

# CS 260: Foundations of Data Science

Prof. Thao Nguyen

Fall 2024



**HVERFORD**  
COLLEGE

# Outline for today

- Midterm 2 Review
  - PCA
  - Naïve Bayes
  - Logistic regression and cross entropy

# Outline for today

- Midterm 2 Review
  - PCA
  - Naïve Bayes
  - Logistic regression and cross entropy

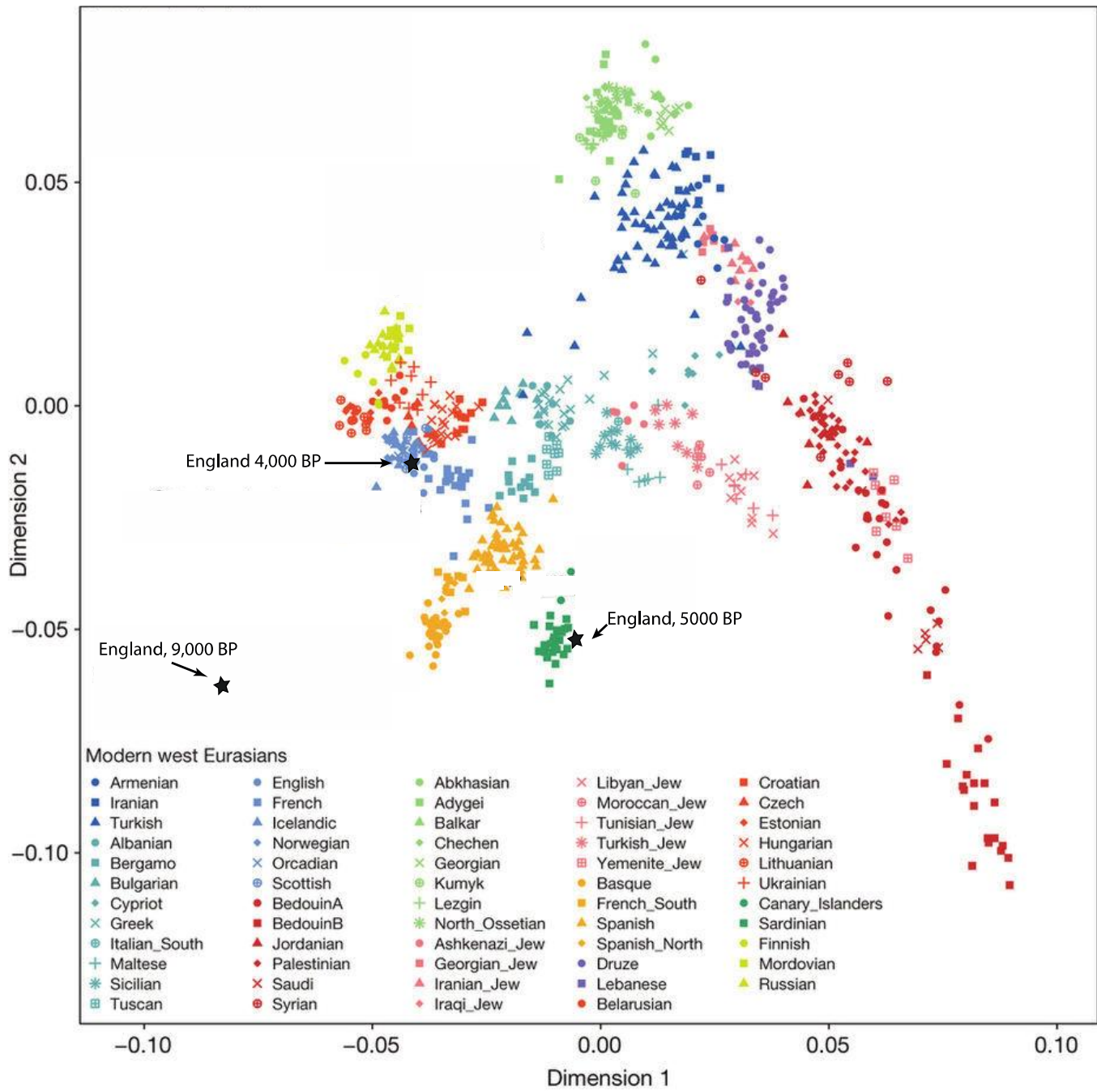
# From the study guide

## 6. Data Visualization

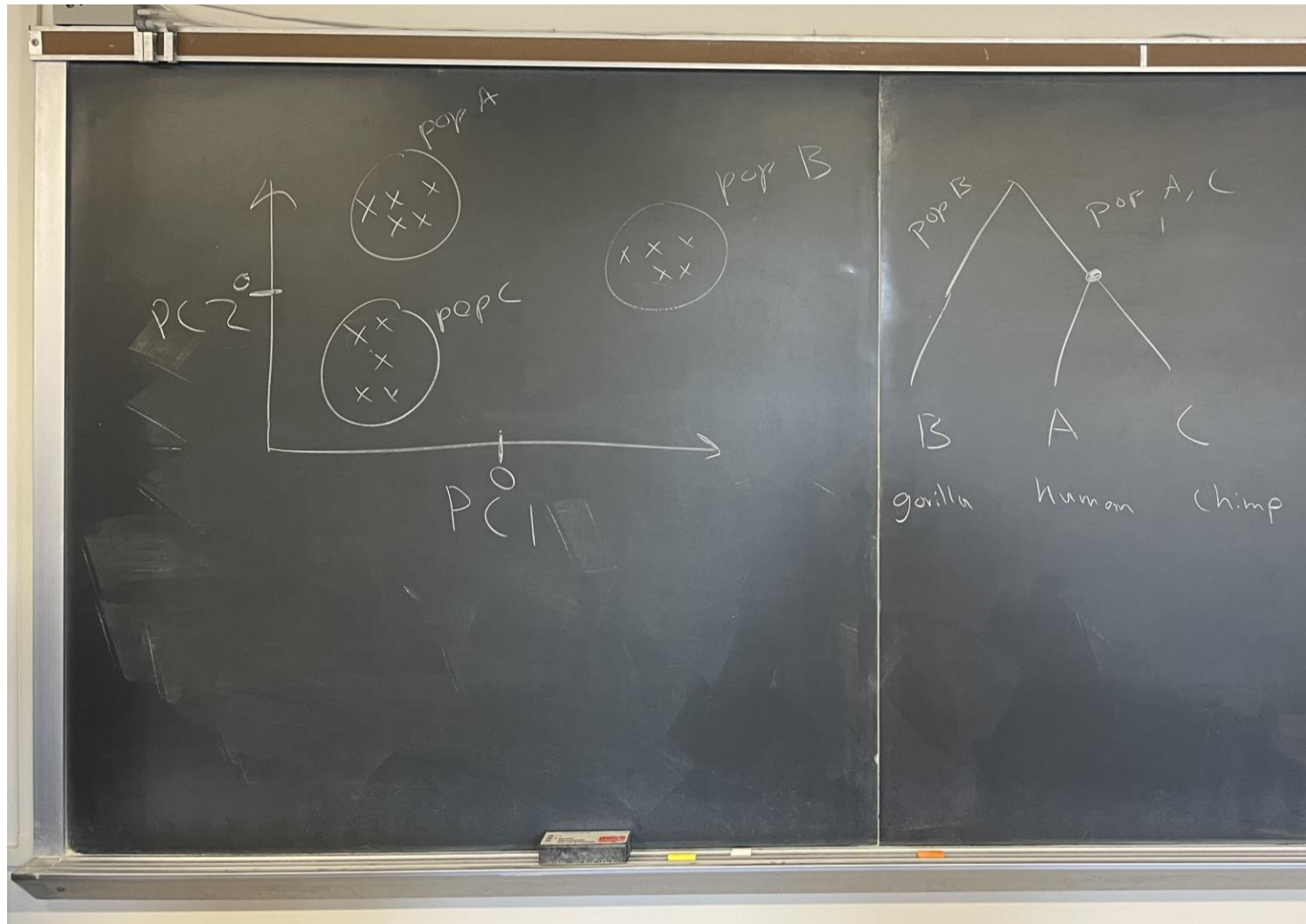
- Best ways of visualizing **discrete** vs. **continuous** data
- How to choose colors; idea of **sequential**, **diverging**, or **qualitative** color schemes
- How to make color schemes color-blind and black/white printing friendly
- Idea of **principal component analysis (PCA)** as a way to accomplish **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

# Principal Component Analysis (PCA)

- Transforms  $p$ -dimensional data so that the new first dimension explains as much of the variation as possible, the new second explains as much of the remaining variation as possible, and so on
- PCA is a linear transformation
- Typically, we look at the first few dimensions of the transformed data as a means of dimensionality reduction and visualization
- PCA is often used for:
  - Data visualization
  - Infer qualitative relationships between groups



# PCA “classic” genetics example



# Outline for today

- Midterm 2 Review
  - PCA
  - Naïve Bayes
  - Logistic regression and cross entropy



# From the study guide

## 2. Naive Bayes

- 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

# Bayes' Theorem

- $P(A,B) = P(A | B)P(B)$
- $P(A,B) = P(B | A)P(A)$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

# Independence

- Independence:  $P(A,B) = P(A)P(B)$
- Conditional independence:  $P(A | B,C) = P(A | C)$

not always true!

↓  
Naïve Bayes  
assumption

# Naïve Bayes Model

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

# Naïve Bayes Prediction

$$\hat{y} = \arg \max_{k \in \{1, 2, \dots, K\}} p(y = k) \prod_{j=1}^p p(x_j | y = k).$$

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 today

- Midterm 2 Review
  - PCA
  - Naïve Bayes
  - Logistic regression and cross entropy

# From the study guide

## 5. Logistic Regression

- 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



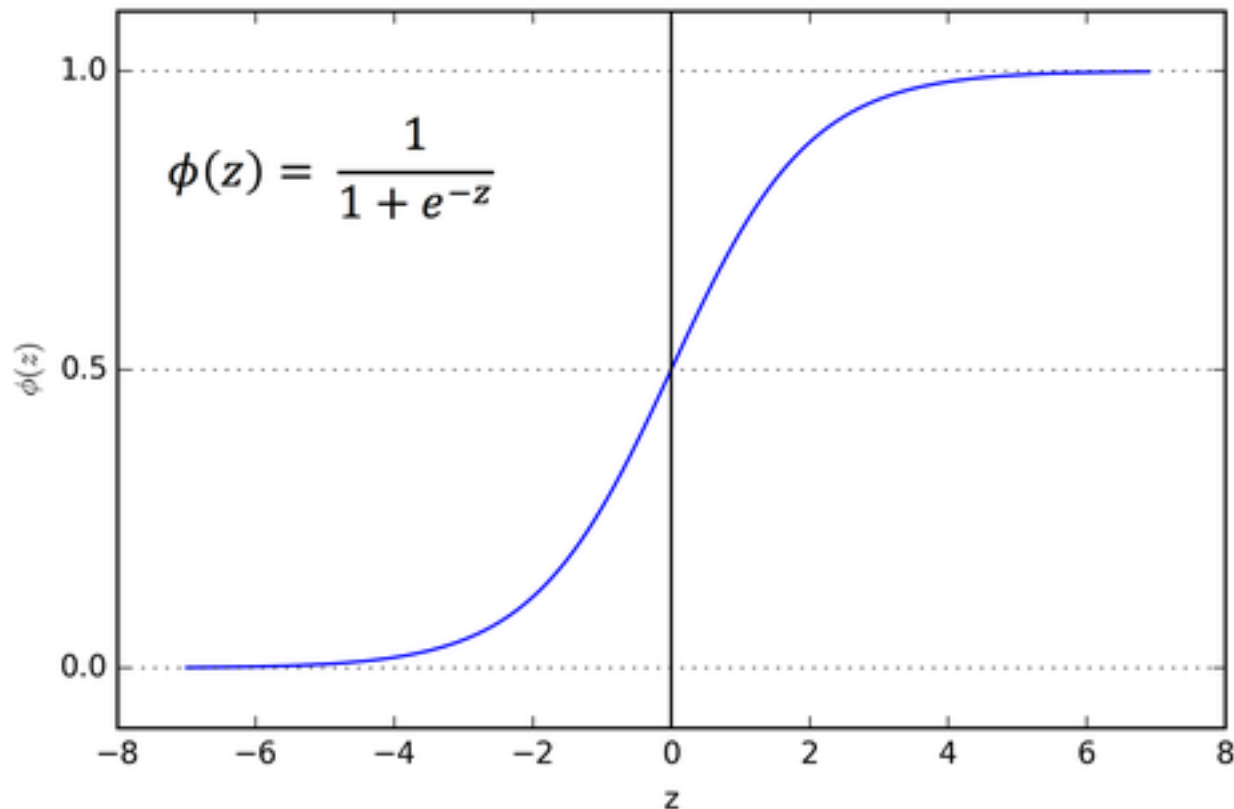
# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

# Logistic (sigmoid) function

Transforms a continuous real number into a range of (0, 1)



# Logistic Regression

- Binary classification  $y \in \{0,1\}$
- Model will be

$$h_{\vec{w}}(\vec{x}) = p(y = 1|\vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$$

- Classification (already have  $\vec{w}$ )

$$\text{if } \vec{w} \cdot \vec{x} \geq 0 \Rightarrow \hat{y} = 1$$

$$\vec{w} \cdot \vec{x} < 0 \Rightarrow \hat{y} = 0$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

# 3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient of cost wrt single data point  $\mathbf{x}_i$

$$\nabla J_{\mathbf{x}_i}(\mathbf{w}) = (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)\mathbf{x}_i$$

# Stochastic Gradient Descent for Logistic Regression (binary classification)

set  $\vec{w} = \vec{0}$

while cost  $J(\vec{w})$  is still changing:

shuffle data points

for  $i = 1, \dots, n$ :

$$\vec{w} \leftarrow \vec{w} - \alpha \underbrace{\nabla_{\vec{x}_i} J(\vec{w})}_{\text{derivative of } J(\vec{w}) \text{ wrt } x_i}$$

store  $J(\vec{w})$

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!*