

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Spring 2025



HAVERFORD
COLLEGE

- **Lab 8** due **TONIGHT!**
- Review today and Tuesday (in class and lab)
- **Exam** next Thursday (week from today)
- **Sophomore CS majors meeting TODAY**
 - 3:30-4:30pm in H109

Outline

- Midterm 2 Review
 - Revisit confusion matrices
 - PCA (linear transformation + interpretation)
 - Central Limit Theorem
 - Entropy vs. classification error
 - Logistic regression and cross entropy
 - Naïve Bayes

Outline

- Midterm 2 Review
 - Revisit confusion matrices
 - PCA (linear transformation + interpretation)
 - Central Limit Theorem
 - Entropy vs. classification error
 - Logistic regression and cross entropy
 - Naïve Bayes

Confusion matrix with more classes

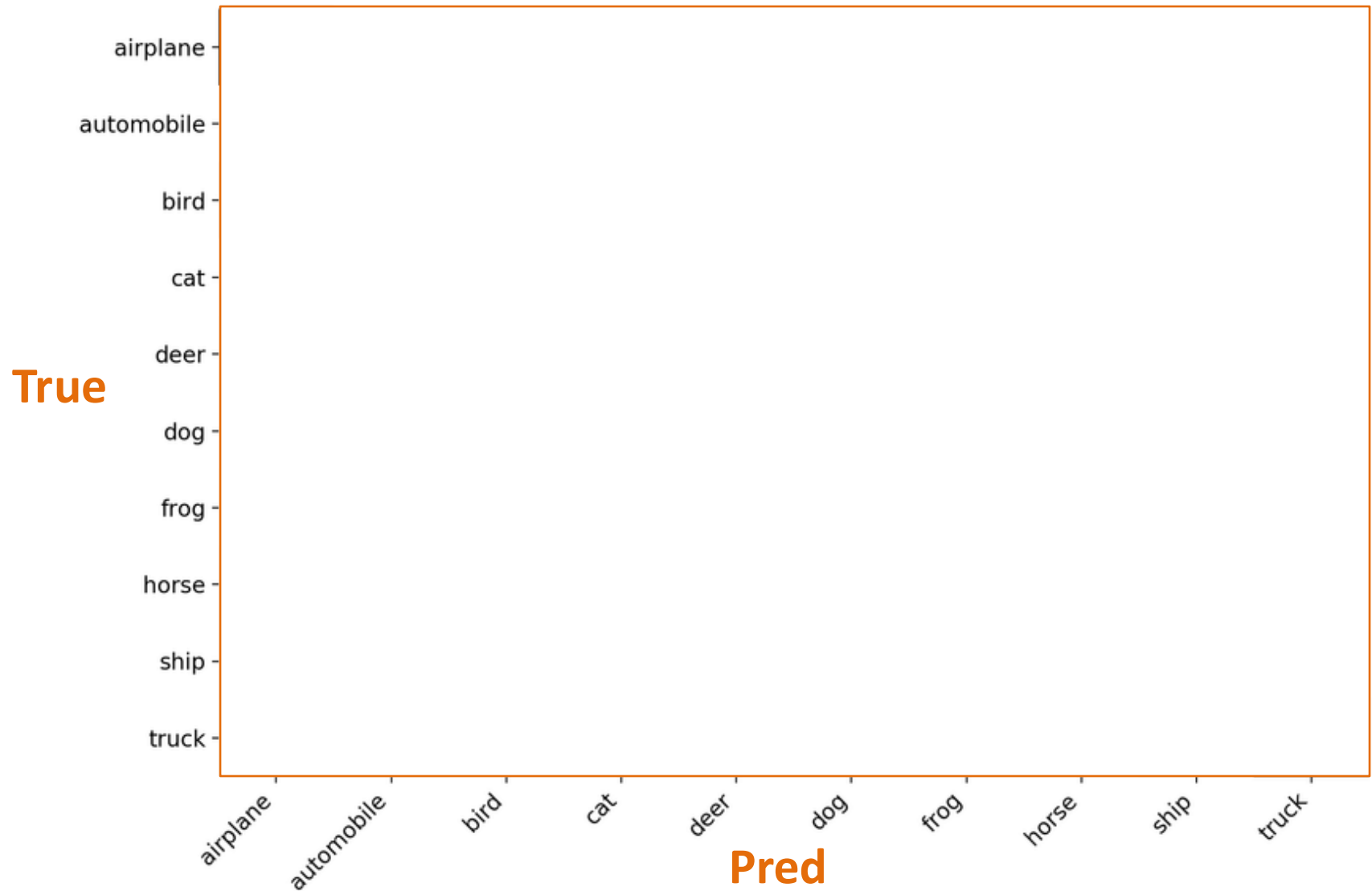
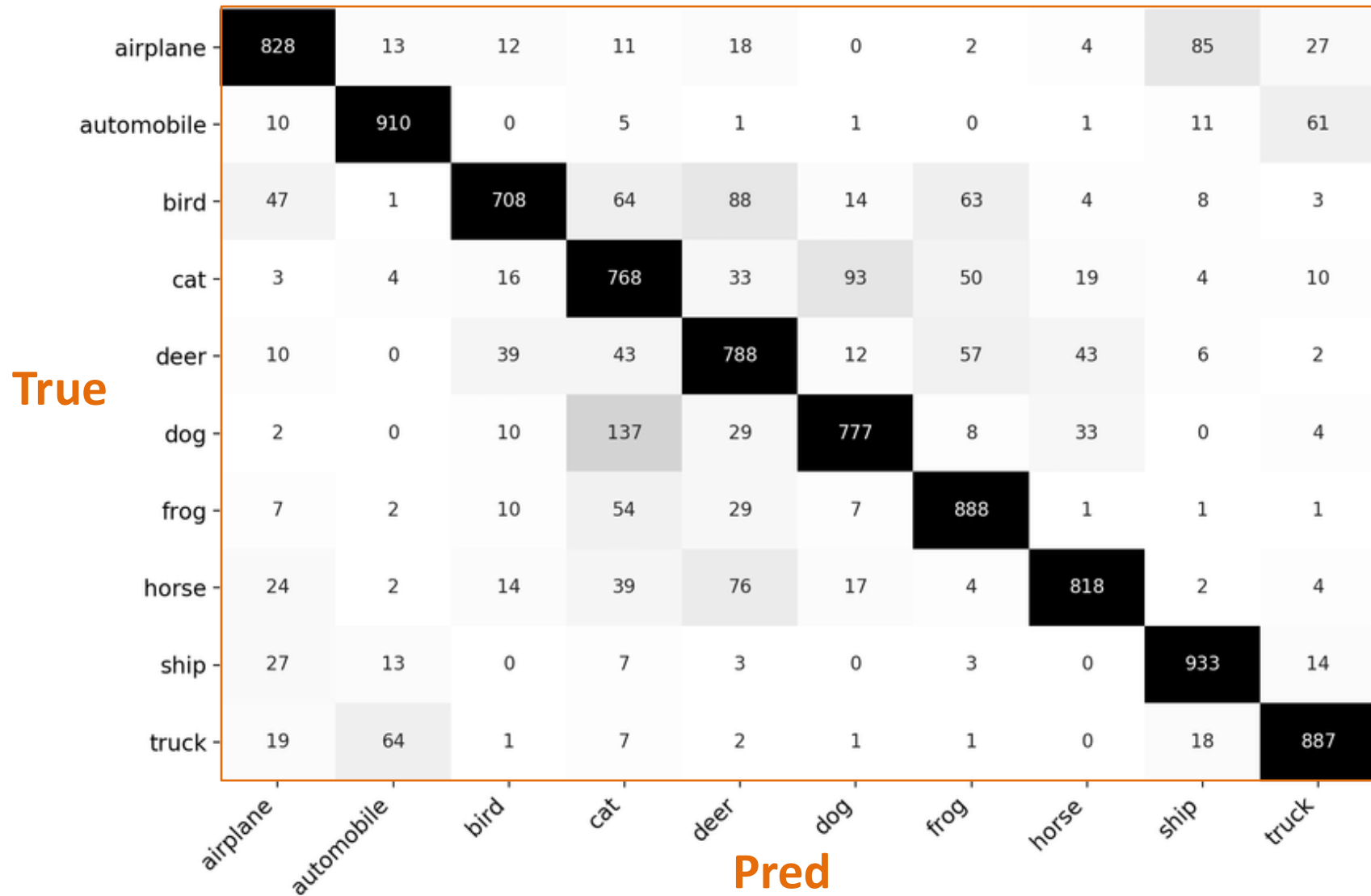


Figure by: Qun Liu (confusion matrix on cifar-10 dataset)

Confusion matrix with more classes



Confusion matrices with just two classes don't have to be “positive” and “negative”

- Example: male and female
 - No “positive” and “negative” class
 - ROC curve not appropriate

Confusion matrices

$cm = np.zeros((K, K))$

for ex in test:

$true = ex.label$

$pred = model.classify(ex.features)$

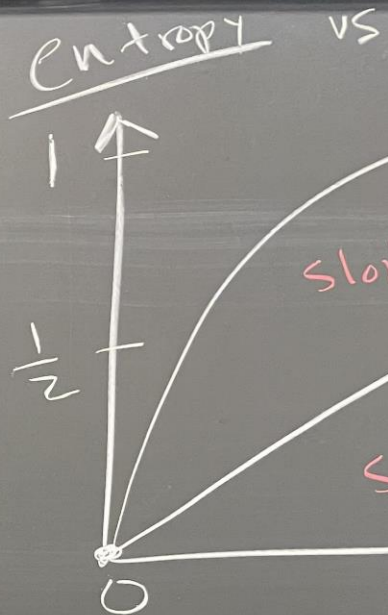
$cm[true, pred] += 1$

i.e. ③

i.e. ①

3, 1

classes



info gain

From the study guide

6. Data Visualization, Dimensionality Reduction, and Unsupervised Learning

- Best ways of visualizing **discrete** vs. **continuous** data
- How to choose colors; idea of **sequential**, **diverging**, or **qualitative** color schemes
- How to make color schemes color-blind and black/white printing friendly
- Idea of **principal component analysis (PCA)** as a way to accomplish **dimensionality reduction**
- Using dimensionality reduction to visualize high-dimensional data
- Details of the PCA algorithm (except computing eigenvalues and eigenvectors)
- Runtime of PCA
- Genealogical interpretation of PCA plots for genetic data
- Basic idea of **tSNE** as an alternative to PCA
- ~~• Idea of **clustering** as a form of **unsupervised learning**, for example **K means**~~

PCA creates linear combinations of features

PCA

X

$\begin{bmatrix} x_{i1}, x_{i2}, \dots, x_{ip} \end{bmatrix}$

$n \times p$

$\begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \\ \vdots \\ w_{p1} \end{bmatrix}$

$p \times 1$

$n \times 1$

$\begin{bmatrix} * \end{bmatrix}$

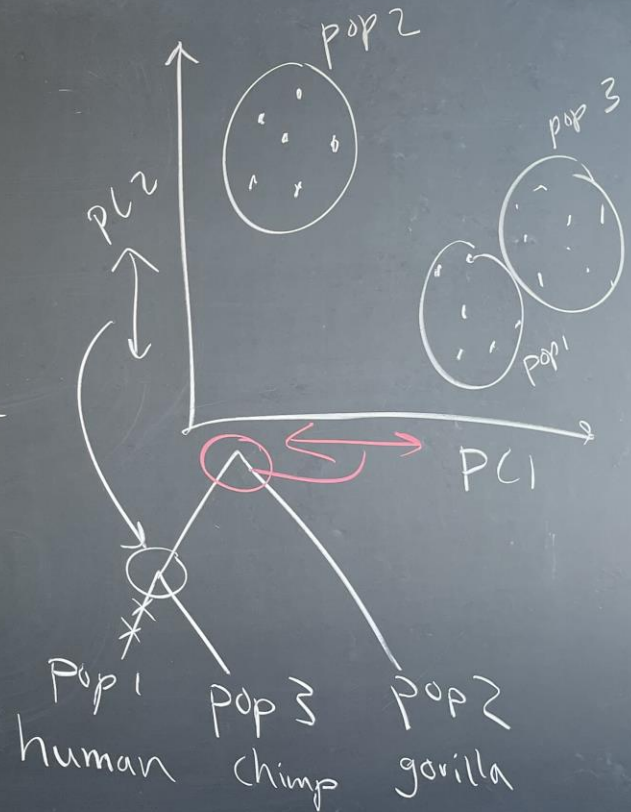
$n \times 1$

i th example

$$x_{i1} \cdot w_{11} + x_{i2} \cdot w_{21} + \dots + x_{ip} \cdot w_{p1}$$
$$= \vec{w}^{(1)} \cdot \vec{x}_i$$

first eigenvector

dot product



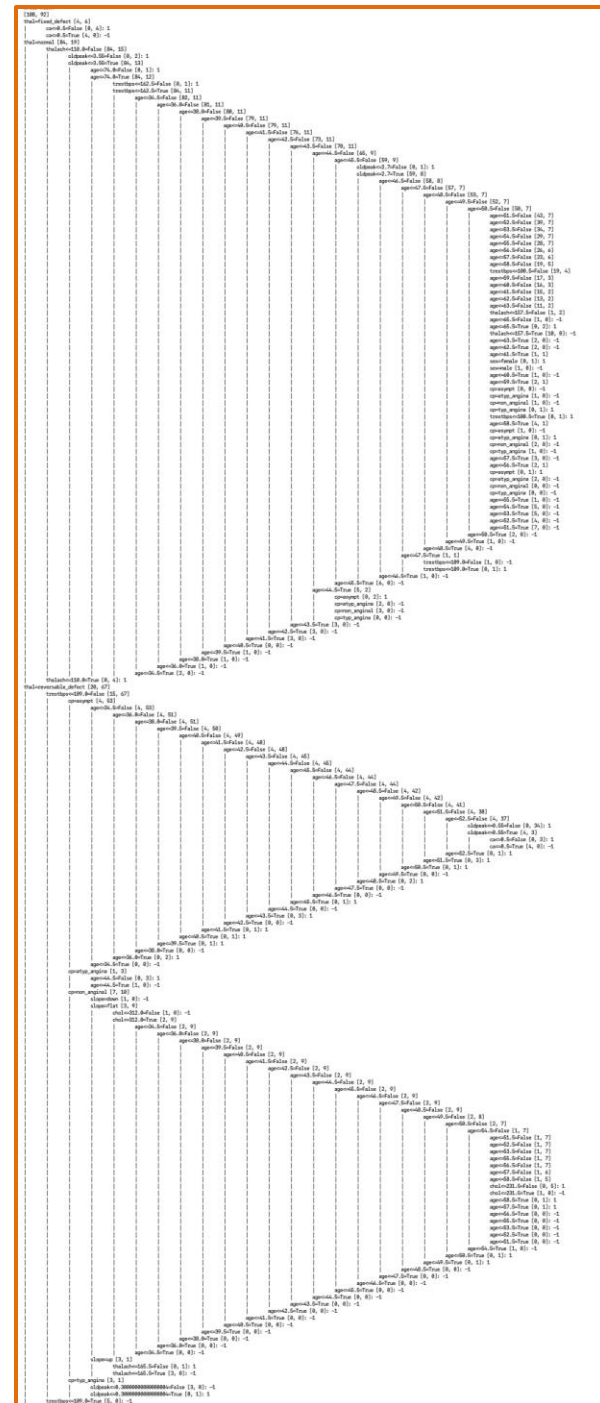
From the study guide

4. Information Theory

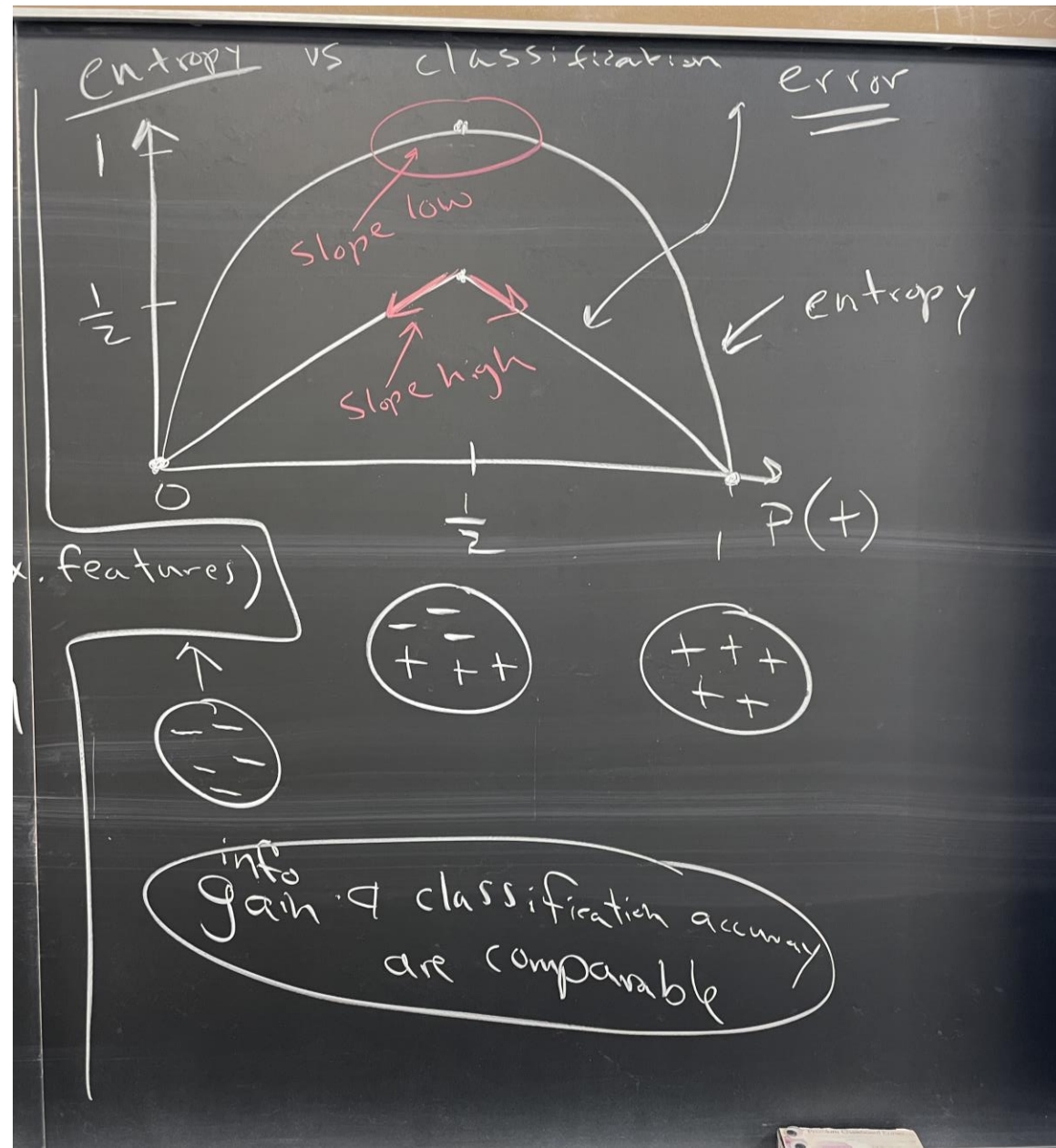
- Conceptual idea of [entropy](#) as well as formal definition
- [Shannon encoding](#) (and decoding), plus how to use entropy to compute average number of bits needed to send one piece of information
- Use of [conditional entropy](#) and [information gain](#) to choose best features
- Comparison with classification accuracy as a way to choose best features
- How to transform continuous features into binary features? (see Handout 14)

Decision trees from entropy (info gain) vs. classification error!

```
[108, 92]
thal=fixed_defect [4, 6]
|
| ca<=0.5=False [0, 6]: 1
| ca<=0.5=True [4, 0]: -1
|
| thal=normal [84, 19]
| |
| | thalach<=110.0=False [84, 15]
| | |
| | | age<=55.5=False [28, 11]
| | | |
| | | | chol<=248.5=False [14, 10]
| | | | |
| | | | | sex=female [13, 3]
| | | | | |
| | | | | | cp=asympt [3, 3]
| | | | | | |
| | | | | | | age<=57.5=False [1, 3]
| | | | | | | |
| | | | | | | | chol<=337.5=False [1, 0]: -1
| | | | | | | | |
| | | | | | | | | chol<=337.5=True [0, 3]: 1
| | | | | | | | |
| | | | | | | | | age<=57.5=True [2, 0]: -1
| | | | | | | | |
| | | | | | | | | cp=atyp_angina [2, 0]: -1
| | | | | | | | | cp=non_anginal [7, 0]: -1
| | | | | | | | | cp=typ_angina [1, 0]: -1
| | | | | | |
| | | | | | | sex=male [1, 7]
| | | | | | | |
| | | | | | | | age<=65.5=False [1, 2]
| | | | | | | | |
| | | | | | | | | age<=66.5=False [0, 2]: 1
| | | | | | | | | |
| | | | | | | | | | age<=66.5=True [1, 0]: -1
| | | | | | | | | |
| | | | | | | | | | age<=65.5=True [0, 5]: 1
| | | | | | |
| | | | | | | chol<=248.5=True [14, 1]
| | | | | | | |
| | | | | | | | oldpeak<=2.7=False [0, 1]: 1
| | | | | | | | |
| | | | | | | | | oldpeak<=2.7=True [14, 0]: -1
| | | | | | |
| | | | | | | age<=55.5=True [56, 4]
| | | | | | | |
| | | | | | | | trestbps<=113.5=False [47, 1]
| | | | | | | | |
| | | | | | | | | oldpeak<=3.55=False [0, 1]: 1
| | | | | | | | | |
| | | | | | | | | | oldpeak<=3.55=True [47, 0]: -1
| | | | | | | | |
| | | | | | | | | trestbps<=113.5=True [9, 3]
| | | | | | | | | |
| | | | | | | | | | oldpeak<=0.05=False [6, 0]: -1
| | | | | | | | | | |
| | | | | | | | | | | oldpeak<=0.05=True [3, 3]
| | | | | | | | | | |
| | | | | | | | | | | cp=asympt [0, 2]: 1
| | | | | | | | | | | |
| | | | | | | | | | | | cp=atyp_angina [2, 0]: -1
| | | | | | | | | | | | cp=non_anginal [1, 1]
| | | | | | | | | | | | |
| | | | | | | | | | | | | age<=41.5=False [0, 1]: 1
| | | | | | | | | | | | | |
| | | | | | | | | | | | | | age<=41.5=True [1, 0]: -1
| | | | | | | | | | | |
| | | | | | | | | | | | cp=typ_angina [0, 0]: -1
| | | | | | |
| | | | | | | thalach<=110.0=True [0, 4]: 1
| | | | | |
| | | | | thal=reversible_defect [20, 67]
| | | | | |
| | | | | | cp=asympt [5, 53]
| | | | | | |
| | | | | | | oldpeak<=0.55=False [0, 43]: 1
| | | | | | | |
| | | | | | | | oldpeak<=0.55=True [5, 10]
| | | | | | | | |
| | | | | | | | | chol<=237.5=False [0, 8]: 1
| | | | | | | | | |
| | | | | | | | | | chol<=237.5=True [5, 2]
| | | | | | | | | | |
| | | | | | | | | | | chol<=179.5=False [4, 0]: -1
| | | | | | | | | | | |
| | | | | | | | | | | | chol<=179.5=True [1, 2]
| | | | | | | | | | | | |
| | | | | | | | | | | | | age<=59.5=False [1, 0]: -1
| | | | | | | | | | | | | |
| | | | | | | | | | | | | | age<=59.5=True [0, 2]: 1
| | | | | | |
| | | | | | | cp=atyp_angina [3, 3]
| | | | | | | |
| | | | | | | | age<=46.5=False [1, 3]
| | | | | | | | |
| | | | | | | | | trestbps<=109.0=False [0, 3]: 1
| | | | | | | | | |
| | | | | | | | | | trestbps<=109.0=True [1, 0]: -1
| | | | | | | | |
| | | | | | | | | age<=46.5=True [2, 0]: -1
| | | | | | |
| | | | | | | cp=non_anginal [9, 10]
| | | | | | | |
| | | | | | | | oldpeak<=1.85=False [0, 5]: 1
| | | | | | | | |
| | | | | | | | | oldpeak<=1.85=True [9, 5]
| | | | | | | | | |
| | | | | | | | | | trestbps<=121.0=False [3, 5]
| | | | | | | | | | |
| | | | | | | | | | | chol<=232.5=False [0, 4]: 1
| | | | | | | | | | | |
| | | | | | | | | | | | chol<=232.5=True [3, 1]
| | | | | | | | | | | | |
| | | | | | | | | | | | | trestbps<=128.5=False [3, 0]: -1
| | | | | | | | | | | | | |
| | | | | | | | | | | | | | trestbps<=128.5=True [0, 1]: 1
| | | | | | | | | | |
| | | | | | | | | | | trestbps<=121.0=True [6, 0]: -1
| | | | | | |
| | | | | | | cp=typ_angina [3, 1]
| | | | | | | |
| | | | | | | | oldpeak<=0.30000000000000004=False [3, 0]: -1
| | | | | | | | |
| | | | | | | | | oldpeak<=0.30000000000000004=True [0, 1]: 1
```



Entropy vs. classification error



Midterm Practice Exam: pg 1 and 2

On your own for ~15 min, then partners

① step 1: get the data ✓ $O(1)$

step 2: subtract off mean $O(np)$

$$O(np) + O(np) \quad \left[\begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \right]_{n \times p}$$

step 3: cov of each pair of features

$$\sum_{i=1}^n (x_{if} - \bar{x}_f)(x_{ig} - \bar{x}_g) \quad \left. \begin{array}{l} \text{2 features} \\ \text{ } \end{array} \right\} O(n)$$

p^2 pairs of features

$$\Rightarrow O(p^2 n)$$

step 4: eigenvalues + eigenvectors of A
 \Rightarrow cubic

$$O(p^3)$$

step 5: transform data

$$\underbrace{(n \times p)}_X \times \underbrace{(p \times r)}_{W_r} = \underbrace{(n \times r)}_{T_r} \Rightarrow O(npr)$$

step 6: plot!

Overall runtime of PCA

$$\cancel{O(np)} + O(np^2) + O(p^3) + \cancel{O(npr)}$$

$r \leq p$

$$\Rightarrow \boxed{O(np^2 + p^3)}$$

(2) $E[Y] = \sum_Y Y P(Y)$

(a)

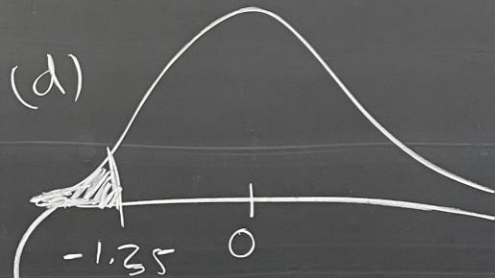
$$= 0 \cdot \frac{1}{8} + 1 \cdot \frac{1}{8} + 2 \cdot \frac{1}{4} + 3 \cdot \frac{1}{2}$$

$$= \boxed{2.125} = \mu$$

(b) $\text{Var}(Y) = E[(Y - \mu)^2]$

$$= (0 - 2.125)^2 \cdot \frac{1}{8} + \dots = \boxed{1.109} = \sigma^2$$

(c) $Z = \frac{\bar{Y}_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} = \frac{1.9 - 2.125}{\sqrt{\frac{1.109}{40}}} = \boxed{-1.39}$



area = p-value = $\boxed{0.08833} > 0.05$

fail to reject H_0

(3) (a) age \Rightarrow $\boxed{\frac{15}{99}}$

old peak $=$ $\boxed{\frac{13}{99}}$ } choose as best feature

(b) $H(y) = - \left(\frac{84}{99} \log_2 \frac{84}{99} + \frac{15}{99} \log_2 \frac{15}{99} \right)$

$w_1 + x_{i2} w_2 y = -1$

$= 0.61$

(c) $H(y | \text{old peak})$

$H(y)$

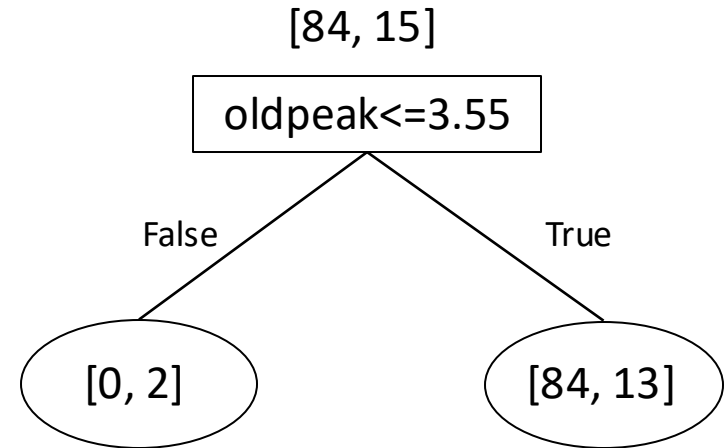
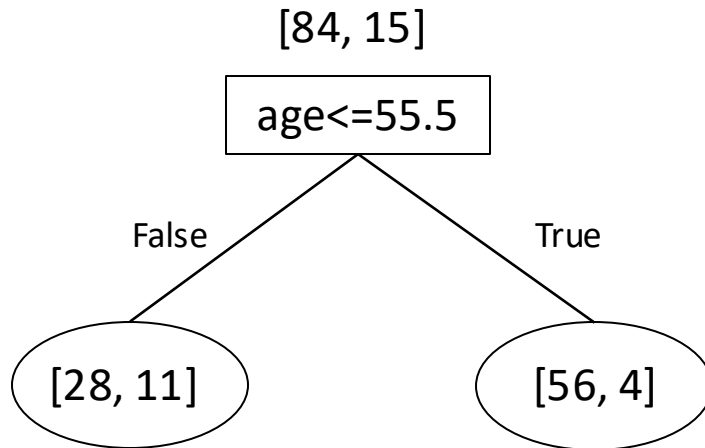
$=$

$$(C) \quad H(Y | \text{oldpeak}) = \frac{2}{99} H(Y | \text{oldpeak} = F) \\ + \frac{97}{99} H(Y | \text{oldpeak} = T)$$

$$H(Y | \text{oldpeak} = T)$$

$$= - \left(\frac{84}{97} \log_2 \frac{84}{97} + \frac{13}{97} \log_2 \frac{13}{97} \right)$$

One feature models (decision stumps): information gain vs. classification error



$$H(Y) = 0.6136190195993708$$

$$H(Y | \text{age} \leq 55.5) = 0.5522480910534322$$

$$H(Y | \text{oldpeak} \leq 3.55) = 0.5568804630596093$$

=> Age feature
produces more
information gain!