# Dis 3: HW1 review

Ziwen Adapted from slides made by Pu Sun

#### Announcements

#### About Discussion

- Odd-numbered discussions will be presented during the 3-4pm TA-led discussion sections. For example, D3 will be presented this week by Ziwen on Tuesday from 3-4 and by Claudio on Thursday 3-4.

- Even-numbered discussions will be presented during the 12-1:30 lecture time.

- HW1 Due Date: Jan 17 at 11:30 pm.
- About Final Project Teams

Ziwen will be releasing this assignment by this Thursday, and **you will have until next Thursday to finalize your groups** 

Groups will be self-selected and limited to 8-9 students

#### libraries

- Pandas
  - A very famous and useful data analyze and visualize tool
  - We will use it for loading data and some calculations
- Scikit-learn
  - Simple and efficient tools for predictive data analysis
- NumPy
  - The fundamental package for scientific computing with Python

# Visual tools

- Motplotlib
- Seaborn
- Pandas
- Plotly

#### Pearson correlation

- It only calculating linear correlation.
- It will give a value [-1, 1].
- Value close to 0 means lower correlations



															10	ř
CRIM ·	i	-0.2	0.41	-0.056	0.42	-0.22	0.35	-0.38		0.58	0.29	-0.39	0.46	-0.39		Į.
ZN	-0.2	1	-0.53	-0.043	-0.52	0.31	-0.57		-0.31	-0.31	-0.39	0.18	-0.41	0.36	- 0.8	3
INDUS -	0.41	-0.53		0.063		-0.39		-0.71	0.6		0.38	-0.36	0.6	-0.48		
CHAS -	-0.056	-0.043	0.063	1	0.091	0.091	0.087	-0.099	-0.0074	-0.036	-0.12	0.049	-0.054	0.18	- 0.6	5
NOX ·	0.42	-0.52	0.76	0.091	1	-0.3	0.73	-0.77	0.61	0.67	0.19	-0.38	0.59	-0.43	- 0.4	4
RM -	-0.22	0.31	-0.39	0.091	-0.3	1	-0.24	0.21	-0.21	-0.29	-0.36	0.13	-0.61	0.7		
AGE ·	0.35	-0.57		0.087		-0.24	1	-0.75	0.46	0.51	0.26	-0.27	0.6	-0.38	- 0.2	2
DIS	-0.38		-0.71	-0.099	-0.77	0.21	-0.75		-0.49	-0.53	-0.23	0.29	-0.5	0.25	- 0 (	0
RAD ·	0.63	-0.31	0.6	-0.0074	0.61	-0.21	0.46	-0.49	1	0.91	0.46	0.44	0.49	-0.38	- 0.0	,
TAX	0.58	-0.31		-0.036		-0.29	0.51	-0.53	0.91		0.46		0.54	-0.47	0	1.2
TRATIO ·	0.29	-0.39	0.38	-0.12	0.19	-0.36	0.26	-0.23	0.46	0.46	1	-0.18	0.37	-0.51		
B·	-0.39	0.18	-0.36	0.049	-0.38	0.13	-0.27	0.29	-0.44		-0.18	1	-0.37	0.33	0	1.4
LSTAT ·	0.46	-0.41	0.6	-0.054	0.59	-0.61	0.6	-0.5	0.49	0.54	0.37	-0.37	1	-0.74	0	).6
MEDV	-0.39	0.36		0.18			-0.38	0.25	-0.38		-0.51	0.33	-0.74	1		
	CRIM	z'n	INDUS	CHAS	NÓX	ВM	AGE	Dis	RÁD	TAX	PTRATIC	b	LSTAT	MEDV		

# Pair plot (a.k.a. scatter plot matrix)



#### Linear and Polynomial regression

$$egin{aligned} y_i &= eta_0 + eta_1 x_i + eta_2 x_i^2 + \cdots + eta_m x_i^m + arepsilon_i \ (i = 1, 2, \dots, n) \ egin{aligned} & y_1 \ & y_2 \ & y_3 \ & \vdots \ & y_n \end{aligned} &= egin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \ 1 & x_2 & x_2^2 & \cdots & x_2^m \ 1 & x_3 & x_3^2 & \cdots & x_3^m \ & \vdots & \vdots & \vdots & \ddots & \vdots \ 1 & x_n & x_n^2 & \cdots & x_n^m \end{aligned} \end{bmatrix} egin{bmatrix} eta_0 \ eta_1 \ eta_2 \ & \vdots \ eta_n \end{aligned} &+ egin{bmatrix} arepsilon_1 \ arepsilon_2 \ arepsilon_3 \ arepsilon_1 \ arepsilon_2 \ arepsilon_1 \ arepsilon_2 \ arepsilon_2 \ arepsilon_2 \ arepsilon_1 \ arepsilon_2 \$$



#### SSE and MSE and variance

SSE: sum of squared error MSE: mean of squared error



sum of the errors of all samples

$$ext{MSE} = rac{1}{n}\sum_{i=1}^n {\left(Y_i - \hat{Y_i}
ight)^2}.$$

### Overfitting and Underfitting



## Overfitting and Underfitting

- Ideal case: low bias (in training) and low variance
- Overfitting: low bias (in training) and high variance
- Underfitting: High bias (in training and test datasets)

#### Outliers



#### Outliers





#### OC SVM





 $W = (X^! X)^{"\#} X^! y$ 

OLS Example