

# Calibration and bias

COMS 4771 Fall 2025

**Predicting conditional probabilities**

### Example #1: Click prediction for online ads

- ▶  $X$  = features of (user, advertisement) pair
- ▶  $Y$  = indicator that user will click on ad
- ▶  $\Pr(Y = 1 \mid X = x)$  is almost always near zero, but useful to know this probability, e.g., to compare ads, estimate revenue

### Example #2: Selective prediction

- ▶ If  $\Pr(Y = 1 \mid X = x) \approx \Pr(Y = 0 \mid X = x)$ , then perhaps should abstain from making a prediction

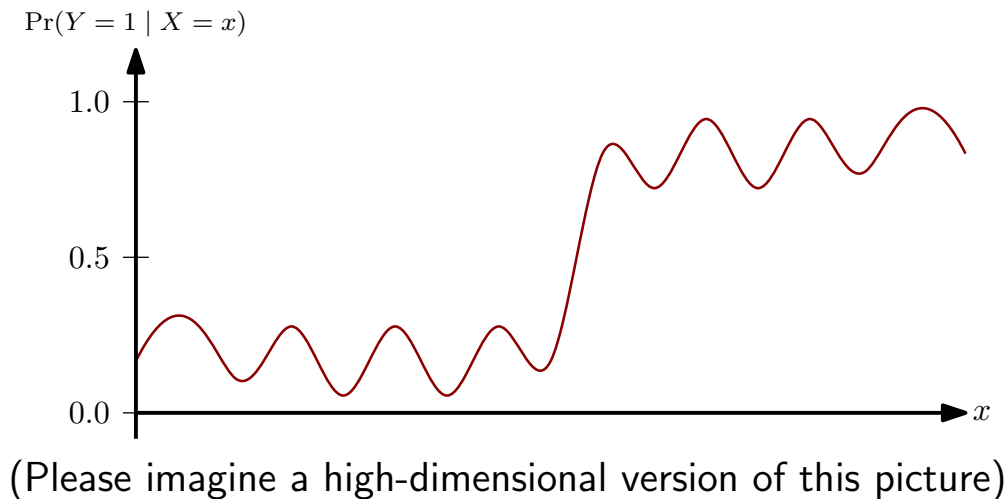
1 / 30

	Estimates $\Pr(Y = 1 \mid X = x)$
nearest neighbors	?
decision trees	?
generative models	✓
logistic regression	✓
Perceptron	no

2 / 30

### Caution:

- ▶ Prediction/estimate of (conditional) probability is still a prediction
  - ▶ Some are accurate, some are inaccurate
  - ▶ Same goes for anything derived from these predictions
- ▶ At least as hard as learning to classify, and can be arbitrarily harder



3 / 30

Ultimately, need to validate accuracy of predictions of (conditional) probabilities

- ▶ **Indirect measures:** logarithmic loss (a.k.a. cross entropy loss), squared error, ...
- ▶ Recall: for binary classification, function  $\hat{p}: \mathcal{X} \rightarrow [0, 1]$  that minimizes  $\mathbb{E}[\text{loss}_{\log}(\hat{p}(X), Y)]$  is

$$\hat{p}(x) = \Pr(Y = 1 \mid X = x)$$

(where  $\text{loss}_{\log}$  is logarithmic loss)

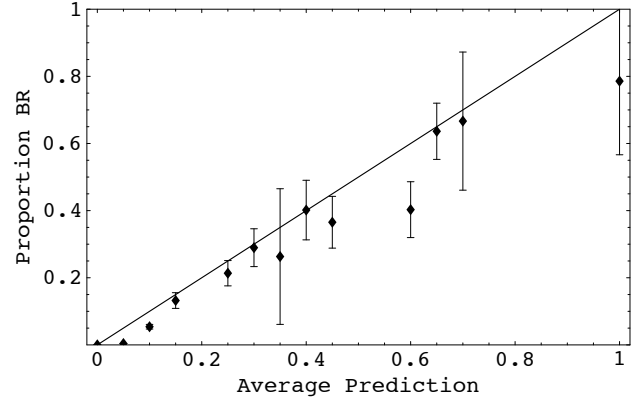
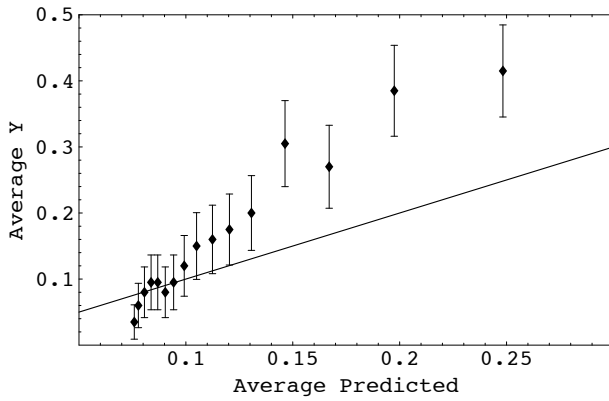
- ▶ Average case measure of the quality of predictions
- ▶ But how accurate is any individual prediction?

4 / 30

## Calibration

Prediction  $\hat{p}(x)$  of  $\Pr(Y = 1 \mid X = x)$  is (approximately) calibrated if

$$\Pr(Y = 1 \mid \hat{p}(X) = p) \approx p \quad \text{for all } p \in [0, 1]$$



(Foster & Stine, 2004, "Building a Predictive Model for Bankruptcy")

6 / 30

Expected calibration error of  $\hat{p}$  (assuming  $\text{range}(\hat{p})$  is finite set  $\mathcal{P} \subset [0, 1]$ ):

$$\sum_{p \in \mathcal{P}} |\Pr(Y = 1 \wedge \hat{p}(X) = p) - p \times \Pr(\hat{p}(X) = p)|$$

Possible to estimate this from test data if  $\mathcal{P}$  is not too large:

$$\sum_{p \in \mathcal{P}} \frac{1}{n} \left| \sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1 \wedge \hat{p}(x^{(i)}) = p\} - p \sum_{i=1}^n \mathbb{1}\{\hat{p}(x^{(i)}) = p\} \right|$$

7 / 30

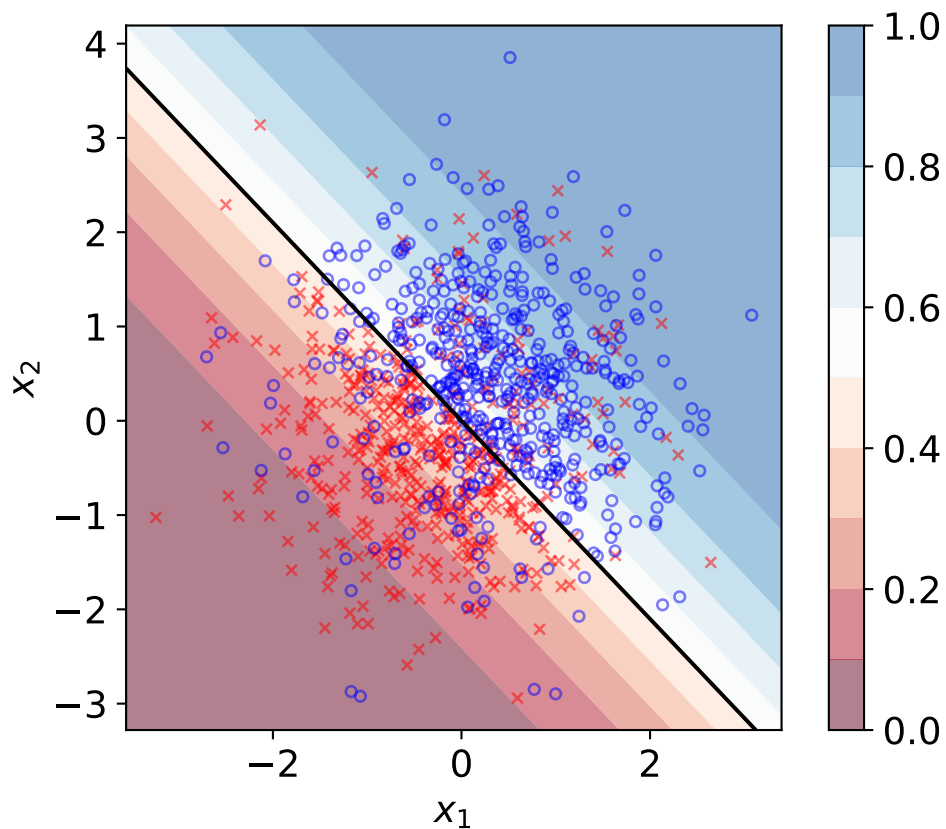
Synthetic example:  $X = (X_1, X_2) \sim N(0, I)$ , and

$$\Pr(Y = 1 \mid X = x) = p^*(x) = \begin{cases} 0.8 & \text{if } x_1 + x_2 > 0 \\ 0.2 & \text{otherwise} \end{cases}$$

Fit logistic regression model to 1000 training examples using MLE

- Error rate is 20.3%, which is nearly optimal
- However, expected calibration error of  $\hat{p}$  is 0.13

8 / 30



9 / 30

## Calibrating conditional probability predictions

Suppose you have real-valued “score” function  $s: \mathbb{R}^d \rightarrow \mathbb{R}$

	Possible score $s(x)$
$k$ -nearest neighbors	
decision trees	
generative models	est. of $\Pr(Y = 1 \mid X = x)$
logistic regression	est. of $\Pr(Y = 1 \mid X = x)$
Perceptron	
(many other possibilities)	

**Goal:** obtain approximately calibrated predictor  $\hat{p}(x)$  of  $\Pr(Y = 1 \mid X = x)$

### (Histogram) binning:

- ▶ Sort  $s(x)$  from training/validation data into  $T$  bins
- ▶ Determine  $T - 1$  boundary values between the bins
- ▶ Let  $\hat{p}^{(i)}$  be estimate of  $\Pr(Y = 1 \mid s(x) \in \text{bin } i)$
- ▶ Then define

$$\hat{p}(x) = \begin{cases} \hat{p}^{(1)} & \text{if } s(x) \text{ falls in bin 1} \\ \hat{p}^{(2)} & \text{if } s(x) \text{ falls in bin 2} \\ \vdots & \\ \hat{p}^{(T)} & \text{if } s(x) \text{ falls in bin } T \end{cases}$$

11 / 30

### **How can this possibly work?**

- ▶ Key idea: score function turns problem into one with only a single feature
- ▶ No curse of dimension to worry about

12 / 30



Synthetic example:  $X = (X_1, X_2) \sim N(0, I)$ , and

$$\Pr(Y = 1 \mid X = x) = p^*(x) = \begin{cases} 0.8 & \text{if } x_1 + x_2 > 0 \\ 0.2 & \text{otherwise} \end{cases}$$

Fit logistic regression model to 1000 training examples using MLE

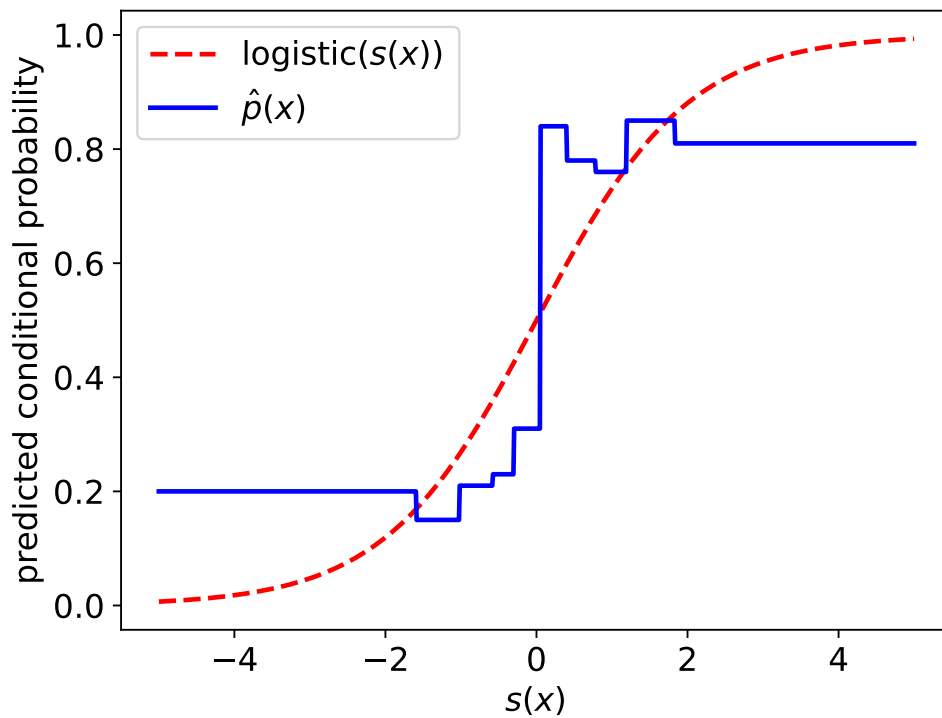
- **Apply binning to**  $s(x) = \hat{w}^\top x$  (with  $T = 10$  bins)
- Expected calibration error: 0.043 (down from 0.13)

13 / 30

Final predictor  $\hat{p}(x)$ :

range of $s(x)$	$\hat{p}(x)$
$s(x) < -1.591$	0.200
$-1.591 \leq s(x) < -1.024$	0.150
$-1.024 \leq s(x) < -0.578$	0.210
$-0.578 \leq s(x) < -0.296$	0.230
$-0.296 \leq s(x) < 0.055$	0.310
$0.055 \leq s(x) < 0.398$	0.840
$0.398 \leq s(x) < 0.777$	0.780
$0.777 \leq s(x) < 1.194$	0.760
$1.194 \leq s(x) < 1.835$	0.850
$1.835 \leq s(x)$	0.810

14 / 30



15 / 30

- Popular way to improve binning: enforce monotonicity (e.g., if you believe  $\Pr(Y = 1 \mid s(x))$  is monotone in  $s(x)$ )
- Caution: a  $\hat{p}$  with low expected calibration error does not necessarily give an accurate predict of  $Y$  from  $X$ 
  - Only gives an accurate predictor of  $Y$  from  $s(X)$
  - But perhaps  $s(X)$  is constant!
  - In this case, suffices to predict the constant  $\Pr(Y = 1)$

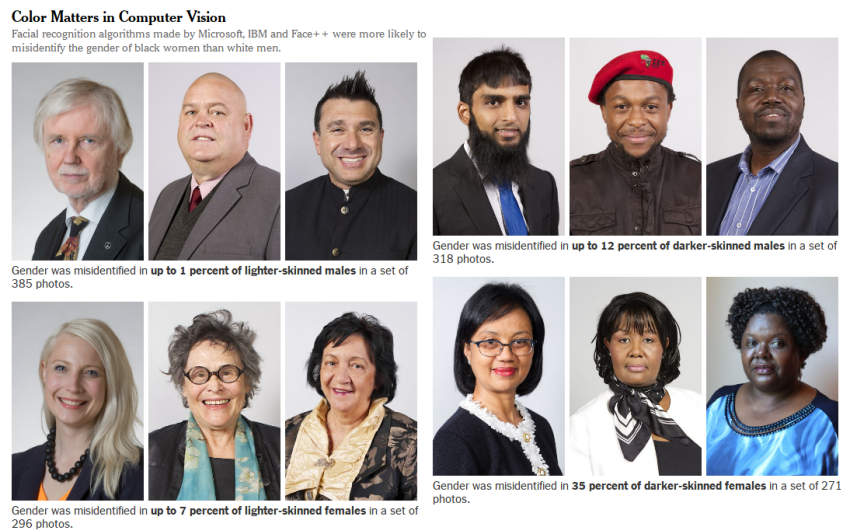
16 / 30

## Calibration versus equalizing error rates

- ▶ Increasing use of predictive models in real-world applications (e.g., admissions, hiring, criminal justice)
- ▶ Do they offer “fair treatment” to individuals/groups?

Well-known example: **“Gender shades”** study (Buolamwini and Gebru, 2018)

- ▶ **Task:** predict gender from image of face
- ▶ **Major finding:** some commercial facial analysis software were less accurate for images of darker-skinned female individuals than for images of lighter-skinned male individuals



18 / 30

ProPublica **“Machine Bias”** study (Angwin et al, 2016)

- ▶ Judge needs to decide whether or not an arrested defendant should be released while awaiting trial
- ▶ Predictive model (“COMPAS”) predicts whether or not defendant will commit (violent) crime if released
- ▶ Study based data from Broward County, Florida argued that COMPAS treated black defendants unfairly in a certain sense

19 / 30

Setup for ProPublica study (highly simplified)

- ▶  $X$ : feature vector specific to arrested defendant
- ▶  $A$ : group membership attribute (e.g., race, sex, age; could be part of  $X$ )
- ▶  $Y$ : outcome to predict (e.g., “will re-offend if released”)
- ▶  $\hat{Y} = f_{\text{COMPAS}}(X)$ : prediction of  $Y$  based on  $X$
- ▶ For simplicity, assume  $A, Y, \hat{Y}$  are all  $\{0, 1\}$ -valued

20 / 30

### Types of errors:

- ▶ False positive rate:  $\text{FPR} = \Pr(\hat{Y} = 1 \mid Y = 0)$
- ▶ False negative rate:  $\text{FNR} = \Pr(\hat{Y} = 0 \mid Y = 1)$
- ▶ Per-group FPR and FNR: for each  $a \in \{0, 1\}$ ,

$$\text{FPR}_a = \Pr(\hat{Y} = 1 \mid Y = 0, A = a)$$

$$\text{FNR}_a = \Pr(\hat{Y} = 0 \mid Y = 1, A = a)$$

Equalized odds: require that  $\text{FPR}_0 \approx \text{FPR}_1$  and  $\text{FNR}_0 \approx \text{FNR}_1$

- ▶ No group incurs errors (either type) at a higher rate than the other

21 / 30

**ProPublica found:** COMPAS software is very far from offering “equalized odds”

- ▶  $\text{FPR}_0 = 45\%$ ,  $\text{FPR}_1 = 23\%$
- ▶  $\text{FNR}_0 = 27\%$ ,  $\text{FNR}_1 = 48\%$

22 / 30

**Response from Northpointe** (creator of COMPAS)

- ▶  $f_{\text{COMPAS}}(x) = \mathbb{1}\{\hat{p}(x) > t\}$  where  $\hat{p}(x)$  is prediction of  $\Pr(Y = 1 \mid X = x)$ , and  $t$  is some suitable threshold parameter
- ▶  $\hat{p}$  approximately-calibrated, and also approximately-calibrated **for each group**

$$\Pr(Y = 1 \mid \hat{p}(X) = p, A = 0) \approx \Pr(Y = 1 \mid \hat{p}(X) = p, A = 1) \approx p$$

- ▶ So  $\hat{p}$  has same probabilistic semantics for each group

23 / 30

**Theorem** (Chouldechova; Kleinberg-Mullainathan-Raghavan): Unless

$$\Pr(Y = 1 \mid A = 0) = \Pr(Y = 1 \mid A = 1) \quad \text{or} \\ \text{FPR}_0 = \text{FPR}_1 = \text{FNR}_0 = \text{FNR}_1 = 0,$$

it is impossible to simultaneously satisfy all of the following:

1.  $\text{FPR}_0 = \text{FPR}_1$
2.  $\text{FNR}_0 = \text{FNR}_1$
3.  $\hat{Y}$  is marginally calibrated for group  $A = 0$
4.  $\hat{Y}$  is marginally calibrated for group  $A = 1$

(Marginal calibration for group  $a$ :  $\Pr(\hat{Y} = 1 \mid A = a) = \Pr(Y = 1 \mid A = a)$ )

## Distribution shift

Distribution shift (a.k.a. train/test mismatch, sample selection bias):

- ▶ Training data is sample from source distribution
- ▶ Care about (average) performance on data from target distribution
- ▶ Distribution shift: source  $\neq$  target

**Example:** care about applying facial analysis software to images from general US population, but only train on images of light-skinned males

- ▶ Hardly any reason to expect things to work well . . .
- ▶ . . . unless you are “testing” only on images of light-skinned males

#### Color Matters in Computer Vision

Facial recognition algorithms made by Microsoft, IBM and Face++ were more likely to misidentify the gender of black women than white men.



Gender was misidentified in **up to 1 percent** of **lighter-skinned males** in a set of 385 photos.



Gender was misidentified in **up to 12 percent** of **darker-skinned males** in a set of 318 photos.



Gender was misidentified in **up to 7 percent** of **lighter-skinned females** in a set of 296 photos.



Gender was misidentified in **35 percent** of **darker-skinned females** in a set of 271 photos.



In many applications, training data is “dataset of convenience”

- ▶ Use whatever data you can get

All methods for addressing distribution shift require

- ▶ Either a lot of domain knowledge,
- ▶ Or additional data from target distribution
- ▶ (Often need both)

27 / 30

**Example:** re-weighting data

- ▶ Suppose you notice that, in training data,

$$\Pr(A = 0) \ll \Pr(A = 1)$$

But you know that in target distribution,  $A = 0$  and  $A = 1$  equally often

- ▶ Use an importance weight of

$$\frac{1}{2 \Pr(A = a)}$$

for every example with  $A = a$  in (empirical) expectation computations

- ▶ **Critical assumption:** conditional distribution of  $(X, Y)$  given  $A$  is the same in source and target; only marginal distribution of  $A$  differs

28 / 30

## Importance-weighted test error rate

- ▶ Test data  $(\tilde{X}^{(1)}, \tilde{Y}^{(1)}, \tilde{A}^{(1)}), \dots, (\tilde{X}^{(m)}, \tilde{Y}^{(m)}, \tilde{A}^{(m)}) \stackrel{\text{i.i.d.}}{\sim} (X, Y, A)$ , from source distribution
- ▶ Define  $p_a = \Pr(A = a)$  for each  $a \in \{0, 1\}$
- ▶ Weighted test error rate:

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{f(\tilde{X}^{(i)}) \neq \tilde{Y}^{(i)}\} \times \frac{1}{2p_{\tilde{A}^{(i)}}}$$

29 / 30

## Expected value of importance-weighted test error rate:

$$\mathbb{E} \left[ \mathbb{1}\{f(X) \neq Y\} \times \frac{1}{2p_A} \right] =$$

---

30 / 30