

# CSCI567 Machine Learning (Fall 2017)

Prof. Fei Sha

U of Southern California

Lecture on Aug. 29, 2017

# Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear regression

# Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear regression

# Administrative stuff

- Please take Syllabus Quiz as well as supplying your Github handle
- Please download the course virtual machine from [http://bytes.usc.edu/files/cs103/install/student-vm\\_2017.ova](http://bytes.usc.edu/files/cs103/install/student-vm_2017.ova).
- Please make sure you attach your @usc email address to your Github handle.

# Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear regression

# Multi-class classification

## Classify data into one of the multiple categories

- Input (feature vectors):  $\mathbf{x} \in \mathbb{R}^D$
- Output (label):  $y \in [C] = \{1, 2, \dots, C\}$
- Learning goal:  $y = f(\mathbf{x})$

## Special case: binary classification

- Number of classes:  $C = 2$
- Labels:  $\{0, 1\}$  or  $\{-1, +1\}$

# Tuning hyperparameter/Model Selection by using a validation dataset

## Training data (set)

- N samples/instances:  $\mathcal{D}^{\text{TRAIN}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- They are used for learning  $f(\cdot)$

## Test (evaluation) data

- M samples/instances:  $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- They are used for assessing how well  $f(\cdot)$  will do in predicting an unseen  $\mathbf{x} \notin \mathcal{D}^{\text{TRAIN}}$

## Development (or validation) data

- L samples/instances:  $\mathcal{D}^{\text{DEV}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_L, y_L)\}$
- They are used to optimize hyperparameter(s).

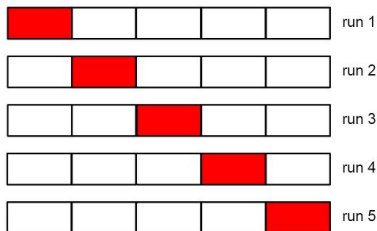
Training data, validation and test data should *not* overlap!

# Cross-validation

## What if we do not have validation data?

- We split the training data into  $S$  equal parts.
- We use each part *in turn* as a validation dataset and use the others as a training dataset.
- We choose the hyperparameter such that *on average*, the model performing the best

$S = 5$ : 5-fold cross validation



*Special case:* when  $S = N$ , this will be leave-one-out.



# Outline

- 1 Administration
- 2 Review of Last Lecture
- 3 Linear regression
  - Motivation
  - Algorithm
  - Univariate solution
  - Multivariate solution in matrix form
  - Computational and numerical optimization
  - Some practical considerations

# Regression

## Predicting a continuous outcome variable

- Predicting a company's future stock price using its past and existing financial information
- Predicting the amount of rain fall
- Predicting ...

# Regression

## Predicting a continuous outcome variable

- Predicting a company's future stock price using its past and existing financial information
- Predicting the amount of rain fall
- Predicting ...

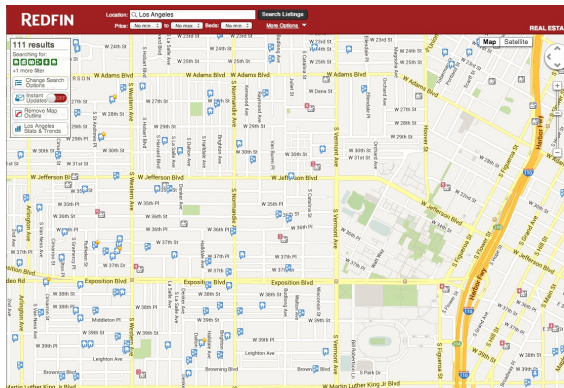
## Key difference from classification

- We measure *prediction errors* differently.
- This will lead to quite different learning models and algorithms.

Ex: be a savvy purchaser by predicting the sale price of a house

Retrieve historical sales records

(This is our training data)



# Features used to predict

[Overview](#)
[Property Details](#)
[Tour Insights](#)
[Property History](#)
[Public Records](#)
[Activity](#)
[Schools](#)

Five unit apartment complex within 2 blocks of USC campus, Gate #6. Great for students (most student leases have parents as guarantors). Most USC students live off campus, so housing units like this are always fully leased. Situated on a gated, corner lot, and across from an elementary school, this complex was recently renovated, and has in-unit laundry hook ups, wall-unit AC, and 12 parking spaces. It is within a DPS (Department of Public Safety) and Campus Cruiser patrolled area. This is a great income generating property, not to be missed!

Property Type: Multi-Family  
Style: Two Level, Low Rise  
Community: Downtown Los Angeles  
County: [Los Angeles](#)  
MLS#: 22176741

## Property Details for 3620 South BUDLONG, Los Angeles, CA 90007

Details provided by i-Tech MLS and may not match the public record. [Learn More](#)

### Interior Features

#### Kitchen Information

- Remodeled
- Oven, Range

#### Multi-Unit Information

- Community Features
  - Units in Complex (Total): 5

#### Multi-Family Information

- # Leased: 5
- # of Buildings: 1
- Owner Pays Water
- Tenant Pays Electricity, Tenant Pays Gas

#### Unit 1 Information

- # of Beds: 2
- # of Baths: 1
- Unfurnished
- Monthly Rent: \$1,700

#### Laundry Information

- Inside Laundry

#### Heating & Cooling

- Wall Cooling Unit(s)

#### Unit 2 Information

- # of Beds: 3
- # of Baths: 1
- Unfurnished
- Monthly Rent: \$2,250

#### Unit 3 Information

- Unfurnished

#### Unit 4 Information

- # of Beds: 3
- # of Baths: 1
- Unfurnished

- Monthly Rent: \$2,350

#### Unit 5 Information

- # of Beds: 3
- # of Baths: 2
- Unfurnished
- Monthly Rent: \$2,325

#### Unit 6 Information

- # of Beds: 3
- # of Baths: 1
- Monthly Rent: \$2,250

### Property / Lot Details

#### Property Features

- Automatic Gate, Card/Code Access

#### Lot Information

- Lot Size (Sq. Ft.): 9,649
- Lot Size (Acres): 0.2215
- Lot Size Source: Public Records

- Automatic Gate, Lawn, Sidewalks
- Corner Lot, Near Public Transit

#### Property Information

- Updated/Remodeled
- Square Footage Source: Public Records

- Tax Parcel Number: 5040017019

### Parking / Garage, Exterior Features, Utilities & Financing

#### Parking Information

- # of Parking Spaces (Total): 12
- Parking Space
- Garage

#### Building Information

- Total Floors: 2

#### Utility Information

- Green Certification Rating: 0.00
- Green Location: Transportation, Walkability
- Green Walk Score: 0
- Green Year Certified: 0

#### Financial Information

- Capitalization Rate (%): 6.25
- Actual Annual Gross Rent: \$126,331
- Gross Rent Multiplier: 11.29

### Location Details, Misc. Information & Listing Information

#### Location Information

- Cross Streets: W 38th Pl

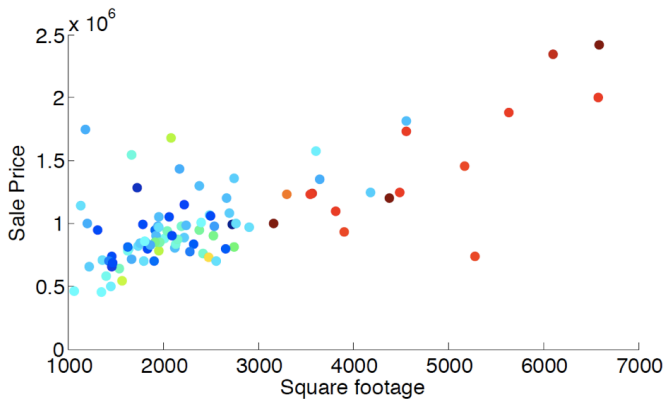
#### Expense Information

- Operating: \$37,664

#### Listing Information

- Listing Terms: Cash, Cash To Existing Loan
- Buyer Financing: Cash

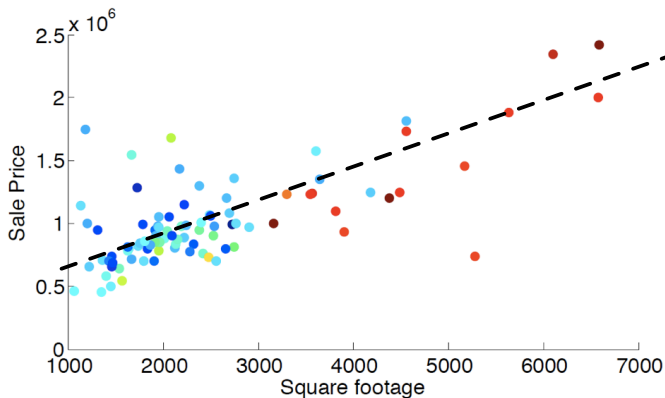
# Correlation between square footage and sale price



(Unlike the Fisher's flower classification example, the colors of the dots in this scatterplot do not mean anything.)

# Possibly linear relationship

Sale price  $\approx$  price\_per\_sqft  $\times$  square\_footage  $+$  fixed\_expense



# How to learn the unknown parameters?

**training data** (past sales record)

sqft	sale price
2000	800K
2100	907K
1100	312K
5500	2,600K
...	...



# Reduce prediction error

## How to measure errors?

- The classification error (got it *right* or *wrong*) is *not appropriate* for continuous outcomes.
- We can look at the *absolute* difference:  $|\text{prediction} - \text{sale price}|$

However, for simplicity, we look at the *squared* errors:  
 $(\text{prediction} - \text{sale price})^2$

sqft	sale price	prediction	error	squared error
2000	810K	720K	90K	8100
2100	907K	800K	107K	$107^2$
1100	312K	350K	38K	$38^2$
5500	2,600K	2,600K	0	0
...	...			

# Minimize squared errors

## Our model

Sale price = price\_per\_sqft  $\times$  square\_footage + fixed\_expense + unexplainable\_stuff

## Training data

sqft	sale price	prediction	error	squared error
2000	810K	720K	90K	8100
2100	907K	800K	107K	$107^2$
1100	312K	350K	38K	$38^2$
5500	2,600K	2,600K	0	0
...	...			
Total				$8100 + 107^2 + 38^2 + 0 + \dots$

## Aim

Adjust price\_per\_sqft and fixed\_expense such that the sum of the squared error is minimized — i.e., the residual/remaining unexplainable\_stuff is minimized.

# Linear regression

## Setup

- Input:  $\mathbf{x} \in \mathbb{R}^D$  (covariates, predictors, features, etc)
- Output:  $y \in \mathbb{R}$  (responses, targets, outcomes, outputs, etc)
- Training data:  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$   
We will use  $x_{nd}$  representing the  $d$ th dimension of the  $n$ th sample  $\mathbf{x}_n$
- Model:  $f : \mathbf{x} \rightarrow y$ , with  $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$ , with  $^T$  standing for vector transpose.  
 $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$  is called *weights*, *parameters*, or *parameter vector*.  $w_0$  is called *bias*.

People also sometimes call  $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$  parameters too!  
And sometimes, people use  $\mathbf{w}$  to mean  $\tilde{\mathbf{w}}$ !

*So please pay attention to contexts when you read papers, textbooks, or assigned reading material.*

# Goal

## Minimize prediction error as much as possible

- Residual Sum of Squares (RSS)

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2$$

- Other definitions of errors are also possible  
We will see an example very soon.

# A simple case: $x$ is just one-dimensional

Our errors are

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + w_1 x_n)]^2$$

Identify stationary points, by taking derivative with respect to parameters, and setting to zeroes

$$\begin{cases} \frac{\partial RSS(\tilde{\mathbf{w}})}{\partial w_0} = 0 \\ \frac{\partial RSS(\tilde{\mathbf{w}})}{\partial w_1} = 0 \end{cases} \Rightarrow \begin{pmatrix} \sum_n 1 & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{pmatrix}$$

# Derivation on the blackboard

# Solution when $x$ is one-dimensional

**Least mean square (LMS) solution (minimizing residual sum of errors)**

$$\begin{pmatrix} \sum_n 1 & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{pmatrix}$$
$$\rightarrow \begin{pmatrix} w_0^{LMS} \\ w_1^{LMS} \end{pmatrix} = \begin{pmatrix} \sum_n 1 & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{pmatrix}$$

**NB.** We sometimes call it least square solutions (LSE) too.

# LMS when $\mathbf{x}$ is D-dimensional

## $RSS(\tilde{\mathbf{w}})$ in matrix form

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{\mathbf{x}} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{\mathbf{w}} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$



# LMS when $x$ is D-dimensional

## $RSS(\tilde{w})$ in matrix form

$$RSS(\tilde{w}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{w}^T \tilde{x}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{x} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{w} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

which leads to

$$RSS(\tilde{w}) = \sum_n (y_n - \tilde{w}^T \tilde{x}_n)(y_n - \tilde{x}_n^T \tilde{w})$$

# LMS when $x$ is D-dimensional

## $RSS(\tilde{w})$ in matrix form

$$RSS(\tilde{w}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{w}^T \tilde{x}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{x} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{w} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

which leads to

$$\begin{aligned} RSS(\tilde{w}) &= \sum_n (y_n - \tilde{w}^T \tilde{x}_n)(y_n - \tilde{x}_n^T \tilde{w}) \\ &= \sum_n \tilde{w}^T \tilde{x}_n \tilde{x}_n^T \tilde{w} - 2y_n \tilde{x}_n^T \tilde{w} + \text{const.} \end{aligned}$$

# LMS when $x$ is D-dimensional

## $RSS(\tilde{w})$ in matrix form

$$RSS(\tilde{w}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{w}^T \tilde{x}_n]^2$$

where we have redefined some variables (by augmenting)

$$\tilde{x} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^T, \quad \tilde{w} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^T$$

which leads to

$$\begin{aligned} RSS(\tilde{w}) &= \sum_n (y_n - \tilde{w}^T \tilde{x}_n)(y_n - \tilde{x}_n^T \tilde{w}) \\ &= \sum_n \tilde{w}^T \tilde{x}_n \tilde{x}_n^T \tilde{w} - 2y_n \tilde{x}_n^T \tilde{w} + \text{const.} \\ &= \left\{ \tilde{w}^T \left( \sum_n \tilde{x}_n \tilde{x}_n^T \right) \tilde{w} - 2 \left( \sum_n y_n \tilde{x}_n^T \right) \tilde{w} \right\} + \text{const.} \end{aligned}$$

# $RSS(\tilde{\mathbf{w}})$ in new notations

## Design matrix and target vector

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times D}, \quad \tilde{\mathbf{X}} = (\mathbf{1} \quad \mathbf{X}) \in \mathbb{R}^{N \times (D+1)}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

# $RSS(\tilde{\mathbf{w}})$ in new notations

## Design matrix and target vector

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times D}, \quad \tilde{\mathbf{X}} = (\mathbf{1} \quad \mathbf{X}) \in \mathbb{R}^{N \times (D+1)}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

## Compact expression

$$RSS(\tilde{\mathbf{w}}) = \left\{ \tilde{\mathbf{w}}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left( \tilde{\mathbf{X}}^T \mathbf{y} \right)^T \tilde{\mathbf{w}} \right\} + \text{const}$$

# Solution in matrix form

## Normal equation

Take derivative with respect to  $\tilde{\mathbf{w}}$

$$\frac{\partial RSS(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} \propto \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \mathbf{w} - \tilde{\mathbf{X}}^T \mathbf{y} = 0$$

This leads to the least-mean-square (LMS) solution

$$\tilde{\mathbf{w}}^{LMS} = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

# Solution in matrix form

## Normal equation

Take derivative with respect to  $\tilde{\mathbf{w}}$

$$\frac{\partial RSS(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} \propto \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \mathbf{w} - \tilde{\mathbf{X}}^T \mathbf{y} = 0$$

This leads to the least-mean-square (LMS) solution

$$\tilde{\mathbf{w}}^{LMS} = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

## Verify the solution when $D = 1$

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \cdots & \cdots \\ 1 & x_N \end{pmatrix} = \begin{pmatrix} \sum_n 1 & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix}$$

*For those who are familiar with this step, you can look up the formula in The Matrix Cookbook*

# Mini-Summary

- Linear regression is the *linear combination of features*.  
 $f : \mathbf{x} \rightarrow y$ , with  $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^T \mathbf{x}$
- If we minimize residual sum squares as our learning objective, we get a *closed-form solution of parameters*.



# Computational complexity

## Bottleneck of computing the solution

$$\tilde{w} = \left( \tilde{X}^T \tilde{X} \right)^{-1} \tilde{X} y$$

is to invert the matrix  $\tilde{X}^T \tilde{X} \in \mathbb{R}^{(D+1) \times (D+1)}$

## How many operations do we need?

- Roughly, on the order of  $O((D+1)^3)$
- Impractical for very large  $D$

We will look at some ideas of addressing this issue later

What if  $\tilde{X}^T \tilde{X}$  is not invertible

Can you think of any reasons why that could happen?

# What if $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ is not invertible

Can you think of any reasons why that could happen?

$N < D + 1$ . **Intuitively, not enough data to estimate all the parameters.**

$D + 1$  unknown but we have only  $N$  training samples

# What if $\tilde{X}^T \tilde{X}$ is not invertible

Can you think of any reasons why that could happen?

$N < D + 1$ . Intuitively, not enough data to estimate all the parameters.

$D + 1$  unknown but we have only  $N$  training samples

**Example:**  $D = 1$ ,  $N = 0$ , or  $1$ , ie, the following “empty” training dataset

sqft	sale price	prediction	error	squared error
1000	2000			

# How about the following?

$$D = 1, N = 2$$

sqft	sale price	prediction	error	squared error
1000	2000			
1000	2000			

# How about the following?

$$D = 1, N = 2$$

sqft	sale price	prediction	error	squared error
1000	2000			
1000	2000			

We still cannot determine the model (uniquely), even now  $N \geq D + 1$ :

- Sale price = sqft  $\times$  2
- Sale price = sqft  $\times$  1 + 1000
- ...

Namely, we need *informative* training data.

# Challenge

Can you summarize those bad scenarios, as illustrated before, with more concise statements about the relationship between training data and the unknown?

*This will be left as a homework exercise.*

# How to solve this problem?

**Intuition:** what does a non-invertible  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  mean?

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbf{U}^T \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r & 0 \\ 0 & \cdots & \cdots & 0 & 0 \end{bmatrix} \mathbf{U}$$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_r > 0$  and  $r < D + 1$ .  $\mathbf{U}$  is unitary matrix.

*Discussion section will talk a bit more about this: this linear algebra step is called eigendecomposition.*



# How to solve this problem?

**Intuition:** what does a non-invertible  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  mean?

$$(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} = \mathbf{U}^T \begin{bmatrix} \lambda_1^{-1} & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2^{-1} & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r^{-1} & 0 \\ 0 & \cdots & \cdots & 0 & \frac{1}{0} \end{bmatrix} \mathbf{U}$$

where  $\frac{1}{0}$  is the issue.

*Discussion section will talk a bit more about this: this linear algebra step is called eigendecomposition.*

# Fix the problem

## Adding something positive

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I} = \mathbf{U}^T \begin{bmatrix} \lambda_1 + \lambda & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 + \lambda & 0 & \cdots & 0 \\ 0 & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & \cdots & \lambda_r + \lambda & 0 \\ 0 & \cdots & \cdots & 0 & \lambda \end{bmatrix} \mathbf{U}$$

where  $\lambda > 0$  and  $\mathbf{I}$  is the identity matrix

*Later, we will justify why this is a sensible thing for us to do*

# Fix the problem

Now we can invert

$$(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I})^{-1} = \mathbf{U}^T \begin{bmatrix} (\lambda_1 + \lambda)^{-1} & 0 & 0 & \dots & 0 \\ 0 & (\lambda_2 + \lambda)^{-1} & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots & 0 \\ 0 & \dots & \dots & (\lambda_r + \lambda)^{-1} & 0 \\ 0 & \dots & \dots & 0 & \frac{1}{\lambda} \end{bmatrix} \mathbf{U}$$

and the solution is

$$\tilde{\mathbf{w}}^{LMS} = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda \mathbf{I} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

*Note that this solution is not the LMS solution to the original problem where the matrix  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  is not invertible.*

# How to choose $\lambda$ ?

Again,  $\lambda$  is a *hyperparameter*, to be distinguished from  $w$ .

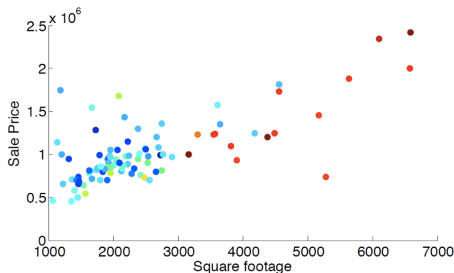
- Use validation or cross-validation
- Other approaches such as Bayesian linear regression — we will describe them briefly later if we have time

# Brain teaser for Linear Regression

**What if  $D = 0$ , ie, not using any predictor/features?**

# What the model looks like when $D = 0$

Sale price = fixed\_expense + unexplainable\_stuff, namely  $f(x) = w_0$



So this is a horizontal line...But where this line should be vertically (ie, what is  $w_0$ )?

# Intuition: the average of the all the sale prices in the training data

From  $D = 1$

$$\begin{pmatrix} \sum_n 1 & \sum_n x_n \\ \sum_n x_n & \sum_n x_n^2 \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} = \begin{pmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{pmatrix}$$

to  $D = 0$

$$\sum_n 1 \times w_0 = \sum_n y_n \rightarrow w_0 = \frac{1}{N} \sum_n y_n$$

In other words, when we say “The housing in City A is more expensive than it in City B”, we are just referring to comparing the bias terms  $w_0$ , ie, the average price in each city, without taking into consideration other features of each individual property.

# linear regression versus nearest neighbors

## Parametric versus non-parametric

- Parametric

The size of the model does *not grow* with respect to the size of the training dataset  $N$ .

In linear regression, there are  $D + 1$  parameters, irrelevant to how many training instances we have.

- Non-parametric

The size of the model *grows* with respect to the size of the training dataset.

In nearest neighbor classification, the training dataset itself needs to be kept in order to make prediction. Thus, the size of the model is the size of the training dataset.

Non-parametric does *not* mean *parameter-less*. It just means the number of parameters is a function of the training dataset.