# CSE 158 – Lecture 2
Web Mining and Recommender Systems

Supervised learning – Regression

# Supervised versus unsupervised learning

**Learning** approaches attempt to **model data** in order to solve a problem

**Unsupervised learning** approaches find patterns/relationships/structure in data, but **are not** optimized to solve a particular predictive task

**Supervised learning** aims to directly model the relationship between input and output variables, so that the output variables can be predicted accurately given the input
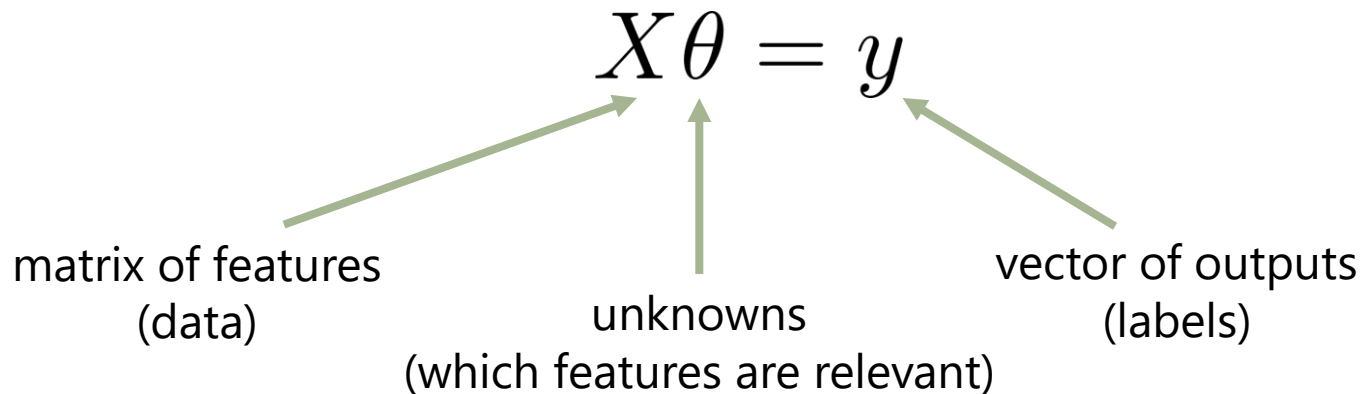
**Regression** is one of the simplest supervised learning approaches to learn relationships  between input variables (features) and output variables (predictions)

**Linear regression** assumes a predictor of the form

$$y_i = x_i \cdot \theta$$

$$X\theta = y$$

matrix of features
(data)

unknowns
(which features are relevant)

vector of outputs
(labels)

(or $Ax = b$ if you prefer)

# Linear regression

**Linear regression** assumes a predictor of the form

$$X\theta = y$$

**Q:** Solve for theta

**A:** $\theta = (X^T X)^{-1} X^T y$

# Example 1

How do preferences toward certain beers vary with age?

# Example 1



**Beers:**

**Ratings/reviews:**

**User profiles:**

# Example 1

50,000 reviews are available on
http://jmcauley.ucsd.edu/cse158/data/beer/beer_50000.json
(see course webpage)

See also – non-alcoholic beers:
http://jmcauley.ucsd.edu/cse158/data/beer/non-alcoholic-beer.json

# Example 1

# Real-valued features

How do preferences toward certain beers vary with age?
How about **ABV**?

$$\text{rating} = \theta_0 + \theta_1 \, age$$

$$\theta \cdot x$$

$$[\theta_0, \theta_1] \cdot [1, age]$$

(code for all examples is on http://jmcauley.ucsd.edu/cse158/code/week1.py)

# Example 1

## Preferences vs **ABV**

$$\text{ratings} = \theta_0 + \theta_1 ABV + \theta_2 ABV^2$$



ratings

3.f

ABV

$$\sum_i \left( y_i - x_i \theta \right)^2 \qquad \theta \cdot x$$

$$\left( 1, ABV, ABV^2, \dots \right)$$

Example 2

# Categorical features

$$\theta = [1, \text{ is male}, \text{ is female}]$$

How do beer preferences vary as a function of **gender**?

$$\text{rating } = \theta_0 + \theta_1 [\text{if male}]$$
$$+ \theta_2 [\text{if female}]$$

$$x_i = [1, 0, 1] \text{ if female}$$
$$= [1, 1, 0] \text{ if male}$$

male   female

(code for all examples is on http://jmcauley.ucsd.edu/cse158/code/week1.py)

# Linearly dependent features

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 5 & 3 & 2 \\ 3 & 3 & 0 \\ 2 & 0 & 2 \end{bmatrix} \begin{matrix} a+b \\ b \\ a \end{matrix}$$

$$rating = 3 + 2(if\ female) + 1(if\ male)$$
$$= 100 - 95(if\ female) - 96(if\ male)$$

# Linearly dependent features

$$\text{rating} = \theta_0 + \theta_1 \; (\text{if female})$$

$$\theta_0 = \text{male rating}$$

$$\theta_1 = \text{how much higher is female rating}$$

How would you build a feature to represent the **month**, and the impact it has on people's rating behavior?

$$rating = \theta_0 + \theta_1 1hf(month)$$

rating

if month=M[j] j = 1...and j-1... M=11

# Exercise

$$\text{rating} = \theta_0 + \theta_1 [is\ jan] + \theta_2 (is\ feb)$$
$$\ldots \quad \theta_{11} [is\ Nov]$$

$$x_i = [1, 0, 0, 0, 1, 0\ 0 \ldots]$$

# What does the data actually look like?

Season vs. rating (overall)

# CSE 158 – Lecture 2
Web Mining and Recommender Systems

Regression Diagnostics

## **Mean-squared error** (MSE)

$$\frac{1}{N}\|y - X\theta\|_2^2$$
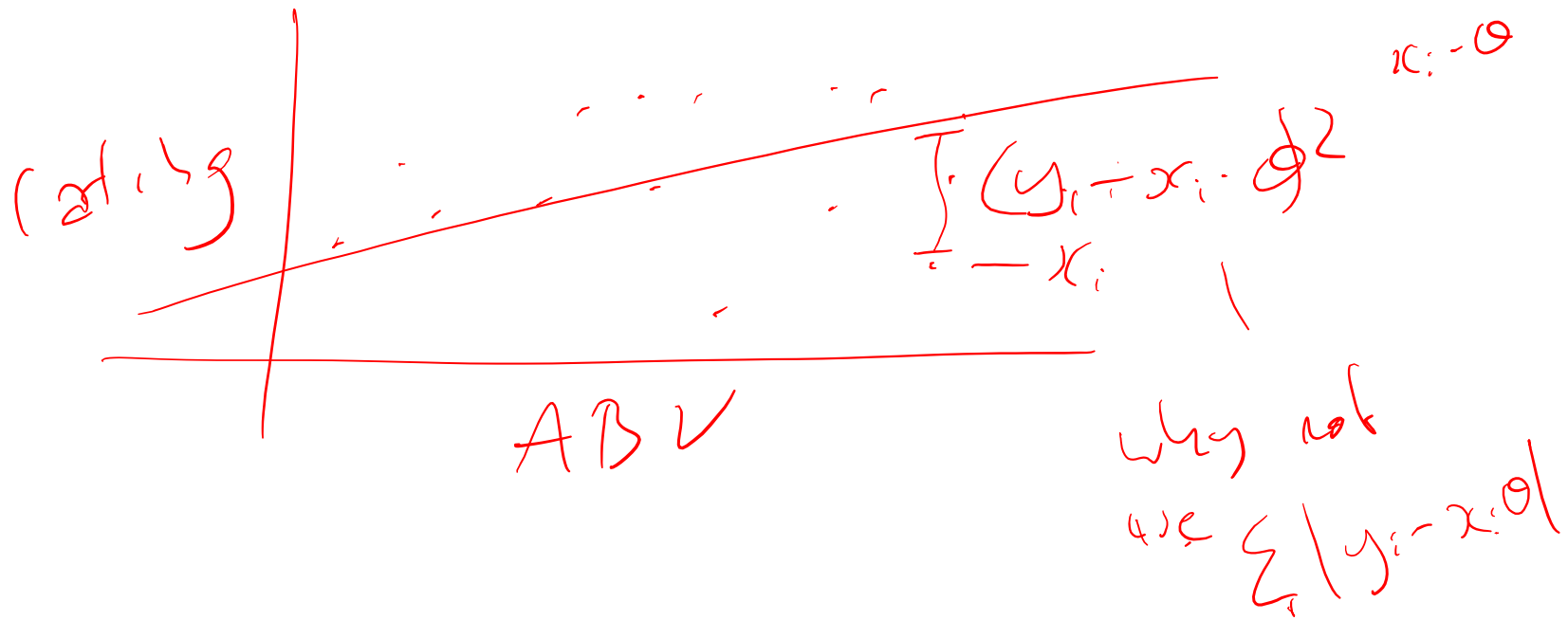
$$= \frac{1}{N}\sum_{i=1}^{N}(y_i - X_i \cdot \theta)^2$$

$$\|x\|_2^2 = \sum_i x_i^2$$

$$\|x\|_1 = \sum_i |x_i|$$

**Q:** Why MSE (and not mean-absolute-error or something else)

# Regression diagnostics

Small errors = common

large errors = <u>very</u> uncommon

$y_i - x_i \theta$

$\mathcal{N}(0, \sigma)$

$y_i = $ prediction error

$= x_i \cdot \theta + \mathcal{N}(0, \sigma)$

$$P_\theta(y \mid X) = \prod_i \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{(y_i - x_i \cdot \theta)^2}{2\sigma^2}}$$

$$\max_\theta P_\theta(y/X) = \prod_i e^{-(y_i - x_i \theta)^2}$$

$$= \min_\theta \sum_i (y_i - x_i \theta)^2$$

# Coefficient of determination

**Q:** How low does the MSE have to be before it's "low enough"?
**A:** It depends! The MSE is proportional to the **variance** of the data

# Regression diagnostics

## Coefficient of determination
## (R^2 statistic)

Mean:

$$\bar{y} = \frac{1}{N} \sum_i y_i$$

Variance:

$$var(y) = \frac{1}{N} \sum_i (y_i - \bar{y})^2$$

MSE:

$$mse = \frac{1}{N} \sum_i (y_i - x_i \cdot \theta)^2$$

# **Coefficient of determination**
## (R^2 statistic)

$$FVU(f) = \frac{MSE(f)}{Var(y)}$$

(FVU = fraction of variance unexplained)

$FVU(f) = 1$ $\longrightarrow$ Trivial predictor

$FVU(f) = 0$ $\longrightarrow$ Perfect predictor

## **Coefficient of determination**
(R^2 statistic)

$$R^2 = 1 - FVU(f) = 1 - \frac{MSE(f)}{Var(y)}$$

R^2  = 0  $\longrightarrow$  Trivial predictor
R^2  = 1  $\longrightarrow$  Perfect predictor

**Q:** But can't we get an R^2 of 1 (MSE of 0) just by throwing in enough random features?

**A:** Yes! This is why MSE and R^2 should always be evaluated on data that **wasn't** used to train the model

A good model is one that **generalizes to new data**

When a model performs well on **training** data but doesn't generalize, we are said to be **overfitting**

When a model performs well on **training** data but doesn't generalize, we are said to be **overfitting**

**Q:** What can be done to avoid overfitting?

# Occam's razor

"Among competing hypotheses, the one with the fewest assumptions should be selected"

$$X\theta = y$$

"hypothesis"

**Q:** What is a "complex" versus a "simple" hypothesis?

# Occam's razor

$$\text{rating} = \Theta_0 + \Theta_1 ABU + \Theta_2 ABU^2 + \dots$$

$$\Theta_{(1)}$$

$$\Theta_{(2)}$$

$$\Theta_{(3)}$$

"less complex"

?

# Occam's razor

**A1:** A "simple" model is one where theta has few non-zero parameters
(only a few features are relevant)

**A2:** A "simple" model is one where theta is almost uniform
(few features are significantly more relevant than others)

# Occam's razor

**A1:** A "simple" model is one where theta has few non-zero parameters $\longrightarrow$ $\|\theta\|_1$ is small

$$\sum_i |\theta_i|$$

**A2:** A "simple" model is one where theta is almost uniform $\longrightarrow$ $\|\theta\|_2$ is small

$$\sum_i \theta_i^2$$

# "Proof"

$$height = \Theta_0 + \Theta_1 age + \Theta_2 shoesize$$

$$\frac{\mid}{a \quad ss}$$

$$\Theta_{(1)}$$

$$\frac{\mid \quad \mid}{a \quad ss}$$

$$\Theta_{(2)}$$

$$\|\Theta_{(1)}\|_1 = \|\Theta_{(2)}\|_1$$

$$\|\Theta_{(1)}\|_2^2 = \|\Theta_{(2)}\|_2^2$$

**Regularization** is the process of penalizing model complexity during training

$$\arg\min_\theta = \frac{1}{N}\|y - X\theta\|_2^2 + \lambda\|\theta\|_2^2$$

MSE

(l2) model complexity

**Regularization** is the process of penalizing model complexity during training

$$\arg\min_\theta = \frac{1}{N}\|y - X\theta\|_2^2 + \lambda\|\theta\|_2^2$$

How much should we trade-off accuracy versus complexity?

# Optimizing the (regularized) model

$$\arg\min_\theta = \frac{1}{N}\|y - X\theta\|_2^2 + \lambda\|\theta\|_2^2$$

$$\underbrace{\phantom{\frac{1}{N}\|y - X\theta\|_2^2 + \lambda\|\theta\|_2^2}}_{f(\theta)}$$

- Could look for a closed form solution as we did before
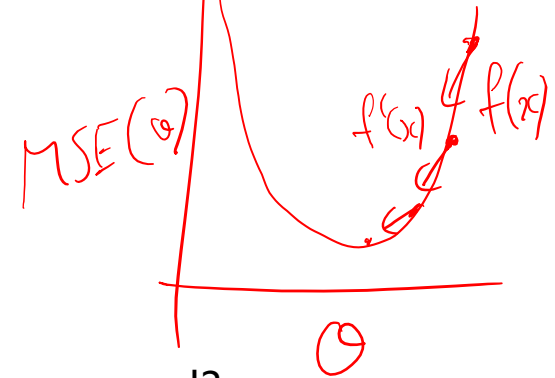- Or, we can try to solve using **gradient descent**

# Gradient descent:

1. Initialize $\theta$ at random
2. While (not converged) do
   $$\theta := \theta - \alpha f'(\theta)$$

All sorts of annoying issues:
- How to initialize theta?
- How to determine when the process has converged?
- How to set the step size alpha

These aren't really the point of this class though

# Optimizing the (regularized) model

$$f(\theta) = \frac{1}{N}\|y - X\theta\|_2^2 + \lambda\|\theta\|_2^2$$

$$\frac{\partial f}{\partial \theta_k} \ ?$$

$$\left( \sum_i \left( y_i - x_i \cdot \theta \right)^2 + \lambda \sum_j \theta_j^2 \right)$$

$$= \sum_i 2 \cdot \theta_k \left( y_i - x_i \cdot \theta \right) + 2\lambda X \theta_k$$

$$\sum_i -2 x_{ik} \left( y_i - x_i \cdot \theta \right) + 2\theta_k$$

# Gradient descent in scipy:

(code for all examples is on http://jmcauley.ucsd.edu/cse158/code/week1.py)

(see "ridge regression" in the "sklearn" module)

# Model selection

$$\arg\min_\theta = \frac{1}{N}\|y - X\theta\|_2^2 + \lambda\|\theta\|_2^2$$

How much should we trade-off accuracy versus complexity?

Each value of lambda generates a different model. **Q:** How do we select which one is the best?

# Model selection

How to select which model is best?

**A1:** The one with the lowest training error?

**A2:** The one with the lowest test error?

We need a **third** sample of the data that is not used for training or testing

# A **validation set** is constructed to "tune" the model's parameters

- Training set: used to **optimize the model's parameters**
- Test set: used to report how well we expect the model to perform on **unseen data**
- Validation set: used to **tune** any model parameters that are not directly optimized

# A few "theorems" about training, validation, and test sets

- The training error **increases** as lambda **increases**
- The validation and test error are at least as large as the training error (assuming infinitely large random partitions)
- The validation/test error will usually have a "sweet spot" between under- and over-fitting

# Model selection

# Summary of Week 1: Regression

- Linear regression and least-squares
- (a little bit of) feature design
- Overfitting and regularization
- Gradient descent
- Training, validation, and testing
- Model selection

# Homework

Homework is **available** on the course webpage
http://cseweb.ucsd.edu/classes/fa17/cse158-a/files/homework1.pdf

Please submit it at the beginning of the **week 3** lecture (Oct 16)

All submissions should be made as **pdf files on gradescope**

# Questions?