



Seeds of Change

Root causes of algorithmic unfairness, and a path forward

go/seeds-of-change

mmitchellai@

Google Confidential

A quick recap

**Training data are
collected and
annotated**

**Model is trained and
evaluated**

**Media are filtered,
ranked, aggregated,
or generated**

Input

Output

Training data are
collected and
annotated

Model is trained and
evaluated

Media are filtered,
ranked, aggregated,
or generated

Input

Output

Training data are
collected and
annotated

Model is trained and
evaluated

Media are filtered,
ranked, aggregated,
or generated

Input

Output

Training data are
collected and
annotated

Model is trained and
evaluated

Media are filtered,
ranked, aggregated,
or generated

Input

Output

At a high level, where is
unfairness creeping in?

Within the data

Reporting bias

Selection bias

Overgeneralization bias

Out-group homogeneity bias

Stereotypical bias

Historical Unfairness

Implicit associations

Implicit stereotypes

Prejudice

Group Attribution error

Halo effect

Within the data

Reporting bias
Selection bias
Overgeneralization bias
Out-group homogeneity bias
Stereotypical bias
Historical Unfairness
Implicit associations
Implicit stereotypes
Prejudice
Group Attribution error
Halo effect

Data collection and annotation

Sampling error
Non-sampling error
Insensitivity to sample size
Correspondence bias
In-group bias
Bias blind spot
Confirmation bias
Subjective validation
Experimenter's bias
Choice-supportive bias
Neglect of probability
Anecdotal fallacy
Illusion of validity
Automation bias

Within the data

Reporting bias
Selection bias
Overgeneralization bias
Out-group homogeneity bias
Stereotypical bias
Historical Unfairness
Implicit associations
Implicit stereotypes
Prejudice
Group Attribution error
Halo effect

Data collection and annotation

Sampling error
Non-sampling error
Insensitivity to sample size
Correspondence bias
In-group bias
Bias blind spot
Confirmation bias
Subjective validation
Experimenter's bias
Choice-supportive bias
Neglect of probability
Anecdotal fallacy
Illusion of validity
Automation bias

Training and evaluation

Evaluation metric
Features
Objective Function
Model architecture
Variables
Tasks
Hyperparameters

Within the data

Reporting bias
Selection bias
Overgeneralization bias
Out-group homogeneity bias
Stereotypical bias
Historical Unfairness
Implicit associations
Implicit stereotypes
Prejudice
Group Attribution error
Halo effect

Data collection and annotation

Sampling error
Non-sampling error
Insensitivity to sample size
Correspondence bias
In-group bias
Bias blind spot
Confirmation bias

Subjective validation
Experimenter's bias
Choice-supportive bias
Neglect of probability
Anecdotal fallacy
Illusion of validity
Automation bias

Training and evaluation

Evaluation metric
Features
Objective Function
Model architecture
Variables
Tasks
Hyperparameters

We use data to estimate how
likely different things are

Stereotypical bias

A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

How could this be?



A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

How could this be?



A man and his son are in a terrible accident and are rushed to the hospital in critical care.

The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"

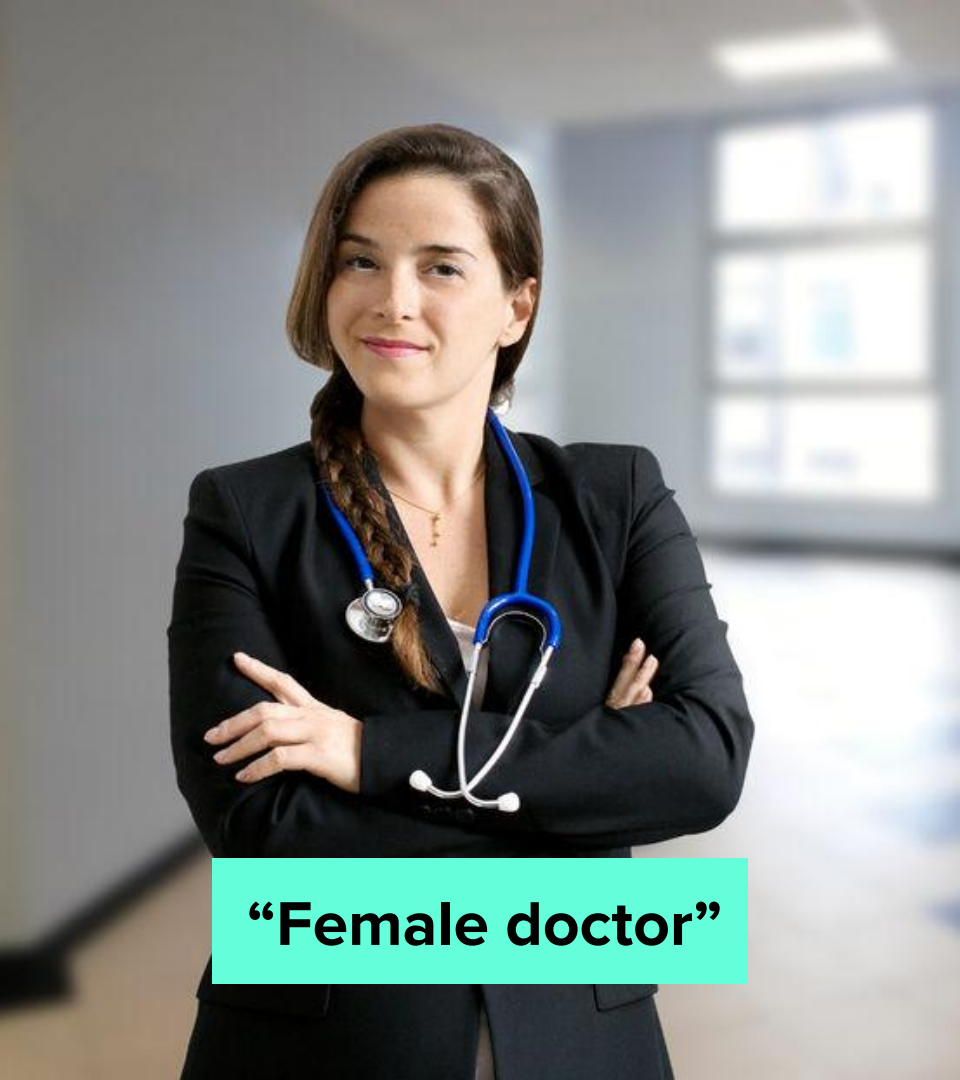
How could this be?



“Female doctor”



“Doctor”



“Female doctor”

The majority of test subjects overlooked the possibility that the doctor is a she—including men, women, and self-described feminists.

[Wapman & Belle, Boston University](#)

Reporting bias



"male surgeon"



All

Images

Videos

News

Shopping

More

Settings

Tools

About 89,800 results (0.28 seconds)



"female surgeon"



All

Images

Videos

News

Shopping

More

Settings

Tools

About 199,000 results (0.53 seconds)



"male surgeon"



All

Images

Videos

News

Shopping

More

Settings

Tools

About 89,800 results (0.28 seconds)

**Real-world diversity
among surgeons**

81%

Male

19%

Female



"female surgeon"



All

Images

Videos

News

Shopping

More

Settings

Tools

About 199,000 results (0.53 seconds)

SOURCE

[Statistics on the Number of Women Surgeons in the United States](#)

World learning from text

Gordon and Van Durme, 2013

Word	Frequency in corpus
“spoke”	11,577,917
“laughed”	3,904,519
“murdered”	2,834,529
“inhaled”	984,613
“breathed”	725,034
“hugged”	610,040
“blinked”	390,692
“exhale”	168,985

World learning from text

Gordon and Van Durme, 2013

Word	Frequency in corpus
“spoke”	11,577,917
“laughed”	3,904,519
“murdered”	2,834,529
“inhaled”	984,613
“breathed”	725,034
“hugged”	610,040
“blinked”	390,692
“exhale”	168,985

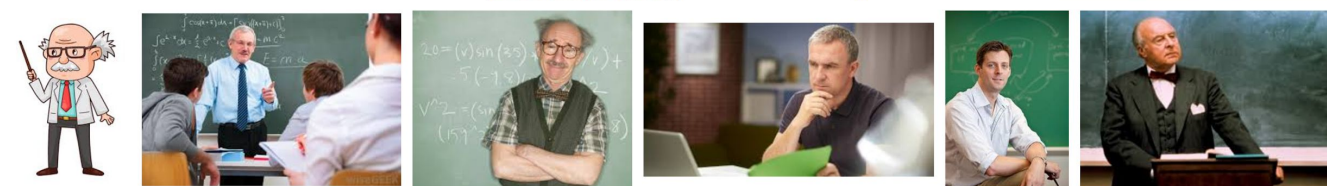
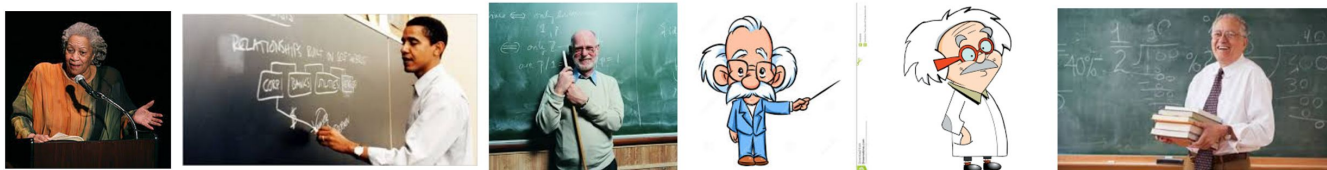
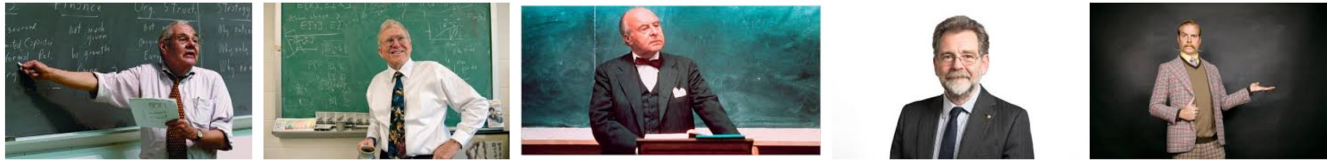
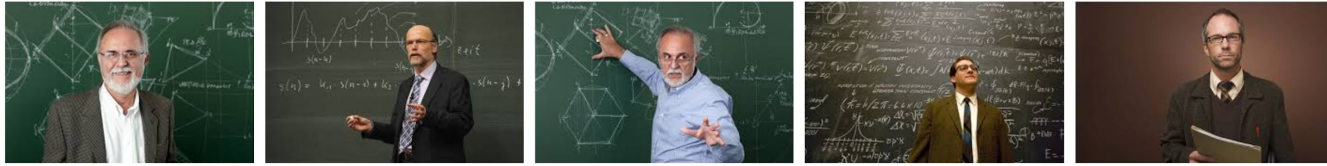
Top results show historical unfairness,
implicit associations, and implicit
stereotypes reflected in **Reporting Bias**

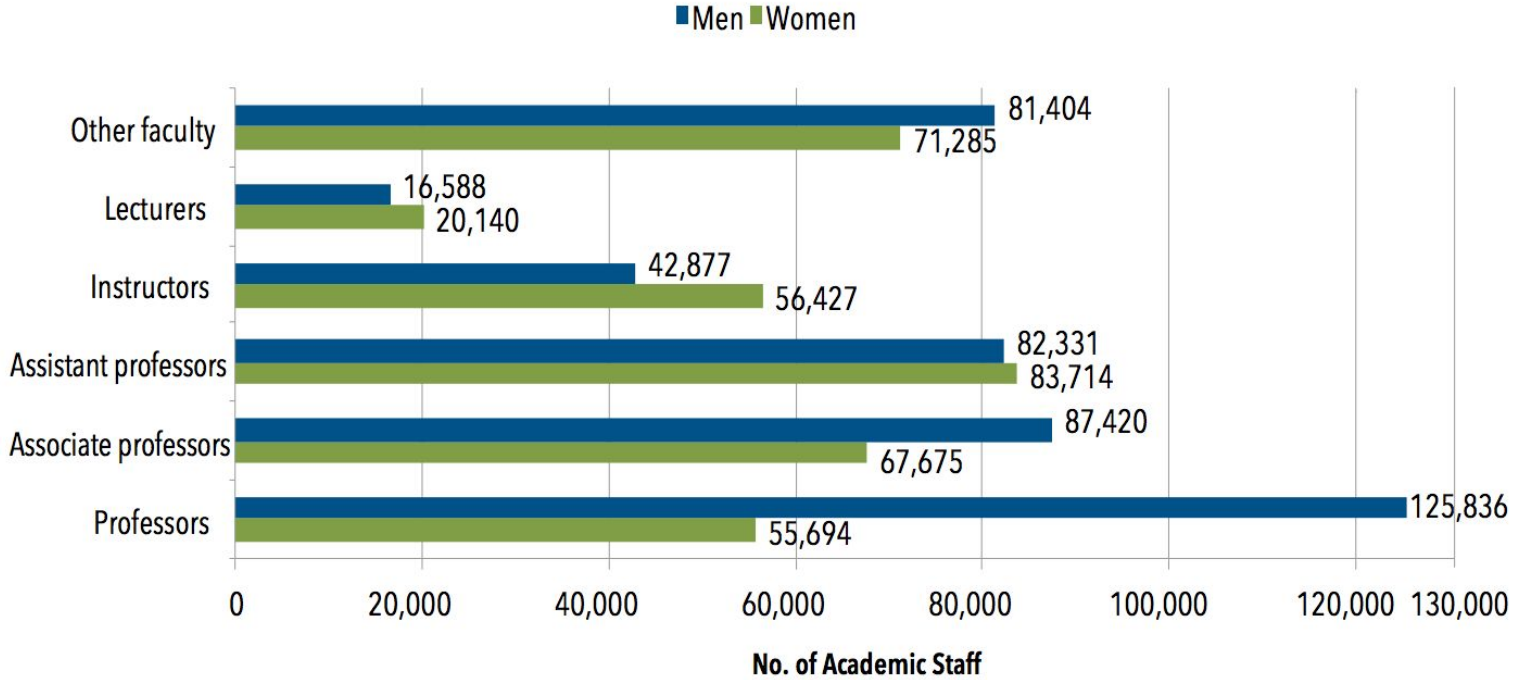
We tend to mention and share things that are outside of our expectation of day-to-day norms; ignoring the things that “go without saying”.

**Training data are
collected and
annotated**

**Media are filtered,
ranked, aggregated,
or generated**

- hot female
- android
- male
- baby
- african american
- indian
- chinese
- japanese
- university
- college
- classroom
- lab
- concord hospital
- cartoon
- meme

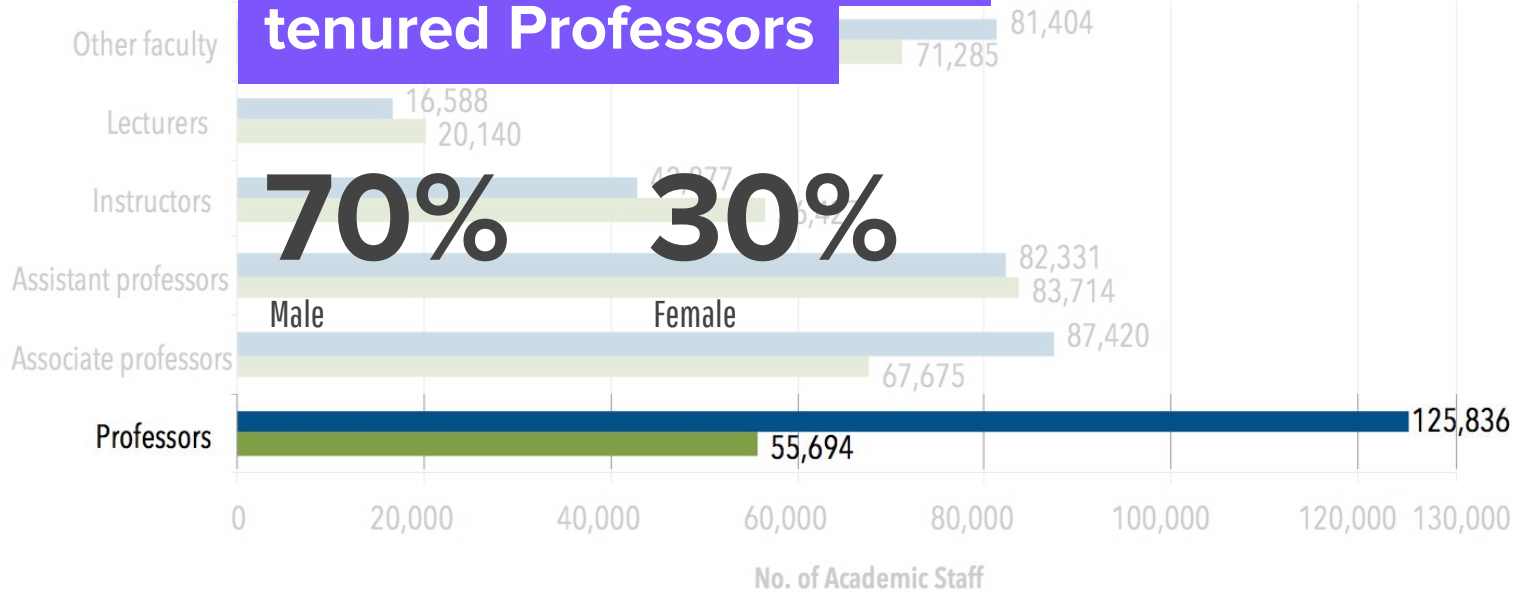




SOURCE

Johnson, Heather L. 2016. Pipelines, Pathways, and Institutional Leadership: An Update on the Status of Women in Higher Education. Washington, DC: American Council on Education

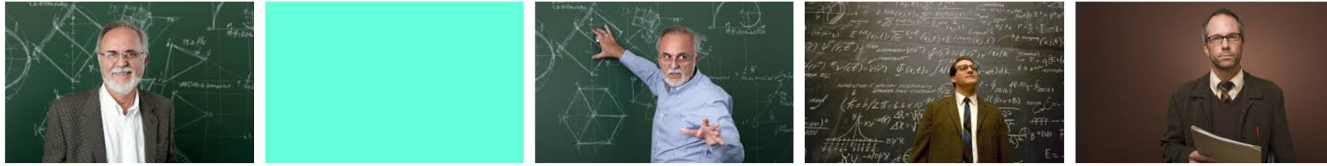
Gender diversity among tenured Professors



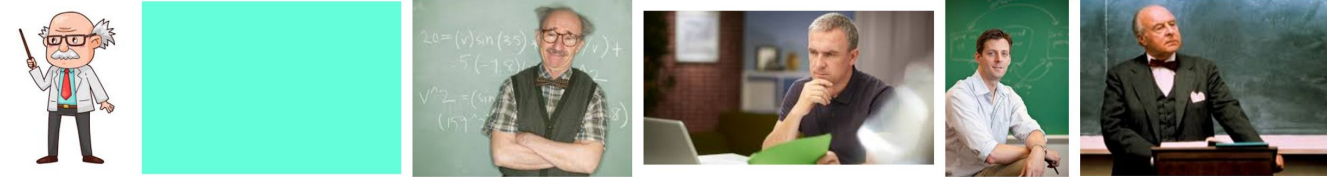
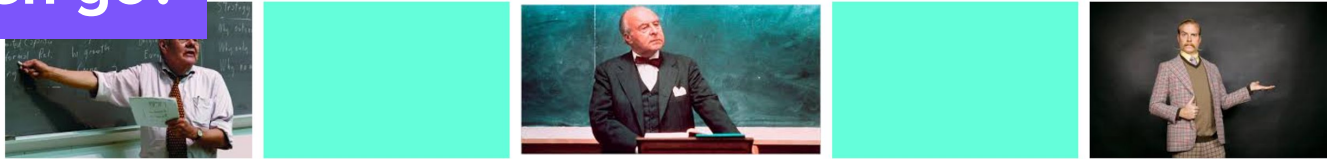
SOURCE

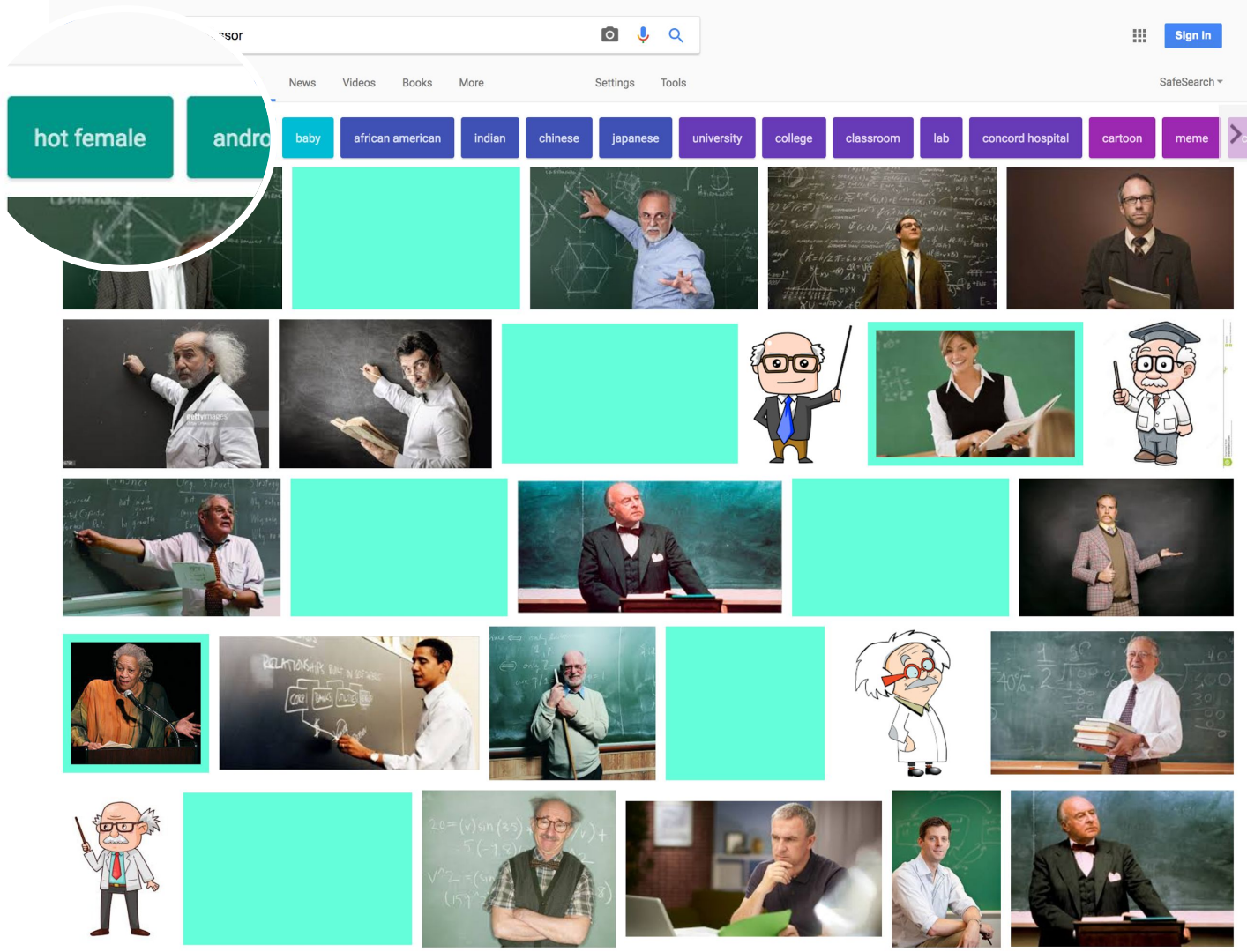
Johnson, Heather L. 2016. Pipelines, Pathways, and Institutional Leadership: An Update on the Status of Women in Higher Education. Washington, DC: American Council on Education

- hot female
- android
- male
- baby
- african american
- indian
- chinese
- japanese
- university
- college
- classroom
- lab
- concord hospital
- cartoon
- meme



Where did the women go?





**Training data are
collected and
annotated**

**Model is trained and
evaluated**

**Media are filtered,
ranked, aggregated,
or generated**

INSIGHT: EVALUATION METRIC

The Confusion Matrix

Evaluation Metric Insights: The Confusion Matrix

Predictions

References

Evaluation Metric Insights: The Confusion Matrix

Predictions

References

Create for each (subgroup, prediction) pair.
Compare across subgroups.

Evaluation Metric Insights: The Confusion Matrix

Predictions

References

Create for each (subgroup, prediction) pair.
Compare across subgroups.

Example: women, face detection
men, face detection

Evaluation Metric Insights: The Confusion Matrix

		Predictions	
		Positive	Negative
References	Positive		
	Negative		

Evaluation Metric Insights: The Confusion Matrix

		Predictions	
		Positive	Negative
References	Positive	Reference says something exists Model predicts it True Positives	Reference says something exists Model doesn't predict it False Negatives <i>Type II Error</i>
	Negative	Reference says something doesn't exist Model predicts it False Positives <i>Type I error</i>	Reference says something doesn't exist Model doesn't predict it True Negatives

Evaluation Metric Insights: The Confusion Matrix

		Predictions	
		Positive	Negative
References	Positive	Reference says something exists Model predicts it True Positives	Reference says something exists Model doesn't predict it False Negatives <i>Type II Error</i>
	Negative	Reference says something doesn't exist Model predicts it False Positives <i>Type I error</i>	Reference says something doesn't exist Model doesn't predict it True Negatives

The Problem Areas

Evaluation Metric Insights: The Confusion Matrix

		Predictions		Calculate
		Positive	Negative	
References	Positive	Reference says something exists Model predicts it True Positives	Reference says something exists Model doesn't predict it False Negatives <i>Type II Error</i>	True Positive Rate/ Sensitivity/ Recall False Negative Rate/ Miss Rate
	Negative	Reference says something doesn't exist Model predicts it False Positives <i>Type I error</i>	Reference says something doesn't exist Model doesn't predict it True Negatives	False Positive Rate/ Fallout True Negative Rate/ Specificity
		Precision / Positive Predictive Value, False Discovery Rate	Negative Predictive Value, False Omission Rate	LR+, LR-


Evaluation Metric Insights: The Confusion Matrix

		Predictions		Calculate
		Positive	Negative	
References	Positive	Reference says something exists Model predicts it True Positives	Reference says something exists Model doesn't predict it False Negatives <i>Type II Error</i>	True Positive Rate/ Sensitivity/ Recall False Negative Rate/ Miss Rate
	Negative	Reference says something doesn't exist Model predicts it False Positives <i>Type I error</i>	Reference says something doesn't exist Model doesn't predict it True Negatives	False Positive Rate/ Fallout True Negative Rate/ Specificity
		Precision / Positive Predictive Value, False Discovery Rate	Negative Predictive Value, False Omission Rate	LR+, LR-

Evaluation Metric Insights: The Confusion Matrix

		Predictions		
		Positive	Negative	Calculate
References	Positive	Reference says something exists Model predicts it True Positives	Reference says something exists Model doesn't predict it False Negatives <i>Type II Error</i>	True Positive Rate/ Sensitivity/ Recall False Negative Rate/ Miss Rate
	Negative	Reference says something doesn't exist Model predicts it False Positives <i>Type I error</i>	Reference says something doesn't exist Model doesn't predict it True Negatives	False Positive Rate/ Fallout True Negative Rate/ Specificity
		Precision / Positive Predictive Value, False Discovery Rate	Negative Predictive Value, False Omission Rate	LR+, LR-


Evaluation Metric: Error trade-offs



You're pregnant

False Positive

(Type I error)



You're not pregnant

False Negative

(Type II Error)

Error trade-offs

Real World Example:



Error trade-offs

Real World Example:

- Project working with clinicians for mental health



Error trade-offs

Real World Example:

- Project working with clinicians for mental health
- Trying to detect **suicide risk**



Error trade-offs

Real World Example:

- Project working with clinicians for mental health
- Trying to detect **suicide risk**
- For patient trust (and sanity), important not to have **False Positives**



Error trade-offs

Real World Example:

- Project working with clinicians for mental health
- Trying to detect **suicide risk**
- For patient trust (and sanity), important not to have **False Positives**
 - Predicting suicide risk when there is not a risk



Error trade-offs

Real World Example:

- Project working with clinicians for mental health
- Trying to detect suicide risk
- For patient trust (and sanity), important not to have False Positives
 - Predicting suicide risk when there is not a risk
- Prioritize **True Positive Rate** at a low **False Positive Rate**



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

TOOL: EVALUATION METRICS

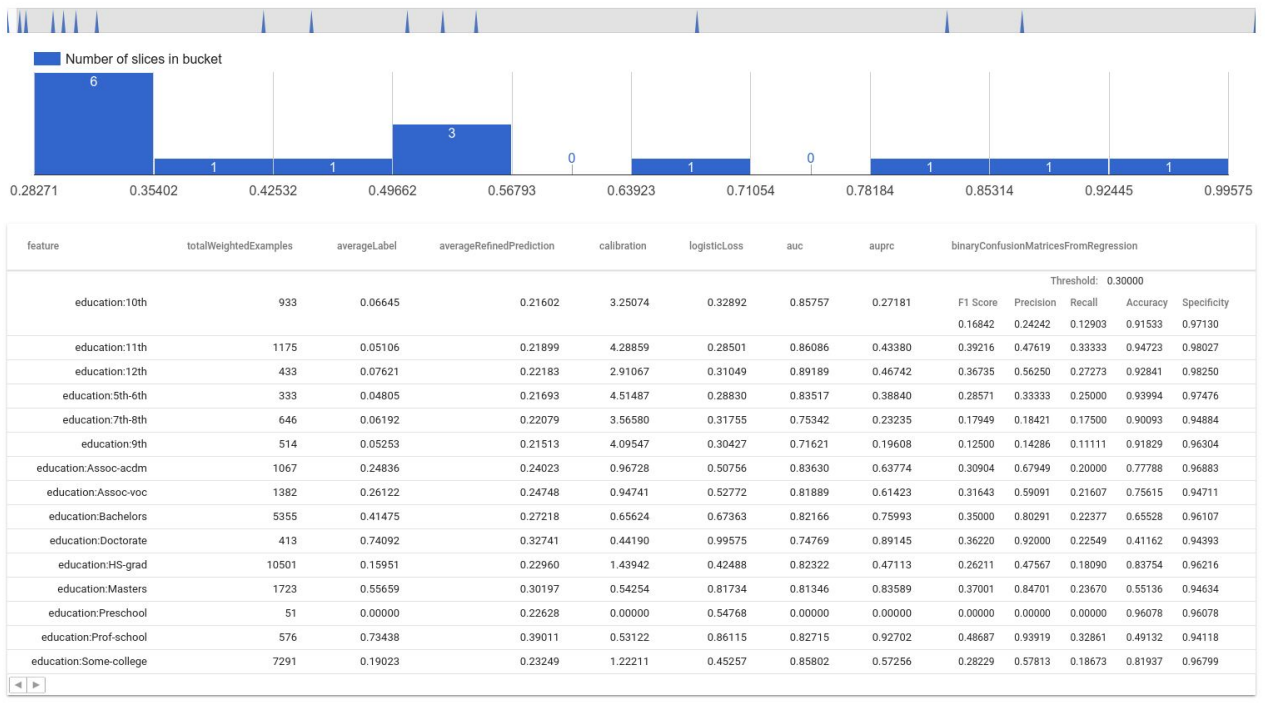
Lantern: Guided Model Analysis

Model Evaluation
+
Data slicing
=
Better Understanding of
Disproportionate
Outcomes



**Colab
Start**

go/lantern-eval-colab



INSIGHT: FEATURES

Word embeddings

Common ML Feature: Word Embeddings

Word embeddings represent each word as a vector.

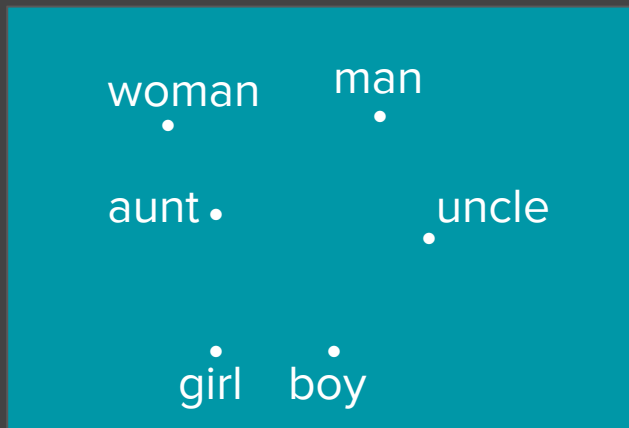


Common ML Feature: Word Embeddings

Word embeddings represent each word as a vector.



Allows us to calculate similarity between words.



Common ML Feature: Word Embeddings

Word embeddings represent each word as a vector.

Similarities between embeddings can be found using **cosine distance**:

$$\cos(\vec{\text{man}}, \vec{\text{woman}}) = \frac{\vec{\text{man}} \cdot \vec{\text{woman}}}{\|\vec{\text{man}}\| \cdot \|\vec{\text{woman}}\|}$$

Common ML Feature: Word Embeddings

Word embeddings represent each word as a vector.

Similarities between embeddings can be found using cosine distance.

Similarities between the **difference** between vectors can also be calculated using cosine distance.

$$\begin{array}{l} \rightarrow \quad \rightarrow \quad \rightarrow \\ \mathbf{g} = \text{man} - \text{woman} \\ \rightarrow \quad \rightarrow \quad \rightarrow \\ \mathbf{r} = \text{king} - \text{queen} \end{array} \quad \cos(\mathbf{g}, \mathbf{r}) = \frac{\begin{array}{l} \rightarrow \quad \rightarrow \\ \mathbf{g} \cdot \mathbf{r} \end{array}}{\begin{array}{l} \rightarrow \quad \rightarrow \\ \|\mathbf{g}\| \cdot \|\mathbf{r}\| \end{array}}$$

[Bolukbasi, Tolga; Chang, Kai-Wei; Zou, James; Saligrama, Venkatesh; Kalai, Adam \(2016\). "Man is to Computer Programmer as Woman is to Homemaker?: Debiasing Word Embeddings". Proceedings of NIPS.](#)

Common ML Feature: Word Embeddings

Word embeddings represent each word as a vector.

Similarities between embeddings can be found using cosine distance.

Similarities between the difference between vectors can also be calculated using cosine distance.

This can show us roughly equivalent relationships between words.

→ → → →
man - woman ≈ king - queen

Common ML Feature: Word Embeddings

Word embeddings represent each word as a vector.

Similarities between embeddings can be found using cosine distance.

Similarities between the difference between vectors can also be calculated using cosine distance.

This can show us roughly equivalent relationships between words ... including unfairness.

→ → → →
man - woman ≈ king - queen
→ → → →
man - woman ≈ computer programmer - homemaker

[Bolukbasi, Tolga; Chang, Kai-Wei; Zou, James; Saligrama, Venkatesh; Kalai, Adam \(2016\). "Man is to Computer Programmer as Woman is to Homemaker?: Debiasing Word Embeddings". Proceedings of NIPS.](#)

Potential Solution: **Debias** your embeddings

High-Level:

1. Calculate the representation of a concept, like “gender”, using word embeddings.
2. Subtract this representation from learned word embeddings.
3. Use a hyperparameter to define how much this subtraction effects the embedding.

[Link to Code](#)

[Bolukbasi, Tolga; Chang, Kai-Wei; Zou, James; Saligrama, Venkatesh; Kalai, Adam \(2016\). “Man is to Computer Programmer as Woman is to Homemaker?: Debiasing Word Embeddings”. Proceedings of NIPS.](#)

TECHNIQUE: EMBEDDINGS

Embeddings with Tensorflow

Embeddings
reveal words used in
similar contexts within
your dataset.



**Colab
Start**

[go/tf-embedding-colab](https://go.tf-embedding-colab)

Id	L2 Distance↑	L2 Norm	Adjust	Word
2000	0.000000	1.000000	Remove	<u>teacher</u>
1736	0.707160	1.000000	Add	<u>teachers</u>
44702	0.732374	1.000000	Add	<u>guidance counselor</u>
6229	0.740699	1.000000	Add	<u>elementary</u>
105512	0.791613	1.000000	Add	<u>paraprofessional</u>
371401	0.795801	1.000000	Add	<u>paraeducator</u>
13229	0.798513	1.000000	Add	<u>Teacher</u>
931	0.829719	1.000000	Add	<u>student</u>
198	0.833520	1.000000	Add	<u>school</u>
4825	0.837015	1.000000	Add	<u>classroom</u>

[https://g3doc.corp.google.com/
engedu/ml/mldays/g3doc/embeddings_demo.md](https://g3doc.corp.google.com/engedu/ml/mldays/g3doc/embeddings_demo.md)

Embeddings Demo

THE JOURNEY CONTINUES

Fairness-Relevant Tools

Google-internal

go/mlx <input type="checkbox"/>	Suite of tools useful for different aspects of fairness/bias. Some key tools also listed below.
go/tfx <input type="checkbox"/> Codelab	Computes statistics over data for visualization and example validation; anomaly detection; etc.
go/mlx tools <input type="checkbox"/>	Great list of tools to help visualize different aspects of your model.
go/mlx-lantern <input type="checkbox"/> Codelab	Computes evaluation metrics and loss for slices of your data with visualization. Interested in adding further support relevant to fairness in particular. Use with go/tfx or Sibyl .
go/ml-dash <input type="checkbox"/>	Compare metrics; visualize loss over time; etc.
go/wide-n-deep <input type="checkbox"/>	Combine the benefits of wide models and deep models (deep learning).
go/multitask <input type="checkbox"/>	Support multitask (multi-headed) learning. Predicting several tasks at once can be useful for the tasks to mutually benefit one another.
go/glassbox <input type="checkbox"/>	Interpretable machine learning.
go/bias	Report biased Google products.

Google-internal

Embedding Projector 📖	View how different strings of text pattern with other strings in a high-dimensional space.
go/mledu-in-embeddings 📖	View word relationships in embedding space.
Rank Lab ☐ Recipes & Best Practices	Supports feature ablation experiments, shuffling.
Fast Feature Ablation ☐	Fast Feature Ablation (FFA) adapts the feature ablation process cpop/jpg developed for SmartASS to an implementation suitable for Tensorflow and TF.Learn specifically.
Chain ☐☐ Codelab	Provides easy handling for moving from detection to evaluation. Includes a face attribute client: Age/Gender/UHS estimates (common in semantic scene understanding).
Affective Computing ☐☐	Label images for affective states, emotions, etc.
VSEval ☐☐ Codelab	Flexible infrastructure to acquire, store, and share high-quality ground truth, as well as by offering insightful statistics and visualization tools to support such research.
Learning Arbiter ☐☐ Codelab	The Arbiter Perception Eval system is in development! It aims to be a modular service oriented ecosystem built to ease up the evaluation of machine perception models.

Thanks!

dsculley@

mmitchellai@

[ML Fairness](#)

Machine Learning, Subgroup Discovery

[go/ml-fairness-tools](#)

[go/ml-fairness-metrics](#)

References

[Benton, Adrian; Mitchell, Margaret; Hovy, Dirk \(2017\). "Multi-task learning for Mental Health Conditions with Limited Social Media Data". Proceedings of EACL.](#)

[Bolukbasi, Tolga; Chang, Kai-Wei; Zou, James; Saligrama, Venkatesh; Kalai, Adam \(2016\). "Man is to Computer Programmer as Woman is to Homemaker?: Debiasing Word Embeddings". Proceedings of NIPS.](#)

[Gordon, Jonathan; Van Durme, Benjamin \(2013\). "Reporting Bias and Knowledge Acquisition". Proceedings of the 2013 workshop on Automated knowledge base construction.](#)

[Kay, Matthew; Matuszek, Cynthia; Munson, Sean A. \(2015\). "Unequal Representation and Gender Stereotypes in Image Search Results for Occupations". Proceedings of CHI.](#)

[Misra, Ishan; Girshick, Ross; Mitchell, Margaret; Zitnick, Larry \(2016\). "Seeing through the Human Reporting Bias: Visual Classifiers from Noisy Human-Centric Labels". Proceedings of CVPR.](#)

[Wapman, Mikaela; Belle, Deborah \(2014\). "A Riddle Reveals Depth of Gender Bias". Boston University. As reported by Barlow, Rich. BU Today. <https://www.bu.edu/today/2014/bu-research-riddle-reveals-the-depth-of-gender-bias/>](#)

KDD Tutorial: http://francescobonchi.com/algorithmic_bias_tutorial.html

THE JOURNEY CONTINUES

Additional Slides

INSIGHT: TASKS

**Leverage multiple tasks to
improve performance across
different subgroups**

[go/tf-multitask](https://go.tf-multitask)

Motivation from “The Karate Kid”

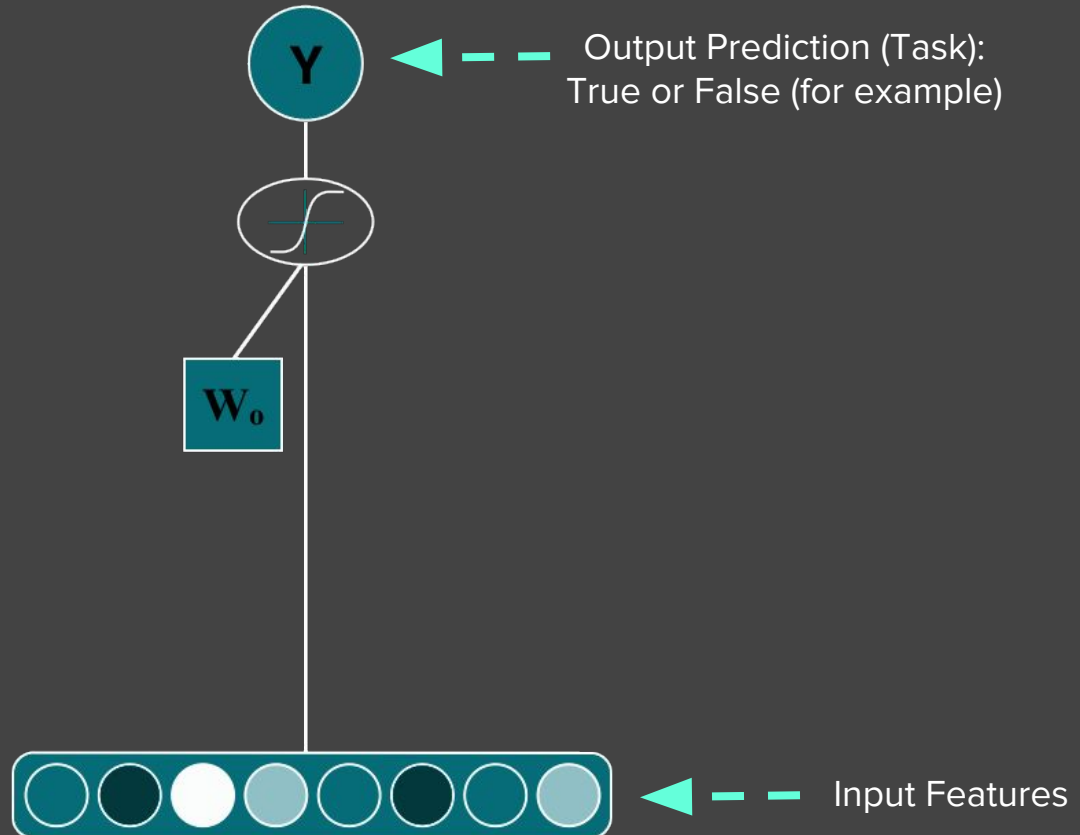


Single-task Learners
(STL)

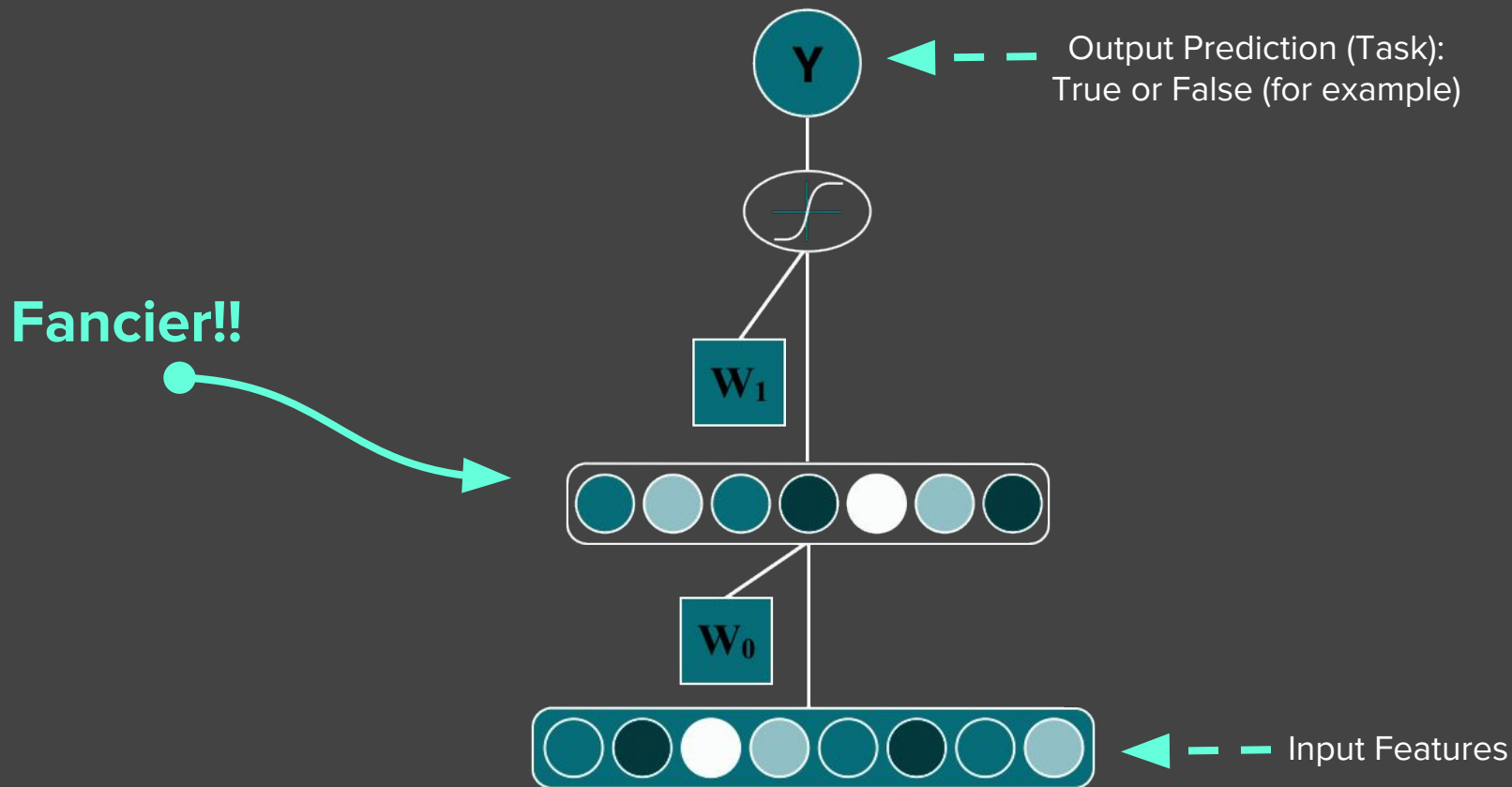


Multitask Learner
(MTL)

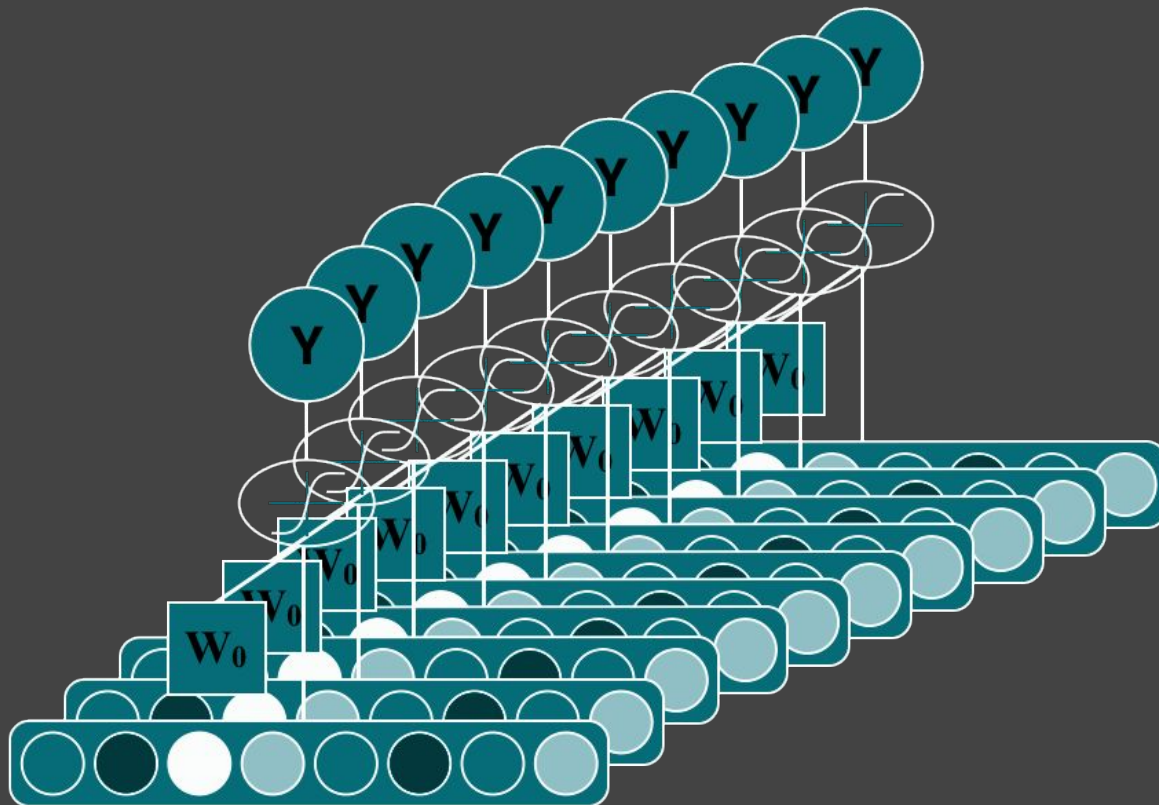
Single-Task: Logistic Regression



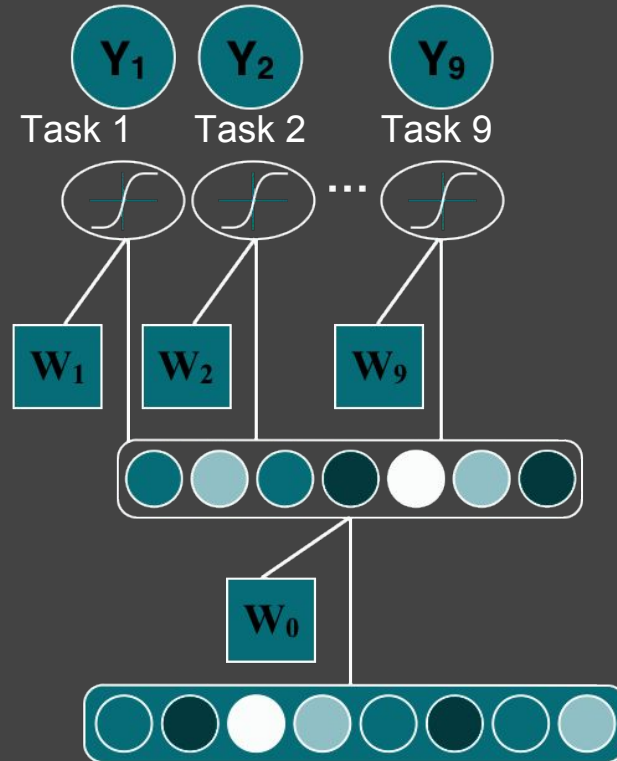
Single-Task: Deep Learning



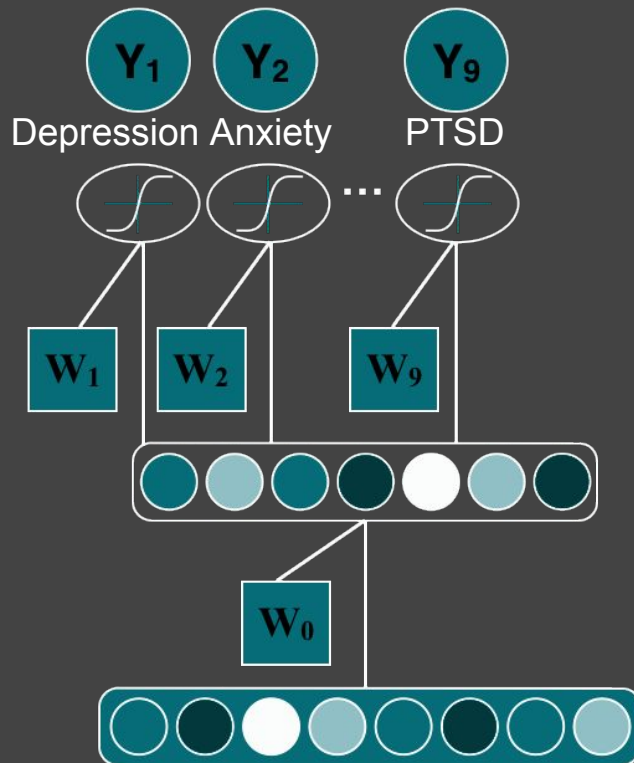
Multiple Tasks with Basic Logistic Regression



Multiple Tasks + Deep Learning: Multi-task Learning



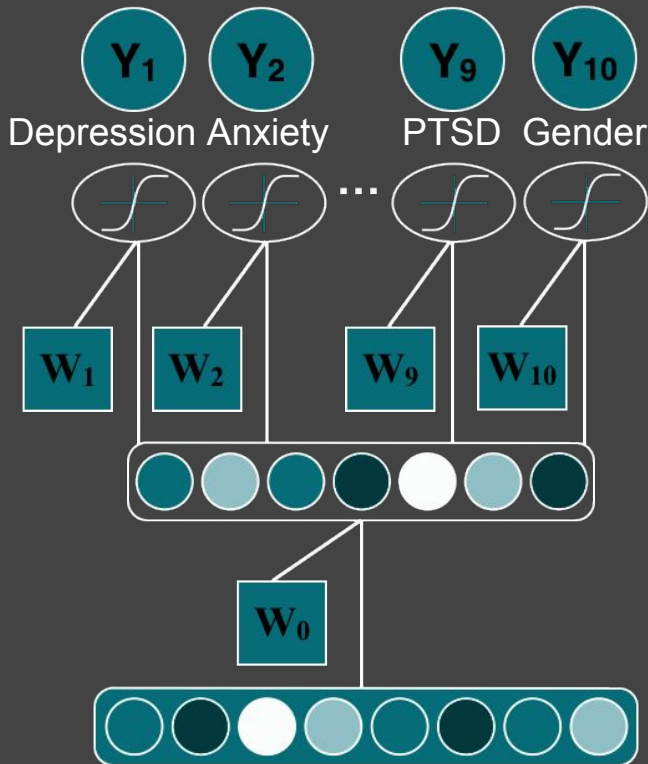
Multiple Tasks + Deep Learning: Multi-task Learning Example



Task	N
Neurotypicality	4791
Anxiety	2407
Depression	1400
Suicide attempt	1208
Eating disorder	749
Schizophrenia	349
Panic disorder	263
PTSD	248
Bipolar disorder	191
All	9611

} <5% positive examples

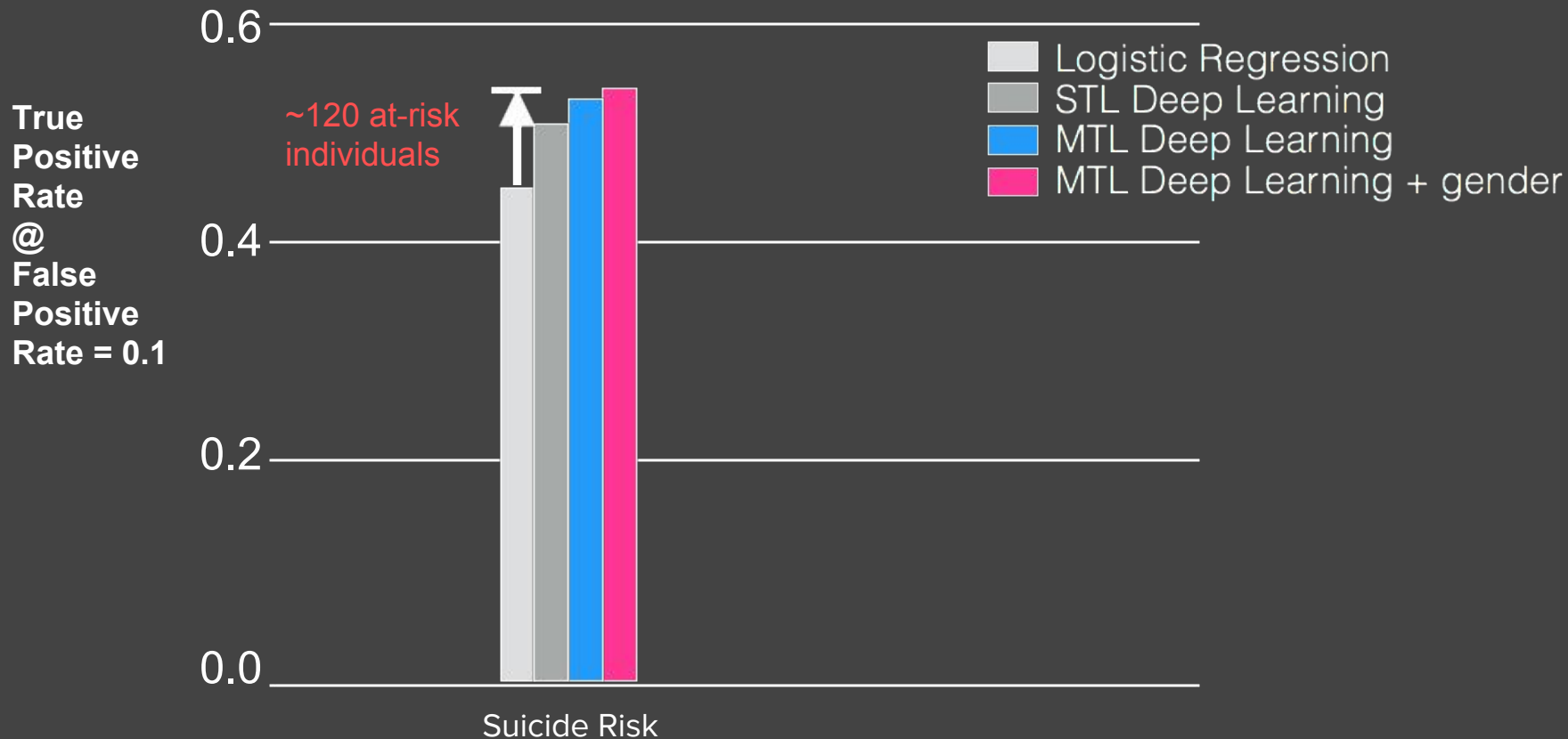
Multiple Tasks + Deep Learning: Multi-task Learning Example



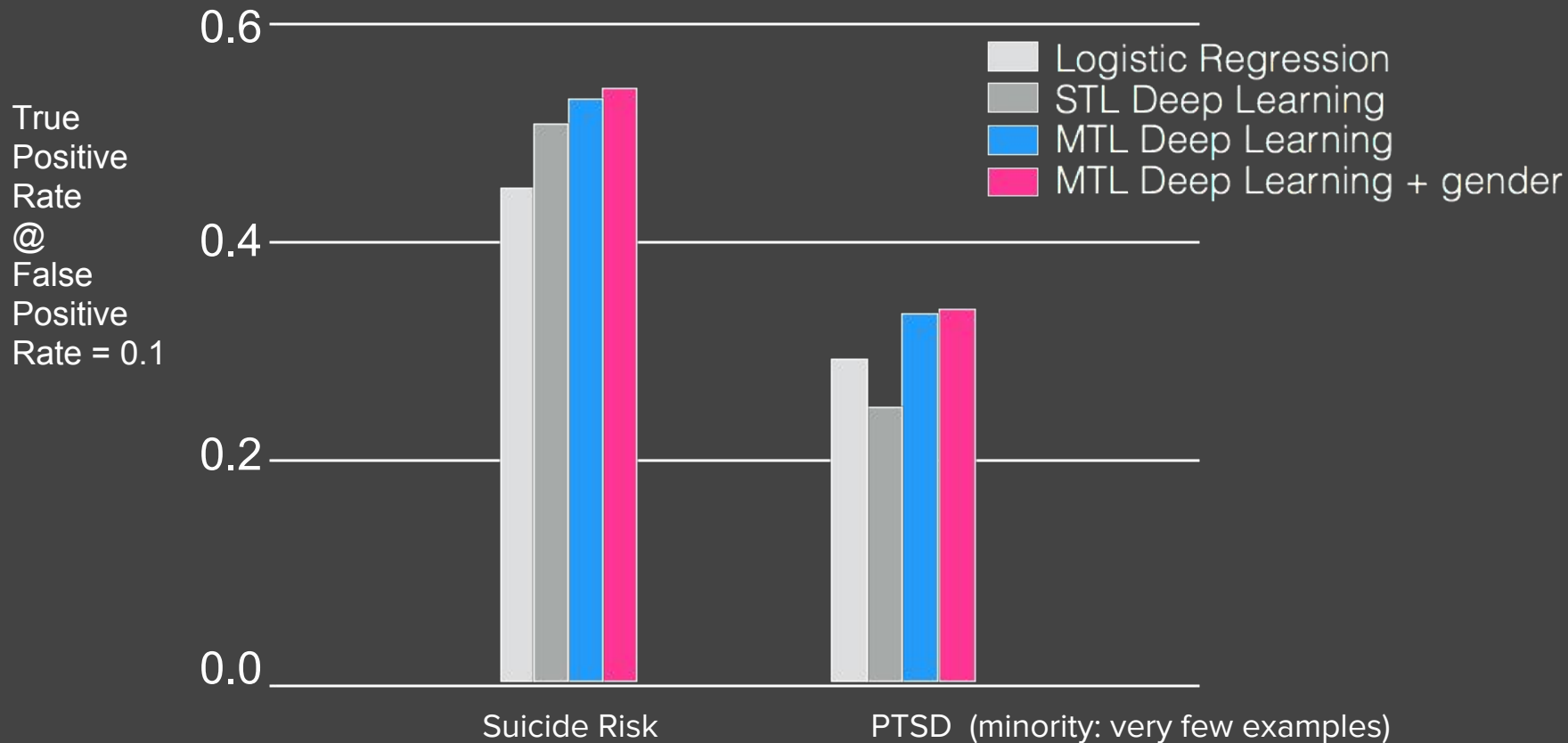
Task	N
Gender	1101
Neurotypicality	4791
Anxiety	2407
Depression	1400
Suicide attempt	1208
Eating disorder	749
Schizophrenia	349
Panic disorder	263
PTSD	248
Bipolar disorder	191
All	9611

} <5% positive examples

Improved Performance across Subgroups



Improved Performance across Subgroups



Lantern: Guided Model Analysis, including Multi-Task!

Includes offline model evaluations, computation of metrics on different slices of the data

feature	auc	auprc	averageLabel	averageRefinedPrediction	binaryConfusionMatricesFromRegression			
					Threshold: 0.75000			
age:19	0.66929	0.70809	0.55309	0.55243	F1 Score	Precision	Recall	Accuracy
					0.40000	0.10000	0.30000	0.20000
age:18	0.68247	0.73486	0.62338	0.46560	0.41000	0.11000	0.31000	0.21000
age:20	0.66872	0.70736	0.55309	0.56765	0.43000	0.13000	0.33000	0.23000
age:22	0.71525	0.77450	0.62338	0.46510	0.44000	0.14000	0.34000	0.24000

[go/mlx-lantern](https://go.mlx-lantern)

[Source Document for Multi-Task Models](#)

INSIGHT: OBJECTIVE FUNCTION

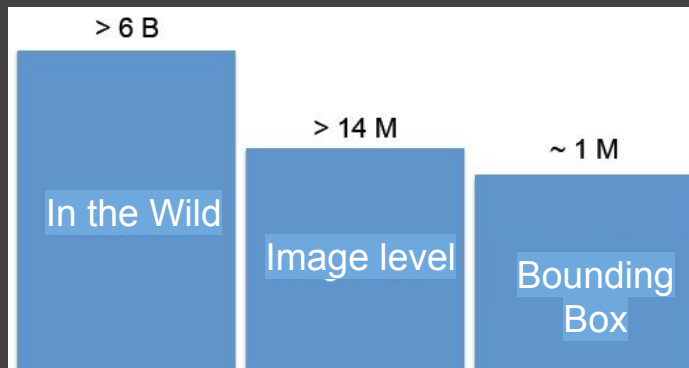
**Visual presence +
Relevance**

Data data everywhere ...

Facebook

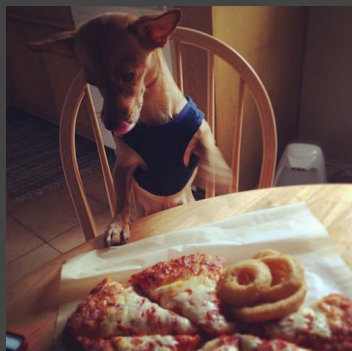
300 Million
images uploaded
everyday

flickr

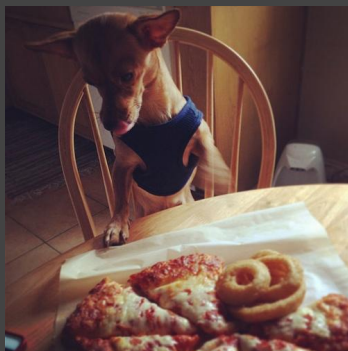


YouTube™

100 hours of video
every minute



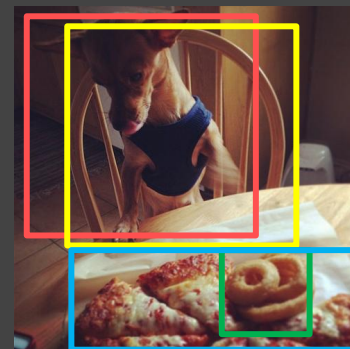
#dog #hungry



OMG Frodo is sitting
eating pizza and donuts.



dog, chair, pizza, donut



dog, chair, pizza, donut

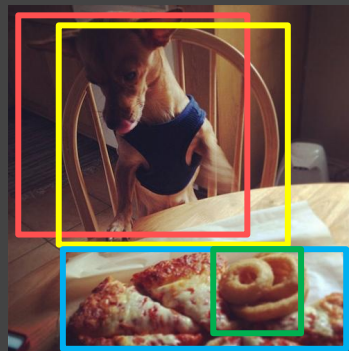
Data data everywhere ...

But not many labels to train

Exhaustively annotated data is expensive

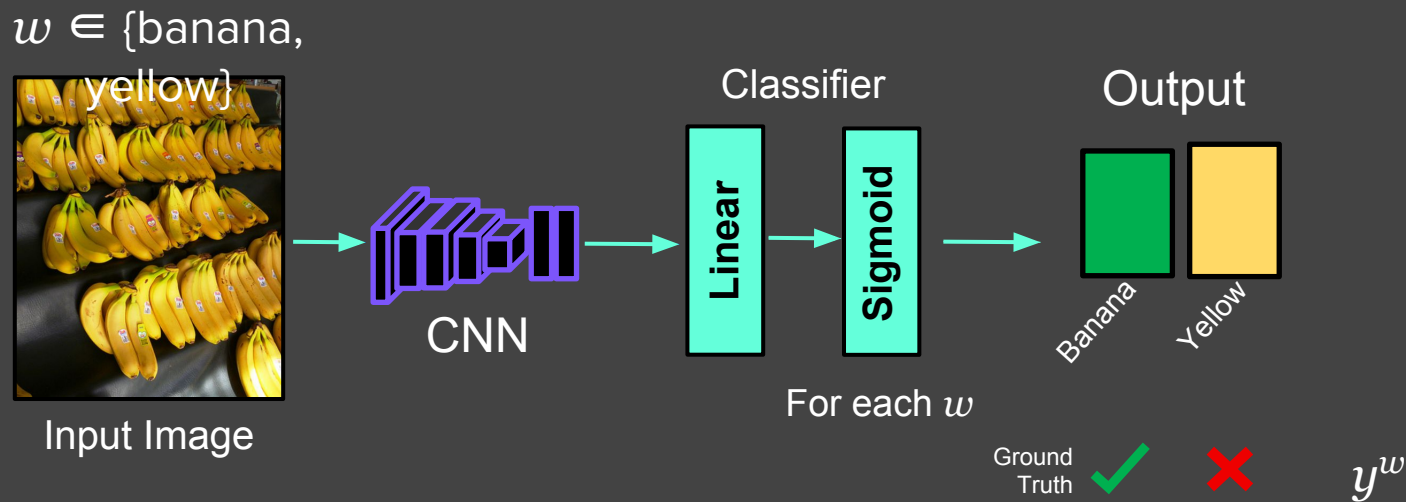


dog, chair, pizza, donut



*dog, chihuahua, brown, chair, table, wall,
space heater, pizza, greasy, donut 1, donut 2,
pizza slice 1, pizza slice 2...*

Simple Image Classification



“Gold standard” Annotation: Human-biased label $y^w \in \{0, 1\}$

Prediction $h^w(y^w|I)$

Factoring in Reporting Bias: Idea

- A human-biased prediction h can be factored into two terms

Factoring in Reporting Bias: Idea

- A human-biased prediction h can be factored into two terms
 - Visual presence v – *Is the concept **visually present**?*



$w \in \{\text{banana},$
 $\text{yellow}\}$ ✓ ✓

Factoring in Reporting Bias: Idea

- A human-biased prediction h can be factored into two terms
 - Visual presence v – *Is the concept **visually present**?*
 - Relevance r – *Is the concept **relevant** for a human?*



$w \in \{\text{banana},$
 $\text{yellow}\}$ ❌

Factoring in Reporting Bias: Idea

- A human-biased prediction h can be factored into two terms
 - Visual presence v – *Is the concept **visually present**?*
 - Relevance r – *Is the concept **relevant** for a human?*

$$h = f(r, v)$$



Factoring in Reporting Bias: Idea

- A human-biased prediction h can be factored into two terms
 - Visual presence v – *Is the concept **visually present**?*
 - Relevance r – *Is the concept **relevant** for a human?*

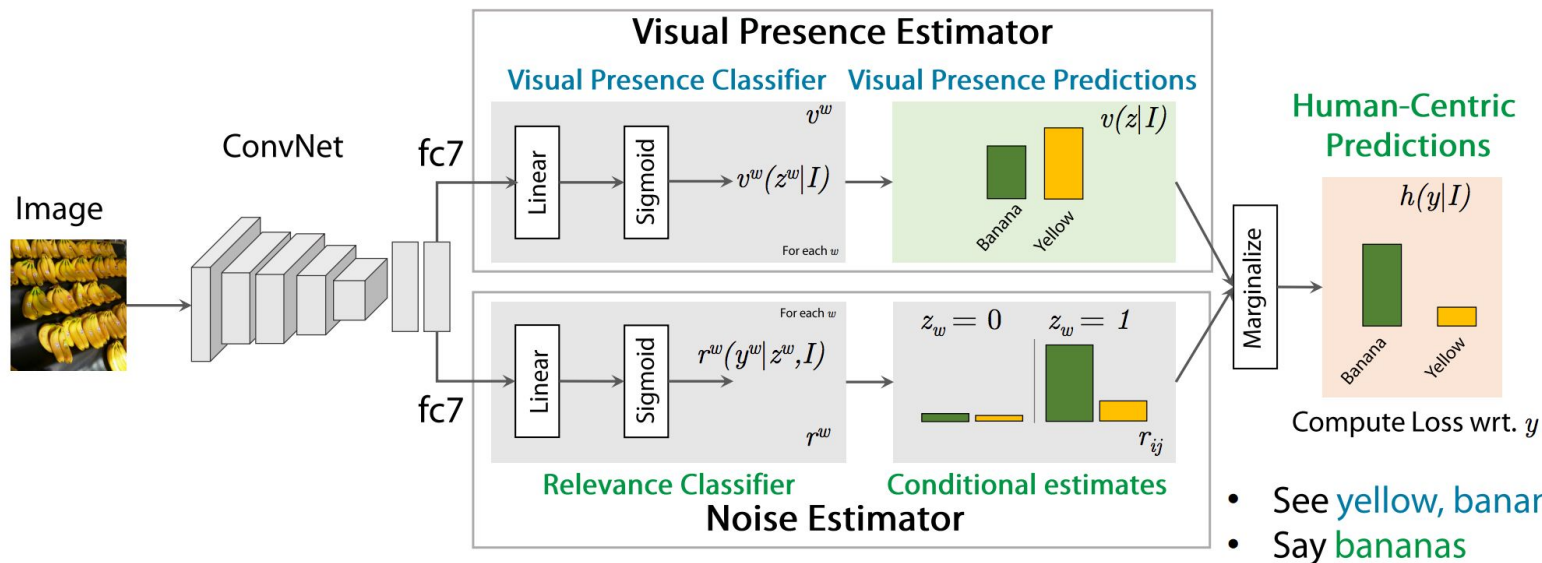


Given visual presence, is concept **relevant**? Is concept **present**?

$$h(y|I) = \sum_{j \in \{0,1\}} r(y|z = j, I) v(z = j|I)$$

	Label	Prediction
Visually correct ground truth (Unknown)	z	v
Available ground truth (human-biased)	y	h

End-to-End Approach



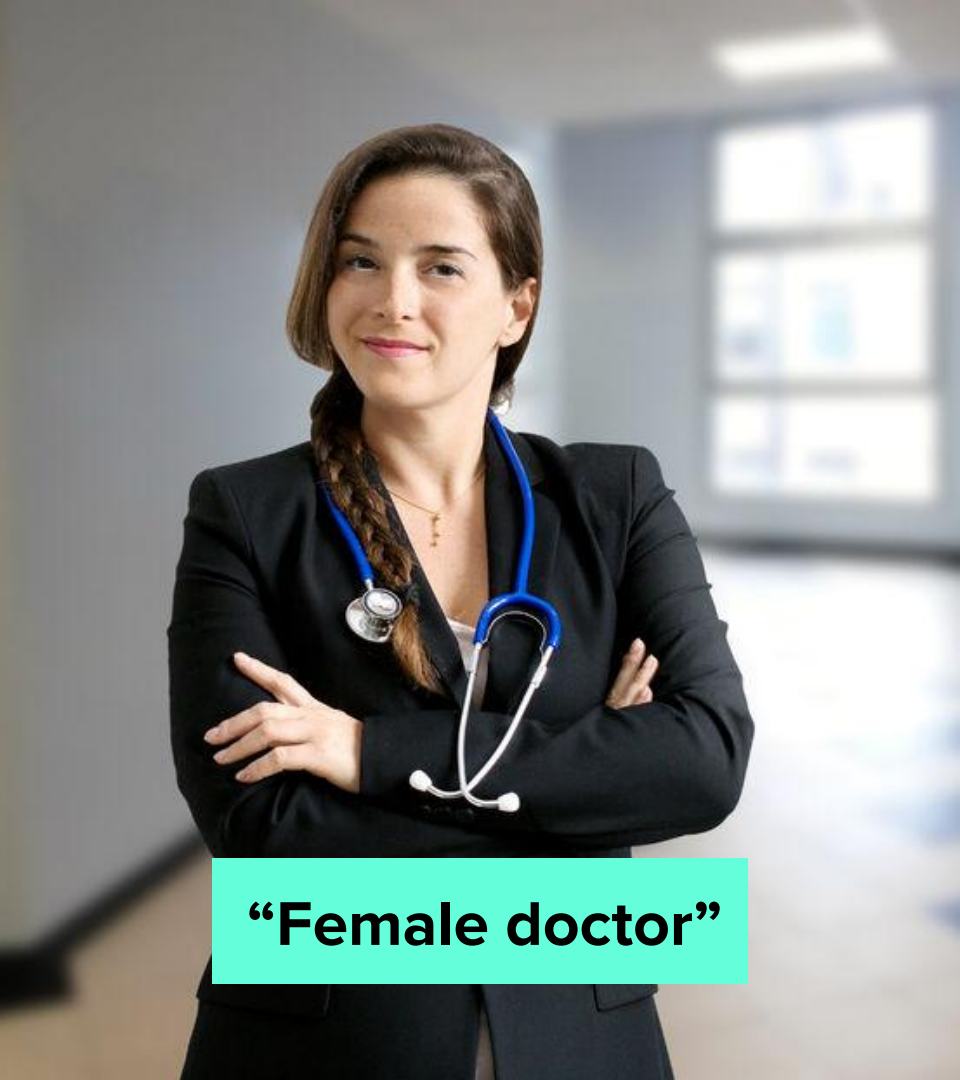
$$\text{Marginalize: } h(y|I) = \sum_{j \in \{0,1\}} r(y|z=j, I) v(z=j|I)$$

SOURCE

Misra, Ishan; Girshick, Ross; Mitchell, Margaret; Zitnick, Larry (2016). "Seeing through the Human Reporting Bias: Visual Classifiers from Noisy Human-Centric Labels". Proceedings of CVPR.



“Male doctor”



“Female doctor”

Thanks!

mmitchellai@

[ML Fairness](#)

Machine Learning, Subgroup Discovery

[go/ml-fairness-tools](#)

[go/ml-fairness-metrics](#)

