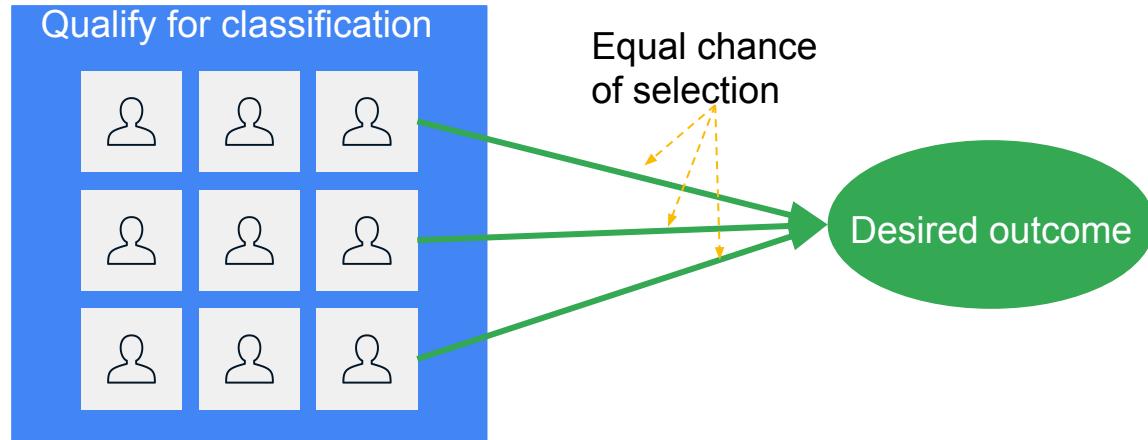
Machine learning and human bias

Evaluating metrics with inclusion for your ML system

Equality of opportunity

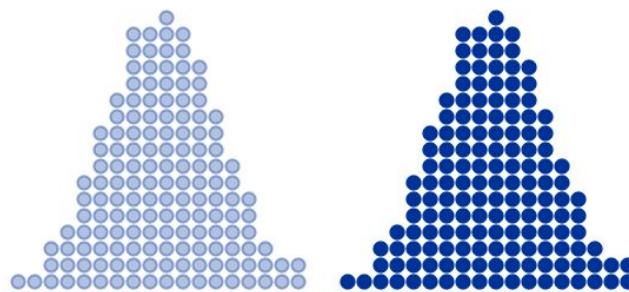
How to find errors in your dataset using Facets

The Equality of Opportunity approach strives to give individuals an equal chance of desired outcome



A toy classifier to predict who will pay back their loan involves two populations that might overlap

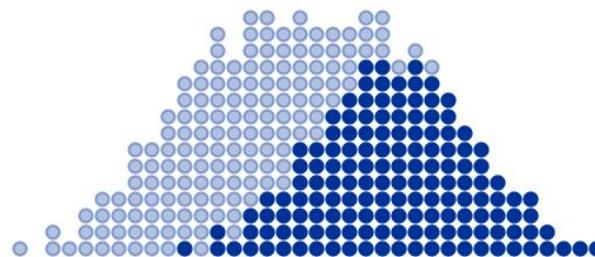
0 10 20 30 40 50 60 70 80 90 100



would default on loan

would pay back loan

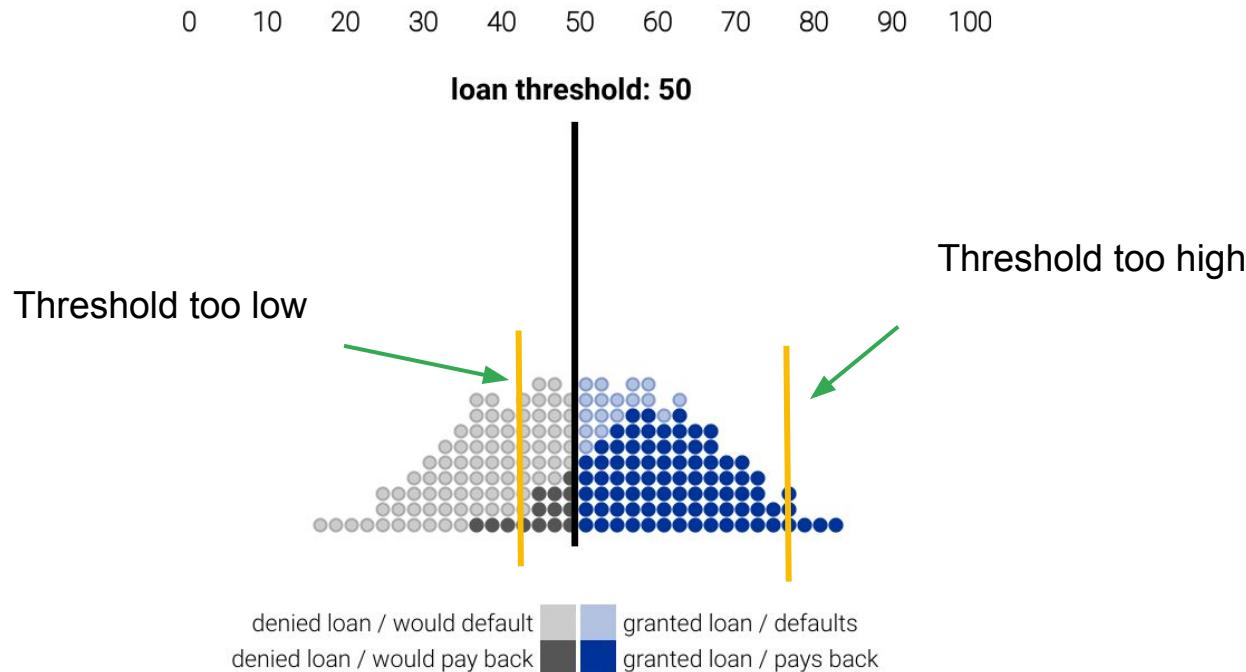
0 10 20 30 40 50 60 70 80 90 100



would default on loan

would pay back loan

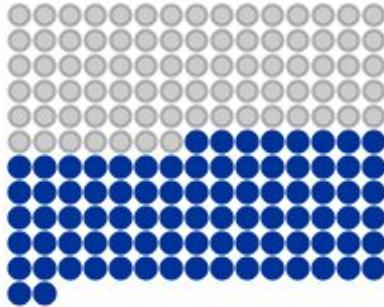
Picking a credit score threshold involves a tradeoff



The impact of a threshold on credit score is evaluated based on its impact on customers and on loan repayment

Correct 84%

loans granted to paying applicants and denied to defaulters



Incorrect 16%

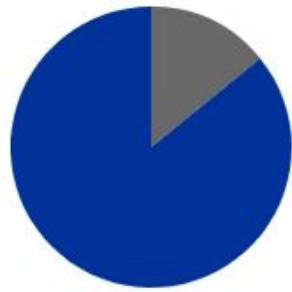
loans denied to paying applicants and granted to defaulters



Simulating the impact of a threshold on profit

True Positive Rate 86%

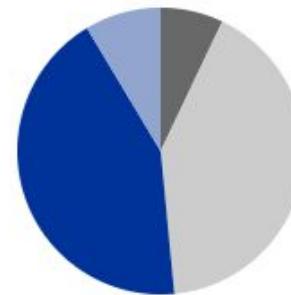
percentage of paying
applications getting loans



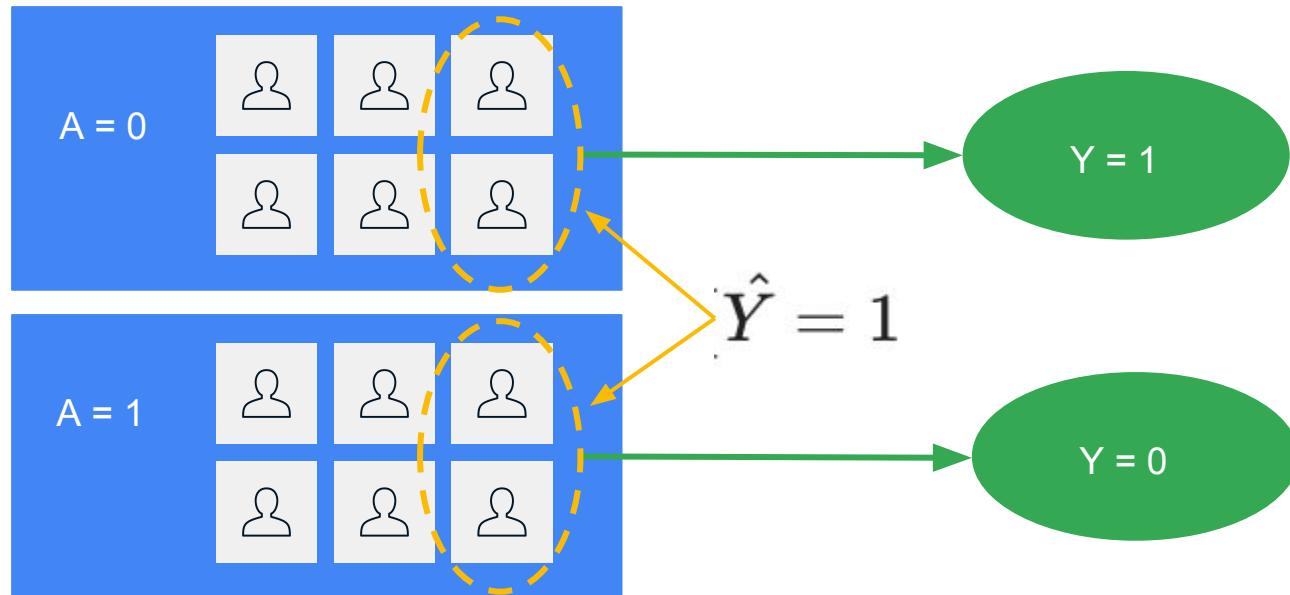
Profit: 13600

Positive Rate 52%

percentage of all
applications getting loans



Classification and Discrimination must obey the Equality of Opportunity Principle



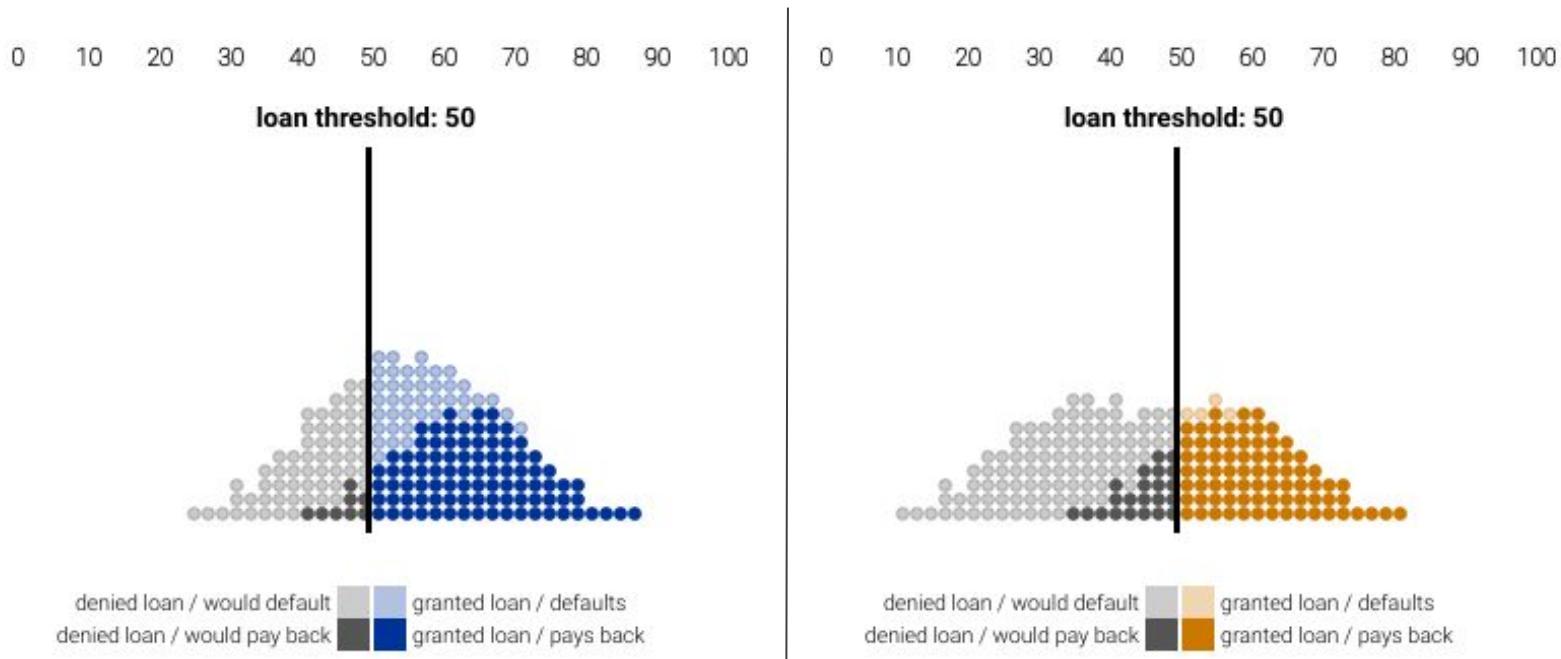
Course 1: How Google does ML

Module 4: Human-centered ML

Lesson Title: **Simulating Decisions**

Format: Screencast

Simulating decisions with no constraints can lead to unequal distribution

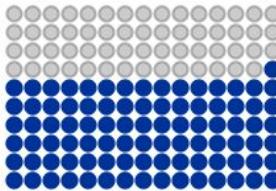


A successful loan makes \$300
An unsuccessful loan costs \$700
Credit scores are between 0 - 100

Simulating decisions with no constraints can lead to unequal distribution

Total profit = 19600

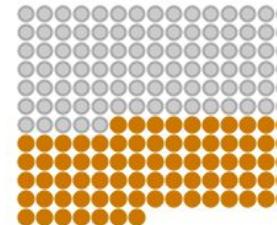
Correct 76%
loans granted to paying applicants and denied to defaulters



Incorrect 24%
loans denied to paying applicants and granted to defaulters



Correct 87%
loans granted to paying applicants and denied to defaulters



Incorrect 13%
loans denied to paying applicants and granted to defaulters

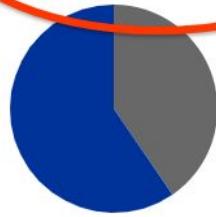


Threshold

- Credit Score of 50 for Blue Group
- Credit Score of 50 for Orange Group

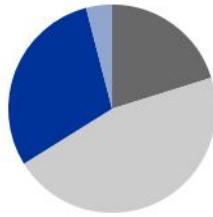
Simulating decisions for max profit result in unequal standards

True Positive Rate 60%
percentage of paying
applications getting loans

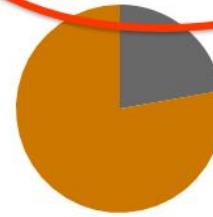


Profit: 12100

Positive Rate 34%
percentage of all
applications getting loans

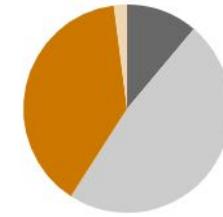


True Positive Rate 78%
percentage of paying
applications getting loans



Profit: 20300

Positive Rate 41%
percentage of all
applications getting loans

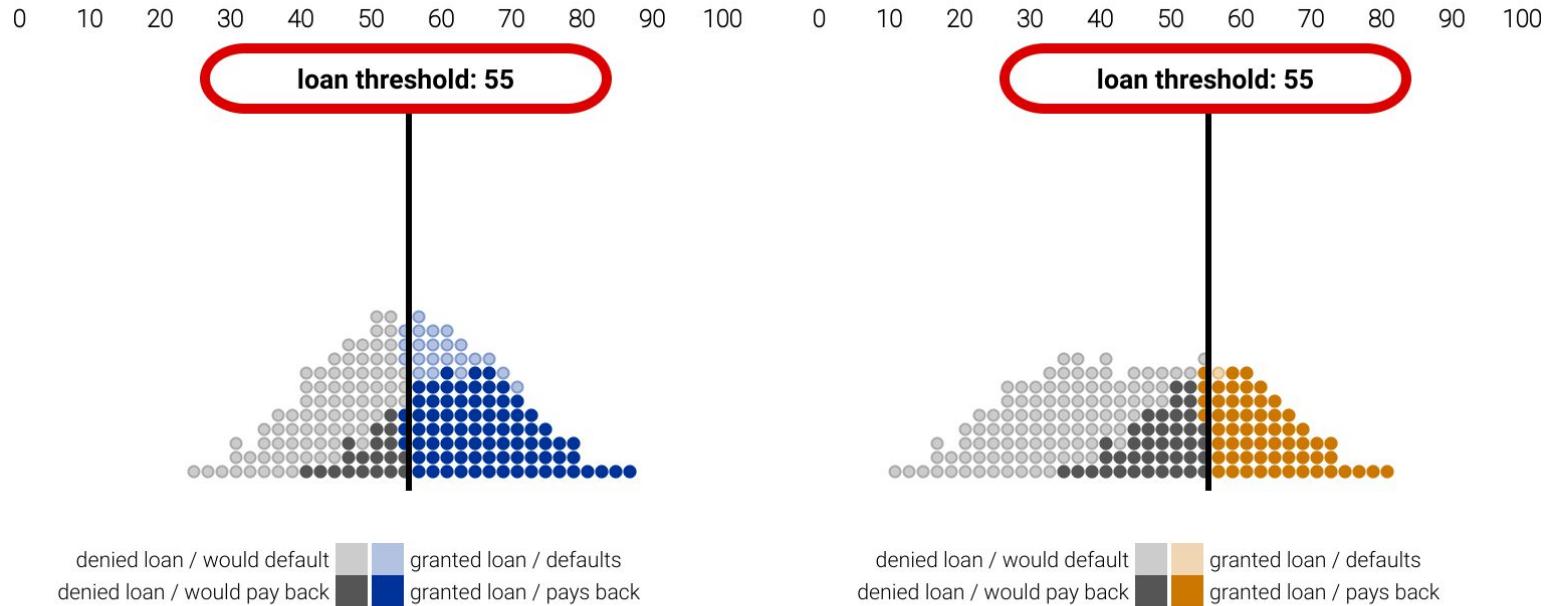


Threshold

- Credit Score of 61 for Blue Group
- Credit Score of 50 for Orange Group

Total profit: 32400

Simulating decisions with group unaware holds everyone to the same standard, which can be unfair to some groups

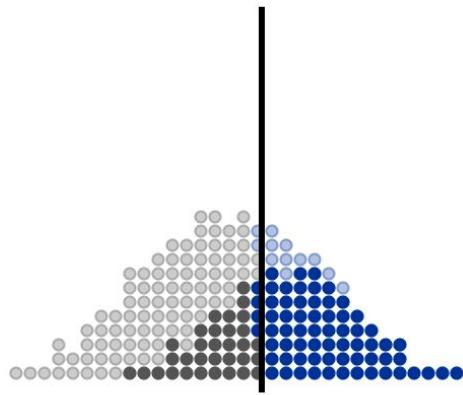


Total profit = 25600

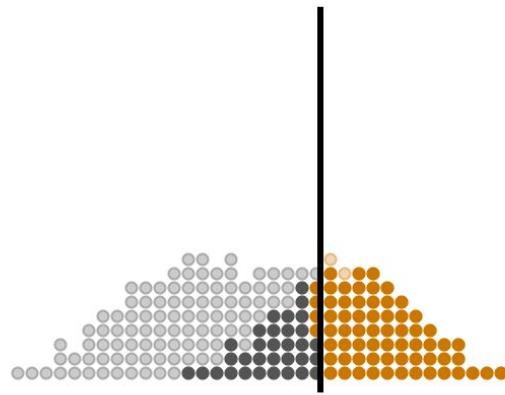
Simulating decisions equal opportunity results in an identical true positive rate for all groups

0 10 20 30 40 50 60 70 80 90 100 0 10 20 30 40 50 60 70 80 90 100

loan threshold: 59



loan threshold: 53



denied loan / would default granted loan / defaults
denied loan / would pay back granted loan / pays back

denied loan / would default granted loan / defaults
denied loan / would pay back granted loan / pays back

Total profit = 30400

Course 1: How Google does ML

Module 4: Human-centered ML

**Lesson Title: Finding Errors in Your Dataset
Using Facets**

Format: Screencast

Agenda

Machine learning and human bias

Evaluating metrics with inclusion for your ML system

Equality of opportunity

How to find errors in your dataset using Facets

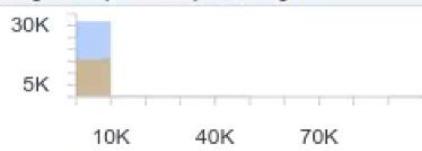
<https://research.googleblog.com/2017/07/facets-open-source-visualization-tool.html>

Facets gives users a quick understanding of the distribution of values across features of their datasets

Sort by Non-uniformity Reverse order Feature search

Features: int(6) string(9)

train test

| Numeric Features (6) | | | | | | | | Chart to show |
|----------------------|---------|----------|----------|--------|-------|--------|-------|---|
| count | missing | mean | std dev | zeros | min | median | max | Standard |
| Capital Gain | | | | | | | | <input type="checkbox"/> log <input type="checkbox"/> expand <input type="checkbox"/> percentages |
| 32.6k | 0% | 1,077.65 | 7,385.29 | 91.67% | 0 | 0 | 100k |  |
| 16.3k | 0% | 1,081.91 | 7,583.94 | 91.87% | 0 | 0 | 100k |  |
| Capital Loss | | | | | | | |  |
| 32.6k | 0% | 87.3 | 402.96 | 95.33% | 0 | 0 | 4,356 |  |
| 16.3k | 0% | 87.9 | 403.11 | 95.31% | 0 | 0 | 3,770 | |
| fnlwgt | | | | | | | | |
| 32.6k | 0% | 190k | 106k | 0% | 12.3k | 178k | 1.48M |  |
| 16.3k | 0% | 189k | 106k | 0% | 13.5k | 178k | 1.49M | |

In Facets features are sorted by non-uniformity, with the feature with the most non-uniform distribution at the top

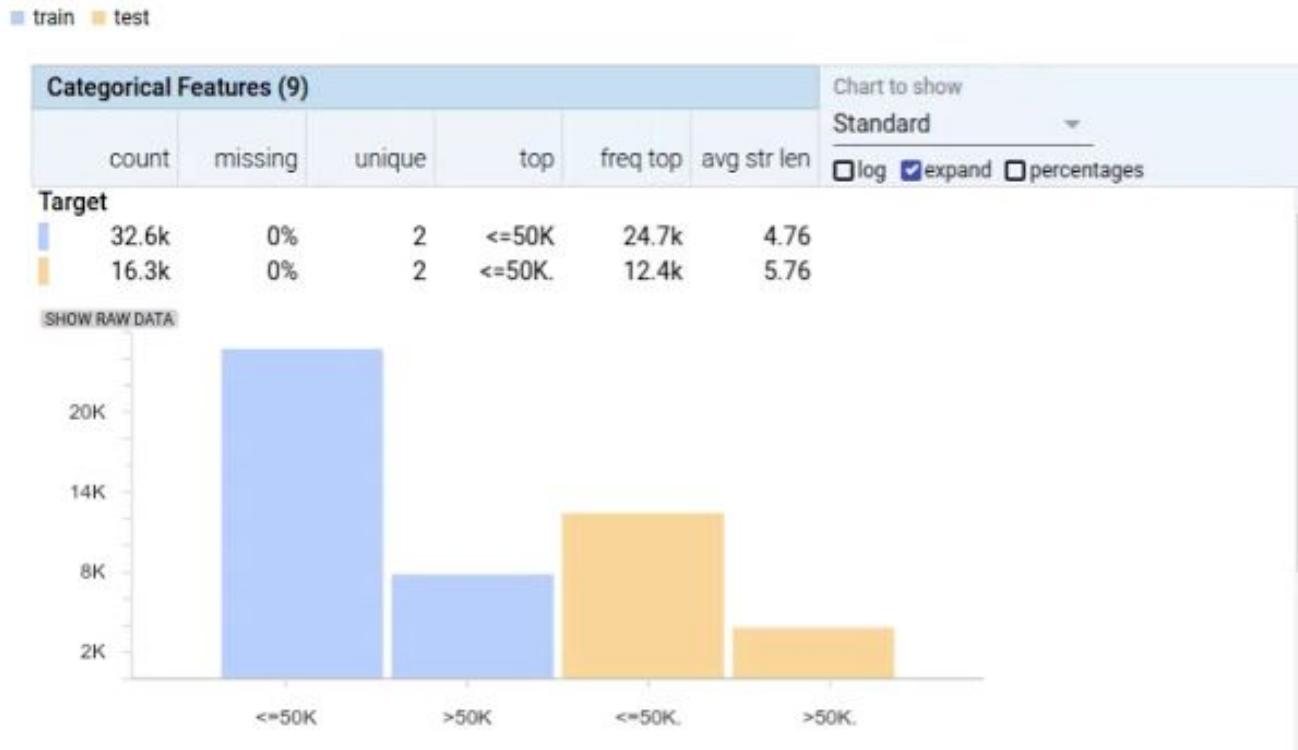
Sort by Non-uniformity Reverse order Feature search

Features: int(6) string(9)

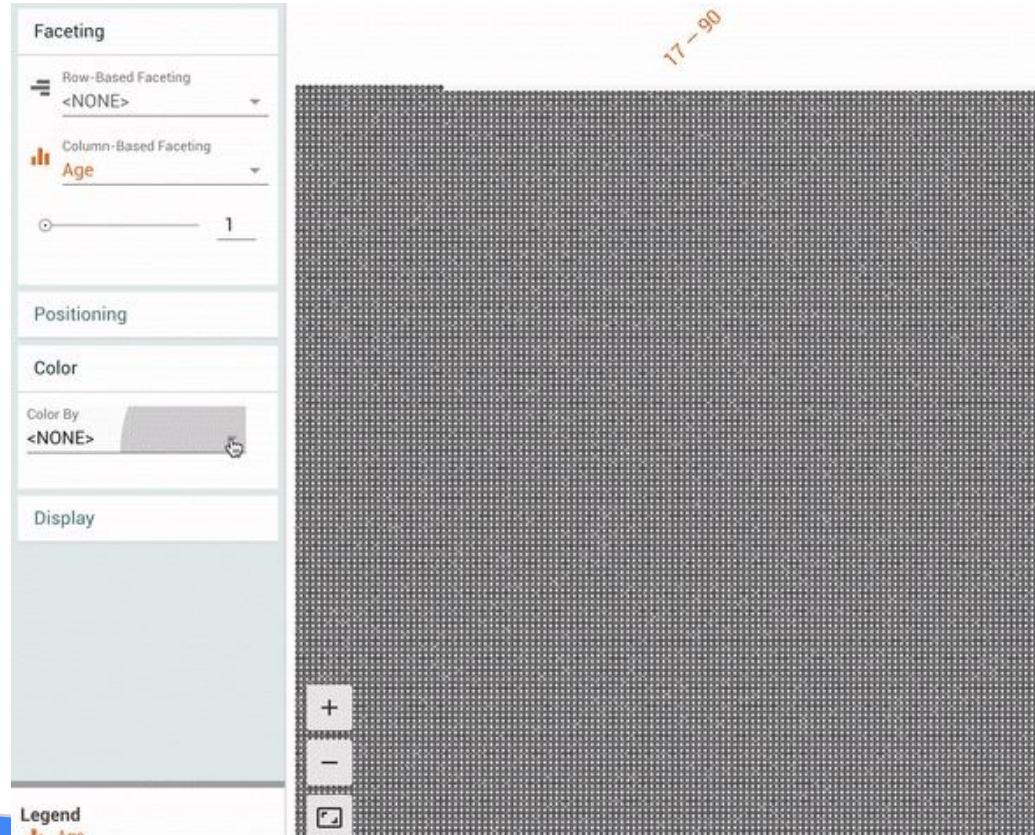
train test

| Numeric Features (6) | | | | | | | | Chart to show | |
|----------------------|-------|---------|----------|----------|---------------|-----|--------|---------------|---|
| | count | missing | mean | std dev | zeros | min | median | max | Standard |
| Capital Gain | | | | | | | | | <input type="checkbox"/> log <input type="checkbox"/> expand <input type="checkbox"/> percentages |
| | 32.6k | 0% | 1,077.65 | 7,385.29 | 91.67% | 0 | 0 | 100k |  |
| | 16.3k | 0% | 1,081.91 | 7,583.94 | 91.87% | 0 | 0 | 100k |  |
| Capital Loss | | | | | | | | | |
| | 32.6k | 0% | 87.3 | 402.96 | 95.33% | 0 | 0 | 4,356 | |
| | 16.3k | 0% | 87.9 | 403.11 | 95.31% | 0 | 0 | 3,770 | |

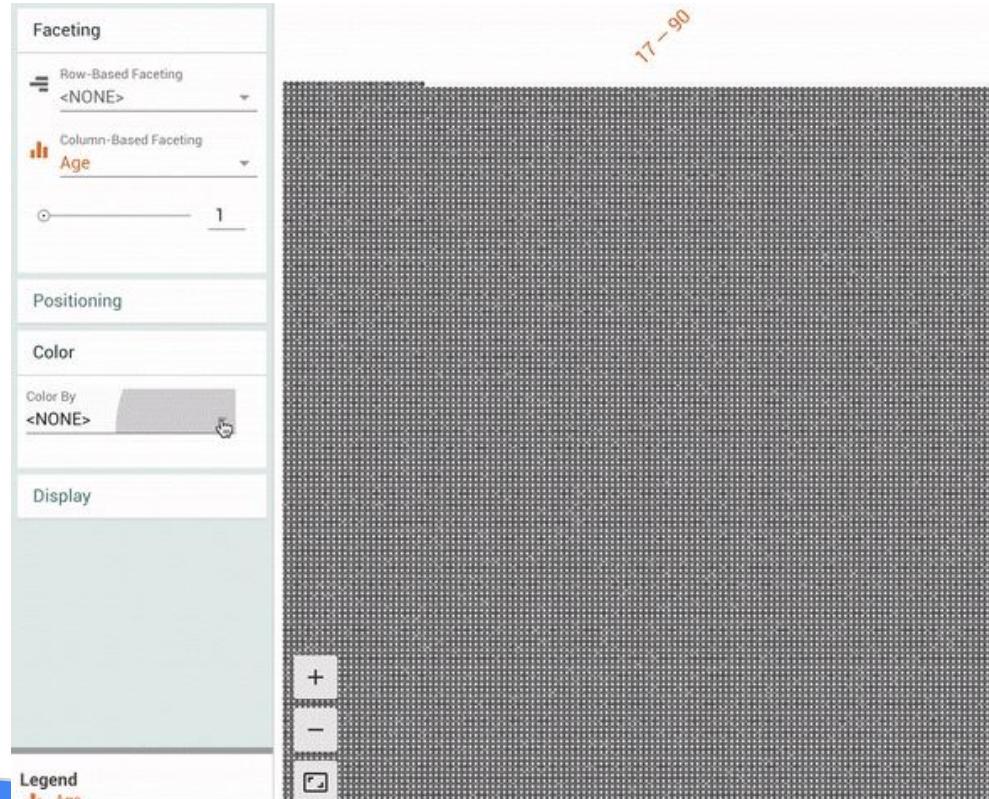
Facets features are sorted by distribution distance



Facets Dive provides an easy-to-customize, intuitive interface



Color the data points by one feature, then facet in another

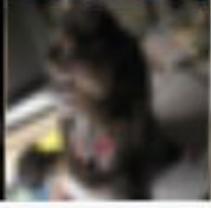
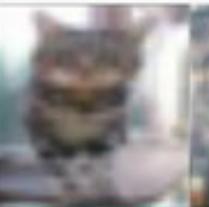


Explore CIFAR-10 for Errors using Facets Dive





Exploring a dataset for Errors



cloud.google.com

