# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023

HAVERFORD
COLLEGE

# Write one midterm question or topic on your notecard

Will answer later in class today or on Piazza

# Admin

- **Exam** due in class on Tuesday
  - Do not open the exam until you're ready to start
  - Time limit: **3 hours**
  - Resources: hand-written study sheet, calculator

- **First candidate talk on Monday!**
  - 4pm tea
  - 4:15pm talk (H109)

# Outline for November 16

- Midterm 2 Review

  – Entropy vs. classification error

  – PCA

  – Naïve Bayes

  – Central Limit Theorem

  – Logistic regression and cross entropy
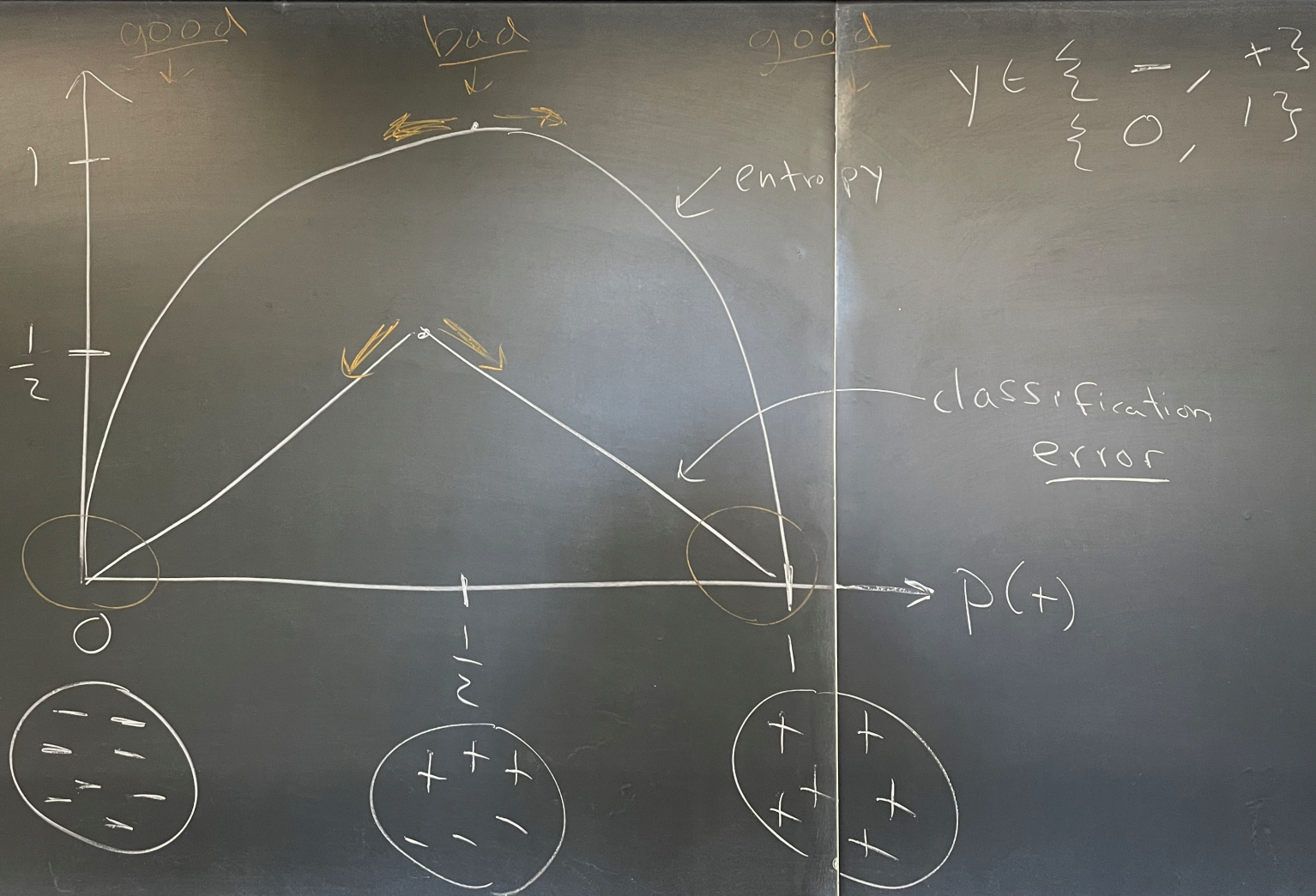
# Outline for November 16

- Midterm 2 Review

  - Entropy vs. classification error
  - PCA
  - Naïve Bayes
  - Central Limit Theorem
  - Logistic regression and cross entropy
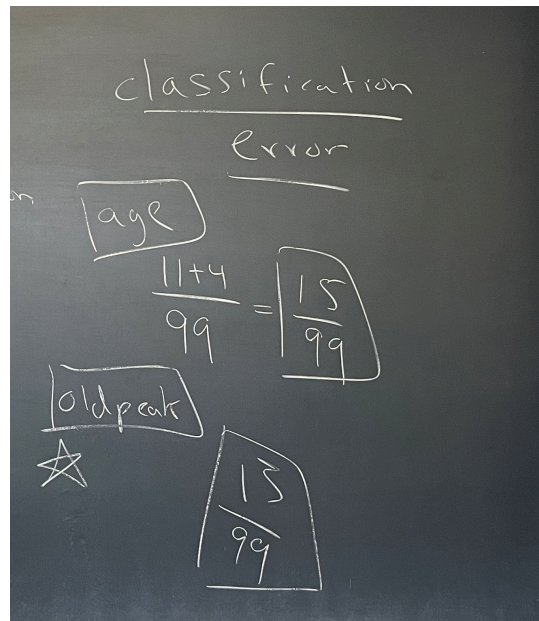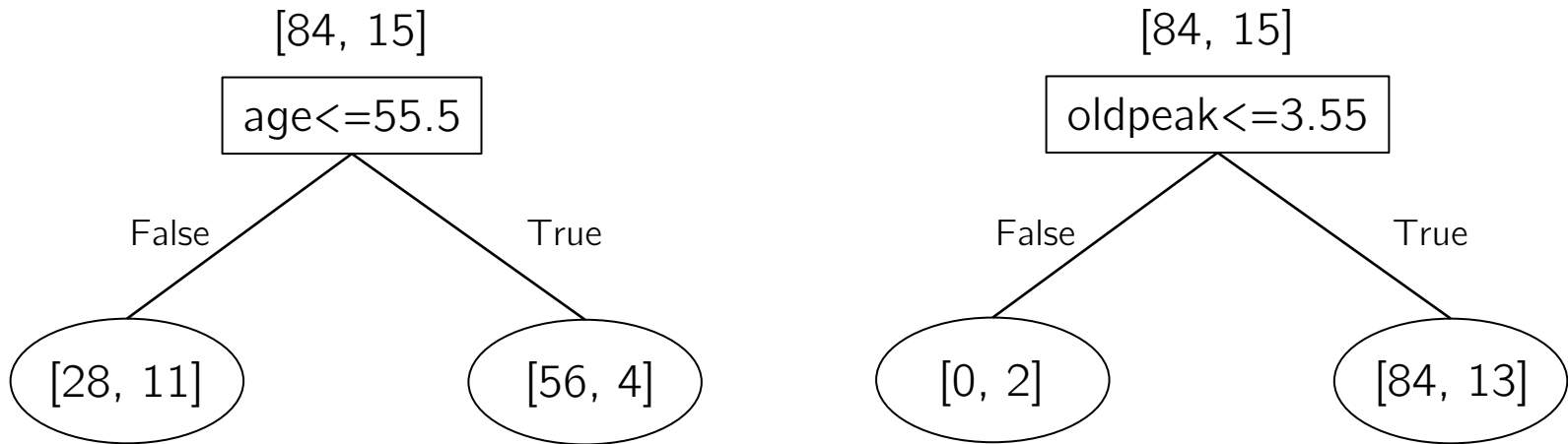
# From the study guide

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)

# Entropy vs. classification error

good        bad        good

$1$

$\frac{1}{2}$

entropy

classification error

$O$        $\frac{1}{2}$        $1$        $p(+)$

$y \in \{-1, +\}$
$\{0, 1\}$

(- - -)        (+ + + - -)        (+ + + + + +)

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

[84, 15]

age<=55.5

False | True

[28, 11]  [56, 4]
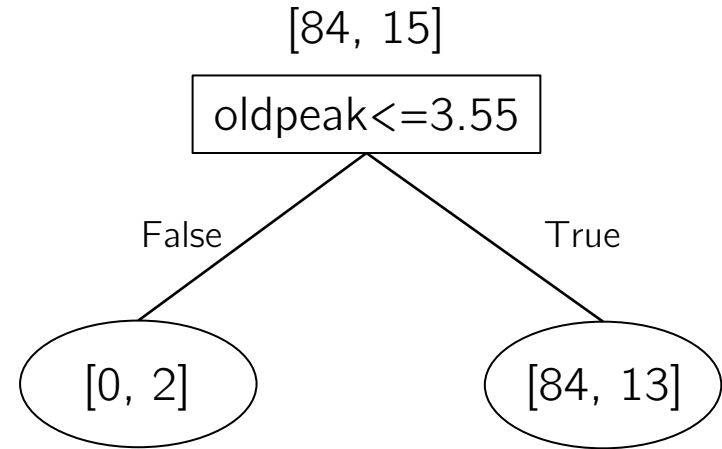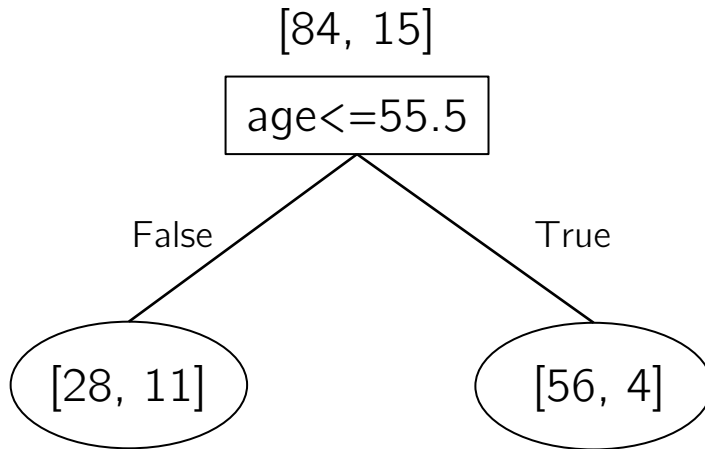
[84, 15]

oldpeak<=3.55

False | True

[0, 2]  [84, 13]

$$H(Y) = -\sum C \; P(y=c) \log_2 P(y=c) \qquad y \in \{-1, +1\}$$

$$C \; Evals(y)$$

$$H(Y) = -\left( \frac{84}{99} \log_2 \frac{84}{99} + \frac{15}{99} \log_2 \frac{15}{99} \right) = 0.61$$

$$H(Y \mid oldpeak) = \frac{2}{99} \underset{\nearrow 0}{H(Y \mid oldpeak = F)} + \frac{97}{99} H(Y \mid oldpeak = T)$$

$$H(Y \mid oldpeak = T) = -\left( \frac{84}{97} \log \frac{84}{97} + \frac{13}{97} \log \frac{13}{97} \right)$$

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

[84, 15]

age<=55.5

False         True

[28, 11]        [56, 4]

[84, 15]

oldpeak<=3.55

False         True

[0, 2]        [84, 13]

$H(Y) = 0.6136190195993708$

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

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

=> Age feature produces more information gain!

# 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=reversable_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
```

# Outline for November 16

- Midterm 2 Review

  - Entropy vs. classification error
  - PCA
  - Naïve Bayes
  - Central Limit Theorem
  - Logistic regression and cross entropy

# From the study guide

6. <u>Data Visualization</u>

- 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

$X$ — varies a lot!

$\vec{W}$   $W$

$w_1 = 2$

$w_2 = 5$

$w_3 = 0$

$w_4 = -1$

$X_{i1}$   $X_{i2}$   $X_{i3}$   $X_{i4}$

$n \times p$

$\uparrow$

$4$

$\Rightarrow$ not important

$n \times 2$

$P \times 2$

$\bigstar = W_1 X_{i1} + W_2 X_{i2} + W_3 X_{i3} + W_4 X_{i4}$

$= \boxed{\vec{W} \cdot \vec{X}_i}$

PCA is a linear transformation.

Modern west Eurasians

| | | | | |
|---|---|---|---|---|
| ● Armenian | ● English | ● Abkhasian | ✕ Libyan_Jew | ■ Croatian |
| ■ Iranian | ■ French | ■ Adygei | ⊕ Moroccan_Jew | ▲ Czech |
| ▲ Turkish | ▲ Icelandic | ▲ Balkar | + Tunisian_Jew | ◆ Estonian |
| ▲ Albanian | ◆ Norwegian | ◆ Chechen | ✳ Turkish_Jew | ✕ Hungarian |
| ■ Bergamo | ✕ Orcadian | ✕ Georgian | ⊞ Yemenite_Jew | ⊕ Lithuanian |
| ▲ Bulgarian | ⊕ Scottish | ⊕ Kumyk | ● Basque | + Ukrainian |
| ◆ Cypriot | ● BedouinA | + Lezgin | ■ French_South | ● Canary_Islanders |
| ✕ Greek | ■ BedouinB | ✳ North_Ossetian | ▲ Spanish | ■ Sardinian |
| ⊕ Italian_South | ▲ Jordanian | ● Ashkenazi_Jew | ◆ Spanish_North | ● Finnish |
| + Maltese | ◆ Palestinian | ■ Georgian_Jew | ● Druze | ■ Mordovian |
| ✳ Sicilian | ✕ Saudi | ▲ Iranian_Jew | ■ Lebanese | ▲ Russian |
| ⊞ Tuscan | ⊕ Syrian | ◆ Iraqi_Jew | ● Belarusian | |

# PCA "classic" genetics example

# Outline for November 16

- Midterm 2 Review

  - Entropy vs. classification error
  - PCA
  - Naïve Bayes
  - Central Limit Theorem
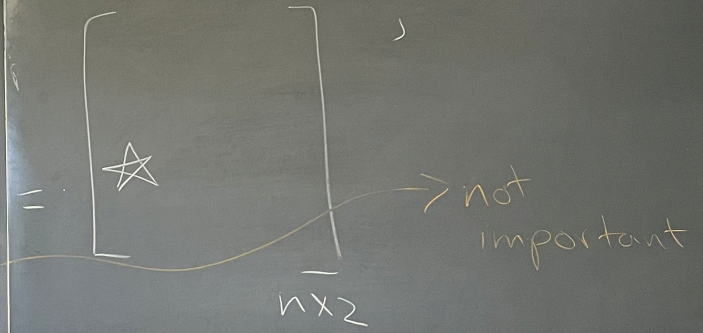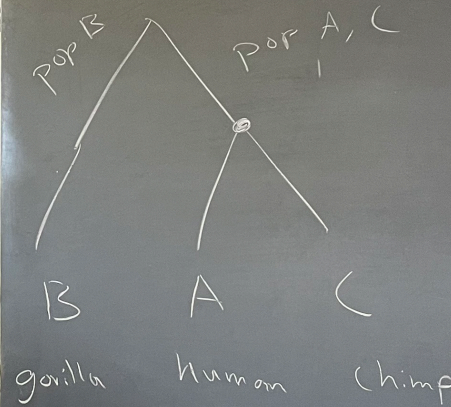  - Logistic regression and cross entropy

# From the study guide

2. <u>Naive Bayes</u>

- Bayes rule in data science: identify and explain the evidence, prior, posterior, likelihood.
- Derivation of the Naive Bayes model for $p(y = k|\vec{x})$ (via the Naive Bayes assumption).
- How do we estimate the probabilities of a Naive Bayes model?
- Laplace counts (motivation, application details)
- How can we predict the label of a new example after fitting a Naive Bayes model?
- What types of features/label do we currently require for Naive Bayes?
- How Naive Bayes can be implemented using dictionaries in Python

# Naïve Bayes assumption



Bayes $\quad P(A,B) = P(A)P(B|A)$

independence $\quad P(A,B) = P(A)P(B)$ ← not always true!

conditional independence $\quad P(A|B,C) = P(A|C)$ ←

likelihood $\quad P(x_1, x_2, x_3 | y) = P(x_1 | y) P(x_2, x_3 | \cancel{x_1}, y)$

$$P(x_1 | y) P(x_2 | y) P(x_3 | \cancel{x_2}, y)$$

$$\leftarrow \prod_{j=1}^{P} P(x_j | y)$$

# Naïve Bayes Model

$$p(y = k | \boldsymbol{x}) \propto p(y = k) \prod_{j=1}^{p} p(x_j | y = k).$$

# Naïve Bayes Prediction

$$\hat{y} = \underset{k \in \{1, 2, \cdots, K\}}{\arg \max} \; p(y = k) \prod_{j=1}^{p} p(x_j | y = k).$$

# Estimating prior: p(y=k)

$$\theta_k = \frac{N_k + 1}{n + K}$$

# Estimating likelihood: $p(x_j = v \mid y = k)$

$$\theta_{k,j,v} = \frac{N_{k,j,v} + 1}{N_k + |f_j|}$$

# Outline for November 16

- Midterm 2 Review

  - Entropy vs. classification error
  - PCA
  - Naïve Bayes
  - Central Limit Theorem
  - Logistic regression and cross entropy

# From the study guide

7. <u>Statistics</u>

- Motivation for studying statistics and hypothesis testing
- Probability distributions (discrete vs. continuous)
- Computing (theoretical) expected value and variance for discrete distributions
- Sample mean and sample variance
- Central limit theorem (CLT) and application in cases where the mean/variance are known
- Computation and interpretation of Z-scores and p-values
- Null vs. alternative hypotheses; when to reject the null hypothesis; significance level $\alpha$
- Using randomized trials and permutation testing to obtain more precise p-values
- Idea of a t-test as a way to test differences in means (not details)
- Bootstrap: sampling from our data with replacement (usually keeping $n$ the same)
- How to use bootstrapping to obtain confidence intervals
- Bagging (Bootstrap Aggregation): create a classifier for each bootstrapped training dataset
- Idea of using an ensemble of classifiers (ideally with low bias) to reduce variance
- To test, let each classifier in the ensemble "vote"

# Bootstrap demo

See video tutorial on Piazza!

# Outline for November 16

- Midterm 2 Review

  - Entropy vs. classification error
  - PCA
  - Naïve Bayes
  - Central Limit Theorem
  - Logistic regression and cross entropy

# From the study guide

5. <u>Logistic Regression</u>

- Motivation for logistic regression; our model is a logistic function that takes in $\vec{w} \cdot \vec{x}$
- Logistic regression creates a *linear* decision boundary (visualize for $p = 1$).
- In logistic regression our cost is the negative log likelihood (don't need to derive)
- Intuition/visualization of the cost function (and relationship to cross entropy)
- Idea of SGD for logistic regression, relationship to linear regression

# For each method/approach, is X continuous or discrete? What about y?

- Linear regression
- Polynomial regression
- Decision trees/stumps
- ROC curve as an evaluation metric
- Naïve bayes
- Logistic regression
- Entropy and information gain
- PCA

*Think about offline!*

# Notecards: will post responses on Piazza!