

# Bias and Fairness

Some ML issues in the real world

CMSC 422

Slides adapted from Prof. CARPUAT

# Many Cars Tone Deaf To Women's Voices

Female voices pose a bigger challenge for voice-activated technology than men's voices



Women have a tougher time using voice-command systems I...

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

May 23, 2016

## Prediction Fails Differently for Black Defendants

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

*Overall, Northpointe's assessment tool correctly predicts recidivism 61 percent of the time. But blacks are almost twice as likely as whites to be labeled a higher risk but not actually re-offend. It makes the opposite mistake among whites: They are much more likely than blacks to be labeled lower risk but go on to commit other crimes. (Source: ProPublica analysis of data from Broward County, Fla.)*

# Recall: Formal Definition of Binary Classification (from CIML)

## TASK: BINARY CLASSIFICATION

*Given:*

1. An input space  $\mathcal{X}$
2. An unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \{-1, +1\}$

*Compute:* A function  $f$  minimizing:  $\mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) \neq y]$

# Train/Test Mismatch

- When working with real world data, training sample
  - reflects human biases
  - is influenced by practical concerns
    - e.g., what kind of data is easy to obtain
- Train/test distribution mismatch is frequent issue
  - aka covariate shift, sample selection bias, domain adaptation

# Domain Adaptation

- What does it mean for 2 distributions to be related?
- When 2 distributions are related how can we build models that effectively share information between them?

# Unsupervised adaptation

- **Goal:** learn a classifier  $f$  that achieves low expected loss under new distribution  $\mathcal{D}^{new}$
- Given labeled training data from old distribution  $\mathcal{D}^{old}$   $(x_1, y_1), \dots, (x_N, y_N)$
- And unlabeled examples from new distribution  $\mathcal{D}^{new}$ :  $z_1, \dots, z_M$

# Relation between test loss in new domain and old domain

$$\text{test loss} \tag{8.1}$$

$$= \mathbb{E}_{(x,y) \sim \mathcal{D}^{\text{new}}} [\ell(y, f(x))] \tag{8.2} \quad \text{definition}$$

$$= \sum_{(x,y)} \mathcal{D}^{\text{new}}(x,y) \ell(y, f(x)) \tag{8.3} \quad \text{expand expectation}$$

$$= \sum_{(x,y)} \mathcal{D}^{\text{new}}(x,y) \frac{\mathcal{D}^{\text{old}}(x,y)}{\mathcal{D}^{\text{old}}(x,y)} \ell(y, f(x)) \tag{8.4} \quad \text{times one}$$

$$= \sum_{(x,y)} \mathcal{D}^{\text{old}}(x,y) \frac{\mathcal{D}^{\text{new}}(x,y)}{\mathcal{D}^{\text{old}}(x,y)} \ell(y, f(x)) \tag{8.5} \quad \text{rearrange}$$

$$= \mathbb{E}_{(x,y) \sim \mathcal{D}^{\text{old}}} \left[ \frac{\mathcal{D}^{\text{new}}(x,y)}{\mathcal{D}^{\text{old}}(x,y)} \ell(y, f(x)) \right] \tag{8.6} \quad \text{definition}$$

# How can we estimate the ratio between $D_{\text{new}}$ and $D_{\text{old}}$ ?

Fixed base distribution

$S$  = selection variable

$$\frac{D^{\text{new}}(x, y)}{D^{\text{old}}(x, y)} = \frac{\frac{1}{Z^{\text{new}}} \mathcal{D}^{\text{base}}(x, y) p(s = 0 | x)}{\frac{1}{Z^{\text{old}}} \mathcal{D}^{\text{base}}(x, y) p(s = 1 | x)} \quad \text{definition (8.9)}$$

$$= \frac{\frac{1}{Z^{\text{new}}} p(s = 0 | x)}{\frac{1}{Z^{\text{old}}} p(s = 1 | x)} \quad \text{cancel base (8.10)}$$

$$= Z \frac{p(s = 0 | x)}{p(s = 1 | x)} \quad \text{consolidate (8.11)}$$

$$= Z \frac{1 - p(s = 1 | x)}{p(s = 1 | x)} \quad \text{binary selection (8.12)}$$

$$= Z \left[ \frac{1}{p(s = 1 | x)} - 1 \right] \quad \text{rearrange (8.13)}$$

We can estimate  $P(s=1 | x)$  using a binary classifier!

---

**Algorithm 23** SELECTIONADAPTATION( $\langle(\mathbf{x}_n, y_n)\rangle_{n=1}^N, \langle\mathbf{z}_m\rangle_{m=1}^M, \mathcal{A}$ )

---

- 1:  $D^{dist} \leftarrow \langle(\mathbf{x}_n, +1)\rangle_{n=1}^N \cup \langle(\mathbf{z}_m, -1)\rangle_{m=1}^M$  // assemble data for distinguishing  
// between old and new distributions
  - 2:  $\hat{p} \leftarrow$  train logistic regression on  $D^{dist}$
  - 3:  $D^{weighted} \leftarrow \langle(\mathbf{x}_n, y_n, \frac{1}{\hat{p}(\mathbf{x}_n)} - 1)\rangle_{n=1}^N$  // assemble weight classification  
// data using selector
  - 4: **return**  $\mathcal{A}(D^{weighted})$  // train classifier
-

# Supervised adaptation

- **Goal:** learn a classifier  $f$  that achieves low expected loss under new distribution  $\mathcal{D}^{new}$
- Given labeled training data from old distribution  
 $\mathcal{D}^{old} : \langle \mathbf{x}_n^{(old)}, y_n^{(old)} \rangle_{n=1}^N$
- And labeled examples from new distribution  
 $\mathcal{D}^{new} : \langle \mathbf{x}_m^{(new)}, y_m^{(new)} \rangle_{m=1}^M$

# One solution: feature augmentation

- Map inputs to a new augmented representation

	shared	old-only	new-only	
$\mathbf{x}_n^{(\text{old})}$	$\mathbf{x}_n^{(\text{old})}$	$\mathbf{x}_n^{(\text{old})}$	$\underbrace{0, 0, \dots, 0}_{D\text{-many}}$	$\rangle$
$\mathbf{x}_m^{(\text{new})}$	$\mathbf{x}_m^{(\text{new})}$	$\underbrace{0, 0, \dots, 0}_{D\text{-many}}$	$\mathbf{x}_m^{(\text{new})}$	$\rangle$

# One solution: feature augmentation

- Transform  $D_{old}$  and  $D_{new}$  training examples
- Train a classifier on new representations
- Done!

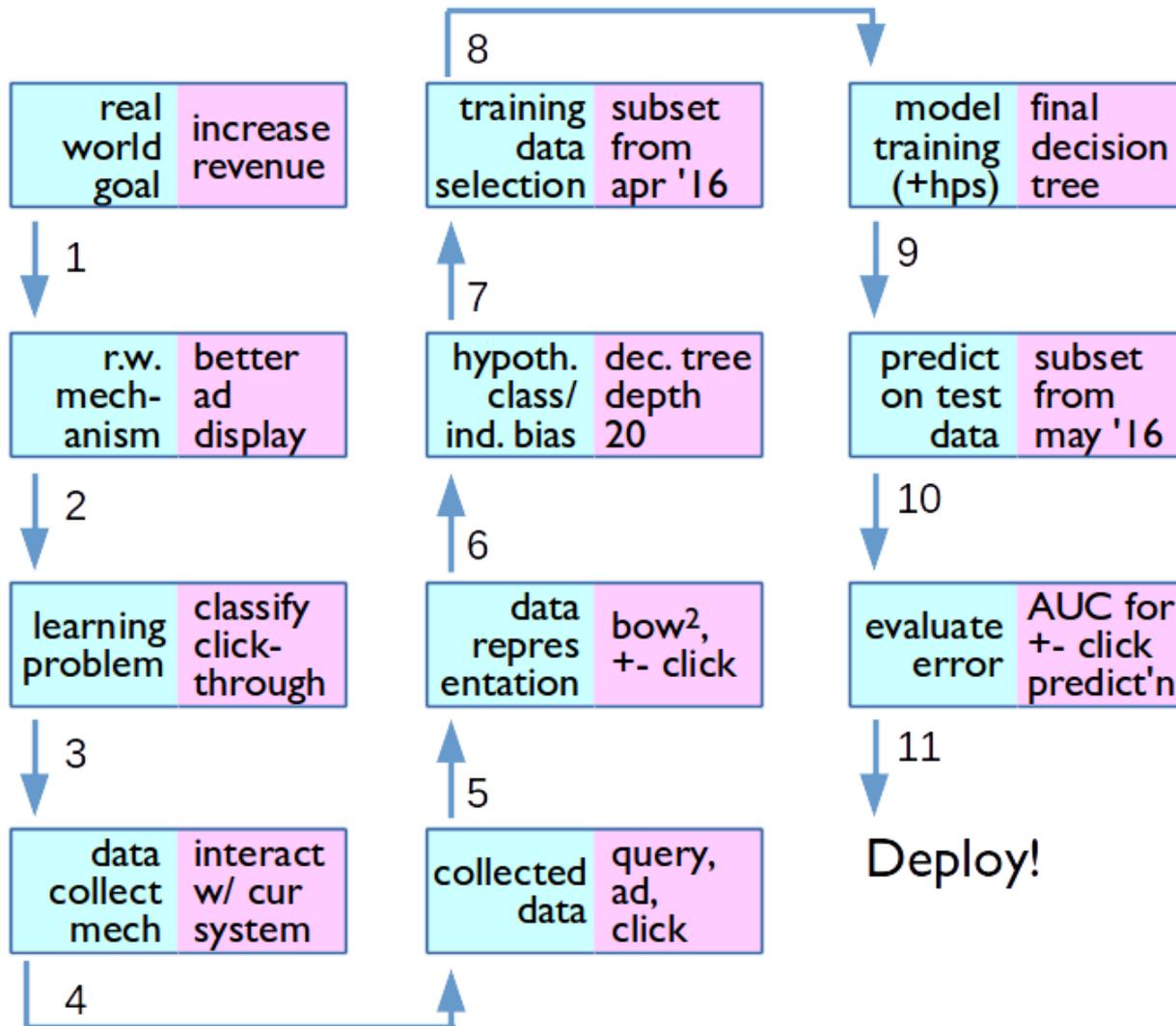
# One solution: feature augmentation

- Adding instance weighting might be useful if  $N \gg M$
- Most effective when distributions are “not too close but not too far”
  - In practice, always try “old only”, “new only”, “union of old and new” as well!

**Theorem 9** (Unsupervised Adaptation Bound). *Given a fixed representation and a fixed hypothesis space  $\mathcal{F}$ , let  $f \in \mathcal{F}$  and let  $\epsilon^{(best)} = \min_{f^* \in \mathcal{F}} \frac{1}{2} [\epsilon^{(old)}(f^*) + \epsilon^{(new)}(f^*)]$ , then, for all  $f \in \mathcal{F}$ :*

$$\underbrace{\epsilon^{(new)}(f)}_{\text{error on } \mathcal{D}^{new}} \leq \underbrace{\epsilon^{(old)}(f)}_{\text{error on } \mathcal{D}^{old}} + \underbrace{\epsilon^{(best)}}_{\text{minimal avg error}} + \underbrace{d_{\mathcal{A}}(\mathcal{D}^{old}, \mathcal{D}^{new})}_{\text{distance}} \quad (8.27)$$

# Typical Design Process for an ML Application



# Bias is pervasive

- Bias in the labeling
- Sample selection bias
- Bias in choice of labels
- Bias in features or model structure
- Bias in loss function
- Deployed systems create feedback loops

# ACM Code of Ethics

“To minimize the possibility of indirectly harming others, computing professionals must minimize malfunctions by following generally accepted standards for system design and testing. Furthermore, it is often necessary to assess the social consequences of systems to project the likelihood of any serious harm to others. If system features are misrepresented to users, coworkers, or supervisors, the individual computing professional is responsible for any resulting injury.”

# Bias and how to deal with it

- Train/test mismatch
- Unsupervised adaptation
- Supervised adaptation