

Slides adapted from Prof Carpuat and Duraiswami



# CMSC 422 Introduction to Machine Learning

## **Lecture 12 Bias and Fairness**

Furong Huang / [furongh@cs.umd.edu](mailto:furongh@cs.umd.edu)

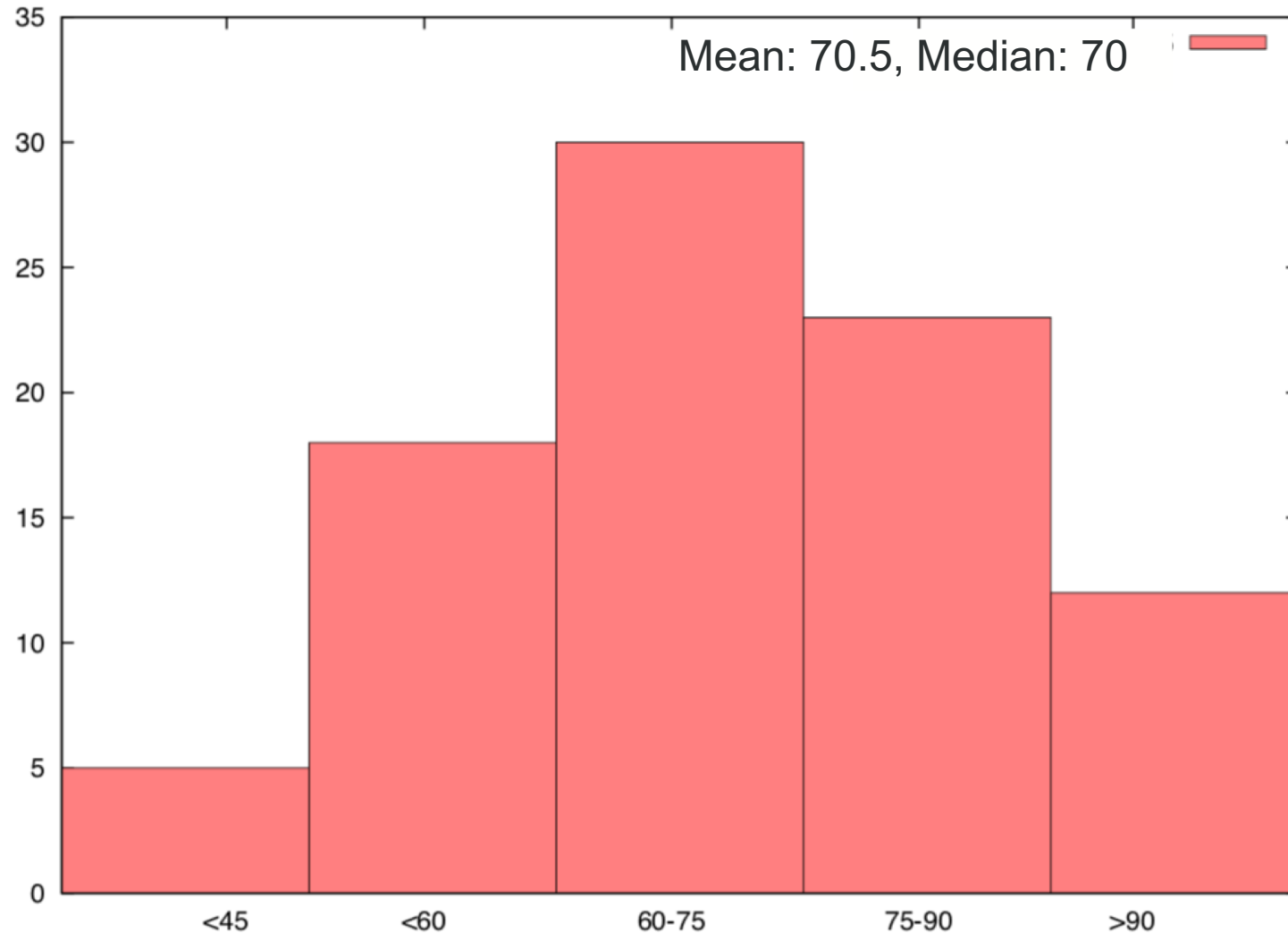


## **Some ML issues in the real world**



UNIVERSITY OF  
MARYLAND

# Midterm – grade distribution (histogram)



# Midterm – makeup policy

---

- Resubmit your answer to question(s)
  - Up to **8 points** to make up
  - **I trust** you, don't seek help from others
  - For T/F problems, give **detailed** justification/proof
  - Grading will be **very strict**, you will get 0 points if the justification is flawed
- Detailed printout of justifications/proofs
  - Your name, your session ID, your UID
  - Deadline: March 15 (Thursday), 11:00 am before class.

# Requirements from my first lecture

---

## ➤ Do the reading before class

- already familiar with high-level concepts and mathematical notation by the time you come to class.
- you can get more out of class time by focusing on understanding the **reasoning** and **clarifying** what didn't make sense when you first read it

## ➤ ***Understanding***

- being able to precisely **define** and **manipulate** the **mathematical** concepts
- being able to discuss the intuition behind algorithms with words is **necessary** but **not sufficient**.

**Why?**



# Final Exam – how to prepare

---

- We will assign more questions that require mathematical reasoning in the homework
- You will read the textbook before coming to class
- You will ask questions during the lecture
- You could form study groups if need to review linear algebraic knowledge, logical reasoning techniques

# Midterm – regrading requests

---

- Written only requests
  - List the problems that should be regraded
  - Suggest the points you should be getting
  - Justify why
- Submit a detailed printout
  - Your name, your session ID, your UID
  - Deadline: March 27, 11:00 am before class

# Ranking

---

Canonical example: web search

Given all the documents on the web

For a user query, retrieve relevant documents, ranked from most relevant to least relevant

---



How can we reduce ranking to binary classification?

# Preference function

---

- Given a query  $q$  and documents  $d_i$  and  $d_j$ , the preference function outputs whether
  - $d_i$  should be preferred to  $d_j$
  - Or  $d_j$  should be preferred to  $d_i$
- That's a binary classification problem!

# Specifying the reduction from ranking to binary classification

---

- How to train classifier that predicts preferences?
- How to turn the predicted preferences into a ranking?

---

**Algorithm 16** NAIVERANKTRAIN(*RankingData*, BINARYTRAIN)

---

```
1:  $\mathbf{D} \leftarrow []$ 
2: for  $n = 1$  to  $N$  do
3:   for all  $i, j = 1$  to  $M$  and  $i \neq j$  do
4:     if  $i$  is preferred to  $j$  on query  $n$  then
5:        $\mathbf{D} \leftarrow \mathbf{D} \oplus (x_{nij}, +1)$ 
6:     else if  $j$  is preferred to  $i$  on query  $n$  then
7:        $\mathbf{D} \leftarrow \mathbf{D} \oplus (x_{nij}, -1)$ 
8:     end if
9:   end for
10: end for
11: return BINARYTRAIN( $\mathbf{D}$ )
```

---

Features associated with comparing document  $i$  and document  $j$  for query  $n$

---

**Algorithm 17** NAIVERANKTEST( $f, \hat{x}$ )

---

```
1:  $score \leftarrow \langle 0, 0, \dots, 0 \rangle$  // initialize  $M$ -many scores to zero
2: for all  $i, j = 1$  to  $M$  and  $i \neq j$  do
3:    $y \leftarrow f(\hat{x}_{ij})$  // get predicted ranking of  $i$  and  $j$ 
4:    $score_i \leftarrow score_i + y$ 
5:    $score_j \leftarrow score_j - y$ 
6: end for
7: return ARGSORT( $score$ ) // return queries sorted by score
```

# Naïve approach

---

Works well for bipartite problems

“is this document relevant or not?”

Not ideal for full ranking problems,  
because

Binary preference problems are not all equally  
important

Separates preference function and sorting



# Improving on naïve approach

## TASK: $\omega$ -RANKING

*Given:*

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

*Compute:* A function  $f : \mathcal{X} \rightarrow \Sigma_M$  minimizing:

$$\mathbb{E}_{(x, \sigma) \sim \mathcal{D}} \left[ \sum_{u \neq v} [\sigma_u < \sigma_v] [\hat{\sigma}_v < \hat{\sigma}_u] \omega(\sigma_u, \sigma_v) \right] \quad (5.7)$$

where  $\hat{\sigma} = f(x)$

# Example of cost functions

---

$$\omega(i, j) = \begin{cases} 1 & \text{if } \min\{i, j\} \leq K \text{ and } i \neq j \\ 0 & \text{otherwise} \end{cases}$$

# Resulting Ranking Algorithms

---

**Algorithm 18**  $\text{RANKTRAIN}(\mathbf{D}^{rank}, \omega, \text{BINARYTRAIN})$

---

```
1:  $\mathbf{D}^{bin} \leftarrow [ ]$ 
2: for all  $(x, \sigma) \in \mathbf{D}^{rank}$  do
3:   for all  $u \neq v$  do
4:      $y \leftarrow \text{SIGN}(\sigma_v - \sigma_u)$                                 // y is +1 if u is preferred to v
5:      $w \leftarrow \omega(\sigma_u, \sigma_v)$                             // w is the cost of misclassification
6:      $\mathbf{D}^{bin} \leftarrow \mathbf{D}^{bin} \oplus (y, w, x_{uv})$ 
7:   end for
8: end for
9: return  $\text{BINARYTRAIN}(\mathbf{D}^{bin})$ 
```

---

---

**Algorithm 19**  $\text{RANKTEST}(f, \hat{x}, \text{obj})$ 

---

**Exercise: understand this offline**

```
1: if  $\text{obj}$  contains 0 or 1 elements then
2:   return  $\text{obj}$ 
3: else
4:    $p \leftarrow$  randomly chosen object in  $\text{obj}$  // pick pivot
5:    $\text{left} \leftarrow [ ]$  // elements that seem smaller than  $p$ 
6:    $\text{right} \leftarrow [ ]$  // elements that seem larger than  $p$ 
7:   for all  $u \in \text{obj} \setminus \{p\}$  do
8:      $\hat{y} \leftarrow f(x_{up})$  // what is the probability that  $u$  precedes  $p$ 
9:     if uniform random variable  $< \hat{y}$  then
10:       $\text{left} \leftarrow \text{left} \oplus u$ 
11:     else
12:       $\text{right} \leftarrow \text{right} \oplus u$ 
13:     end if
14:   end for
15:    $\text{left} \leftarrow \text{RANKTEST}(f, \hat{x}, \text{left})$  // sort earlier elements
16:    $\text{right} \leftarrow \text{RANKTEST}(f, \hat{x}, \text{right})$  // sort later elements
17:   return  $\text{left} \oplus \langle p \rangle \oplus \text{right}$ 
18: end if
```

---

# RankTest

---

A probabilistic version of the quicksort algorithm

Only  $O(M \log_2 M)$  calls to  $f$  in expectation

Better error bound than naïve algorithm  
(see CIML for theorem)

# What you should know

---

- What are reductions and why they are useful
- Implement, analyze and prove error bounds of algorithms for
  - Weighted binary classification
  - Multiclass classification (OVA, AVA)
- Understand algorithms for
  - $\omega$  –ranking



## **Bias and Fairness**



UNIVERSITY OF  
MARYLAND

# Word Embeddings Could be Biased

---

Man is to Computer Programmer as Woman is to Homemaker?  
Debiasing Word Embeddings

## Extreme *she* occupations

- |                 |                       |                        |
|-----------------|-----------------------|------------------------|
| 1. homemaker    | 2. nurse              | 3. receptionist        |
| 4. librarian    | 5. socialite          | 6. hairdresser         |
| 7. nanny        | 8. bookkeeper         | 9. stylist             |
| 10. housekeeper | 11. interior designer | 12. guidance counselor |

## Extreme *he* occupations

- |                |                   |                |
|----------------|-------------------|----------------|
| 1. maestro     | 2. skipper        | 3. protege     |
| 4. philosopher | 5. captain        | 6. architect   |
| 7. financier   | 8. warrior        | 9. broadcaster |
| 10. magician   | 11. fighter pilot | 12. boss       |

Figure 1: The most extreme occupations as projected on to the *she*–*he* gender direction on g2vNEWS. Occupations such as *businesswoman*, where gender is suggested by the orthography, were excluded.



# 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.)*

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

# 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 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

---

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

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

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

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

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

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

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

Fixed base  
distribution

$S$  = selection  
variable

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

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

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

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

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

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

# Clarifications

---

- Note  $Z^{\text{new}}$  is the same for all the data points (i.e.,  $\forall(x, y)$ ). It is also a normalization to make sure  $\sum_{(x,y)} \mathcal{D}^{\text{new}}(x, y) = 1$ .

$$Z^{\text{new}} = \sum_{(x,y)} [\mathcal{D}^{\text{base}}(x, y)p(s = 0|x)]$$

- Therefore,  $Z^{\text{new}}$  is not a function of  $(x, y)$ , in fact it is a constant.
- Similarly for  $Z^{\text{old}}$ .
- All examples are drawn from fixed base distribution, some are selected to go into  $\mathcal{D}^{\text{new}}$ , some are selected to go into  $\mathcal{D}^{\text{old}}$  according to a selection variable  $s$ .



**Algorithm 23** SELECTIONADAPTATION( $\langle (x_n, y_n) \rangle_{n=1}^N, \langle z_m \rangle_{m=1}^M, \mathcal{A}$ )

```

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:  $D^{weighted} \leftarrow \left\langle (\mathbf{x}_n, y_n, \frac{1}{\hat{p}(\mathbf{x}_n)} - 1) \right\rangle_{n=1}^N$  // assemble weight classification
// data using selector

4: return  $\mathcal{A}(D^{weighted})$  // train classifier

```

# We will revisit logistic regression!

# 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
$\mathbf{x}_n^{(\text{old})}$	$\mapsto \left\langle \mathbf{x}_n^{(\text{old})} \right.$	$, \mathbf{x}_n^{(\text{old})}$	$, \underbrace{0, 0, \dots, 0}_{D\text{-many}} \rangle$
$\mathbf{x}_m^{(\text{new})}$	$\mapsto \left\langle \mathbf{x}_m^{(\text{new})} \right.$	$, \underbrace{0, 0, \dots, 0}_{D\text{-many}}$	$, \mathbf{x}_m^{(\text{new})} \rangle$

# One solution: feature augmentation

---

- Transform  $D_{\text{old}}$  and  $D_{\text{new}}$  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  $f \in \mathcal{F}$ :*

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

# 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

# Bias and how to deal with it

---

- Train/test mismatch
- Unsupervised adaptation
- Supervised adaptation



# 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.”

<https://www.acm.org/about-acm/acm-code-of-ethics-and-professional-conduct>



UNIVERSITY OF  
MARYLAND

**Furong Huang**

3251 A.V. Williams, College Park, MD 20740

301.405.8010 / [furongh@cs.umd.edu](mailto:furongh@cs.umd.edu)