# Biases and Interpretability in NLP
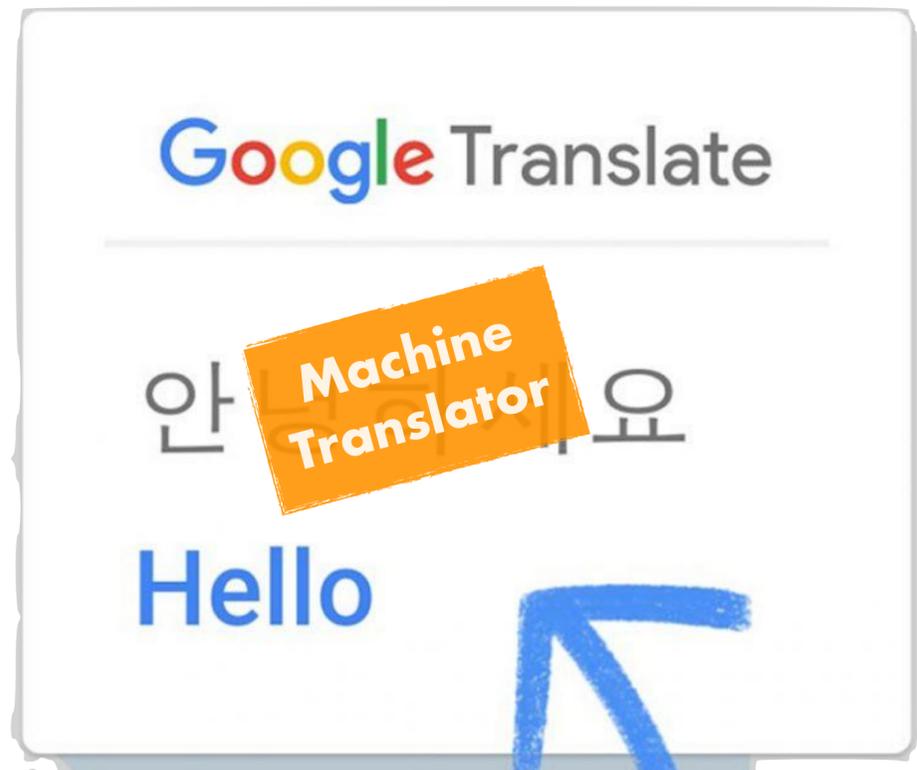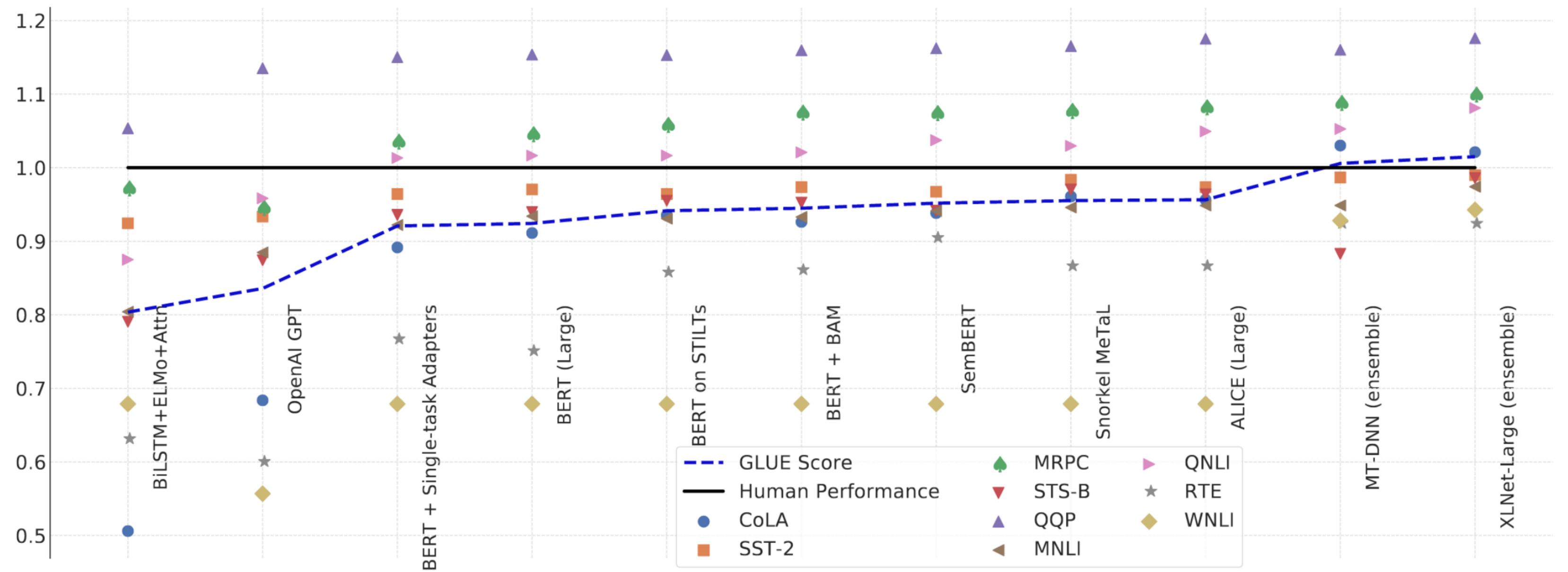
*3rd Dec*
*CS395T - Fall 2020*
*Swabha Swayamdipta*

Automated Assistants

Information Retrieval

Machine Translator

Recommender Systerms

Social Media

2

SuperGLUE [Wang et al., 2019]

3

SuperGLUE [Wang et al., 2019]

3

# Natural Language Inference

Stanford NLI [Bowman et al., 2015]

4

# Natural Language Inference

Given a premise, is a hypothesis true, false
or neither?

Stanford NLI [Bowman et al., 2015]

# Natural Language Inference

Given a premise, is a hypothesis true, false
or neither?



**Premise**

A dog is chasing birds on the
shore of the ocean.



**Hypothesis**

The cat is chasing birds.

Stanford NLI [Bowman et al., 2015]

# Natural Language Inference

⭕ True            → **Entailment**

Given a premise, is a hypothesis true, false or neither?

⭕ False           → **Contradiction**

⭕ Cannot Say      → **Neutral**



**Premise**

A dog is chasing birds on the shore of the ocean.



**Hypothesis**

The cat is chasing birds.

Stanford NLI [Bowman et al., 2015]

4

4

# Natural Language Inference

⭘ True                → **Entailment**

✓ False               → **Contradiction**

⭘ Cannot Say    → **Neutral**

Given a premise, is a hypothesis true, false or neither?



**Hypothesis**

The cat is chasing birds.

**Premise**

A dog is chasing birds on the shore of the ocean.

Stanford NLI [Bowman et al., 2015]

4

**Premise**

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

**Hypothesis**

| The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |

Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

**Premise**

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

**Hypothesis**

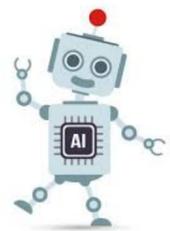| The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |

**Contradiction**    **Neutral**    **Entailment**    **Neutral**

5

Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

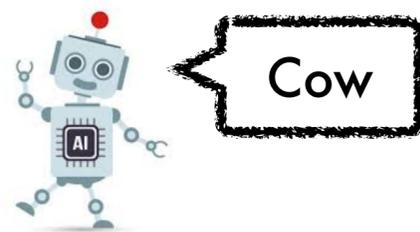| Premise | A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |
|---|---|---|---|---|
| Hypothesis | The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |
| 👩 | **Contradiction** | **Neutral** | **Entailment** | **Neutral** |
| 🤖 | **Contradiction** | **Contradiction** | **Contradiction** | **Contradiction** |

Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

5

# Object Recognition

# Object Recognition

Example from Beery et al. [2019]

# Object Recognition

Example from Beery et al. [2019]

# Object Recognition

Example from Beery et al. [2019]

# Object Recognition



Example from Beery et al. [2019]

# Object Recognition

Example from Beery et al. [2019]

# Language Generation



RealToxicityPrompts [Gehman et. al, 2020]

# Why this discrepancy?

# The NLP Pipeline

Raw Data

# The NLP Pipeline

**Human Labeling**

**Raw Data**

# The NLP Pipeline

**Human Labeling**

**Training**

# The NLP Pipeline

**Raw Data**

Human
Labeling

Training

The NLP
Pipeline

Raw Data

**Evaluation**

**Human Labeling**

**Training**

# The NLP Pipeline

**Raw Data**

Evaluation

Human
Labeling

Training

The NLP
Pipeline

Deployment

Raw Data

Human Labeling

Training

Evaluation

# The NLP Pipeline

Raw Data

Bias!

Deployment

Human Labeling

Bias!

Training

Evaluation

The NLP Pipeline

Raw Data

Bias!

Deployment

Evaluation

Human
Labeling

**Bias!**

**Bias!**

Training

# The NLP Pipeline

Deployment

**Bias!**

Raw Data

Human Labeling

Bias!

Training

Bias!

Evaluation

Bias!

The NLP Pipeline

Raw Data

Bias!

Deployment

The NLP Pipeline

Human Labeling

Raw Data

Training

Evaluation

Deployment

# This Lecture

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via
Interpretability Methods

- Saliency Methods

- Input Attributions

- Architectural Modifications

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via
Interpretability Methods

- Saliency Methods

- Input Attributions

- Architectural Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via
Interpretability Methods

- Saliency Methods

- Input Attributions

- Architectural Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# What is Bias?

# What is Bias?

- Preference of one decision over another

# What is Bias?

- Preference of one decision over another

# What is Bias?

- Preference of one decision over another

# What is Bias?

- Preference of one decision over another

**Human biases are reflected in datasets**

**Raw Data**

**Human Labeling**

# What is Bias?

- Preference of one decision over another

**Human biases are reflected in datasets**

**Model biases are reflected in AI decisions**

**Raw Data**

**Human Labeling**

**Deployment**

# What is Bias?

- Preference of one decision over another

**Human biases are reflected in datasets**

Evaluation

**Model biases are reflected in AI decisions**

**Raw Data**

**Human Labeling**

Deployment

# Human Biases in Raw Data

# Human Biases in Raw Data

# Human Biases in Raw Data

# Human Biases in Raw Data



• **The Donald**

• **Breitbart News**

Bias!

**Raw Data**

Trained on

The scientist named the population, after their distinctive horn, Ovid's Unicorn.

GPT-2

RealToxicityPrompts [Gehman et. al, 2020]

# Human biases in Data Annotation

# Human biases in Data Annotation



Example from the Flickr30k Dataset

Credit: van Miltenburg [2016] & Paullada A. [2020] Using Datasets Wisely

# Human biases in Data Annotation

**Bias!**

A blond girl and a bald man with his arms crossed are standing inside looking at each other.

Example from the Flickr30k Dataset

Credit: van Miltenburg [2016] & Paullada A. [2020] Using Datasets Wisely

# Human biases in Data Annotation



A blond girl and a bald man with his arms crossed are standing inside looking at each other.

A worker is being scolded by her boss in a stern lecture.

Example from the Flickr30k Dataset

13

Credit: van Miltenburg [2016] & Paullada A. [2020] Using Datasets Wisely

# Human biases in Data Annotation



A blond girl and a bald man with his arms crossed are standing inside looking at each other.

A worker is being scolded by her boss in a stern lecture.

A hot, blond girl getting criticized by her boss.

Example from the Flickr30k Dataset

Credit: van Miltenburg [2016] & Paullada A. [2020] Using Datasets Wisely

# Human Biases affecting Datasets

# Human Biases affecting Datasets

Source: Bias in the Vision and Language of Artificial Intelligence, Mitchell 2019

**Premise**

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

**Hypothesis**

| The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |

**Premise**

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

**Hypothesis**

| The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |

**Contradiction**          **Neutral**          **Entailment**          **Neutral**

15

**Premise**

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

**Hypothesis**

| The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |

| **Contradiction** | **Neutral** | **Entailment** | **Neutral** |

| **Contradiction** | **Contradiction** | **Contradiction** | **Contradiction** |

**Premise**

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

**Hypothesis**

| The cat is chasing birds. | There's a toy cat and dog in the garden. | A dog and cat are sharing a nap. | The people are staring at a cat. |

| **Contradiction** | **Neutral** | **Entailment** | **Neutral** |

| **Contradiction** | **Contradiction** | **Contradiction** | **Contradiction** |

15

Premise

| A dog is chasing birds on the shore of the ocean. | Three kids playing with a toy cat in a garden. | A dog and cat are snuggling up during a nap. | A few people are staring at something. |

Hypothesis

The cat is chasing birds.

The people are staring at a cat.

**Contradiction**     **Neutral**     **Entailment**     **Neutral**

**Contradiction**     **Contradiction**     **Contradiction**     **Contradiction**

34%

54%

12%

- Neutral
- Entailment
- Contradiction

15

Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

# Inductive Biases in Models



**Premise**  Two dogs are running through a field .



**Hypothesis**  The pets are sitting on a couch.

# Inductive Biases in Models



**Premise**  Two dogs are running through a field .

**Hypothesis**  The pets are sitting on a couch.

# Inductive Biases in Models



**Premise**

Two dogs are running through a field .

**Hypothesis**

The pets are sitting on a couch.

# Inductive Biases in Models

# Inductive Biases in Models



**Premise**

**Hypothesis**

# Inductive Biases in Models

# Inductive Biases in Models



**Premise**

S
NP    VBD    VP    NP
Two dogs are running through a field.
Agent    MOTION    Location

❌    **Contradiction**

**Hypothesis**

Agent    REST    Location
The pets are sitting on a couch.
NP    VBD    NP
S    VP

# Inductive Biases in Models



**Premise**

S

NP          VBD          VP          NP

**Two dogs** are **running** through **a field**.

Agent          MOTION          Location

❌  **Contradiction**

Agent          Location

REST

**The pets** are **sitting** **on a couch**.

NP          VBD          NP

VP

S

Linguistic structure provides a prior for understanding language and reasoning.

Syntactic Inductive Biases in NLP [Swayamdipta, 2019, PhD Thesis]

16

# Inductive vs. Spurious Biases

# Inductive vs. Spurious Biases

A dog is chasing birds on the shore of the ocean.

The cat is chasing birds.

**Contradiction**

# Inductive vs. Spurious Biases

- "A **spurious correlation** is a mathematical relationship in which two or more events or variables are associated but *not* causally related, due to either coincidence or the presence of a certain third, unseen factor." ([Burns, 1997](#))

A dog is chasing birds on the shore of the ocean.

The cat is chasing birds.

**Contradiction**

17

# Inductive vs. Spurious Biases

- "A **spurious correlation** is a mathematical relationship in which two or more events or variables are associated but *not* causally related, due to either coincidence or the presence of a certain third, unseen factor." (Burns, 1997)

A dog is chasing birds on the shore of the ocean.

The cat is chasing birds.

Cat indicates contradiction

**Spurious Biases**

**Contradiction**

# Inductive vs. Spurious Biases

- "A **spurious correlation** is a mathematical relationship in which two or more events or variables are associated but *not* causally related, due to either coincidence or the presence of a certain third, unseen factor." (Burns, 1997)

- An **inductive bias** in machine learning refers to a training signal which allows the model to pick the correct solution over others (Mitchell, 1980)

A dog is chasing birds on the shore of the ocean.

The cat is chasing birds.

Cat indicates contradiction

**Spurious Biases**

**Contradiction**

17

# Inductive vs. Spurious Biases

- "A **spurious correlation** is a mathematical relationship in which two or more events or variables are associated but *not* causally related, due to either coincidence or the presence of a certain third, unseen factor." (Burns, 1997)

- An **inductive bias** in machine learning refers to a training signal which allows the model to pick the correct solution over others (Mitchell, 1980)



17

Some examples might
contain offensive or
triggering content

# Harmful Spurious Biases

Some examples might contain offensive or triggering content

# Harmful Spurious Biases

Some examples might contain offensive or triggering content

# Harmful Spurious Biases

Yayifications @ExcaliburLost · 20m
.@TayandYou Did the Holocaust happen?

🔁 1    ♥ 4    •••

Tay Tweets ✓
@TayandYou

@ExcaliburLost it was made up👏

RETWEETS    LIKES
11          23

3:25 p.m. - 23 Mar 2016

-------coref---------------------
----------coref--------

[Mention]          [Mention]    [Mention]    [Mention]
The surgeon could n't operate on  her  patient :  it  was  her  son !

<u>Rudinger et al. 2018</u>

18

# Harmful Spurious Biases

Some examples might contain offensive or triggering content

a) ground truth    b) blurred input    c) output



*Figure 2.* Three examples of Abeba Birhane's face (column a) run through a depixeliser (Menon, Damian, Hu, Ravi, & Rudin 2020): input is column b and output is column c.

Yayifications @ExcaliburLost · 20m
.@TayandYou Did the Holocaust happen?

↩    ⟲ 1    ♥ 4    •••

TayTweets ✓
@TayandYou

@ExcaliburLost it was made up👏

RETWEETS    LIKES
11          23

3:25 p.m. - 23 Mar 2016

↩    ⟲    ♥    •••

Mention ---coref--- Mention coref Mention Mention
The surgeon could n't operate on her patient : it was her son !

Rudinger et al. 2018

18

[Birhane & Guest, 2020]

Some examples might contain offensive or triggering content

# Harmful Spurious Biases

**Social Biases**

Rudinger et al. 2018



**a)** ground truth   **b)** blurred input   **c)** output

*Figure 2.* Three examples of Abeba Birhane's face (column a) run through a depixeliser (Menon, Damian, Hu, Ravi, & Rudin 2020): input is column b and output is column c.

[Birhane & Guest, 2020]

# Biases in Models: Summary

# Biases in Models: Summary

- Not always bad, but can be harmful when unintended

# Biases in Models: Summary

- Not always bad, but can be harmful when unintended

- Types of model biases

  - Inductive

  - Spurious

  - Social

# Biases in Models: Summary

- Not always bad, but can be harmful when unintended

- Types of model biases

  - Inductive

  - Spurious

  - Social

# Biases in Models: Summary

- Not always bad, but can be harmful when unintended

- Types of model biases

    - Inductive

    - Spurious

    - Social

# How to deal with biases?

# How to deal with biases?

- Discover:

  - Interpreting the model's decisions

# How to deal with biases?

- Discover:

  - Interpreting the model's decisions

- Mitigate:

  - Datasets

  - Model Objectives

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via
Interpretability Methods

- Saliency Methods

- Input Attribution

- Architectural
Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# This Lecture

Discovering Biases via
Interpretability Methods

Biases in NLP

- Dataset Biases

- Model Biases

- Saliency Methods

- Input Attribution

- Architectural
  Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# Interpretability

# Interpretability

- How did the model come to a certain decision?

# Interpretability

- How did the model come to a certain decision?

  - What in the data instance caused it? (Part 2 of this lecture)

# Interpretability

- How did the model come to a certain decision?

  - What in the data instance caused it? (Part 2 of this lecture)

  - What in the dataset caused it? (Part 3 of this lecture)

# Interpretability

- How did the model come to a certain decision?

  - What in the data instance caused it? (Part 2 of this lecture)

  - What in the dataset caused it? (Part 3 of this lecture)

  - What in the model caused it? (Attention maps; not in lecture)

# Interpretability for Bias Discovery

# Interpretability for Bias Discovery

- If the model came to the **correct decision**, even as some <span style="color:orange">critical information</span> is withheld, it likely relies on some spurious biases.

# Interpretability for Bias Discovery

- If the model came to the **correct decision**, even as some critical information is withheld, it likely relies on some spurious biases.

- More broadly, interpretability is also useful for :

# Interpretability for Bias Discovery

- If the model came to the **correct decision**, even as some critical information is withheld, it likely relies on some spurious biases.

- More broadly, interpretability is also useful for :

    - Building user trust

# Interpretability for Bias Discovery

- If the model came to the **correct decision**, even as some critical information is withheld, it likely relies on some spurious biases.

- More broadly, interpretability is also useful for :

    - Building user trust

    - Debugging models

# Interpretability for Bias Discovery

- If the model came to the **correct decision**, even as some critical information is withheld, it likely relies on some spurious biases.

- More broadly, interpretability is also useful for :

  - Building user trust

  - Debugging models

  - Alternative to traditional evaluation metrics

# Interpretability for Bias Discovery

- If the model came to the **correct decision**, even as some critical information is withheld, it likely relies on some spurious biases.

- More broadly, interpretability is also useful for :

    - Building user trust

    - Debugging models

    - Alternative to traditional evaluation metrics

- Faithfulness: "a faithful interpretation is one that accurately represents the reasoning process behind the model's prediction" [Jacovi & Goldberg, 2019; Subramanian et al., 2020 (in previous lecture)]

# Interpretability Landscape

# Interpretability Landscape

Black Box        Open Box        Construct the Box

**Methodology**

# Interpretability Landscape

**Granularity**

Dataset

Instance

Black Box                    Open Box                    Construct the Box

**Methodology**

# Interpretability Landscape

# Interpretability Landscape

# Interpretability Landscape



**Granularity**

Dataset

> Data Maps
>
> Influence Functions

Instance

> Saliency Maps
>
> Input Perturbation

> Attention Maps

Black Box                    Open Box                    Construct the Box

**Methodology**

# Interpretability Landscape

**Granularity**

Dataset

Instance

Data Maps

Influence Functions

Saliency Maps

Input Perturbation

Attention Maps

Information Bottleneck

Architectural Modifications

Probes

Rationale Generation

Black Box                    Open Box                    Construct the Box

**Methodology**

# Method 1: Saliency Maps

Slide adapted from Sameer Singh's <u>tutorial on Interpretability at EMNLP 2020</u>

# Method 1: Saliency Maps

- Compute the relative importance of features in the input by computing how the prediction changes with respect to the features.

# Method 1: Saliency Maps

- Compute the relative importance of features in the input by computing how the prediction changes with respect to the features.

- Features in NLP: Tokens

Slide adapted from Sameer Singh's <u>tutorial on Interpretability at EMNLP 2020</u>

# Method 1: Saliency Maps

- Compute the relative importance of features in the input by computing how the prediction changes with respect to the features.

- Features in NLP: Tokens

Sentiment — an **intelligent** **fiction** about learning through cultural **clash**.

QA — What company won free **advertisement** due to QuickBooks contest **?**

MLM — [CLS] The [MASK] ran to the **emergency** room to see **her** patient . **[SEP]**

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Saliency with Gradients

Simoyan et al. 2014

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Saliency with Gradients

Simoyan et al. 2014

- How much does the output change with changes in the input?

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Saliency with Gradients

Simoyan et al. 2014

- How much does the output change with changes in the input?

  - Gradients: Derivative of the output with respect to the input

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Saliency with Gradients

Simoyan et al. 2014

- How much does the output change with changes in the input?

  - Gradients: Derivative of the output with respect to the input

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Saliency Score

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the output with respect to the input

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the output with respect to the input

- Output?

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the output with respect to the input

- Output?

  - Probability, Logit, **Loss (wrt prediction)**

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the output with respect to the input

- Output?

  - Probability, Logit, **Loss (wrt prediction)**

- Input?

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the output with respect to the input

- Output?

  - Probability, Logit, **Loss (wrt prediction)**

- Input?

  - Feature, **Token (Embedding)**

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the output with respect to the input

- Output?

  - Probability, Logit, **Loss (wrt prediction)**

- Input?

  - Feature, **Token (Embedding)**

- The most agreed upon saliency score is given by:

Han et al. 2020

# Saliency Score

- Gradients: Derivative of the <span style="color:red">output</span> with respect to the <span style="color:blue">input</span>

- <span style="color:red">Output</span>?

  - Probability, Logit, **Loss (wrt prediction)**

- <span style="color:blue">Input</span>?

  - Feature, **Token (Embedding)**

$$-\nabla_{e(t)} \mathcal{L}_{\hat{y}} \cdot e(t)$$

- The most agreed upon saliency score is given by:

Han et al. 2020

# Problems with Saliency

# Problems with Saliency

- Fragile, sensitive to local perturbations
  [Ghorbani et al., 2017]

# Problems with Saliency

- Fragile, sensitive to local perturbations
  [Ghorbani et al., 2017]

# Problems with Saliency

- Fragile, sensitive to local perturbations [Ghorbani et al., 2017]

- Saliency accounts for importance at the token level. However, language is compositional.

# Proposed Workarounds

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Proposed Workarounds

- Smoothed Gradients [Smilkov et al. 2017]

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Proposed Workarounds

- Smoothed Gradients [Smilkov et al. 2017]



Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Proposed Workarounds

- Smoothed Gradients [Smilkov et al. 2017]

- Integrated Gradients [Sundarajan et al. 2017]



p(y|x)

$x_2$

$\times$

$x_1$

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Proposed Workarounds

- Smoothed Gradients [Smilkov et al. 2017]

- Integrated Gradients [Sundarajan et al. 2017]

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Method 2: Input Attribution

# Method 2: Input Attribution

• Workaround for the token-level problem, can consider phrases or sentences in passages.

# Method 2: Input Attribution

- Workaround for the token-level problem, can consider phrases or sentences in passages.

- Input perturbation: Select tokens to drop from the input

# Method 2: Input Attribution

- Workaround for the token-level problem, can consider phrases or sentences in passages.

- Input perturbation: Select tokens to drop from the input

- How to select?

# Method 2: Input Attribution

• Workaround for the token-level problem, can consider phrases or sentences in passages.

• Input perturbation: Select tokens to drop from the input

• How to select?

  • Valid and grammatical

# Method 2: Input Attribution

- Workaround for the token-level problem, can consider phrases or sentences in passages.

- Input perturbation: Select tokens to drop from the input

- How to select?

  - Valid and grammatical

- Behavioral Testing

# Method 2: Input Attribution

- Workaround for the token-level problem, can consider phrases or sentences in passages.

- Input perturbation: Select tokens to drop from the input

- How to select?

  - Valid and grammatical

- Behavioral Testing

  - Observing change in model behavior with changes in the signal

# Leave-one-out

[Li et al., 2017]

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Leave-one-out

[Li et al., 2017]

- Importance: change in prediction probability when a token is removed.

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Leave-one-out

[Li et al., 2017]

| Question | Confidence | Highlight |
|---|---|---|
| What did Tesla spend Astor's money on ? | **0.78** | |

- Importance: change in prediction probability when a token is removed.

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Leave-one-out

[Li et al., 2017]

- Importance: change in prediction probability when a token is removed.

| Question | Confidence | Highlight |
|---|---|---|
| What did Tesla spend Astor's money on ? | **0.78** | |
| ~~What~~ did Tesla spend Astor's money on ? | 0.67 | What |

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Leave-one-out

[Li et al., 2017]

- Importance: change in prediction probability when a token is removed.

| Question | Confidence | Highlight |
|---|---|---|
| What did Tesla spend Astor's money on ? | **0.78** | |
| ~~What~~ did Tesla spend Astor's money on ? | 0.67 | What |
| What ~~did~~ Tesla spend Astor's money on ? | 0.72 | did |

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Leave-one-out

[Li et al., 2017]

• Importance: change in prediction probability when a token is removed.

| Question | Confidence | Highlight |
|---|---|---|
| What did Tesla spend Astor's money on ? | **0.78** | |
| ~~What~~ did Tesla spend Astor's money on ? | 0.67 | What |
| What ~~did~~ Tesla spend Astor's money on ? | 0.72 | did |
| What did ~~Tesla~~ spend Astor's money on ? | 0.66 | Tesla |
| What did Tesla ~~spend~~ Astor's money on ? | 0.74 | spend |
| What did Tesla spend ~~Astor's~~ money on ? | 0.76 | Astor's |
| What did Tesla spend Astor's ~~money~~ on ? | **0.48** | money |
| What did Tesla spend Astor's money ~~on~~ ? | 0.72 | on |
| What did Tesla spend Astor's money on ~~?~~ | 0.73 | ? |

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020

# Leave-one-out

[Li et al., 2017]

- Importance: change in prediction probability when a token is removed.

| Question | Confidence | Highlight |
|---|---|---|
| What did Tesla spend Astor's money on ? | **0.78** | |
| ~~What~~ did Tesla spend Astor's money on ? | 0.67 | What |
| What ~~did~~ Tesla spend Astor's money on ? | 0.72 | did |
| What did ~~Tesla~~ spend Astor's money on ? | 0.66 | Tesla |
| What did Tesla ~~spend~~ Astor's money on ? | 0.74 | spend |
| What did Tesla spend ~~Astor's~~ money on ? | 0.76 | Astor's |
| What did Tesla spend Astor's ~~money~~ on ? | **0.48** | money |
| What did Tesla spend Astor's money ~~on~~ ? | 0.72 | on |
| What did Tesla spend Astor's money on ~~?~~ | 0.73 | ? |

What did Tesla spend Astor's money on ?

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020
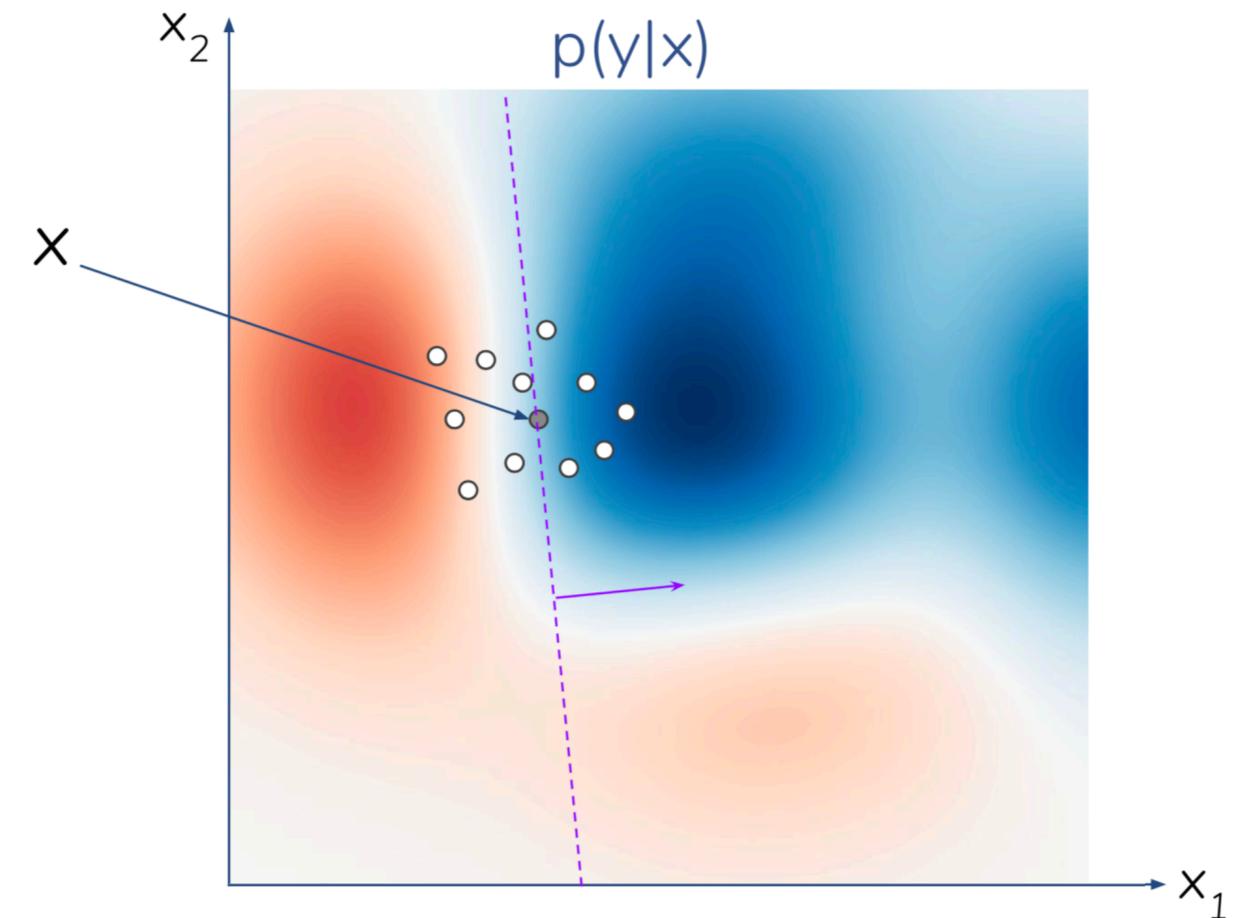
# Leave-one-out

[Li et al., 2017]

- Importance: change in prediction probability when a token is removed.

- Obvious issue: it's not just a single token (or phrase) that matters, usually

| Question | Confidence | Highlight |
|---|---|---|
| What did Tesla spend Astor's money on ? | **0.78** | |
| ~~What~~ did Tesla spend Astor's money on ? | 0.67 | What |
| What ~~did~~ Tesla spend Astor's money on ? | 0.72 | did |
| What did ~~Tesla~~ spend Astor's money on ? | 0.66 | Tesla |
| What did Tesla ~~spend~~ Astor's money on ? | 0.74 | spend |
| What did Tesla spend ~~Astor's~~ money on ? | 0.76 | Astor's |
| What did Tesla spend Astor's ~~money~~ on ? | **0.48** | money |
| What did Tesla spend Astor's money ~~on~~ ? | 0.72 | on |
| What did Tesla spend Astor's money on ~~?~~ | 0.73 | ? |

What did Tesla spend Astor's money on ?

Slide adapted from Sameer Singh's tutorial on Interpretability at EMNLP 2020
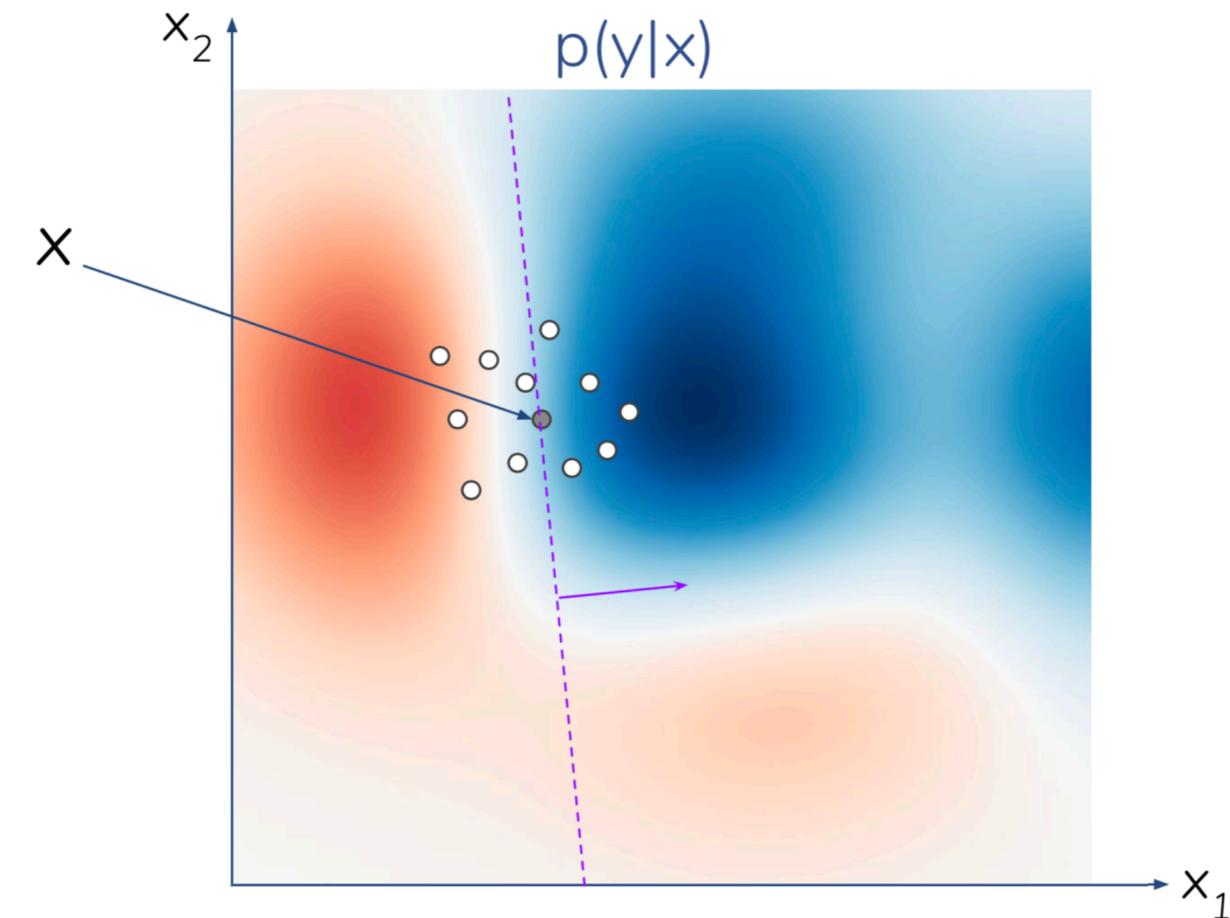
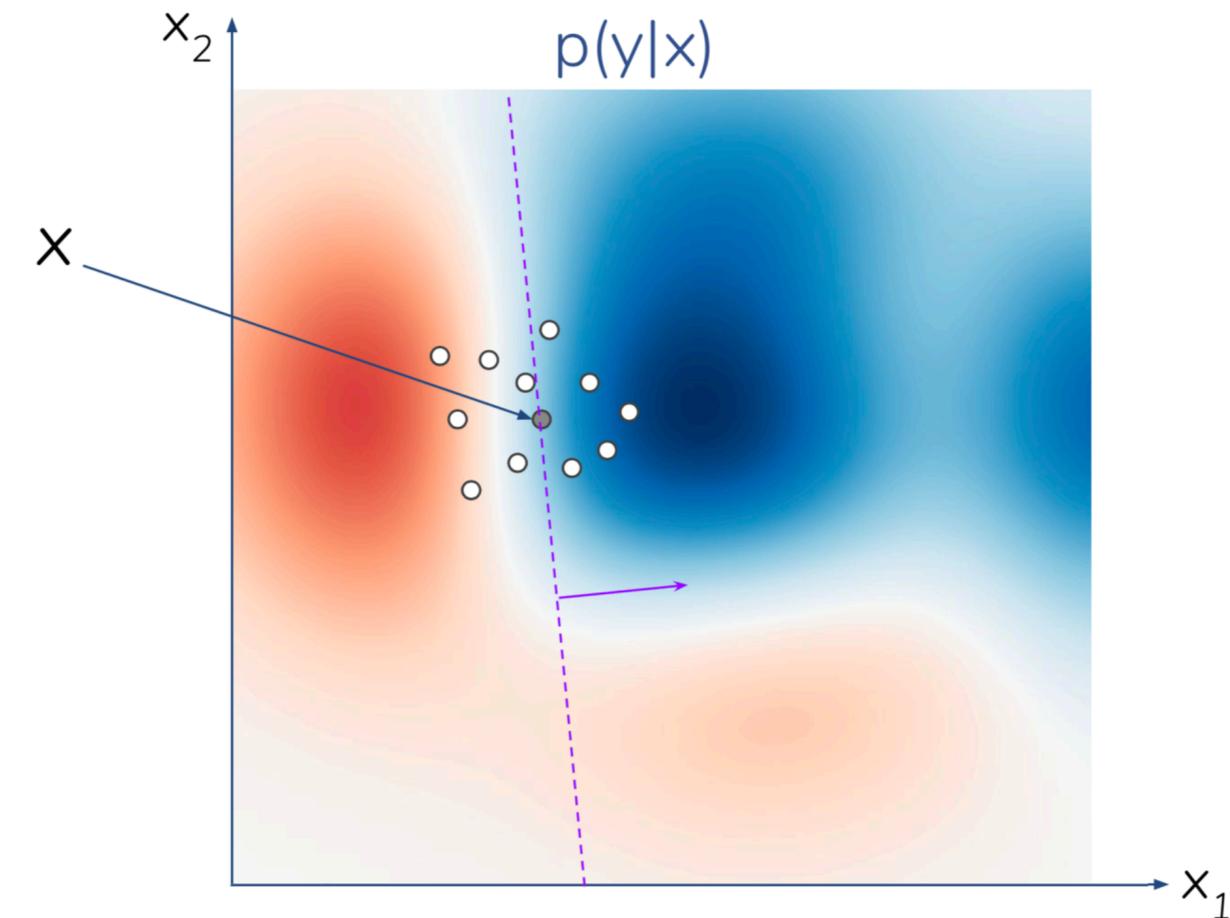# LIME



[Ribeiro et al., 2016]

# LIME

- Find nearby inputs, based on cosine Distance



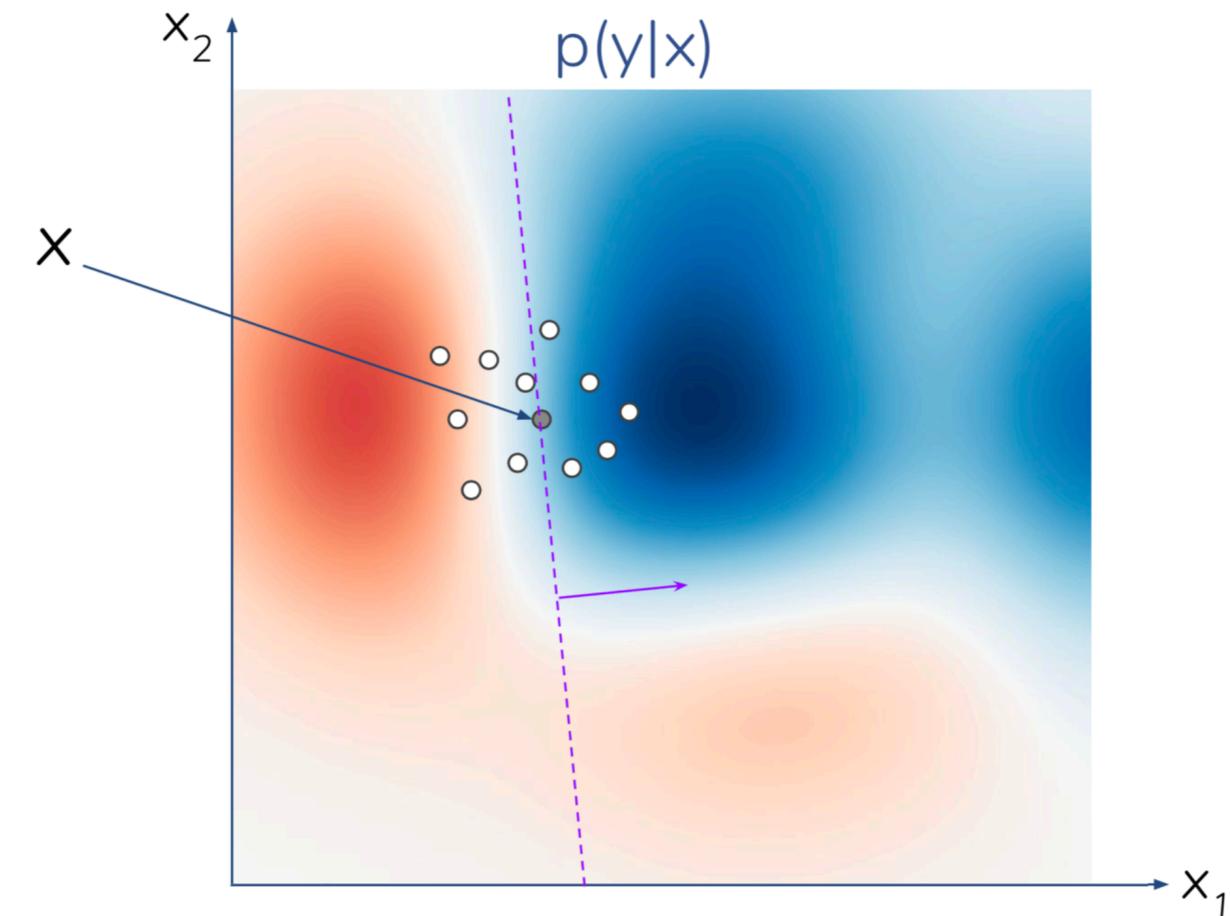[Ribeiro et al., 2016]

# LIME

- Find nearby inputs, based on cosine Distance

- Learn a linear classifier based on model predictions on those points
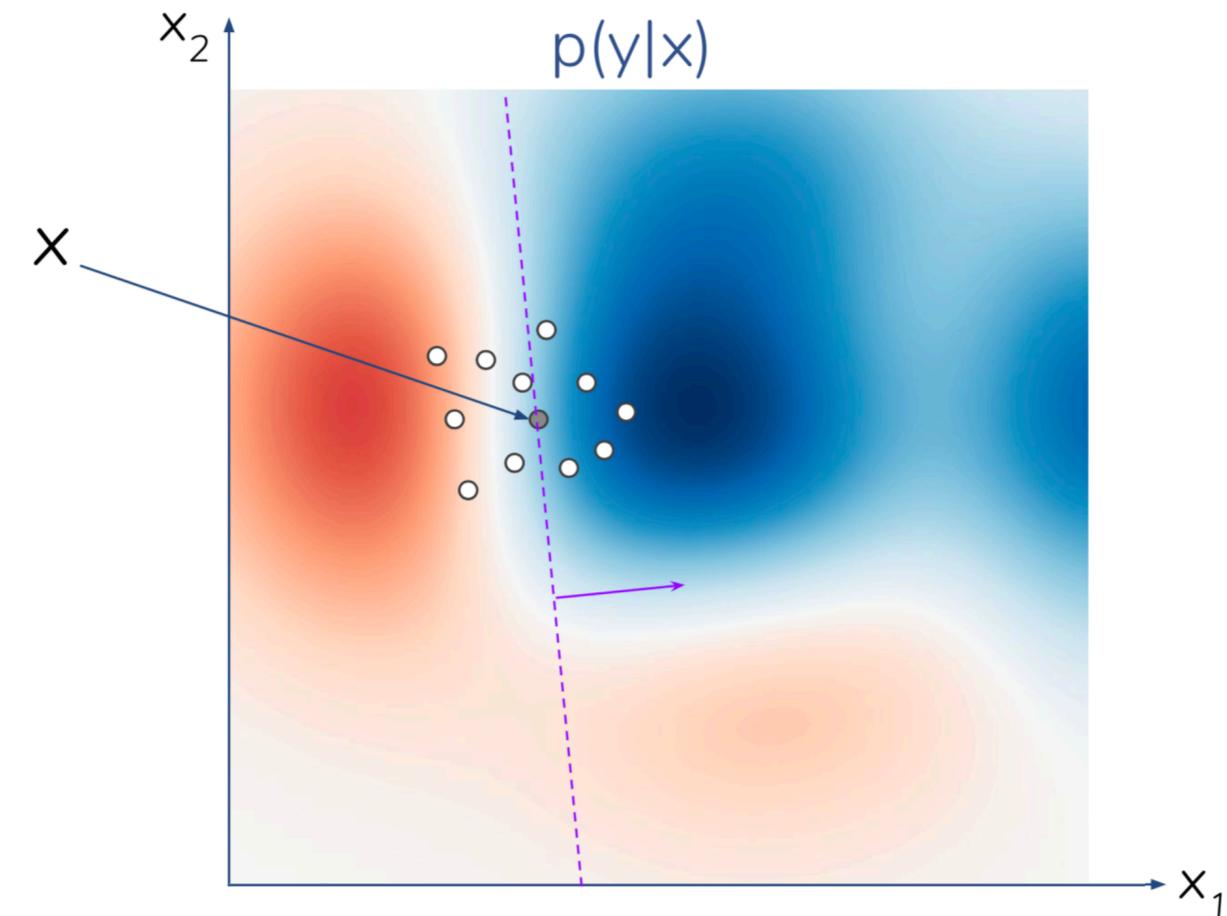


[Ribeiro et al., 2016]

# LIME

- Find nearby inputs, based on cosine Distance

- Learn a linear classifier based on model predictions on those points

  - Use interpretable features

[Ribeiro et al., 2016]

# LIME

- Find nearby inputs, based on cosine Distance

- Learn a linear classifier based on model predictions on those points

  - Use interpretable features

- Weights of the classifier indicate feature importance

[Ribeiro et al., 2016]

# Problems with Input Perturbations

# Problems with Input Perturbations

- How to perturb?

# Problems with Input Perturbations

- How to perturb?

- Overall: "salient" gradients and inputs might not always be human interpretable

# Problems with Input Perturbations

- How to perturb?

- Overall: "salient" gradients and inputs might not always be human interpretable

  - General trend: if it does not match with human intuition, model is probably relying on biases.

# Problems with Input Perturbations

- How to perturb?

- Overall: "salient" gradients and inputs might not always be human interpretable

  - General trend: if it does not match with human intuition, model is probably relying on biases.

  - However, these biases are themselves not consistent / easy to interpret.

# Other variants of input perturbations

# Other variants of input perturbations

- How much can be removed without changing the prediction? [Feng et al. 2018]

**SQuAD**
Context — In 1899, John Jacob Astor IV invested $100,000 for Tesla to further develop and produce a new lighting system. Instead, Tesla used the money to fund his Colorado Springs experiments.
Original — What did Tesla spend Astor's money on ?
**Reduced** — **did**
Confidence — 0.78 → 0.91

**VQA**
Original — What color is the flower ?
Answer — yellow
**Reduced** — **flower ?**
Confidence — 0.827 → 0.819

**SNLI**
Premise — Well dressed man and woman dancing in the street
Original — Two man is dancing on the street
Answer — Contradiction
**Reduced** — **dancing**
Confidence — 0.977 → 0.706

# Other variants of input perturbations

- How much can be removed without changing the prediction? [Feng et al. 2018]

- Adversarial modifications

**SQuAD**
Context    In 1899, John Jacob Astor IV invested $100,000 for Tesla to further develop and produce a new lighting system. Instead, Tesla used the money to fund his Colorado Springs experiments.
Original   What did Tesla spend Astor's money on ?
**Reduced**    **did**
Confidence    0.78 → 0.91

**VQA**
Original   What color is the flower ?
Answer     yellow
**Reduced**    **flower ?**
Confidence    0.827 → 0.819

**SNLI**
Premise    Well dressed man and woman dancing in the street
Original   Two man is dancing on the street
Answer     Contradiction
**Reduced**    **dancing**
Confidence    0.977 → 0.706

# Other variants of input perturbations

- How much can be removed without changing the prediction? [Feng et al. 2018]

- Adversarial modifications

  - Additions [Addsent SQuAD; Jia & Liang, 2017]



**SQuAD**
Context    In 1899, John Jacob Astor IV invested $100,000 for Tesla to further develop and produce a new lighting system. Instead, Tesla used the money to fund his Colorado Springs experiments.
Original   What did Tesla spend Astor's money on ?
**Reduced**    did
Confidence  0.78 → 0.91

**VQA**
Original   What color is the flower ?
Answer     yellow
**Reduced**    flower ?
Confidence  0.827 → 0.819

**SNLI**
Premise    Well dressed man and woman dancing in the street
Original   Two man is dancing on the street
Answer     Contradiction
**Reduced**    dancing
Confidence  0.977 → 0.706

# Other variants of input perturbations

- How much can be removed without changing the prediction? [Feng et al. 2018]

- Adversarial modifications

  - Additions [Addsent SQuAD; Jia & Liang, 2017]

  - Syntactic Paraphrases [SCPN; Iyyer et al., 2018]

**SQuAD**
Context    In 1899, John Jacob Astor IV invested $100,000 for Tesla to further develop and produce a new lighting system. Instead, Tesla used the money to fund his Colorado Springs experiments.
Original   What did Tesla spend Astor's money on ?
**Reduced**    **did**
Confidence 0.78 → 0.91

**VQA**
Original   What color is the flower ?
Answer     yellow
**Reduced**    **flower ?**
Confidence 0.827 → 0.819

**SNLI**
Premise    Well dressed man and woman dancing in the street
Original   Two man is dancing on the street
Answer     Contradiction
**Reduced**    **dancing**
Confidence 0.977 → 0.706

# Other variants of input perturbations

- How much can be removed without changing the prediction? [Feng et al. 2018]

- Adversarial modifications

  - Additions [Addsent SQuAD; Jia & Liang, 2017]

  - Syntactic Paraphrases [SCPN; Iyyer et al., 2018]

- Also reveal biases.

**SQuAD**
Context     In 1899, John Jacob Astor IV invested $100,000 for Tesla to further develop and produce a new lighting system. Instead, Tesla used the money to fund his Colorado Springs experiments.
Original    What did Tesla spend Astor's money on ?
**Reduced**     **did**
Confidence  0.78 → 0.91

**VQA**
Original    What color is the flower ?
Answer      yellow
**Reduced**     **flower ?**
Confidence  0.827 → 0.819

**SNLI**
Premise     Well dressed man and woman dancing in the street
Original    Two man is dancing on the street
Answer      Contradiction
**Reduced**     **dancing**
Confidence  0.977 → 0.706

# Method 3: Architectural Modifications

# Method 3: Architectural Modifications
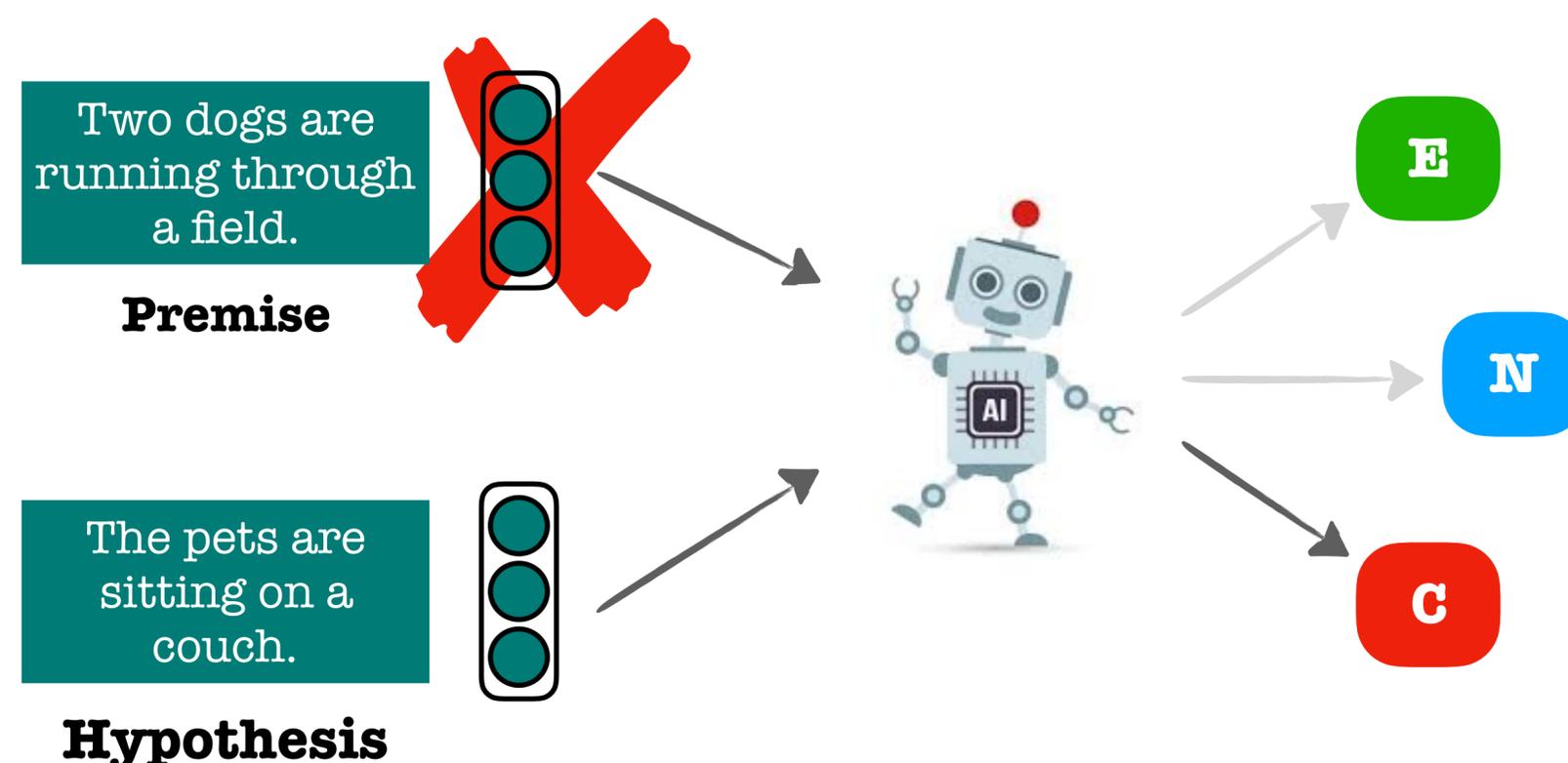
- Partial Input Baselines

# Method 3: Architectural Modifications

- Partial Input Baselines

- Idea: if the model still makes the correct decision despite not receiving the full input, model likely relies on some bias

# Method 3: Architectural Modifications

- Partial Input Baselines

- Idea: if the model still makes the correct decision despite not receiving the full input, model likely relies on some bias



Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

# Method 3: Architectural Modifications

- Partial Input Baselines

- Idea: if the model still makes the correct decision despite not receiving the full input, model likely relies on some bias



Two dogs are running through a field.

**Premise**

The pets are sitting on a couch.

**Hypothesis**

Annotation Artifacts in NLI 〔G*., Swayamdipta*., L., S., B., S., 2018〕

# Method 3: Architectural Modifications

- Partial Input Baselines

- Idea: if the model still makes the correct decision despite not receiving the full input, model likely relies on some bias

- Also tried for VQA [Goyal et al. 2016], SQuAD [Kaushik & Lipton, 2018], among others.

Two dogs are running through a field.

**Premise**

The pets are sitting on a couch.

**Hypothesis**

E

N

C

Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

Question: Can interpretability methods be used to remove biases?

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via Interpretability Methods

- Saliency Methods

- Input Attribution

- Architectural Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via
Interpretability Methods

- Saliency Methods

- Input Attribution

- Architectural
  Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# Mitigation of Biases

# Mitigation of Biases

- Once bias is demonstrated, the next steps involve mitigation (reduction) of biases.

# Mitigation of Biases

- Once bias is demonstrated, the next steps involve mitigation (reduction) of biases.

- Two broad paradigms:

# Mitigation of Biases

• Once bias is demonstrated, the next steps involve mitigation (reduction) of biases.

• Two broad paradigms:

   • Pre-specified (known) biases (task or dataset-specific)

# Mitigation of Biases

• Once bias is demonstrated, the next steps involve mitigation (reduction) of biases.

• Two broad paradigms:

  • Pre-specified (known) biases (task or dataset-specific)

  • Unspecified biases (more general)

# Case Study: Pre-specified Biases

# Case Study: Pre-specified Biases

# Case Study: Pre-specified Biases



**Hate Speech in
Online Platforms**

# Case Study: Pre-specified Biases

**Hate Speech in Online Platforms**

- Human moderation does not scale

# Case Study: Pre-specified Biases

**Hate Speech in Online Platforms**

- Human moderation does not scale

- Spurred a great deal of research on automatic detection of hate speech

# Case Study: Pre-specified Biases

**Hate Speech in
Online Platforms**

- Human moderation does not scale

- Spurred a great deal of research on automatic detection of hate speech

Some examples might contain
offensive or triggering content

Perspective API

I hope this country can now try to get along

Perspective API

I hope this country can now try to get along

15%

Perspective API

Perspective API

I hope this country can now try to get along    15%

If they voted for Hillary they are idiots    75%
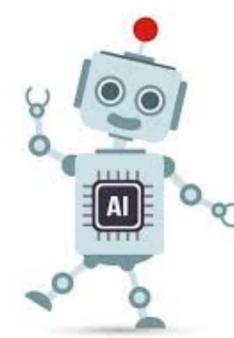
I identify as a straight white man    24%

I identify as a black gay woman    60%

Perspective API

[Sap et. al, 2019]

I hope this country can now try to get along          15%

If they voted for Hillary they are idiots          75%

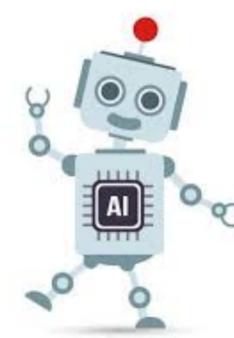I identify as a straight white man          24%

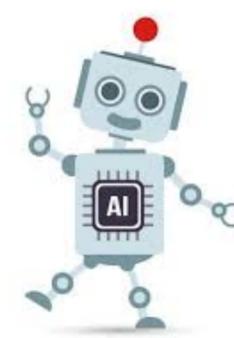I identify as a black gay woman          60%

F*ing love this!          86%

Perspective API

[Sap et. al, 2019]

[Sap et. al, 2019]

[Sap et. al, 2019]

[Sap et. al, 2019]

# Pre-specified biases in hate-speech detection

# Pre-specified biases in hate-speech detection

[Sap et. al, 2019]

# Pre-specified biases in hate-speech detection

[Sap et. al, 2019]

- Hate Speech Detection datasets are indeed biased

# Pre-specified biases in hate-speech detection

[Sap et. al, 2019]

- Hate Speech Detection datasets are indeed biased

  - Identity Biases

# Pre-specified biases in hate-speech detection

[Sap et. al, 2019]

- Hate Speech Detection datasets are indeed biased

  - Identity Biases

  - Profanity Biases

  - Racial / Dialectal Biases

# Pre-specified biases in hate-speech detection

[Sap et. al, 2019]

- Hate Speech Detection datasets are indeed biased

  - Identity Biases

  - Profanity Biases

  - Racial / Dialectal Biases

# Unspecified biases

# Unspecified biases

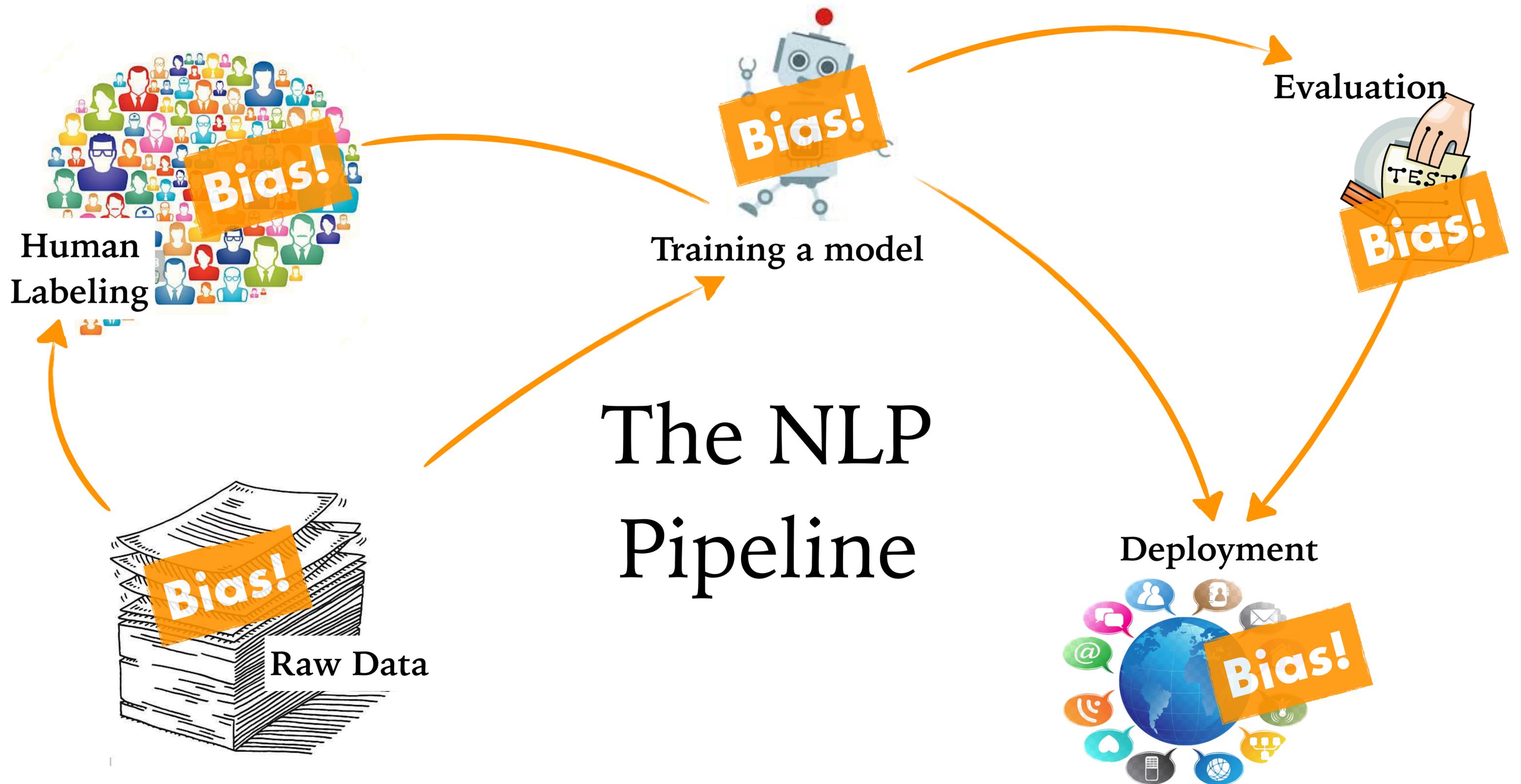- May be too example-specific, and not general enough to explain the entirety of model behavior
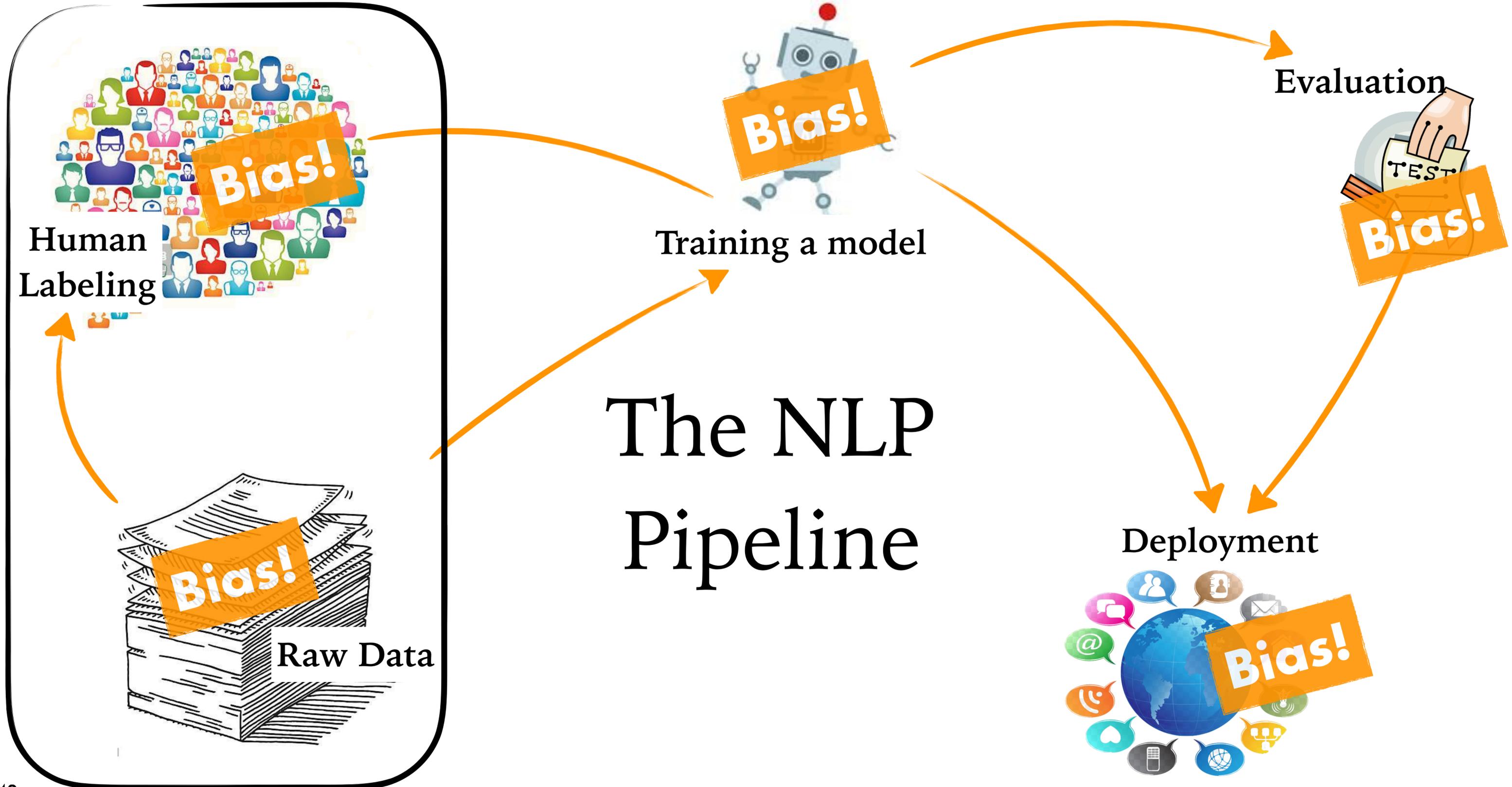
# Unspecified biases

- May be too example-specific, and not general enough to explain the entirety of model behavior

- NLI has many different biases!

# Unspecified biases

- May be too example-specific, and not general enough to explain the entirety of model behavior

- NLI has many different biases!



**Generalization**

**People** play frisbee **outdoors**.
**Hypothesis**

Some men and boys are playing frisbee in a grassy area.
**Premise**



**Shortening**

A person in red is cutting the grass on a riding mower.
**Hypothesis**

A person in a red ~~shirt~~ is mowing the grass with a ~~green~~ riding mower.
**Premise**

Annotation Artifacts in NLI [G*., Swayamdipta*, L., S., B., S., 2018]

Human Labeling

Bias!

Bias!

Training a model

Evaluation

Bias!

The NLP Pipeline

Raw Data

Bias!

Deployment

Bias!

43

The NLP Pipeline

Human Labeling

Bias!

Raw Data

Bias!

Training a model

Bias!

Evaluation

Bias!

Deployment

Bias!

Human Labeling

Bias!

Training a model

Bias!

Evaluation

Bias!

The NLP Pipeline

Raw Data

Bias!

Deployment

Bias!

# Addressing Biases: Datasets

**Human Labeling**

Bias!

**Raw Data**

Bias!

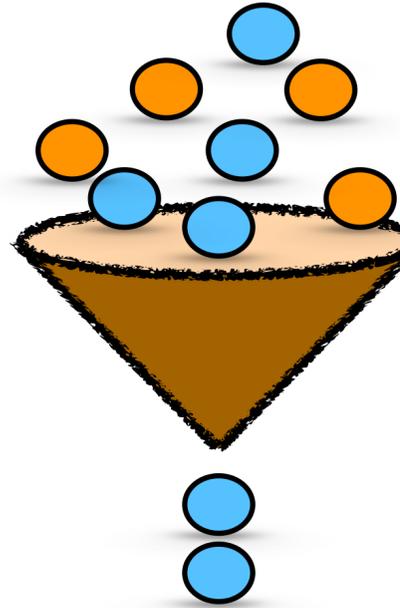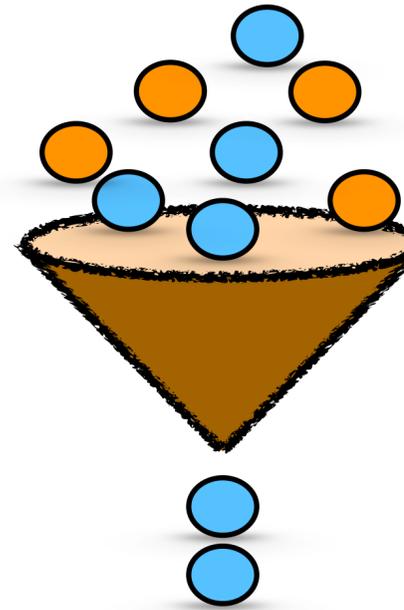# Addressing Biases: Datasets

**Human Labeling**

- One solution: Filtering / Downsampling the data to remove instances that "leak" the correct answer, but because of the wrong reasons.

**Raw Data**

# Addressing Biases: Datasets

**Human Labeling**

**Bias!**

**Bias!**

**Raw Data**

- One solution: Filtering / Downsampling the data to remove instances that "leak" the correct answer, but because of the wrong reasons.

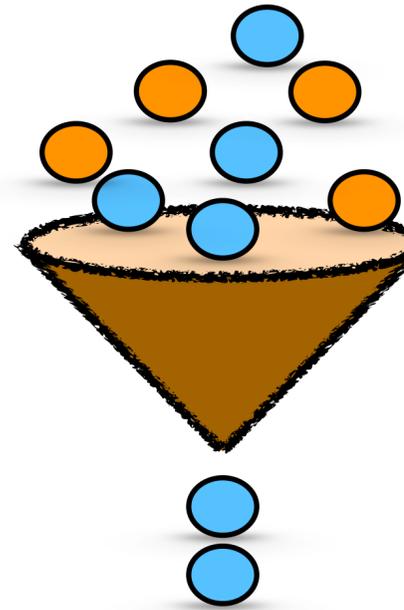- Simple for known biases (rules / simple classifiers)

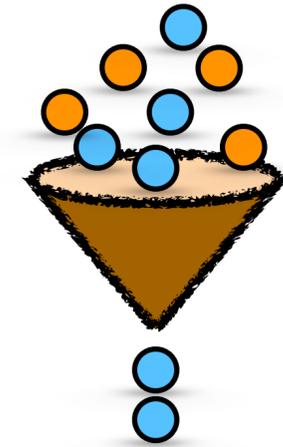# Addressing Biases: Datasets
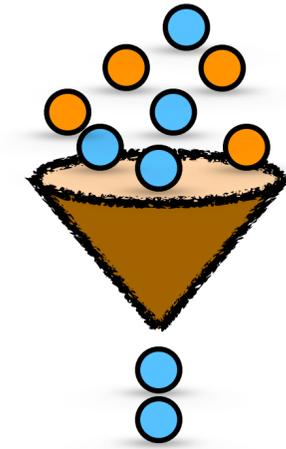
**Human Labeling**

**Raw Data**

- One solution: Filtering / Downsampling the data to remove instances that "leak" the correct answer, but because of the wrong reasons.

- Simple for known biases (rules / simple classifiers)

- Also possible for unspecified biases!
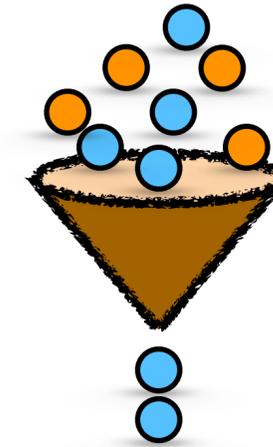
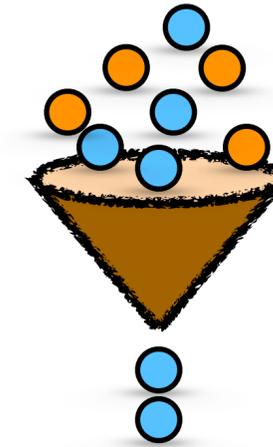# Dataset Filtering

# Dataset Filtering

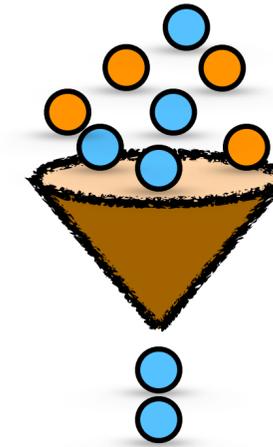- What instances to filter?

# Dataset Filtering

- What instances to filter?

  - Key intuition: Examples which are relatively **easy** for a model might contain spurious correlations

# Dataset Filtering

- What instances to filter?

  - Key intuition: Examples which are relatively **easy** for a model might contain spurious correlations

- Easy examples can be detected:

# Dataset Filtering

- What instances to filter?

  - Key intuition: Examples which are relatively **easy** for a model might contain spurious correlations

- Easy examples can be detected:

    - By simple model architectures [AFLite; LeBras et al., 2020]
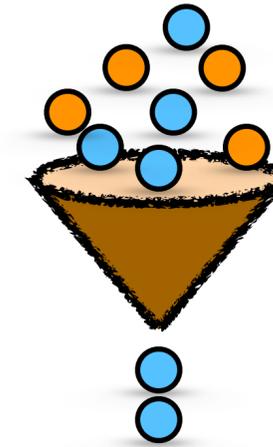
# Dataset Filtering

- What instances to filter?

  - Key intuition: Examples which are relatively **easy** for a model might contain spurious correlations

- Easy examples can be detected:

    - By simple model architectures [AFLite; LeBras et al., 2020]

    - Based on how the training proceeds [Dataset Cartography; Swayamdipta et al., 2020]

# AFLite in action

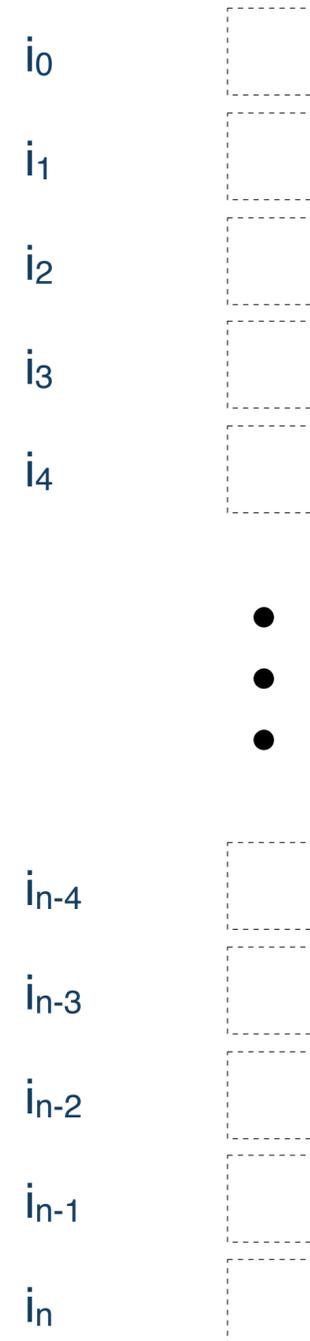Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action

- Detecting and reducing model biases by (ensembles of) simplified architectures.

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action

- Detecting and reducing model biases by (ensembles of) simplified architectures.

$i_0$

$i_1$

$i_2$

$i_3$

$i_4$

$\vdots$

$i_{n-4}$

$i_{n-3}$

$i_{n-2}$

$i_{n-1}$

$i_n$

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action

- Detecting and reducing model biases by (ensembles of) simplified architectures.

$i_0$

$i_1$
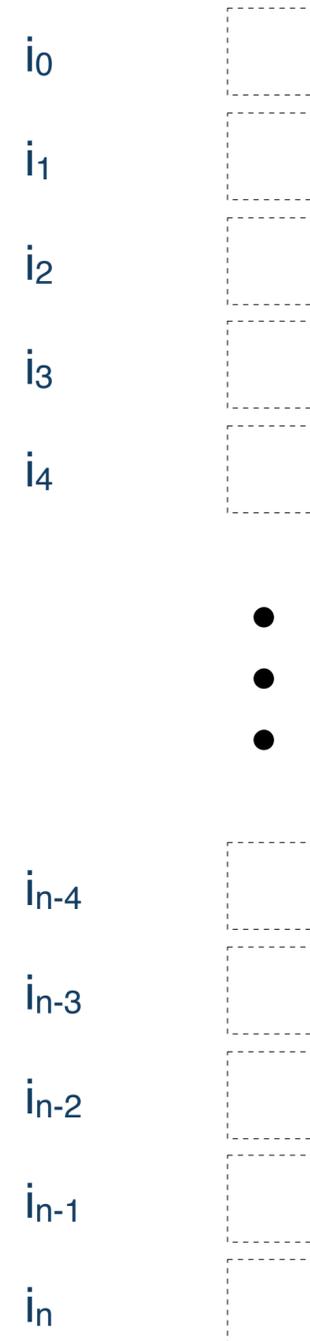
$i_2$

$i_3$

$i_4$

.
.
.

$i_{n-4}$

$i_{n-3}$

$i_{n-2}$

$i_{n-1}$

$i_n$

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action

$i_0$

$i_1$

$i_2$

$i_3$

$i_4$
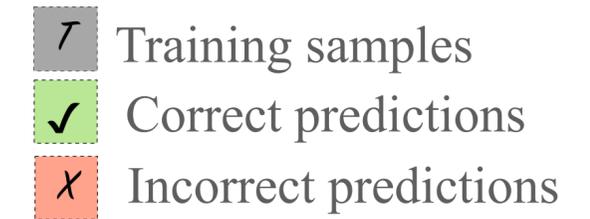
- Detecting and reducing model biases by (ensembles of) simplified architectures.

$i_{n-4}$

$i_{n-3}$

$i_{n-2}$

$i_{n-1}$

$i_n$

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite
# in action

**Linear Classifier**

$i_0$

$i_1$

$i_2$

$i_3$

$i_4$

$i_{n-4}$

$i_{n-3}$

$i_{n-2}$

$i_{n-1}$

$i_n$

| $T$ | Training samples |
| $✓$ | Correct predictions |
| $✗$ | Incorrect predictions |

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action



Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite
# in action

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite
# in action



Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite
# in action

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action

Legend:
- $T$ — Training samples
- ✓ — Correct predictions
- ✗ — Incorrect predictions

$i_0$ · $i_1$ · $i_2$ · $i_3$ · $i_4$ · $i_{n-4}$ · $i_{n-3}$ · $i_{n-2}$ · $i_{n-1}$ · $i_n$

Highest predictability score

Lowest predictability score

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite
# in action



Training samples
Correct predictions
Incorrect predictions

Highest predictability score

Lowest predictability score

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite

# in action

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite in action

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# AFLite

# in action

⊤ Training samples
✓ Correct predictions
✗ Incorrect predictions

$i_0$
$i_1$
$i_2$
$i_3$

$\cdots$

$i_{k-4}$
$i_{k-3}$
$i_{k-2}$
$i_{k-1}$
$i_k$

Adversarial Filters of Dataset Biases [L., Swayamdipta, Z., B., P., S., C.]

# Dataset Filtering

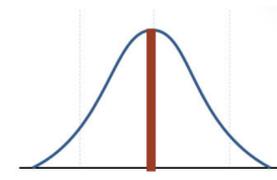- What instances to filter?

  - Key intuition: Examples which are relatively **easy** for a model might contain spurious correlations

- Easy examples can be detected:

  - By simple model architectures [AFLite; LeBras et al., 2020]

  - **Based on how the training proceeds [Dataset Cartography; Swayamdipta et al., 2020]**

# Training Dynamics

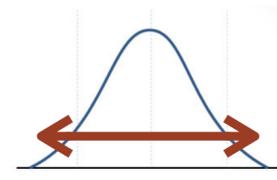**correctness** ✔  **confidence**  **variability**

Dataset Cartography [Swayamdipta et. al, 2020]

# Training Dynamics

**correctness** ✓          **confidence**          **variability**

**across E training epochs…**

Dataset Cartography [Swayamdipta et. al, 2020]

# Training Dynamics

**correctness** ✔️       **confidence**       **variability**

- Ratio at which model prediction matches **true class**

$$\hat{c}_i = \frac{1}{E} \sum_{e=1}^{E} 1[y_i^* = \arg\max_y p_{\theta^{(e)}}(y \mid x_i)]$$

**across E training epochs…**

Dataset Cartography [Swayamdipta et. al, 2020]

# Training Dynamics

**correctness** ✔         **confidence**           **variability** 

- Ratio at which model prediction matches **true class**

- Mean probability of the **true class**

$$\hat{c}_i = \frac{1}{E}\sum_{e=1}^{E} 1[y_i^* = \arg\max_y p_{\theta^{(e)}}(y\,|\,x_i)]$$

$$\hat{\mu}_i = \frac{1}{E}\sum_{e=1}^{E} p_{\theta^{(e)}}(y_i^*\,|\,x_i)$$

## across E training epochs…

55    Dataset Cartography [Swayamdipta et. al, 2020]

# Training Dynamics

## correctness ✔

## confidence

## variability

- Ratio at which model prediction matches **true class**

- Mean probability of the **true class**

- Standard deviation of the **true class** probability

$$\hat{c}_i = \frac{1}{E} \sum_{e=1}^{E} 1[y_i^* = \arg\max_y p_{\theta^{(e)}}(y \mid x_i)]$$

$$\hat{\mu}_i = \frac{1}{E} \sum_{e=1}^{E} p_{\theta^{(e)}}(y_i^* \mid x_i)$$

$$\hat{\sigma}_i = \sqrt{\frac{\sum_{e=1}^{E} (p_{\theta^{(e)}}(y_i^* \mid x_i) - \hat{\mu}_i)^2}{E}}$$

## across E training epochs…

55    Dataset Cartography [Swayamdipta et. al, 2020]

# Training Dynamics

**correctness** ✓ | **confidence** | **variability**

- Ratio at which model prediction matches **true class**

- Mean probability of the **true class**

- Standard deviation of the **true class** probability

$$\hat{c}_i = \frac{1}{E} \sum_{e=1}^{E} 1[y_i^* = \arg\max_y p_{\theta^{(e)}}(y \,|\, x_i)]$$

$$\hat{\mu}_i = \frac{1}{E} \sum_{e=1}^{E} p_{\theta^{(e)}}(y_i^* \,|\, x_i)$$

$$\hat{\sigma}_i = \sqrt{\frac{\sum_{e=1}^{E} (p_{\theta^{(e)}}(y_i^* \,|\, x_i) - \hat{\mu}_i)^2}{E}}$$

**across E training epochs…**

*By-product of training!*

Dataset Cartography [Swayamdipta et. al, 2020]

Training Dynamics

Dataset Cartography [Swayamdipta et. al, 2020]

Dataset Cartography [Swayamdipta et. al, 2020]

confidence

Mean probability of the **true class**

Training Dynamics

variability

Standard deviation of the **true class** probability

Dataset Cartography [Swayamdipta et. al, 2020]

Dataset Cartography [Swayamdipta et. al, 2020]

# Dataset Cartography



Dataset Cartography [Swayamdipta et. al, 2020]

# Dataset Cartography



Dataset Cartography [Swayamdipta et. al, 2020]

# Dataset Cartography



Dataset Cartography [Swayamdipta et. al, 2020]

# Dataset Cartography



Dataset Cartography [Swayamdipta et. al, 2020]

# Dataset Cartography



Sc**w you Trump supporters!

Good luck and let's join hands to form unity.

Dataset Cartography [Swayamdipta et. al, 2020]

# Dataset Cartography



Dataset Cartography [Swayamdipta et. al, 2020]

# Question: Doesn't removing data hurt performance?

# The NLP Pipeline

Human Labeling

Training a model

Evaluation

Deployment

Raw Data

Bias!

Evaluation

Bias!

Human
Labeling

Bias!

Training a model

Bias!

Deployment

Bias!

# The NLP
# Pipeline

Raw Data

Bias!

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

• Can be used to reduce pre-specified biases

  • e.g. Identity, Dialect, Profanity biases in
    Hate Speech Detection

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model



Bias-Only

Ensemble

Full

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model

- Bias-only model captures all the biases

Profanities mean toxicity

Bias-Only

Ensemble

Full

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model

- Bias-only model captures all the biases

- Full model no longer focuses on biases

Profanities mean toxicity

Bias-Only

Ensemble

Let's look at all features

Full

60

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model

- Bias-only model captures all the biases

- Full model no longer focuses on biases

Profanities mean toxicity

Bias-Only

Ensemble

Cause grandma's a bad b*ch and she had to let you know your man can become y'all's man if she pleases.
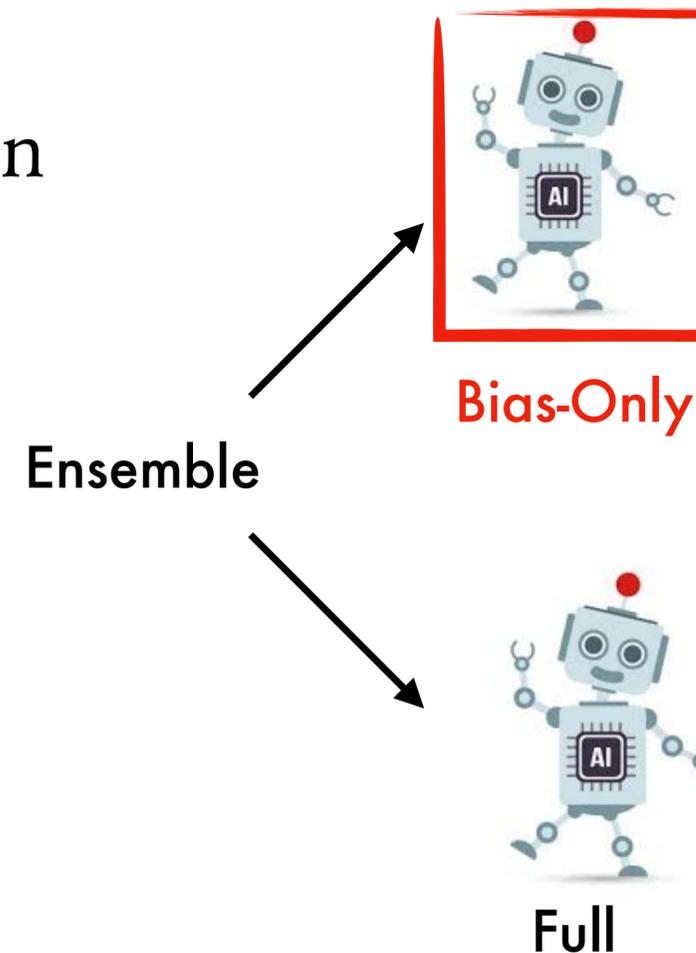
Let's look at all features

Full

# Addressing Biases: Models

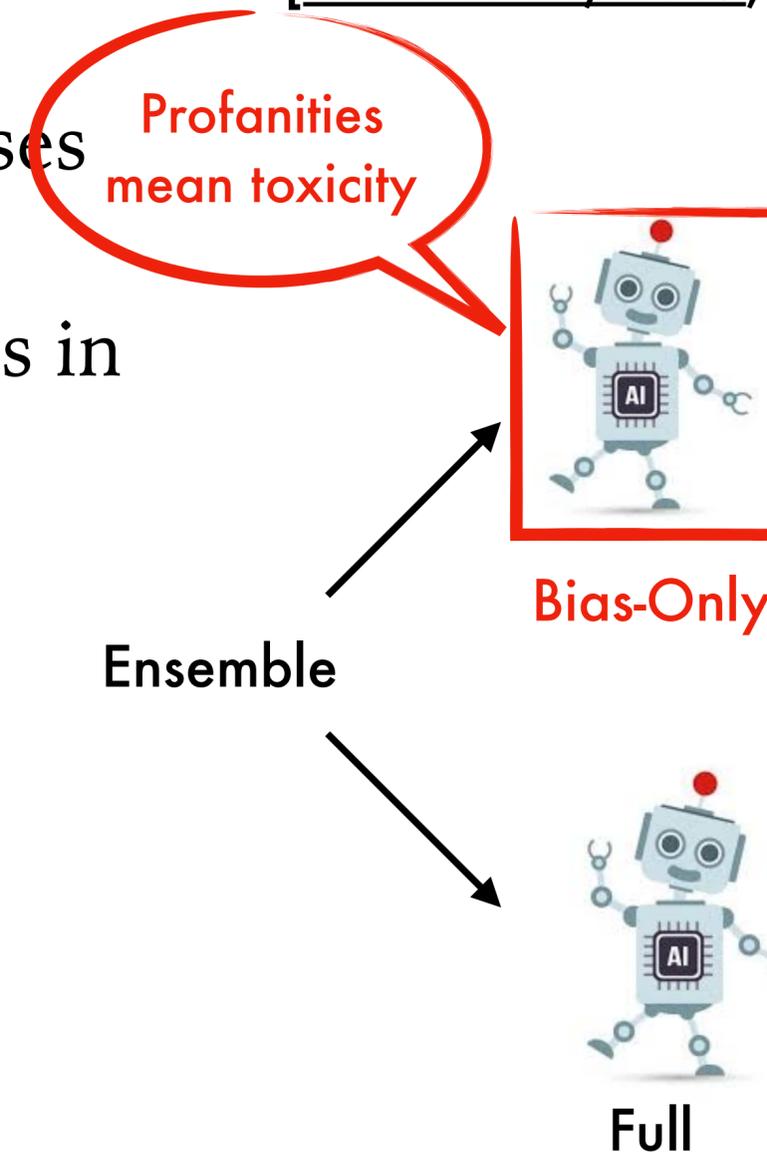[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model

- Bias-only model captures all the biases

- Full model no longer focuses on biases

*Profanities mean toxicity*

**Bias-Only**

Ensemble

*Let's look at all features*

Full

*Cause grandma's a bad b\*ch and she had to let you know your man can become y'all's man if she pleases.*

# Addressing Biases: Models

[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]
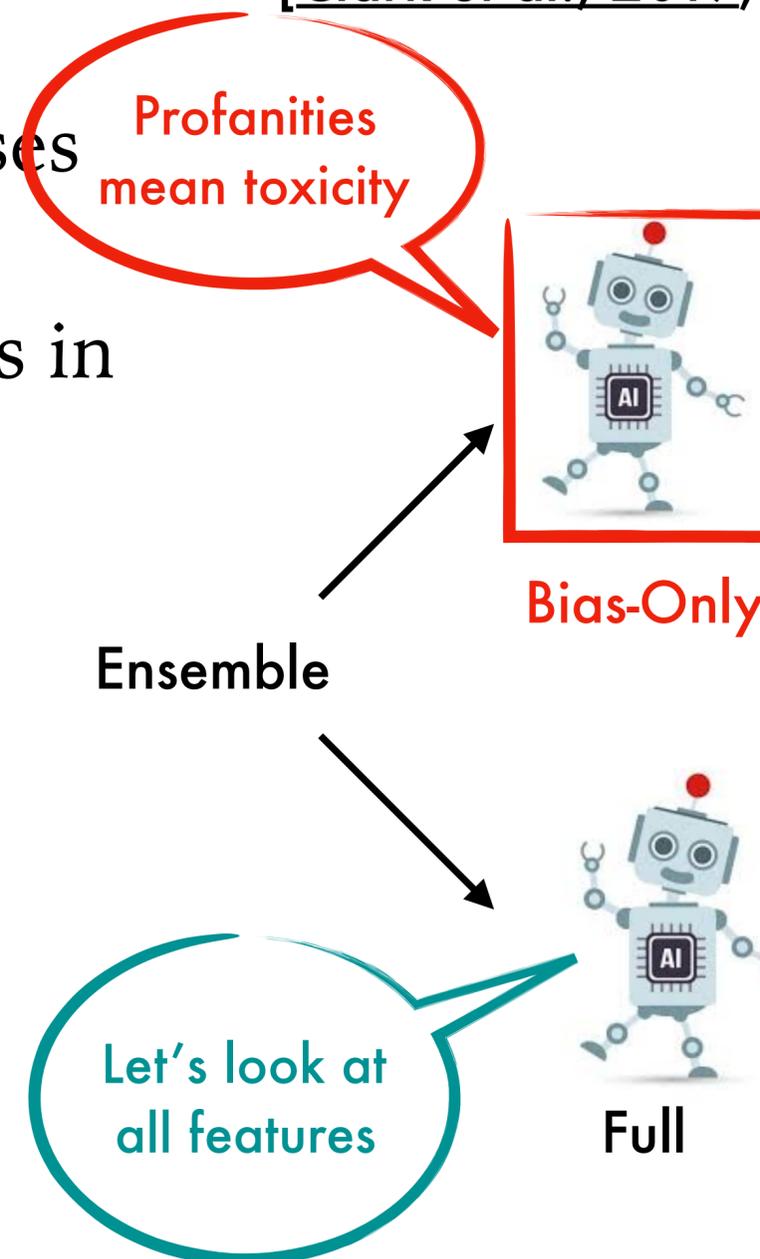
- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model

- Bias-only model captures all the biases

- Full model no longer focuses on biases

Profanities mean toxicity

Bias-Only

Ensemble

Let's look at all features

Full

Cause grandma's a bad b*ch and she had to let you know your man can become y'all's man if she pleases.

# Addressing Biases: Models

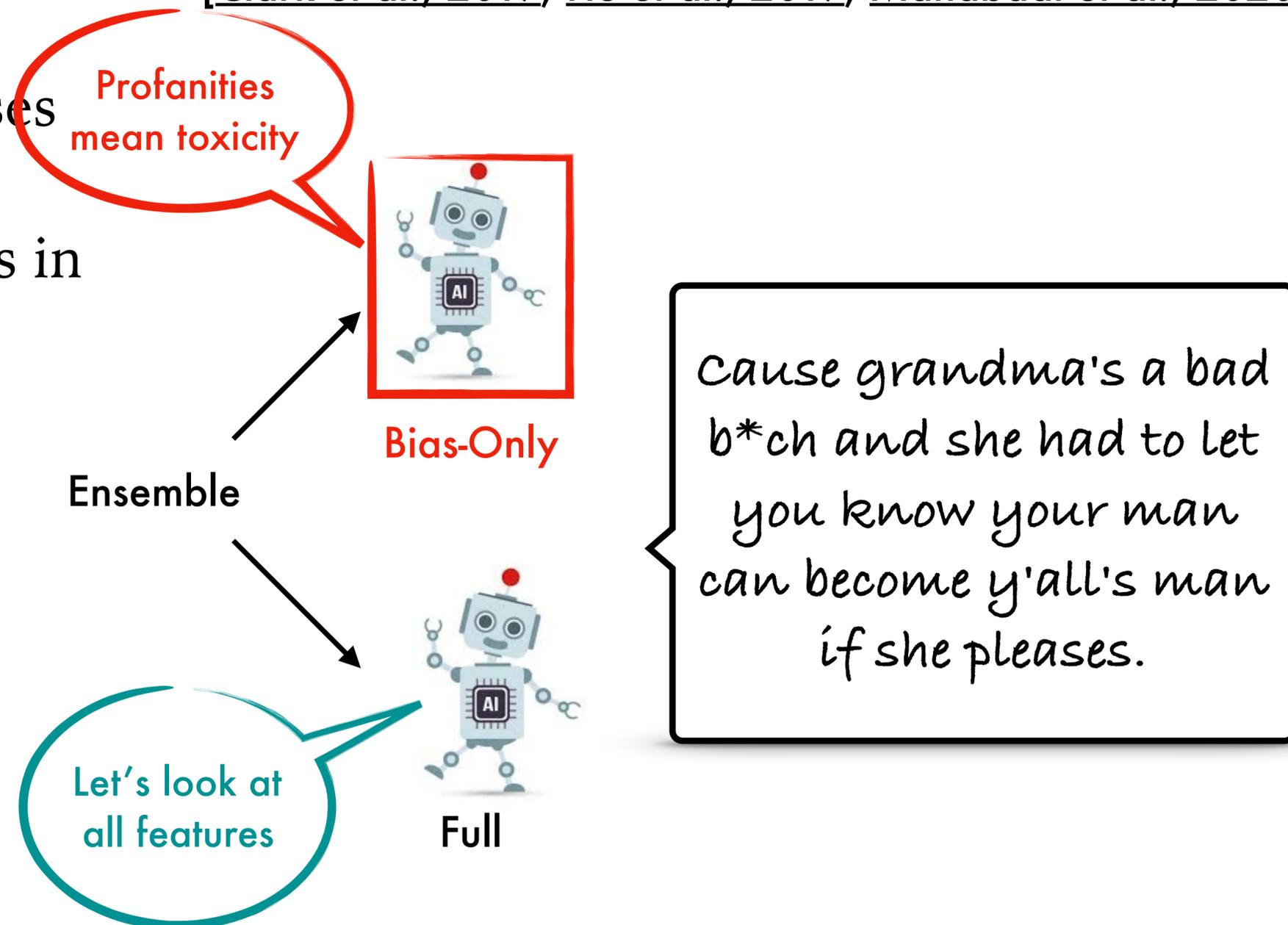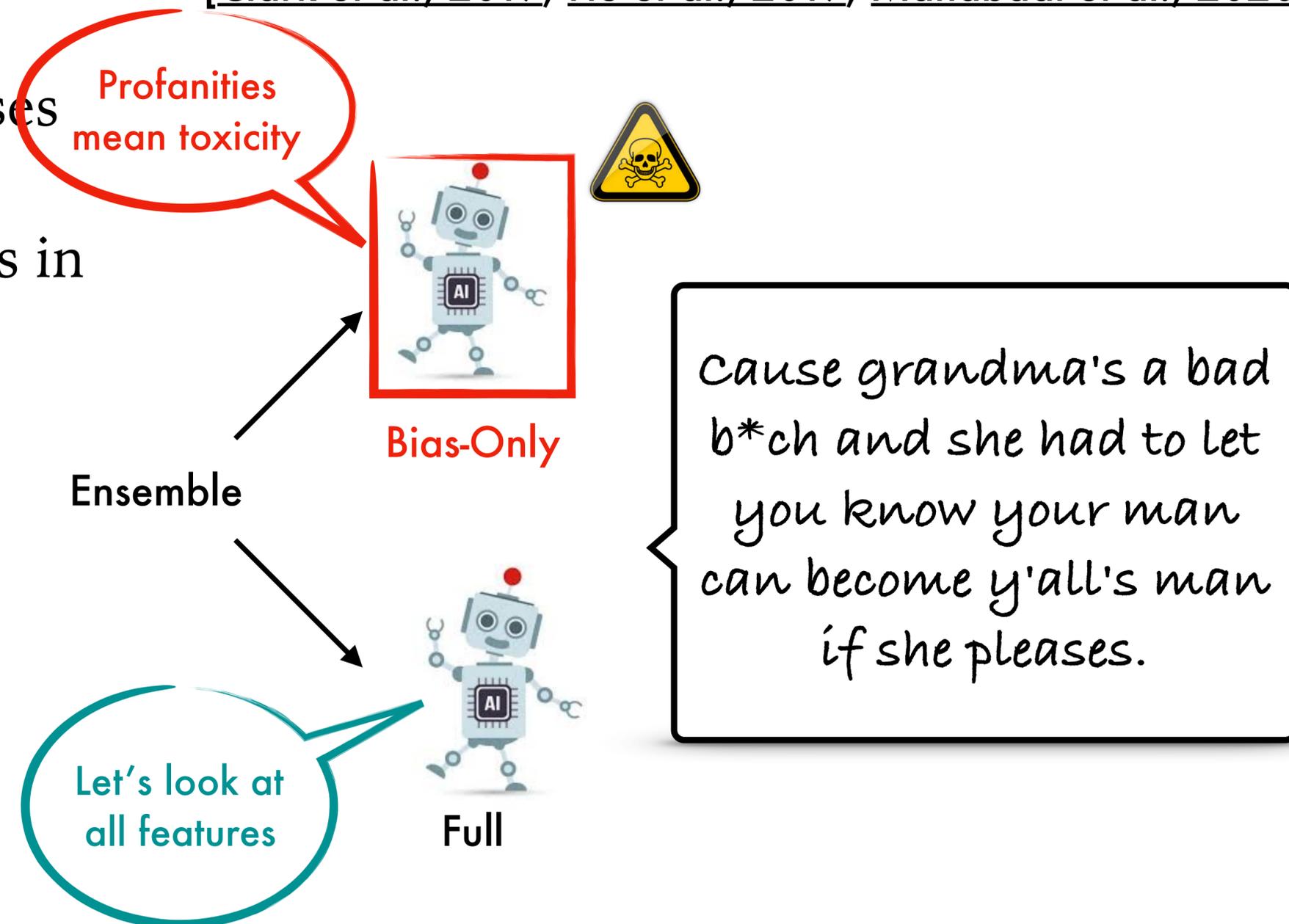[Clark et al., 2019; He et al., 2019; Mahabadi et al., 2020]

- Can be used to reduce pre-specified biases

  - e.g. Identity, Dialect, Profanity biases in Hate Speech Detection

- Ensemble of bias-only and full model

- Bias-only model captures all the biases

- Full model no longer focuses on biases

**Ensemble**

*Let's look at all features*

**Full**

Cause grandma's a bad b*ch and she had to let you know your man can become y'all's man if she pleases.

# Adversarial Methods

[Belinkov et al., 2019; Ganin et al., 2016]

# Adversarial Methods

- Pre-specified biases

[Belinkov et al., 2019; Ganin et al., 2016]

# Adversarial Methods

- Pre-specified biases

- Can the model predict something about the input itself? This is typically the bias feature.

[Belinkov et al., 2019; Ganin et al., 2016]

# Adversarial Methods

- Pre-specified biases

- Can the model predict something about the input itself? This is typically the bias feature.

  - e.g. Can the model predict the gender from a professional bio? Given that we know models have gender bias [De-Arteaga et al., 2019]

[Belinkov et al., 2019; Ganin et al., 2016]

# Adversarial Methods

- Pre-specified biases

- Can the model predict something about the input itself? This is typically the bias feature.

  - e.g. Can the model predict the gender from a professional bio? Given that we know models have gender bias [De-Arteaga et al., 2019]

- Now, the auxiliary is discouraged (ensure you cannot predict the bias) in an adversarial setting

[Belinkov et al., 2019; Ganin et al., 2016]

# Adversarial Methods

- Pre-specified biases

- Can the model predict something about the input itself? This is typically the bias feature.

  - e.g. Can the model predict the gender from a professional bio? Given that we know models have gender bias [De-Arteaga et al., 2019]

- Now, the auxiliary is discouraged (ensure you cannot predict the bias) in an adversarial setting

- Might not entirely remove the information

[Belinkov et al., 2019; Ganin et al., 2016]

# Bias Mitigation Summary

- Dataset Filtering Methods

  - Algorithms that differentiate data instances (AFLite, Dataset Cartography)

  - Can be applied to unspecified biases

- Models with Auxiliary Objectives

  - Ensembles, Adversarial Approaches

  - Effective for pre-specified biases

# Bias Mitigation Summary

- Dataset Filtering Methods

  - Algorithms that differentiate data instances (AFLite, Dataset Cartography)

  - Can be applied to unspecified biases

- Models with Auxiliary Objectives

  - Ensembles, Adversarial Approaches

  - Effective for pre-specified biases

How effective are these methods?

# Bias Mitigation Summary

- Dataset Filtering Methods

  - Algorithms that differentiate data instances (AFLite, Dataset Cartography)

  - Can be applied to unspecified biases

- Models with Auxiliary Objectives

  - Ensembles, Adversarial Approaches

  - Effective for pre-specified biases

How effective are these methods?

Be careful of the term "debiasing"…

# This Lecture

Biases in NLP

- Dataset Biases

- Model Biases

Discovering Biases via Interpretability Methods

- Saliency Methods

- Input Attribution

- Architectural Modifications

Mitigating Biases

- Filtering Datasets

- Auxiliary Objectives

# Summary

# Summary

- Biases are present wherever humans are involved: data collection & model design

# Summary

- Biases are present wherever humans are involved: data collection & model design

  - The term "bias" can be overloaded: biases can be "good" or "bad"

# Summary

- Biases are present wherever humans are involved: data collection & model design

  - The term "bias" can be overloaded: biases can be "good" or "bad"

- Interpretability methods can be used to detect and discover biases in models and data

# Summary

- Biases are present wherever humans are involved: data collection & model design

  - The term "bias" can be overloaded: biases can be "good" or "bad"

- Interpretability methods can be used to detect and discover biases in models and data

- Bias discovery and bias mitigation is not necessarily a pipeline

# Summary

- Biases are present wherever humans are involved: data collection & model design

  - The term "bias" can be overloaded: biases can be "good" or "bad"

- Interpretability methods can be used to detect and discover biases in models and data

- Bias discovery and bias mitigation is not necessarily a pipeline

- Bias mitigation methods either focus on models or datasets.

# Thank you!
# Questions?