

DS 4400

Machine Learning and Data Mining I  
Spring 2021

Alina Oprea  
Associate Professor  
Khoury College of Computer Science  
Northeastern University

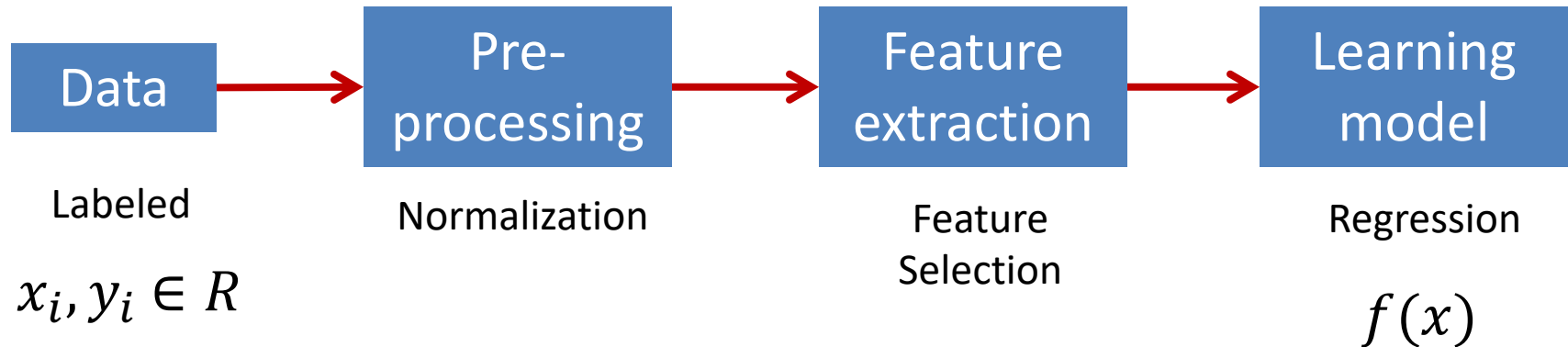
February 2 2021

# Today's Outline

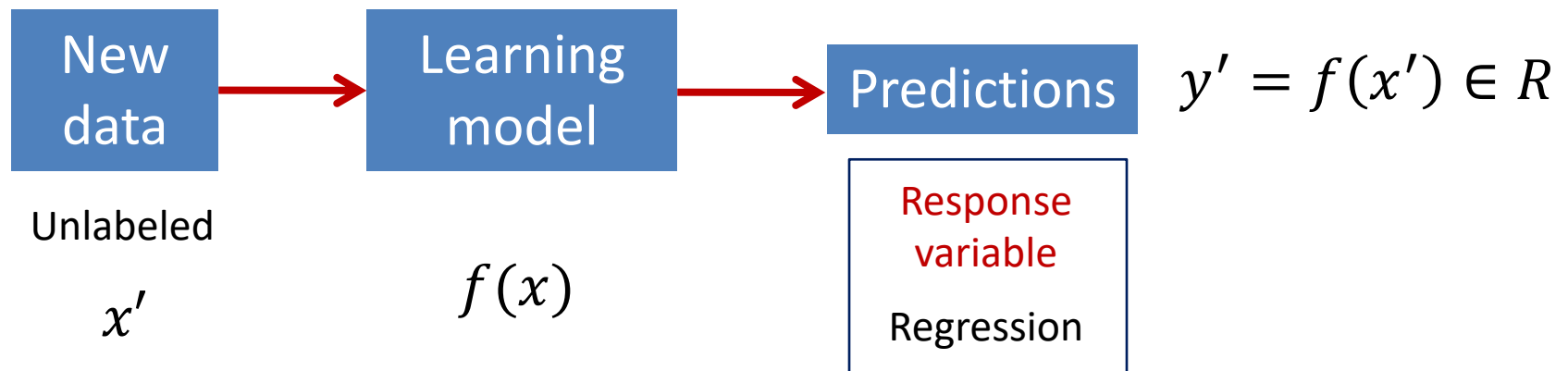
- Announcements
  - First tutorial (recording and notebook) available on Canvas
  - Panda tutorial by Omkar on Wed at 4pm
- Multiple Linear Regression
  - Vector and matrix gradients
  - Closed-form solution derivation
- Lab in Python
  - Simple LR
  - Multiple LR
- Practical issues when training LR models

# Supervised Learning: Regression

## Training



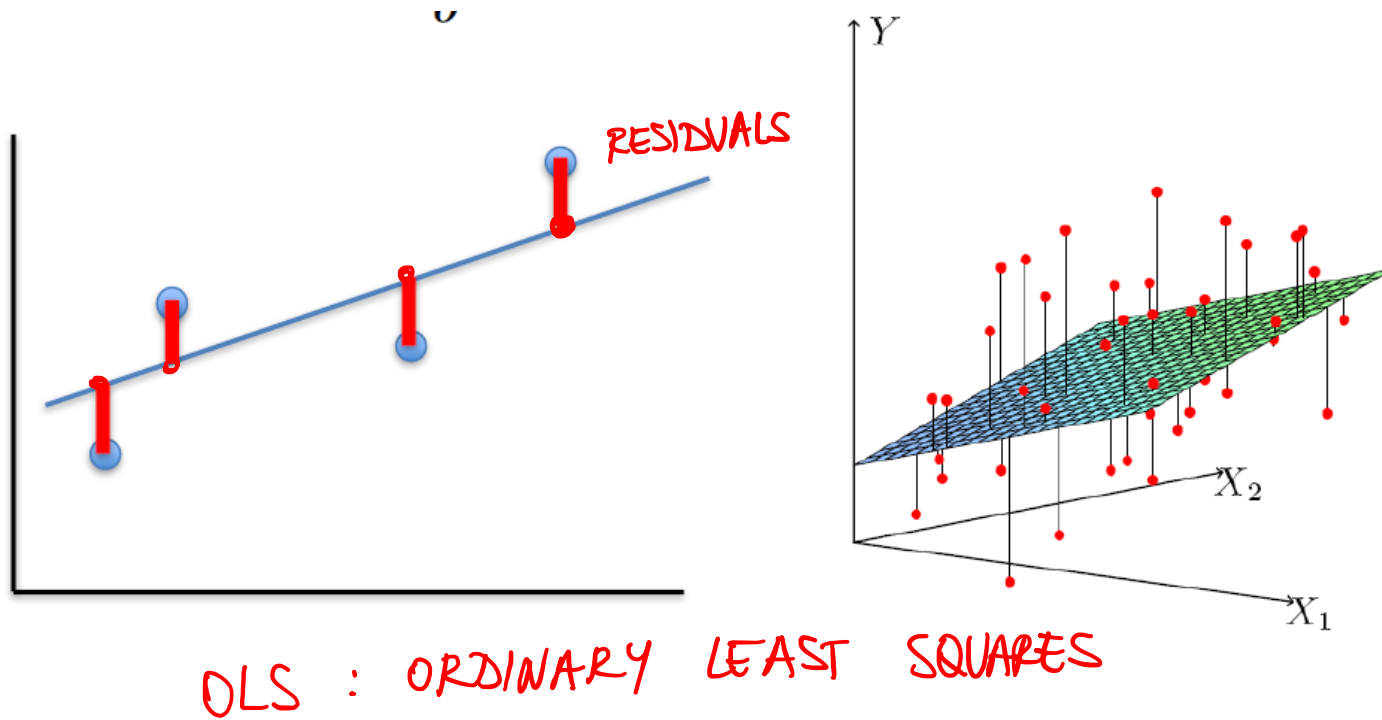
## Testing



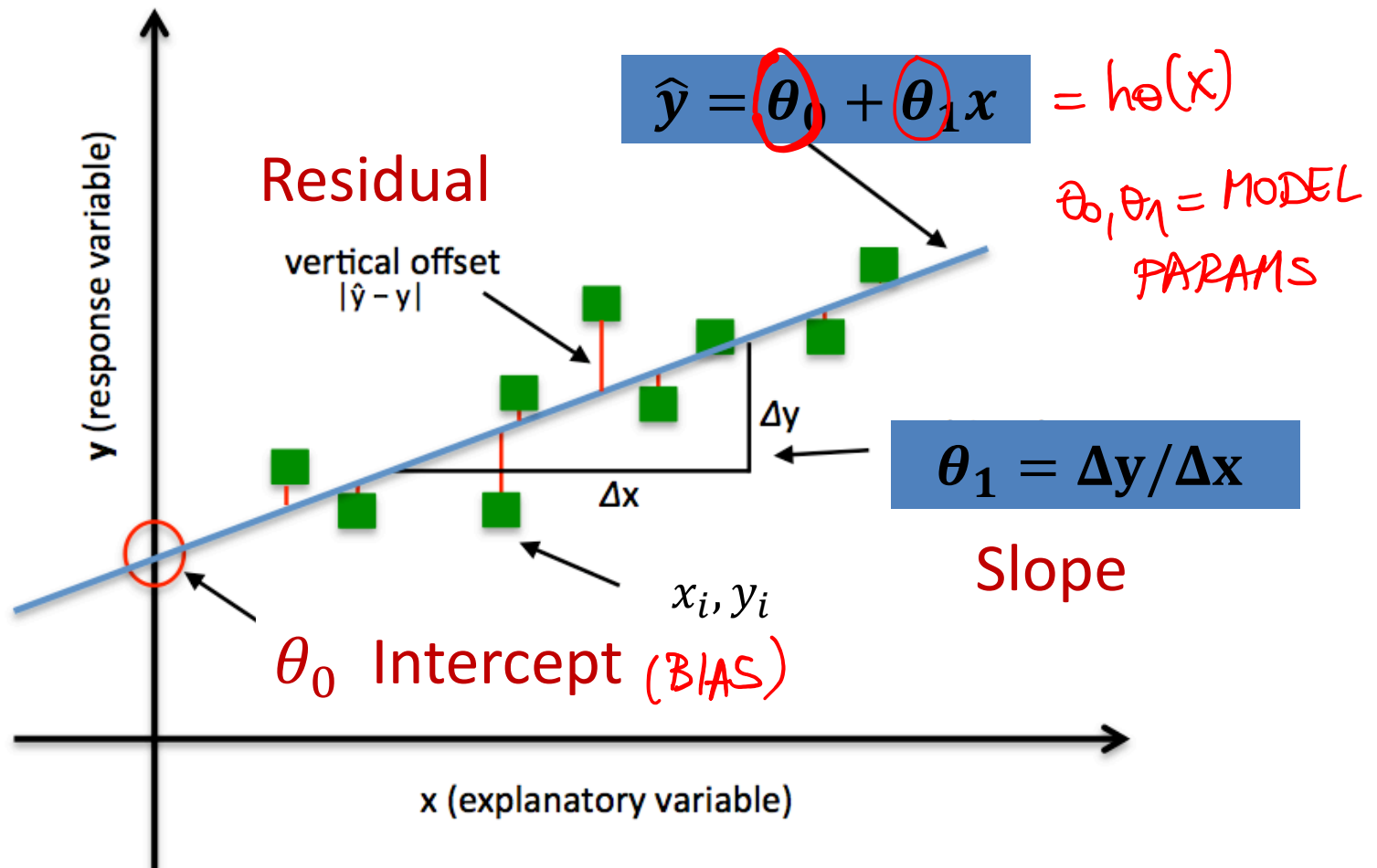
# Least-Squares Linear Regression

$$\text{MSE: } J(\theta) = \frac{1}{N} \sum_{i=1}^N \underbrace{[h_{\theta}(x_i) - y_i]^2}_{\text{RESIDUAL}}$$

$x_i, y_i$  TRAINING DATA  
↓  
FEATURES      RESPONSE VARIABLE



# Interpretation



$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N [h_\theta(x_i) - y_i]^2$$

# Solution for simple linear regression

TRAINING  $i=1, \dots, N$

- Dataset  $x_i \in \mathbb{R}, y_i \in \mathbb{R}, h_{\theta}(x) = \theta_0 + \theta_1 x$

Min •  $J(\theta) = \frac{1}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i)^2$  **MSE / Loss**

$$\left[ \frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{N} \sum_{i=1}^N (\theta_0 + \theta_1 x_i - y_i) = 0 \right.$$

$$\left. \frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{N} \sum_{i=1}^N x_i (\theta_0 + \theta_1 x_i - y_i) = 0 \right]$$

- Solution of min loss

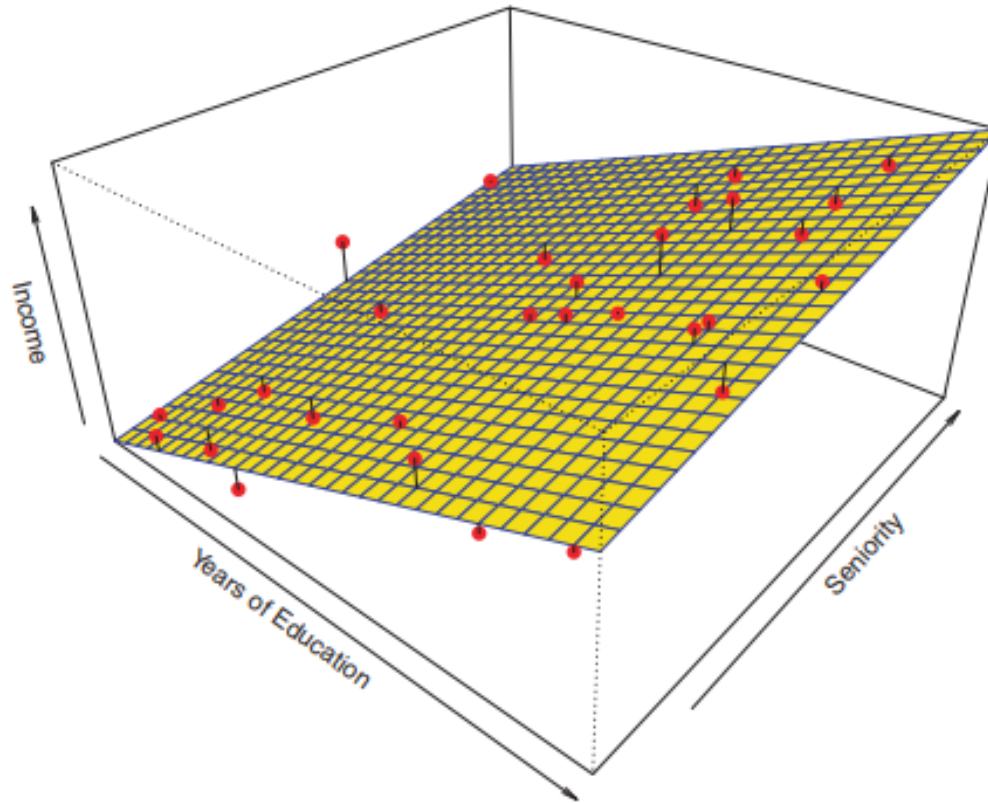
$$-\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$-\theta_1 = \frac{\frac{1}{N-1} \sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{N-1} \sum (x_i - \bar{x})^2} = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

$\text{Var}(x)$

$$\bar{x} = \frac{\sum_{i=1}^N x_i}{N}$$
$$\bar{y} = \frac{\sum_{i=1}^N y_i}{N}$$

# Multiple Linear Regression



- Linear Regression with at least 2 predictors  $d \geq 2$
- Dataset:  $x_i \in R^d, y_i \in R$

# Vector Norms

**Vector norms:** A norm of a vector  $\|x\|$  is informally a measure of the “length” of the vector.

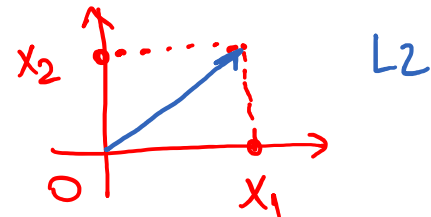
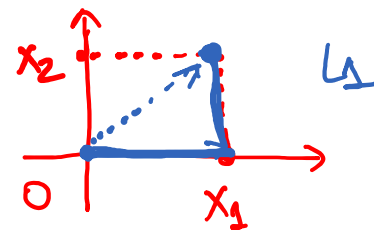
$$\hookrightarrow \|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

– Common norms:  $L_1$ ,  $L_2$  (Euclidean)

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

–  $L_\infty$

$$\|x\|_\infty = \max_i |x_i|$$



EUCLIDEAN

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2}$$



# Vector products

We will use lower case letters for vectors

The elements are referred by  $x_i$ .

- Vector dot (inner) product:

$$x^T y \in \mathbb{R} = \underbrace{\begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}}_{\text{row, size: } 1 \times n} \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\text{column, size: } n \times 1} = \sum_{i=1}^n x_i y_i \in \mathbb{R}$$

- Vector outer product:

$$xy^T \in \mathbb{R}^{m \times n} = \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}}_{\text{column, size: } m \times 1} \underbrace{\begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix}}_{\text{row, size: } 1 \times n} = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}_{\text{size: } m \times n}$$

# Hypothesis Multiple LR

SIMPLE LR:  $h_{\theta}(x) = \theta_0 + \theta_1 x$

MULTIPLE LR:  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$  (\*)

$$X = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

↑ → BIAS

$d = \# \text{ features}$

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$d+1$

MODEL PARAMETER

(\*)  $h_{\theta}(x) = \theta^T x$

DOT PRODUCT OF  $\theta$  and  $x$

# Training data

$N$  training examples

$x_i = (x_{i1}, x_{i2}, \dots, x_{id})$   $d$  features

$y_1, \dots, y_N$  RESPONSE VAR

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1j} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2j} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \dots & x_{Nj} & \dots & x_{Nd} \end{bmatrix}$$

TRAINING EXAMPLE  $i$

FEATURE  $j$

SIZE:  $N \times (d+1)$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

SIZE  $N \times 1$

# Use Vectorization

- $\hat{y}_i = h_{\theta}(x_i) = \sum_{j=0}^d \theta_j x_{ij} = \theta^T x_i$

$\hat{y}_i$  = PREDICTED  
RESPONSE

MATRIX:  $X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{N1} & \dots & x_{Nd} \end{pmatrix}$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$$X\theta = \begin{bmatrix} \theta_0 + \theta_1 x_{11} + \dots + \theta_d x_{1d} \\ \vdots \\ \theta_0 + \theta_1 x_{N1} + \dots + \theta_d x_{Nd} \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \hat{y}$$

size  $N \times 1$

- $h_{\theta}(x) = X\theta$

$$\hat{y} = X\theta$$

VECTORIZED FORM

# Loss function MSE

- For the linear regression cost function:

MSE: 
$$J(\theta) = \frac{1}{N} \sum_{i=1}^N [h_{\theta}(x_i) - y_i]^2$$

$$y - \hat{y} = \begin{bmatrix} y_1 - \hat{y}_1 \\ \vdots \\ y_N - \hat{y}_N \end{bmatrix}$$

$$= \frac{1}{N} \sum_{i=1}^N [\hat{y}_i - y_i]^2$$

$$\|y - \hat{y}\|_2 = \sqrt{(y_1 - \hat{y}_1)^2 + \dots + (y_N - \hat{y}_N)^2}$$

$$= \frac{1}{N} \|\hat{y} - y\|_2^2$$

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

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}$$

# Matrix and vector gradients

- If  $y = f(x)$ ,  $y \in R$  scalar,  $x \in R^n$  vector

$$\frac{\partial f(x)}{\partial x} = \left[ \frac{\partial f(x)}{\partial x_1} \quad \dots \quad \frac{\partial f(x)}{\partial x_n} \right]$$

VECTOR GRADIENT  
ROW VECTOR.

If  $y = f(x)$ ,  $y \in R^m$ ,  $x \in R^n$

$$\frac{\partial f(x)}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

MATRIX GRADIENT  
JACOBIAN MATRIX

SIZE  $m \times n$

# Properties

- 1) • If  $w, x$  are  $(d \times 1)$  vectors,  $\frac{\partial w^T x}{\partial x} = w^T$
- 2) • If  $A: (n \times d) \ x: (d \times 1)$ ,  $\frac{\partial Ax}{\partial x} = A$
- If  $A: (d \times d) \ x: (d \times 1)$ ,  $\frac{\partial x^T Ax}{\partial x} = (A + A^T)x$
- If  $A$  symmetric:  $\frac{\partial x^T Ax}{\partial x} = 2Ax$
- ~~3~~ 5) • If  $x: (d \times 1)$ ,  $\frac{\partial ||x||^2}{\partial x} = 2x^T$

$$\textcircled{1} \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}; \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\bullet \quad w^T x = w_1 x_1 + \dots + w_d x_d \in \mathbb{R}$$

$$\frac{\partial w^T x}{\partial x} = \left[ \frac{\partial w^T x}{\partial x_1} \quad \dots \quad \frac{\partial w^T x}{\partial x_d} \right] = [w_1 \quad w_2 \dots w_d] = w^T$$

$$\textcircled{5} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\|x\|^2 = x_1^2 + \dots + x_d^2 \in \mathbb{R}$$

$$\frac{\partial \|x\|^2}{\partial x} = \left[ \frac{\partial \|x\|^2}{\partial x_1} \quad \dots \quad \frac{\partial \|x\|^2}{\partial x_d} \right] = [2x_1 \quad 2x_2 \dots 2x_d] = 2x^T$$



# Min loss function

– Notice that the solution is when  $\frac{\partial}{\partial \theta} J(\theta) = 0$

→  $X$ : SIZE  $N \times (d+1)$

$\theta$ : SIZE  $(d+1) \times 1$

$y$ : SIZE  $N \times 1$

$$J(\theta) = \frac{1}{N} \|\underbrace{X\theta - y}\|^2 \in \mathbb{R}$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{N} \cdot 2 \cdot (X\theta - y)^T \cdot \frac{\partial (X\theta - y)}{\partial \theta} \quad \text{CHAIN RULE}$$

$$= \frac{2}{N} (X\theta - y)^T \cdot X = 0$$

$$(X\theta - y)^T \cdot X = 0$$

TRANSPOSE

$$\begin{cases} (A^T)^T = A \\ (A \cdot B)^T = B^T A^T \end{cases}$$

$$X^T (X\theta - y) = 0 \Rightarrow (X^T X) \theta = X^T y$$

CLOSED FORM OPTIMAL  
SOLUTION TO MIN MSE

$$\theta = (X^T X)^{-1} \cdot X^T y$$

$(d+1) \times (d+1) \cdot (d+1) \times N, (N \times 1) \rightarrow \text{SIZE } (d+1) \times 1$

INPUT:  $X$  MATRIX TRAINING DATA  
 $y$  VECTOR RESPONSES

OUTPUT:  $\theta = (X^T X)^{-1} X^T y$  MIN MSE  
ASSUMED  $X^T X$  INVERTIBLE

For.  $d=1$ ,  $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$

$X$ : size  $N \times 2$

# How Well Does the Model Fit?

- Residual Sum of Squares

→  $RSS = \sum [R_i]^2 = \sum [y_i - (\theta_0 + \theta_1 x_i)]^2$

- Total Sum of Squares

→  $TSS = \sum [y_i - \bar{y}]^2$   $\text{Var}(Y)$

– Total variance of the response

- Proportion of variability in Y that can be explained using X

–  $R^2 = 1 - \frac{RSS}{TSS} \in [0,1]$   $\uparrow$   $\text{FEATURE}$

- Correlation between feature and response

$$\rho = \text{Corr}(X, Y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}}$$

- For simple regression  $R^2$  is equal to  $\rho^2$
- For multiple LR,  $R^2$  can be used

# Vectorization

- Two options for operations on training data
  - Matrix operations
  - For loops to update individual entries
- Most software packages are highly optimized for matrix operations
  - Python numpy
  - Preferred method!
- Matrix operations are much faster than loops!

# Closed-form solution

- Can obtain  $\theta$  by simply plugging  $X$  and  $y$  into

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

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \dots & x_{id} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \dots & x_{Nd} \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

- If  $X^T X$  is not invertible (i.e., singular), may need to:

– Use pseudo-inverse instead of the inverse

• In python, `numpy.linalg.pinv(a)`

– Remove redundant (not linearly independent) features

– Remove extra features to ensure that  $d \leq n$

A MATRIX,  $G$  PSEUDO-INVERSE

$$AGA = A$$

If  $A$  is invertible

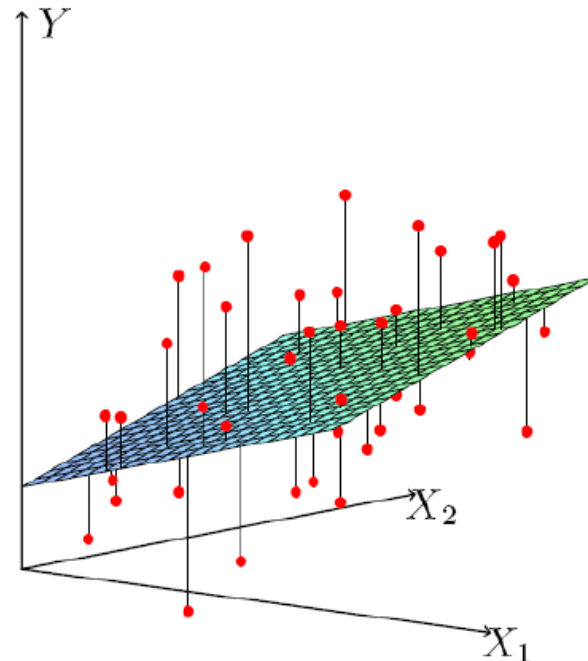
$$\Rightarrow G = A^{-1}$$

USE RIDGE REG.  $\Rightarrow$  ALWAYS INVERTIBLE

# Multiple Linear Regression

- Dataset:  $x_i \in R^d, y_i \in R$
- Hypothesis  $h_\theta(x) = \theta^T x$
- $MSE = \frac{1}{N} \sum (\theta^T x_i - y_i)^2$  Loss / cost

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



# Lab Simple Linear Regression

```
#!/usr/bin/env python

import numpy as np
import matplotlib.pyplot as plt

import pandas as pd
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

from sklearn.linear_model import Ridge, RidgeCV, Lasso, LassoCV

from sklearn.datasets import load_boston
boston_dataset = load_boston()
```

Boston house prediction dataset

# Lab

```
boston = pd.DataFrame(boston_dataset.data, columns=boston_dataset.feature_names)
boston['MEDV'] = boston_dataset.target
boston.head(5)
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

```
print(len(boston))
boston.shape
```

506

```
(506, 14)
```



```

: correlation_matrix = boston.corr().round(2)
  # annot = True to print the values inside the square
  sns.heatmap(data=correlation_matrix, annot=True)

```

: <AxesSubplot:>



# Lab

```
# Simple LR
X = pd.DataFrame(np.c_[boston['RM']], columns = ['RM'])

Y = boston['MEDV']

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state=5)
print(X_train.shape)
print(X_test.shape)
print(Y_train.shape)
print(Y_test.shape)
```

```
(404, 1)
(102, 1)
(404,)
(102,)
```

```
slr = LinearRegression()
slr.fit(X_train, Y_train)
```

# Lab

```
print(slr.intercept_)  
print(slr.coef_)
```

```
-32.839129906011266  
[8.82345634]
```

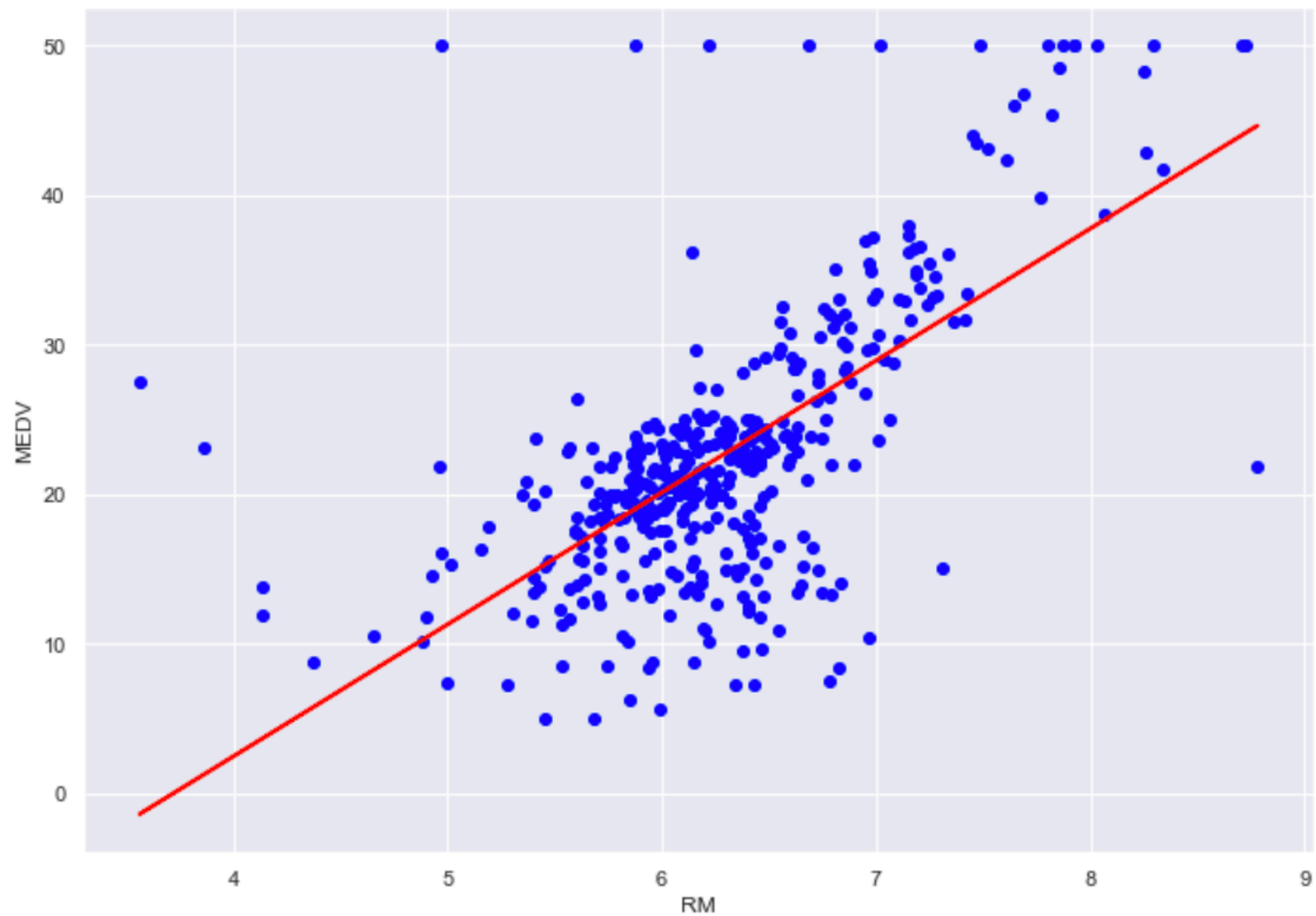
```
Y_train_predict = slr.predict(X_train)  
mse = mean_squared_error(Y_train, Y_train_predict)  
  
print("The model performance for training set")  
print('MSE is {}'.format(mse))
```

```
The model performance for training set  
MSE is 48.612648648611334
```

```
df = pd.DataFrame({'Actual': Y_train, 'Predicted': Y_train_predict})  
df.head()
```

	Actual	Predicted
33	13.1	17.463395
283	50.0	37.069115
418	8.8	19.722200
502	20.6	21.160423
402	12.1	23.666285

```
plt.scatter(X_train, Y_train, color='blue')
plt.plot(X_train, Y_train_predict, color='red', linewidth=2)
plt.xlabel('RM')
plt.ylabel('MEDV')
plt.show()
```



# Multiple LR Lab

```
: # Multiple LR

#X_multi = pd.DataFrame(np.c_[boston['LSTAT'], boston['RM']], columns = ['LSTAT', 'RM'])
X_multi = boston.loc[:, boston.columns != 'MEDV']
Y = boston['MEDV']

X_m_train, X_m_test, Y_m_train, Y_m_test = train_test_split(X_multi, Y, test_size = 0.2, random_state=5)
print(X_m_train.shape)
print(X_m_test.shape)
print(Y_m_train.shape)
print(Y_m_test.shape)

(404, 13)
(102, 13)
(404,)
(102,)

mlr = LinearRegression()
mlr.fit(X_m_train, Y_m_train)

: LinearRegression()
```

# Multiple LR Lab

```
: coeff_df = pd.DataFrame(mlr.coef_, X_m_train.columns, columns=['Coefficient'])  
coeff_df
```

:

	Coefficient
<b>CRIM</b>	-0.130800
<b>ZN</b>	0.049403
<b>INDUS</b>	0.001095
<b>CHAS</b>	2.705366
<b>NOX</b>	-15.957050
<b>RM</b>	3.413973
<b>AGE</b>	0.001119
<b>DIS</b>	-1.493081
<b>RAD</b>	0.364422
<b>TAX</b>	-0.013172
<b>PTRATIO</b>	-0.952370
<b>B</b>	0.011749
<b>LSTAT</b>	-0.594076

# Simple vs Multiple LR

```
print(slr.intercept_)  
print(slr.coef_)
```

```
-32.839129906011266  
[8.82345634]
```

```
Y_train_predict = slr.predict(X_train)  
mse = mean_squared_error(Y_train, Y_train_predict)  
  
print("The model performance for training set")  
print('MSE is {}'.format(mse))
```

```
The model performance for training set  
MSE is 48.612648648611334
```

```
: Y_m_train_predict = mlr.predict(X_m_train)  
mse = mean_squared_error(Y_m_train, Y_m_train_predict)  
  
print("The model performance for training set")  
print("-----")  
print('MSE is {}'.format(mse))  
print("\n")
```

```
The model performance for training set  
-----  
MSE is 22.477090408387635
```

# Simple vs Multiple LR

```
df_m = pd.DataFrame({'Actual': Y_train, 'Predicted simple': Y_train_predict, 'Predicted multi': Y_m_train_predict})  
df_m.head(10)
```

	Actual	Predicted simple	Predicted multi
33	13.1	17.463395	13.828770
283	50.0	37.069115	44.528528
418	8.8	19.722200	3.915991
502	20.6	21.160423	22.377959
402	12.1	23.666285	18.235923
368	50.0	11.013448	25.523748
201	24.1	21.531008	29.439747
310	16.1	11.039918	18.694533
343	23.9	26.242734	27.856463
230	24.3	19.933962	24.644734