

# Assignment Project Exam Help

Machine learning lecture slides

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

Prediction theory

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

- ▶ Statistical model for binary outcomes
- ▶ Plug-in principle and IID model
- ▶ Maximum likelihood estimation
- ▶ Statistical model for binary classification
- ▶
- ▶
- ▶

Assignment Project Exam Help

▶ <https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Statistical model for binary outcomes

- ▶ Example: coin toss
- ▶ Physical model: hard
- ▶ Statistical model: outcome is random
  - ▶ Bernoulli distribution with heads probability  $\theta \in [0, 1]$

Assignment Project Exam Help

<https://eduassistpro.github.io>

- ▶ Goal: correctly predict outcome

Add WeChat edu\_assist\_pr

## Optimal prediction

- ▶ Suppose  $Y \sim \text{Bernoulli}(\theta)$ .
- ▶ Suppose  $\theta$  known.
- ▶ Optimal prediction:

Assignment Project Exam Help

$$\mathbf{1}_{\{\theta > 1/2\}}$$

<https://eduassistpro.github.io>

- ▶ The optimal prediction is incorrect with

Add WeChat edu\_assist\_pr

$$\max_{\hat{y}} \mathbb{E}[\hat{y} - Y]$$

## Learning to make predictions

- ▶ If  $\theta$  unknown:

- ▶ Assume we have data: outcomes of previous coin tosses

- ▶ Data should be related to what we want to predict: same coin is being tossed

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Plug-in principle and IID model

- ▶ Plug-in principle:

- ▶ Estimate unknown(s) based on data (e.g.,  $\theta$ )
- ▶ Plug estimates into formula for optimal prediction

- ▶ <https://eduassistpro.github.io>

- ▶ IID model: Observations & (unseen) observations of  $m$  variables

- ▶ iid: independent and identically distributed

- ▶ Crucial modeling assumption that makes

- ▶ When is the IID assumption not reasonable? ...

- ▶ Parametric statistical model  $\{P_\theta : \theta \in \Theta\}$ 
  - ▶ collection of parameterized probability distributions for data
  - ▶  $\Theta$  is the parameter space
  - ▶ One distribution per parameter value  $\theta \in \Theta$
- ▶

<https://eduassistpro.github.io> (pmf)

for the distribution.

- ▶ What is formula for  $P_\theta(y_1, \dots, y_n)$

Add WeChat edu\_assist\_pro



## Maximum likelihood estimation (1)

- ▶ Likelihood of parameter  $\theta$  (given observed data)

- ▶  $L(\theta) = P_{\theta}(y_1, \dots, y_n)$

- ▶ Maximum likelihood estimator

- ▶ Choose  $\theta$  with highest likelihood



<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Maximum likelihood estimation (2)

- ▶ Coin toss example
  - ▶ Log-likelihood

Assignment Project Exam Help

$$\ln L(\theta) = \sum_i y_i \ln \theta + (1 - y_i) \ln(1 - \theta)$$

<https://eduassistpro.github.io>

Add WeChat  $\hat{\theta}_{MLE} := -$  edu\_assist\_pr

## Back to plug-in principle

- ▶ We are given data  $y_1, \dots, y_n \in \{0, 1\}^n$ , which we model using the IID model from before
- ▶ Obtain estimate  $\hat{\theta}_{\text{MLE}}$  of known  $\theta$  based on  $y_1, \dots, y_n$
- ▶ Plug-in  $\hat{\theta}_{\text{MLE}}$  for  $\theta$  in formula for optimal prediction:

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Analysis of the plug-in prediction (1)

- ▶ How good is the plug-in prediction?
  - ▶ Study behavior under the IID model, where
$$Y_1, \dots, Y_n, Y \sim_{\text{iid}} \text{Bernoulli}(\theta).$$
    - ▶  $Y_1, \dots, Y_n$  are the data we collected
    - ▶  $Y$  is the outcome to predict

<https://eduassistpro.github.io>

worse.

Add WeChat edu\_assist\_pr

## Analysis of the plug-in prediction (2)

► **Theorem:**

$$\Pr(\hat{Y} \neq Y) \leq \min\{\theta, 1 - \theta\} + \frac{1}{2} \cdot |\theta - 0.5| \cdot e^{-2n(\theta - 0.5)^2}.$$

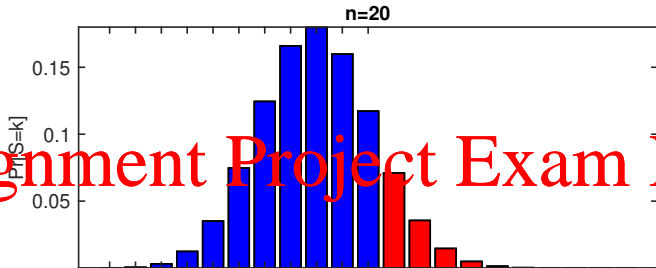
► The first term is the optimal error probability.

► The second term comes from the probability that the  $\hat{\theta}_{\text{MLE}}$  is on the opposite side of  $1/2$  as  $\theta$ .

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

# Assignment Project Exam Help



<https://eduassistpro.github.io>

Figure 1:  $\Pr(S > n/2)$  for  $S \sim$

Add WeChat edu\_assist\_pr

# Assignment Project Exam Help

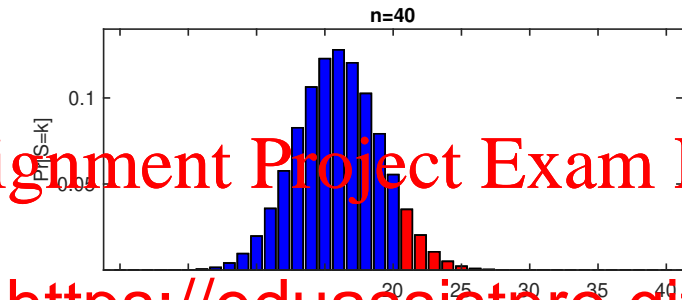
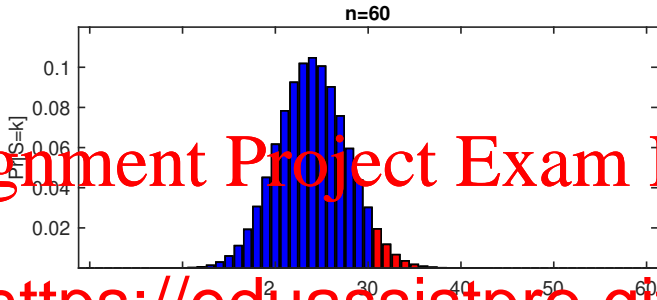


Figure 2:  $\Pr(S > n/2)$  for  $S \sim$

Add WeChat edu\_assist\_pr

Assignment Project Exam Help



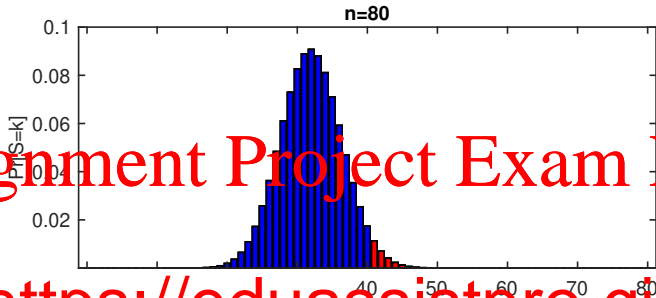
<https://eduassistpro.github.io>

Figure 3:  $\Pr(S > n/2)$  for  $S \sim$

Add WeChat edu\_assist\_pr



Assignment Project Exam Help



<https://eduassistpro.github.io>

Figure 4:  $\Pr(S > n/2)$  for  $S \sim$

Add WeChat edu\_assist\_pr

# Statistical model for labeled data in binary classification

- ▶ Example: spam filtering
- ▶ Labeled example:  $(x, y) \in \mathcal{X} \times \{0, 1\}$
- ▶  $\mathcal{X}$  is input (feature) space,  $\{0, 1\}$  is the output (label) space
- ▶  $\mathcal{X}$  is not necessarily the space of inputs itself (e.g., space of all

▶ <https://eduassistpro.github.io>

- ▶  $X$  has some marginal probability  $d$
- ▶ Conditional probability distribution  
Bernoulli with heads probability
- ▶  $\eta: \mathcal{X} \rightarrow [0, 1]$  is a function, sometimes called regression function or conditional mean function (since  $\mathbb{E}[Y | X = x] = \eta(x)$ ).

## Error rate of a classifier

- ▶ For a classifier  $f: \mathcal{X} \rightarrow \{0, 1\}$ , the error rate of  $f$  (with respect to the distribution of  $(X, Y)$ ) is

Assignment Project Exam Help

$$\text{err}(f) := \Pr(f(X) \neq Y).$$

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

Add WeChat edu\_assist\_pro

which is the same as  $\Pr(f(X) \neq Y)$  if  $(X, Y)$  is uniform over the labeled examples.

- ▶ Caution: This notation  $\text{err}(f)$  does not make explicit the dependence on (the distribution of) the random example  $(X, Y)$ . You will need to determine this from context.

## Conditional expectations (1)

- ▶ Consider any random variables  $A$  and  $B$ .
- ▶ Conditional expectation of  $A$  given  $B$ :
  - ▶ Winter  $\mathbb{E}[A|B]$
  - ▶ A random variable! What is its expectation?
  - ▶ Law of iterated expectations (a.k.a. tower property):

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Conditional expectations (2)

- ▶ Example: roll a fair 6-sided die
  - ▶  $A$  = number shown facing up
  - ▶  $B$  = parity of number shown facing up
  - ▶  $C := \mathbb{E}[A \mid B]$  is random variable with

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

- Optimal classifier (Bayes classifier):

# Assignment Project Exam Help

where  $\eta$  is the conditional mean function

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

- Write error rate as  $\text{err}(f^*) = \Pr[f^*(X) \neq Y]$

- Conditional on  $X$ , probability of misclassification is

- $\min\{\eta(X), 1 - \eta(X)\}$

- So, optimal error rate is

$$\begin{aligned}\text{err}(f^*) &= \mathbb{E}[\mathbf{1}_{\{f^*(X) \neq Y\}}] \\ &= \mathbb{E}[\mathbb{E}[\mathbf{1}_{\{f^*(X) \neq Y\}} \mid X]] \\ &= \mathbb{E}[\min\{\eta(X), 1 - \eta(X)\}].\end{aligned}$$

## Example: spam filtering

- ▶ Suppose input  $x$  is a single (binary) feature, “is email all-caps?”
- ▶ How to interpret “the probability that email is spam given  $x=1$ ?”
- ▶ What does it mean for the Bayes classifier  $f^*$  to be optimal?

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

- ▶ What to do if  $\eta$  is unknown?
  - ▶ Training data:  $(x_1, y_1), \dots, (x_n, y_n)$
  - ▶ Assume data are related to what we want to predict
  - ▶ Let  $Z := (X, Y)$ , and  $Z_i := (X_i, Y_i)$  for  $i = 1, \dots, n$ .
  - ▶ IID model:  $Z_1, \dots, Z_n, Z$  are iid random variables

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr



## Performance of nearest neighbor classifier

- ▶ Study in context of IID model
- ▶ Assume  $\eta(x) \approx \eta(x')$  whenever  $x$  and  $x'$  are close.
  - ▶ This is where the modeling assumption comes in (via choice of distance function)!
- ▶ Let  $(X, Y)$  be the “test” example, and suppose  $(X_i, Y_i)$  is the

▶ <https://eduassistpro.github.io>

$$\eta(X) \approx \eta(X_i).$$

- ▶ Prediction is  $Y_i$ , true label is  $Y$ .
- ▶ Conditional on  $X$  and  $X_i$ , what is prob?
  - ▶  $\eta(X)(1 - \eta(X_i)) + (1 - \eta(X))$
- ▶ Conclusion: expected error rate is
$$\mathbb{E}[\text{err}(\text{NN}_S)] \approx 2 \cdot \mathbb{E}[\eta(X)(1 - \eta(X))] \text{ for large } n$$
  - ▶ Recall that optimal is  $\mathbb{E}[\min\{\eta(X), 1 - \eta(X)\}]$ .
  - ▶ So  $\mathbb{E}[\text{err}(\text{NN}_S)]$  is at most twice optimal.
  - ▶ Never exactly optimal unless  $\eta(x) \in \{0, 1\}$  for all  $x$ .

## Test error rate (1)

- ▶ How to estimate error rate?
- ▶ IID model:

$(X_1, Y_1), \dots, (X_n, Y_n), (X'_1, Y'_1), \dots, (X'_m, Y'_m), (X, Y)$  are iid.

- ▶ Training examples (that you have):  $(X_1, Y_1), \dots, (X_n, Y_n)$

▶ <https://eduassistpro.github.io>

- ▶ Hence, **test examples are independent** (important!)

▶ We would like to estimate  $\text{err}(\hat{f})$

- ▶ Caution: since  $\hat{f}$  depends on training data
- ▶ Convention: When we write  $\text{err}(\hat{f})$  where  $\hat{f}$  is random, we really mean  $\Pr(\hat{f}(X) \neq Y \mid \hat{f})$ .
- ▶ Therefore  $\text{err}(\hat{f})$  is a random variable!

## Test error rate (2)

- ▶ Conditional distribution of  $S := \sum_{i=1}^m \mathbf{1}_{\{\hat{f}(X'_i) \neq Y'_i\}}$  given training data:

▶  $S \mid \text{training data} \sim \text{Binomial}(m, \varepsilon)$  where  $\varepsilon := \text{err}(\hat{f})$

- ▶ By law of large numbers,

<https://eduassistpro.github.io>

Add WeChat [edu\\_assist\\_pro](#)

is close to  $\varepsilon$  when  $m$  is large

- ▶ How accurate is the estimate? Depends on the (conditional) variance!

- ▶  $\text{var}(\frac{1}{m}S \mid \text{training data}) = \frac{\varepsilon(1-\varepsilon)}{m}$

- ▶ Standard deviation is  $\sqrt{\frac{\varepsilon(1-\varepsilon)}{m}}$

► True positive rate (*recall*):  $\Pr(f(X) = 1 \mid Y = 1)$

► False positive rate:  $\Pr(f(X) = 1 \mid Y = 0)$

► Precision:  $\Pr(Y = 1 \mid f(X) = 1)$

► ...

► \_\_\_\_\_

<https://eduassistpro.github.io>

|         |             |                   |
|---------|-------------|-------------------|
| $y = 1$ | $\parallel$ | # false negatives |
|---------|-------------|-------------------|

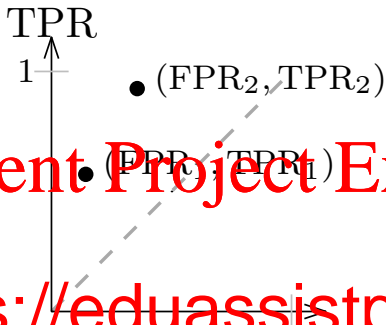
Add WeChat edu\_assist\_pr

- ▶ Receiver operating characteristic (ROC) curve
  - ▶ What points are achievable on the TPR-FPR plane?
  - ▶ Use randomization to combine classifiers

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

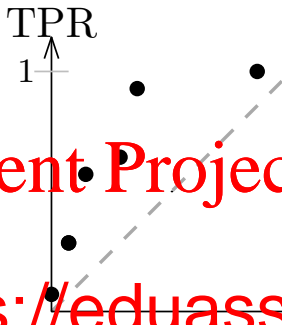


Assignment Project Exam Help

<https://eduassistpro.github.io>

Figure 5: TPR vs FPR plot with two poi

Add WeChat edu\_assist\_pr



Assignment Project Exam Help

<https://eduassistpro.github.io>

Figure 6: TPR vs FPR plot with many po

Add WeChat edu\_assist\_pr

## More than two outcomes

- ▶ What if there are  $K > 2$  possible outcomes?
- ▶ Replace coin with  $K$ -sided die
- ▶ Say  $X$  has a categorical distribution over  $[K] := \{1, \dots, K\}$ ,  
determined probability vector  $\theta = (\theta_1, \dots, \theta_K)$

- ▶ <https://eduassistpro.github.io>

$$\hat{y} := \arg \max_{k \in [K]}$$

Add WeChat edu\_assist\_pro



## Statistical model for multi-class classification

- ▶ Statistical model for labeled examples  $(X, Y)$ , where  $Y$  takes values in  $[K]$

Assignment Project Exam Help

- ▶ Now,  $Y | X = x$  has a categorical distribution with parameter vector  $\eta(x) = (\eta(x)_1, \dots, \eta(x)_K)$

- ▶ Conditional probability function:  $\eta(x)_k := \Pr(Y = k | X = x)$

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

Add WeChat edu\_assist\_pr

## Potential downsides of the IID model

- ▶ Example: Train OCR digit classifier using data from Alice's handwriting, but eventually use on digits written by Bob.

- ▶ What is a better evaluation?

Assignment Project Exam Help

<https://eduassistpro.github.io>

- ▶ What if we want to eventually use on digits written by Alice and Bob?

Add WeChat edu\_assist\_pro