# CSE 255 – Lecture 2

Data Mining and Predictive Analytics

## Supervised learning – Regression

**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
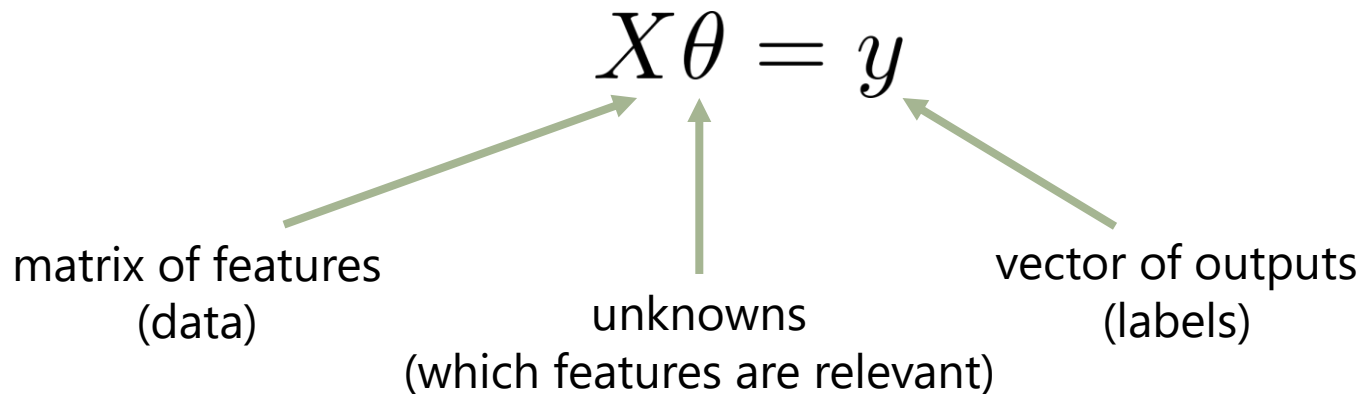
**Regression** is one of the simplest supervised learning approaches to learn relationships  between input variables (features) and output variables (predictions)

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

$$x_i \cdot \theta = y_i$$

$$X\theta = y$$

matrix of features
(data)

unknowns
(which features are relevant)

vector of outputs
(labels)

(or $Ax = b$ if you prefer)

**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?

$$rating = \theta_0 + \theta_1 \, age$$

$$\begin{bmatrix} y_1 \\ y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & age_1 \\ 1 & age_2 \\ 1 & \vdots \\ 1 & age_N \end{bmatrix} \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix}$$

rating

age

# Example 1



**Beers:**

**Ratings/reviews:**

**User profiles:**

# Example 1

50,000 reviews are available on
http://jmcauley.ucsd.edu/cse255/data/beer/beer_50000.json
(see course webpage)

See also – non-alcoholic beers:
http://jmcauley.ucsd.edu/cse255/data/beer/non-alcoholic-beer.json

# Example 1

# Real-valued features

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

(code for all examples is on http://jmcauley.ucsd.edu/cse255/code/week1.py)

# Example 1

## Preferences vs **ABV**



$$\Theta_0 + \Theta_1 \times ABV$$
$$+ \Theta_2 \times ABV^2$$
$$+ \Theta_3 \times ABV^3$$

# Example 1

## Real-valued features

What is the interpretation of:

$$\theta = (3.4, 10e^{-7})$$

$$3.4 \quad + \quad 10e^{-7} \times \text{length seconds}$$

# Example 2

# Categorical features

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

$$r = \Theta_0 + \Theta_1 \times gender$$

$$female = [1] \quad male = [0]$$

$$male = \Theta_0 \quad , \quad female = \Theta_0 + \Theta_1$$

(code for all examples is on http://jmcauley.ucsd.edu/cse255/code/week1.py)

# Linearly dependent features

$Male = [0, 1] \qquad female = [1, 0]$

$Male = \Theta_0 + \Theta_2 \qquad female = \Theta_0 + \Theta_1$

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

$$X^T X = \begin{bmatrix} 7 & 4 & 3 \\ 4 & 4 & 0 \\ 3 & 0 & 3 \end{bmatrix} \begin{matrix} A \\ B \\ A-B \end{matrix}$$
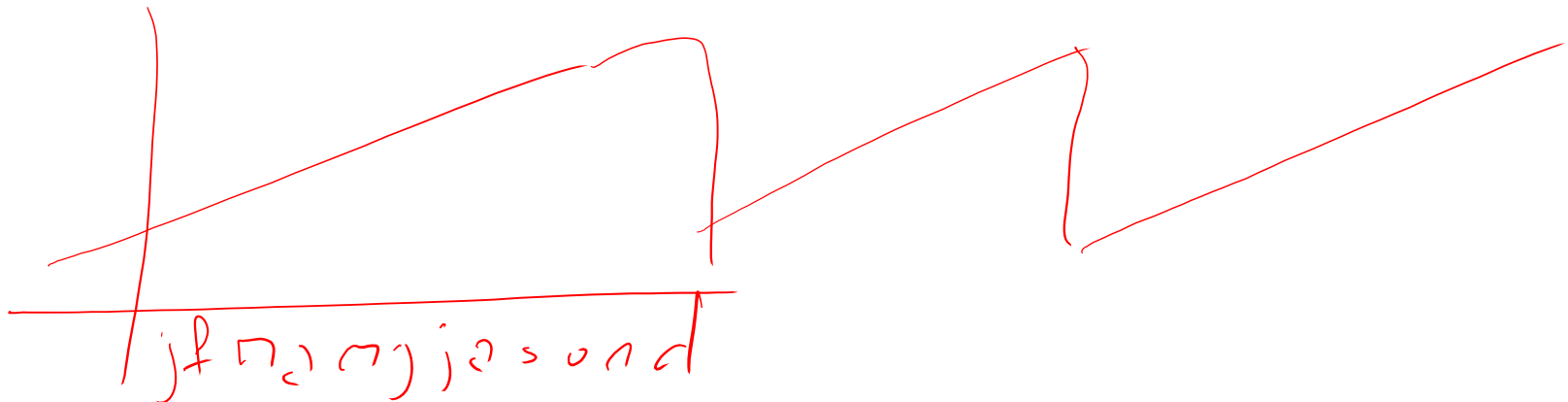
$(X^T X)^{-1} X y$

$r = 2 + 3 [m] + 4 [-f]$

$r = 1 + 4 [m] + 5 [f]$

# Exercise

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

jan = (1)  feb = (2) ...
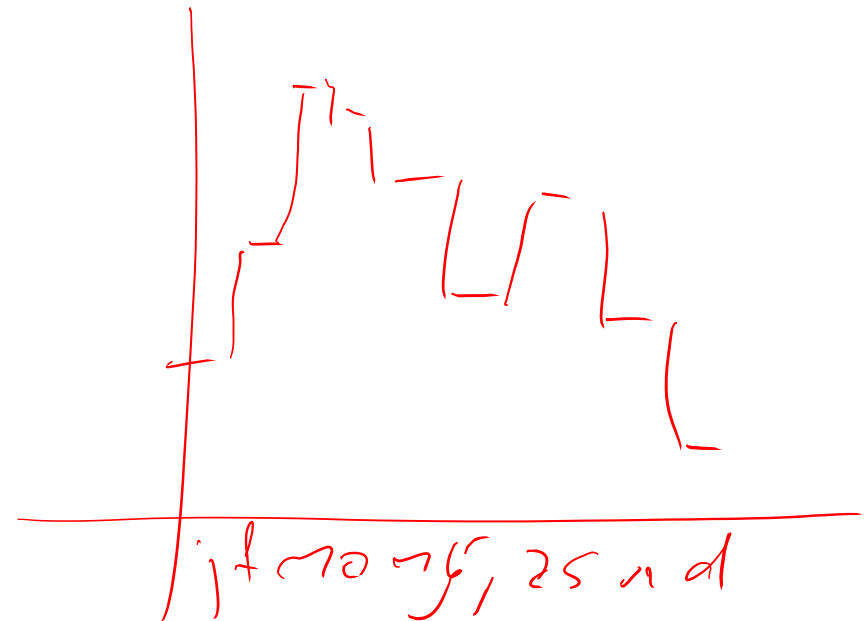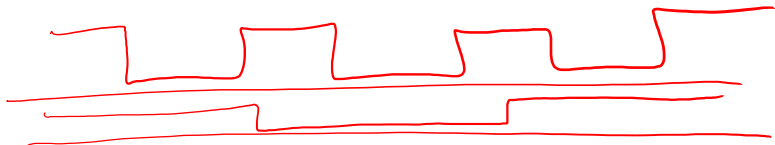
# Exercise

jan [ 1 0 0 0 0 0 0 0 0 0 0 ]

feb [ 0 1 0 0 0 . . . . ]

jan = [ 0 0 0 1 ]
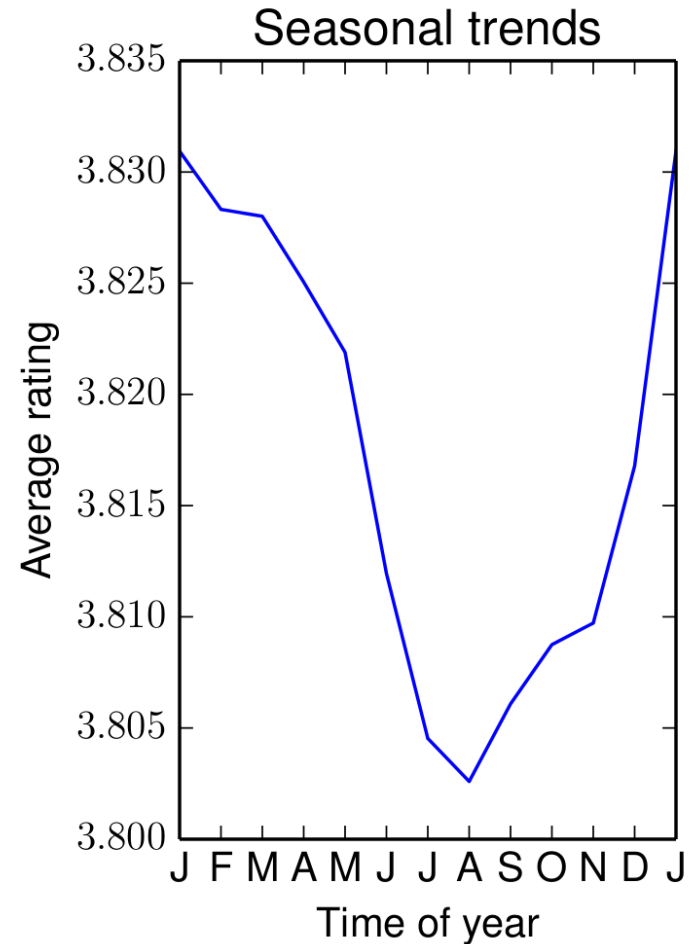
feb = [ 0 0 1 0 ]

j                    d

itmanf, 25 n d

# What does the data actually look like?

Season vs. rating (overall)

# Example 3

## Random features

What happens as we add more and more **random** features?

(code for all examples is on http://jmcauley.ucsd.edu/cse255/code/week1.py)

# CSE 255 – Lecture 2
Data Mining and Predictive Analytics

Regression Diagnostics

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

$$\frac{1}{N}\|y - X\theta\|_2^2$$

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

$$= \frac{1}{N}\sum_i |y_i - X_i \cdot \theta|$$

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

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

$$y_i - x_i \cdot \theta \sim \mathcal{N}(0, \sigma)$$

$$\prod_i \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - x_i \cdot \theta)^2}{2\sigma^2}}$$

$$\sum_i (y_i - x_i \cdot \theta)^2$$

$$(\max)$$
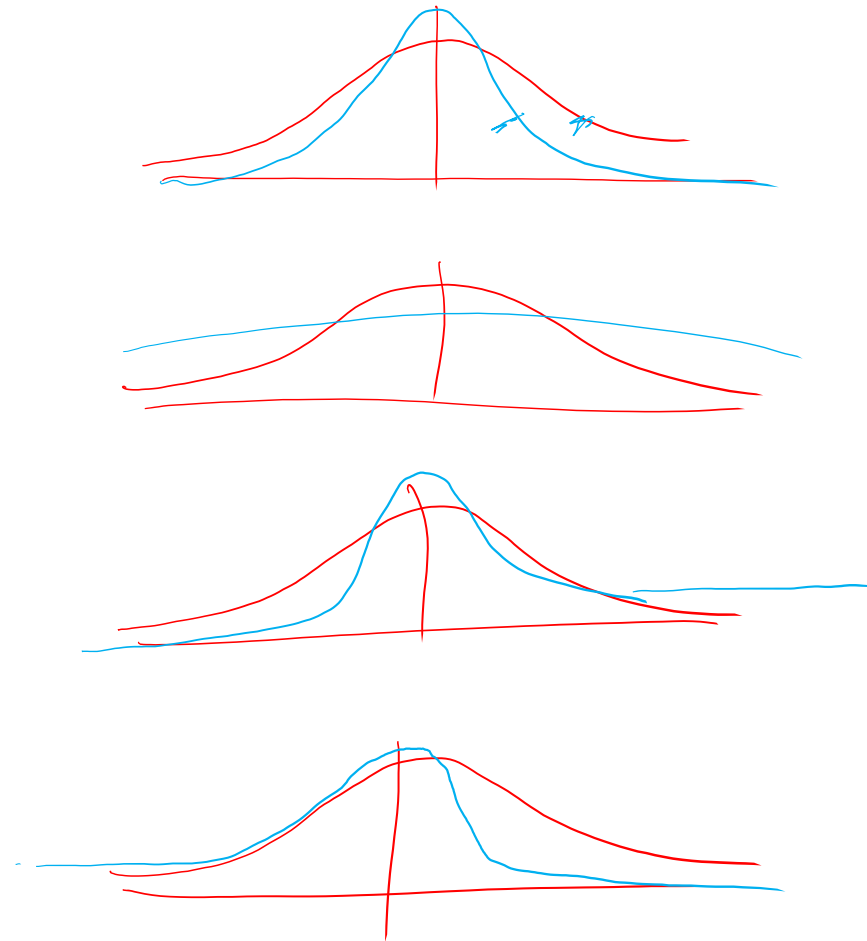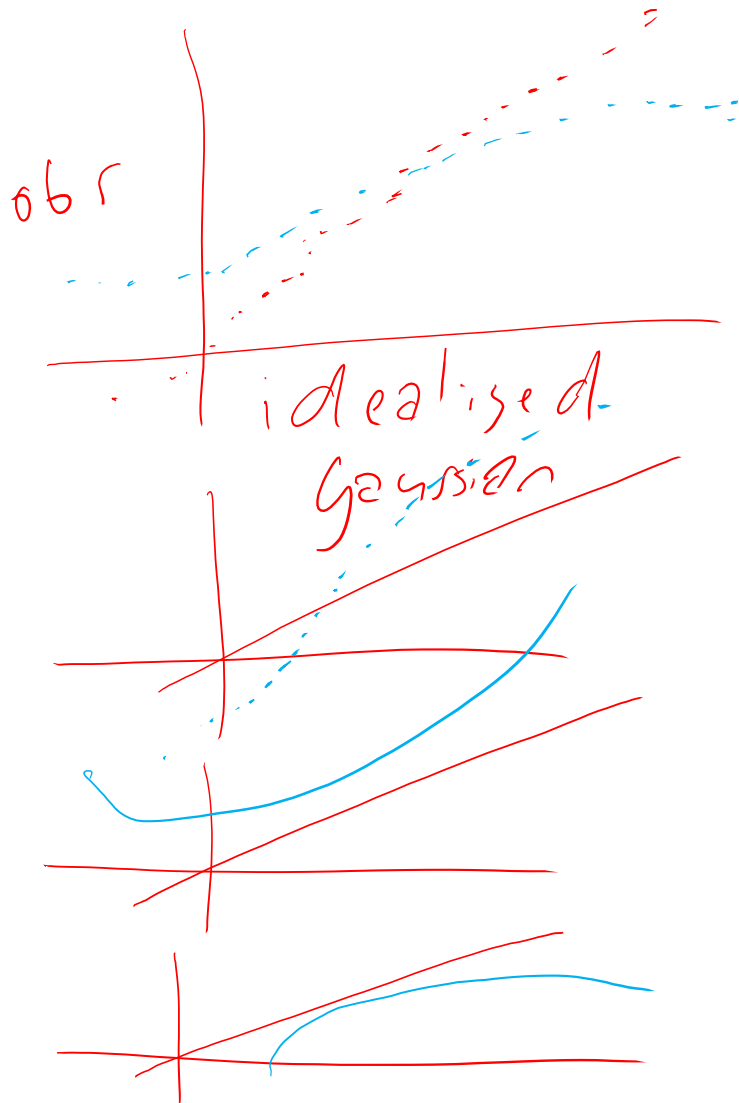
$$(\min)$$

# Regression diagnostics

$$[5.5, -0.5, 3, 1, -2, 4, 0]$$

$$[-2, -0.5, 6, 1, 3, 4, 5.5]$$

$$[-6, -3, -1, 0, 1, 3, 6]$$

# Regression diagnostics



obr

idealised Gaussian

# Regression diagnostics

**Quantile-Quantile (QQ)-plot**



$$\theta_0 + \theta_1 \sqrt{ABU}$$

**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

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

Mean:

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

Variance:

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

MSE:

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

## **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**

**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

**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

$$\|\theta\|_k = \sqrt[k]{\sum_i \theta_i^k}$$

$$\sum_i |\theta_i|$$

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

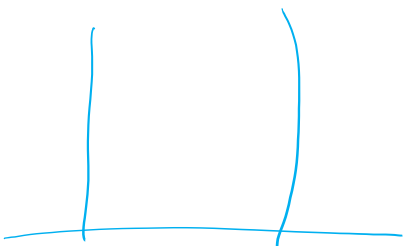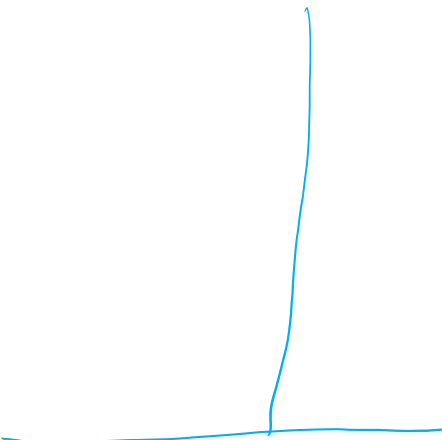**A2:** A "simple" model is one where theta is almost uniform $\longrightarrow$ $\|\theta\|_2^2$ is small

$$\sum_i \theta_i^2$$

# "Proof"

$$\text{height} = \theta_0 + \theta_1 \times \text{weight}$$
$$+ \theta_2 \times \text{shoe size}$$

$\theta_1 \quad \theta_2$

$\theta_1' \quad \theta_2'$

$$\|\theta\|_1 = \|\theta'\|_1$$

$$\|\theta\|_2 > \|\theta'\|_2$$

**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?
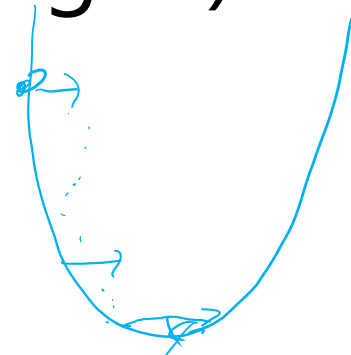
# Optimizing the (regularized) model

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

$$\underbrace{\qquad\qquad\qquad\qquad}_{f(\theta)}$$

- We no longer have a convenient closed-form solution for theta
- Need to resort to some form of approximation algorithm

# 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

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

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

$$f = \frac{1}{N}\sum_i \left(x_i \cdot \theta - y_i\right)^2 + \lambda \sum_k \theta_k^2$$

$$\frac{\partial f}{\partial \theta_k} = \frac{1}{N}\sum_i 2x_{ik}\left(x_i \cdot \theta - y_i\right) + \lambda \cdot 2\theta_k$$

# Gradient descent in scipy:

(code for all examples is on http://jmcauley.ucsd.edu/cse255/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

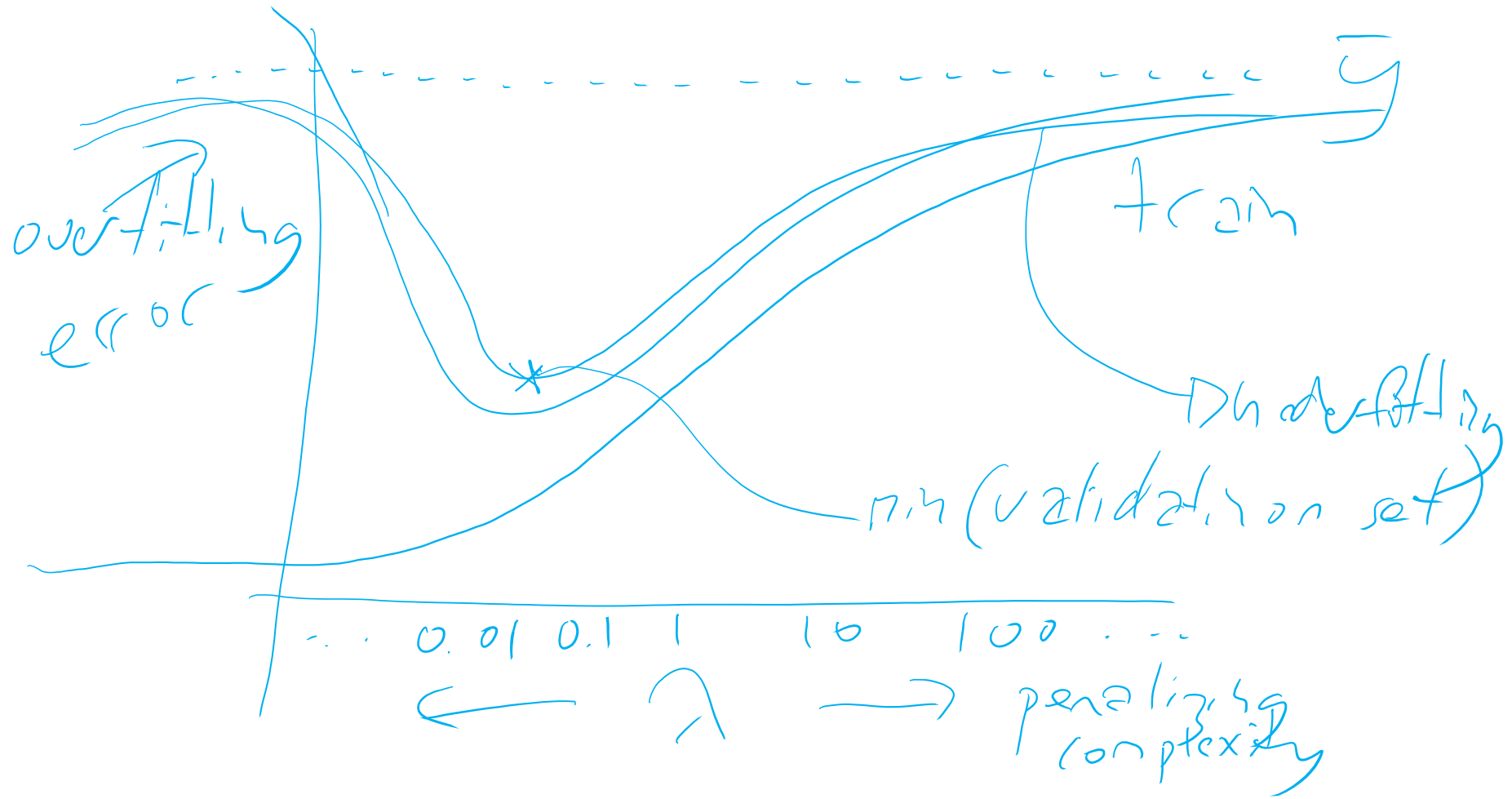# A **validation set** is constructed to "tune" the model's parameters

- 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

overfitting
error

train

Dh underfitting

min (validation set)

... 0.01 0.1 1 10 100 ...

← λ → penalizing complexity

# 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/fa15/cse255-a/files/homework1.pdf

Please submit it at the beginning of the **week 3** lecture (Oct 12)

# Questions?