

# MSCI 546: Advanced Machine Learning

Sirisha Rambhatla

University of Waterloo

Lecture  
Linear Regression

# Outline

- 1 Linear regression
- 2 Reading

# Outline

- 1 Linear regression
- 2 Reading

# Regression

## Predicting a continuous outcome variable using past observations

- Predicting future temperature
- Predicting the amount of rainfall
- Predicting the demand of a product
- Predicting the sale price of a house
- ...

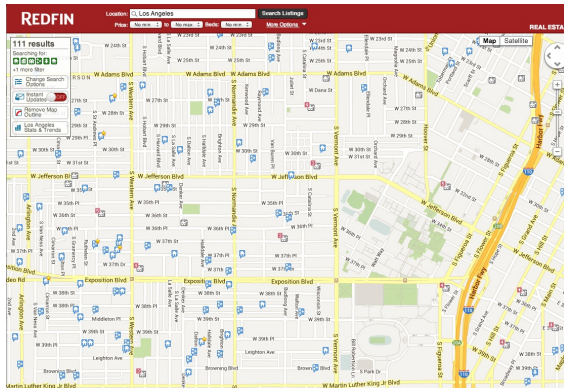
## Key difference from classification

- continuous vs discrete
- measure *prediction errors* differently.
- lead to quite different learning algorithms.

**Linear Regression:** regression with linear models

# Ex: Predicting the sale price of a house

## Retrieve historical sales records (training data)



# Features used to predict

**3620 South BUDLONG**  
Los Angeles, CA 90007  
Status: Closed

**\$1,510,000**  
Last Sold Price

**14**  
Beds

**6**  
Baths

**4,418** Sq. Ft.  
\$342 / Sq. Ft.

Built: 1956 Lot Size: 9,649 Sq. Ft. Sold On: Jul 26, 2013

Overview Property Details Tour Insights Property History Public Records Activity Schools

1 of 12

Five unit apartment complex within 2 blocks of USC campus, Gate #8. 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 quiet, 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  
Community: Downtown Los Angeles  
MLS#: 22176741

Style: Two Level, Low Rise  
County: Los Angeles

## 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

#### Laundry Information

- Inside Laundry

#### Heating & Cooling

- Wall Cooling Unit(s)

### 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

#### 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

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

- Tax Parcel Number: 5040017019

#### Lot Information

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

#### Property Information

- Updated/Remodeled
- Square Footage Source: Public Records

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

#### Parking Information

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

#### Building Information

- Total Floor: 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 98th Pl

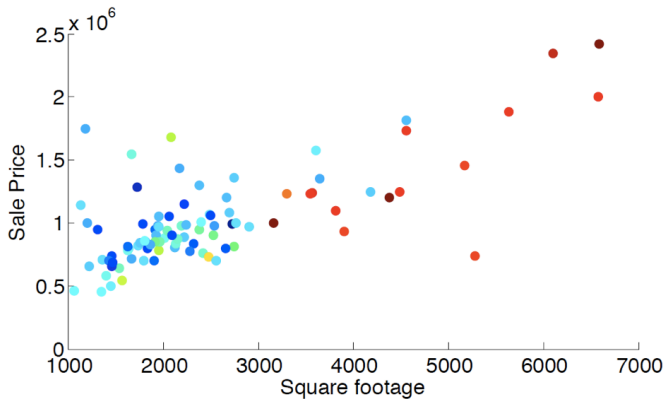
#### Expense Information

- Operating: \$37,664

#### Listing Information

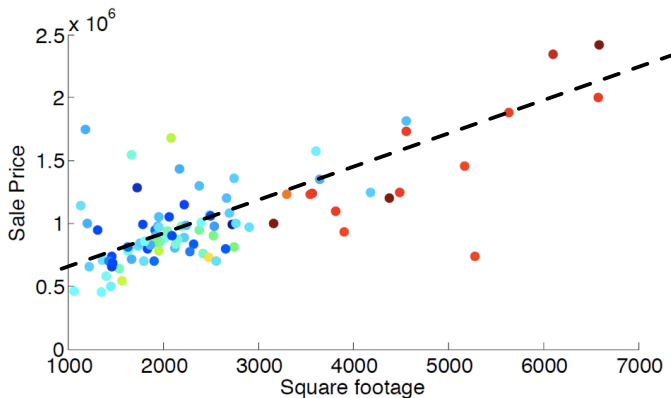
- Listing Terms: Cash, Cash To Existing Loan
- Buyer Financing: Cash

# Correlation between square footage and sale price



# Possibly linear relationship

Sale price  $\approx$  **price\_per\_sqft**  $\times$  square\_footage + **fixed\_expense**  
(*slope*) (*intercept*)





# How to learn the unknown parameters?

## How to measure error for one prediction?

- The classification error (0-1 loss, i.e. *right* or *wrong*) is *inappropriate* for continuous outcomes.
- We can look at
  - *absolute* error:  $|\text{prediction} - \text{sale price}|$
  - or *squared* error:  $(\text{prediction} - \text{sale price})^2$  (**most common**)

**Goal: pick the model (unknown parameters) that minimizes the average/total prediction error, but *on what set*?**

- test set, ideal but we *cannot use test set while training*
- training set ✓

# Formal setup for linear regression

**Input:**  $\mathbf{x} \in \mathbb{R}^D$  (features, covariates, context, predictors, etc)

**Output:**  $y \in \mathbb{R}$  (responses, targets, outcomes, etc)

**Training data:**  $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

**Linear model:**  $f : \mathbb{R}^D \rightarrow \mathbb{R}$ , with  $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^\top \mathbf{x}$   
(superscript  $^\top$  stands for transpose), i.e. a *hyper-plane* parametrized by

- $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^\top$  (weights, weight vector, parameter vector, etc)
- bias  $w_0$

**NOTE:** for notation convenience, very often we

- append 1 to each  $x$  as the first feature:  $\tilde{\mathbf{x}} = [1 \ x_1 \ x_2 \ \dots \ x_D]^\top$
- let  $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \dots \ w_D]^\top$ , a concise representation of all  $D + 1$  parameters. For ease of notation, we will drop the  $(\tilde{\cdot})$  and use  $D$  to subsume the constant term.
- the model becomes simply  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  (but don't forget that the bias term is in there!)

# Goal

Minimize total squared error (note that  $\mathbf{x}_n^\top \mathbf{w} = \mathbf{w}^\top \mathbf{x}_n$ )

- **Residual Sum of Squares** (RSS), a function of  $\mathbf{w}$

$$\text{RSS}(\mathbf{w}) = \sum_n (f(\mathbf{x}_n) - y_n)^2 = \sum_n (\mathbf{x}_n^\top \mathbf{w} - y_n)^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

where

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

# Goal

Minimize total squared error (note that  $\mathbf{x}_n^\top \mathbf{w} = \mathbf{w}^\top \mathbf{x}_n$ )

- **Residual Sum of Squares** (RSS), a function of  $\mathbf{w}$

$$\text{RSS}(\mathbf{w}) = \sum_n (f(\mathbf{x}_n) - y_n)^2 = \sum_n (\mathbf{x}_n^\top \mathbf{w} - y_n)^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

- find  $\mathbf{w}^* = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{argmin}} \text{RSS}(\mathbf{w})$ , i.e. **least (mean) squares solution**  
(more generally called **empirical risk minimizer**)
- in principle can apply any optimization algorithm, but linear regression admits a *closed-form solution*

# General least square solution

## Objective

$$\text{RSS}(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

$$\begin{aligned}\text{RSS}(\mathbf{w}) &= \sum_n (\mathbf{x}_n^\top \mathbf{w} - y_n)^2 = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \\ &= (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{y}^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}\end{aligned}$$

# General least square solution

## Find stationary points (Matrix Calculus)

$$\begin{aligned}\nabla \text{RSS}(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} \left( \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{y}^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \right) \\ &= \left( \mathbf{X}^\top \mathbf{X} \mathbf{w} + \left( \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \right)^\top \right) - (\mathbf{y}^\top \mathbf{X})^\top - \mathbf{X}^\top \mathbf{y} + 0 \\ &= 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y}\end{aligned}$$

Setting the gradient to zero

$$\begin{aligned}2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} &= 0 \implies \mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{y} \\ \mathbf{w}^* &= \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

assuming  $\mathbf{X}^\top \mathbf{X}$  is invertible for now. By convexity  $\mathbf{w}^*$  is the minimizer of RSS.

# Computational complexity

**Bottleneck** of computing

$$\mathbf{w}^* = \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

is to invert the matrix  $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{D \times D}$

- aka *pseudo-inverse*\* denoted by  $(\cdot)^\dagger$ , i.e.  $\mathbf{X}^\dagger = \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top$
- naively need  $O(D^3)$  time
- there are many faster approaches

\*see [https://en.wikipedia.org/wiki/Moore-Penrose\\_inverse](https://en.wikipedia.org/wiki/Moore-Penrose_inverse)

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

## What does that imply?

Recall  $(\mathbf{X}^\top \mathbf{X}) \mathbf{w}^* = \mathbf{X}^\top \mathbf{y}$ . If  $\mathbf{X}^\top \mathbf{X}$  not invertible, this equation aka *Normal Equations* has

- infinitely many solutions ( $\Rightarrow$  infinitely many minimizers)
- This is because *Normal Equations* are always *consistent*<sup>†</sup> meaning a solution *always* exists! It may not be unique though.

<sup>†</sup>See <https://sites.math.washington.edu/~burke/crs/308/LeastSquares.pdf>



# How to resolve this issue?

**Intuition:** what does inverting  $\mathbf{X}^\top \mathbf{X}$  do?

**eigendecomposition:**  $\mathbf{X}^\top \mathbf{X} = \mathbf{U}^\top \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \lambda_D \end{bmatrix} \mathbf{U}$

where  $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_{D+1} \geq 0$  are **eigenvalues**.

**inverse:**  $(\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{U}^\top \begin{bmatrix} \frac{1}{\lambda_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \frac{1}{\lambda_D} \end{bmatrix} \mathbf{U}$

*i.e. just inverse of the eigenvalues*

## How to solve this problem?

Non-invertible  $\Rightarrow$  some eigenvalues are 0.

**One natural fix: add something positive**

$$\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} = \mathbf{U}^\top \begin{bmatrix} \lambda_1 + \lambda & 0 & \cdots & 0 \\ 0 & \lambda_2 + \lambda & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \lambda_D + \lambda \end{bmatrix} \mathbf{U}$$

where  $\lambda > 0$  and  $\mathbf{I}$  is the identity matrix. Now it is invertible:

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} = \mathbf{U}^\top \begin{bmatrix} \frac{1}{\lambda_1 + \lambda} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2 + \lambda} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \frac{1}{\lambda_D + \lambda} \end{bmatrix} \mathbf{U}$$

## Fix the problem

The solution becomes

$$\mathbf{w}^* = \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- not a minimizer of the original RSS

This in fact comes from minimizing **regularized** RSS aka **Ridge Regression**!

$$\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$\lambda$  is a *hyper-parameter*, can be tuned by cross-validation.

# Outline

- 1 Linear regression
- 2 Reading

# Reading

## PiML1

- Chapter 11 – 11.1, 11.2 (11.2.1, 11.2.2), 11.3 (11.3.1)