

# Assignment Project Exam Help

Machine learning lecture slides

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

**Classification III: Classification objectives**

**Assignment Project Exam Help**

**<https://eduassistpro.github.io>**

**Add WeChat edu\_assist\_pr**

- ▶ Scoring functions
- ▶ Cost-sensitive classification
- ▶ Conditional probability estimation
- ▶ Reducing multi-class to binary
- ▶

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Scoring functions in general

- ▶ Statistical model:  $(X, Y) \sim P$  for distribution  $P$  over  $\mathcal{X} \times \{-1, +1\}$

- ▶ Binary classifiers are generally of the form

$$x \mapsto \text{sign}(h(x))$$

<https://eduassistpro.github.io>

where  $\eta(x) = \Pr(Y = +1 \mid X = x)$

- ▶ Use with loss functions like  $\ell_{0/1}$ ,

$$\mathcal{R}(h) = \mathbb{E}[\ell(h(x), Y)]$$

- ▶ Issues to consider:

- ▶ Different types of mistakes have different costs
- ▶ How to get  $\Pr(Y = +1 \mid X = x)$  from  $h(x)$ ?
- ▶ More than two classes

## Cost-sensitive classification

- Cost matrix for different kinds of mistakes (for  $c \in [0, 1]$ )

	$\hat{y} = -1$	$\hat{y} = +1$
$y = -1$	0	$c$
$y = +1$	1	0

- <https://eduassistpro.github.io>

$$\ell^{(c)}(y, \hat{y}) = (\mathbf{1}_{\{y=+1\}}) \cdot (1 - c)$$

Add WeChat: edu\_assist\_pro

- If  $\ell$  is convex in  $\hat{y}$ , then so is  $\ell^{(c)}$
- Cost-sensitive (empirical) risk:

$$\mathcal{R}^{(c)}(h) := \mathbb{E}[\ell^{(c)}(Y, h(X))]$$

$$\hat{\mathcal{R}}^{(c)}(h) := \frac{1}{n} \sum_{i=1}^n \ell^{(c)}(y_i, h(x_i))$$

## Minimizing cost-sensitive risk

- ▶ What is the analogue of Bayes classifier for cost-sensitive (zero-one loss) risk?

- ▶ Let  $\eta(x) = \Pr(Y=1 | X=x)$
- ▶ Fix  $x$ ; what is conditional cost-sensitive risk of predicting  $\hat{y}$ ?

- ▶ <https://eduassistpro.github.io>

$\hat{y} = \begin{cases} +1 & \text{if } \eta(x) \cdot (1 - c) > c \\ -1 & \text{otherwise} \end{cases}$

Add WeChat edu\_assist\_pro

- ▶ So use scoring function  $h(x) = \eta(x) - c$ 
  - ▶ Equivalently, use  $\eta$  as scoring function, but threshold at  $c$  instead of  $1/2$
- ▶ Where does  $c$  come from?

## Example: balanced error rate

- Balanced error rate:  $\text{BER} := \frac{1}{2}\text{FNR} + \frac{1}{2}\text{FPR}$
- Which cost sensitive risk to try to minimize?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

where  $\pi = \Pr(Y = +1)$ .

- Therefore, we want to use the following cost

	$\hat{y} = -$	$\hat{y} = +$
$y = -1$	0	$\frac{1}{1-\pi}$
$y = +1$	$\frac{1}{\pi}$	0

- This corresponds to  $c = \pi$ .

## Importance-weighted risk

- ▶ Perhaps the world tells you how important each example is
- ▶ Statistical model:  $(X, Y, W) \sim P$ 
  - ▶  $W$  is (non-negative) importance weight of example  $(X, Y)$
- ▶ Importance-weighted  $\ell$ -risk of  $h$ :

- ▶ <https://eduassistpro.github.io>

Add WeChat  $\frac{1}{n} \sum_{i=1}^n w_i \cdot \ell(h(x_i), y_i)$  edu\_assist\_pro



## Conditional probability estimation (1)

- ▶ How to get estimate of  $\eta(x) = \Pr(Y = +1 \mid X = x)$ ?
- ▶ Useful if want to know expected cost of a prediction

# Assignment Project Exam Help

$$\mathbb{E}[\ell_{0/1}^{(c)}(Yh(X)) \mid X = x] = \begin{cases} (1 - c) \cdot \eta(x) & \text{if } h(x) \leq 0 \\ c \cdot \eta(x) & \text{if } h(x) > 0 \end{cases}$$

- ▶ <https://eduassistpro.github.io>

$$h(x) = 2\eta(x) - 1$$

Add WeChat edu\_assist\_pro

- ▶ Therefore, given  $h$ , can estimate  $\eta$
- ▶ Recipe:
  - ▶ Find scoring function  $h$  that (approximately) minimizes (empirical) squared loss risk
  - ▶ Construct conditional probability estimate  $\hat{\eta}$  using above formula

## Conditional probability estimation (2)

- ▶ Similar strategy available for logistic loss
- ▶ But not for hinge loss!

Assignment Project Exam Help

- ▶ Hinge loss risk is minimized by  $\hat{h}(x) = \text{sgn}(2\hat{\eta}(x) - 1)$
- ▶ Cannot recover  $\eta$  from  $\hat{h}$
- ▶  $h$  (e.g.,

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Application: Reducing multi-class to binary

- ▶ Multi-class: Conditional probability function is vector-valued function

Assignment Project Exam Help

$$\eta(x) = \begin{bmatrix} \Pr(Y = 1 | X = x) \\ \vdots \end{bmatrix}$$

- ▶ <https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

$$\eta_k(x) = \Pr(Y =$$

- ▶ This can be done by creating  $I_k$  where in problem  $k$ , label is  $1_{\{y\}}$
- ▶ Given the  $K$  learned conditional probability functions  $\hat{\eta}_1, \dots, \hat{\eta}_K$ , we form a final predictor  $\hat{f}$

$$\hat{f}(x) = \arg \max_{k=1, \dots, K} \hat{\eta}_k(x).$$

## When does one-against-all work well?

- ▶ If learned conditional probability functions  $\hat{\eta}_k$  are accurate, then behavior of one-against-all classifier  $\hat{f}$  is similar to optimal classifier

# Assignment Project Exam Help

$$f^*(x) = \arg \max_k \Pr(Y = k \mid X = x).$$

- ▶ <https://eduassistpro.github.io>

$$\text{err}(\hat{f}) \leq \text{err}(f^*) + 2 \cdot \mathbb{E}[m]$$

Add WeChat edu\_assist\_pr

- ▶ Use of predictive models (e.g., in admissions, hiring, criminal justice) has raised concerns about whether they offer “fair treatment” to individuals and/or groups
- ▶ We will focus on group-based fairness

Assignment Project Exam Help

d

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Disparate treatment

- ▶ Often predictive models work better for some groups than for others

- ▶ Example: face recognition (Buolamwini and Gebru, 2018; Lohr, 2018)

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Possible causes of unfairness

- ▶ People deliberately being unfair
- ▶ Disparity in number of available training data for different groups
- ▶ Disparity in usefulness of available features for different groups
- ▶
- ▶

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

- ▶ ProPublica (investigative journalism group) studied a particular predictive model being used to determine “pre-trial detention”
  - ▶ Angwin et al., 2016
  - ▶ Judge needs to decide whether or not an arrested defendant should be released while awaiting trial

<https://eduassistpro.github.io/>

- ▶ Study argued that COMPAS treated black in a certain sense
  - ▶ What sense? How do they make this argu



- ▶ Setup:

- ▶  $X$ : features for individual

- ▶  $A$ : group membership attribute (e.g., race, sex, age, religion)

- ▶  $Y$ : outcome variable to predict (e.g., "will repay loan", "will re-offend")

$A$ ))

- ▶ <https://eduassistpro.github.io>

$(A, Y,$

- ▶ Add WeChat [edu\\_assist\\_pr](#)

Caveat: Often, we don't have access to

## Classification parity

- Fairness criterion: Classification parity

Assignment Project Exam Help

- Sounds reasonable, but easy to satisfy with perverse methods



<https://eduassistpro.github.io>

$(A = 0)$		$\hat{Y} = 0$	$\hat{Y} = 1$		
$Y = 0$		1/2	0		
$Y = 1$		0	1/2		

- For  $A = 0$  people, correctly give loans to people who will repay
- For  $A = 1$  people, give loans randomly (Bernoulli(1/2))
- Satisfies criterion, but bad for  $A = 1$  people

## Equalized odds (1)

- Fairness criterion: Equalized odds

$\Pr(\hat{Y} = 1 | Y = y, A = 0) \approx \Pr(\hat{Y} = 1 | Y = y, A = 1)$   
for both  $y \in \{0, 1\}$ .



<https://eduassistpro.github.io>



Add WeChat edu\_assist\_pro

$(A=0)$		$(A=1)$	
$Y=0$	$Y=1$	$Y=0$	$Y=1$
1/2	0	1/4	1/4

E.g.,  $A = 0$  group has 0% FPR, while  $A = 1$  has 50% FPR.

- Criteria imply constraints on the classifier / scoring function
  - Can try to enforce constraint during training

## Equalized odds (2)

- ▶ ProPublica study:

- ▶ Found that FPR for  $A = 0$  group (black defendants; 45%) was higher than FPR for  $A = 1$  group (white defendants; 23%)

$(A = 0) \parallel \hat{Y} = 0 \mid \hat{Y} = 1$			$(A = 1) \parallel \hat{Y} = 0 \mid \hat{Y} = 1$		
					0.14
					0.21

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr