

JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May 02  
Monday

## 09.00 RESAMPLING AND REGULARIZATION

Resampling = important for model performance and validation. Done in every data science task.

Ex - Cross Validation

11.00

Regularization - used to prevent overfitting and improve performance.

12.00

Ex - Ridge Regression and Lasso

13.00

Resampling -

• Repeatedly draw samples from a training set and refit the model.

15.00

• Gain more information  
• See how model would perform on new data without collecting new data.

16.00

→ Cross validation is widely used method for resampling  
→ It is used to evaluate model's best performance and find best parameters for the model.

## 3 APPROACHES TO DO - CROSS VALIDATION

- Validation Set
- Leave One Out Cross Validation
- K-Fold Cross Validation



03

'22

May

Week 19 Day (123-242)

Tuesday

MAY 2022

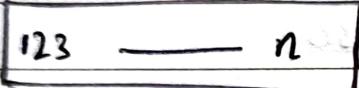
MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
					1	2	3	4	5	6	7	8	
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

## 09.00 - VALIDATION SET

- VALIDATION SET IS THE MOST BASIC APPROACH.

- 10.00 • WE HAVE A DATASET OF  $n^2$  POINTS AND WE RANDOMLY SPLIT DATA SET INTO TRAINING SET AND TEST SET.
- 11.00 • WE FIT THE MODEL ON TRAINING SET AND MAKE THE PREDICTIONS ON TEST SET.

12.00



13.00



14.00

## 15.00 SOME DRAWBACKS -

- TEST ERROR RATE IS VARIABLE DEPENDING ON WHICH OBSERVATIONS WERE IN THE TRAINING AND VALIDATION SET. SINCE WE ARE SPLITTING RANDOMLY.
- ALSO ONLY SMALL SUBSET OF DATA IS USED FOR TRAINING. WHEN IDEALLY WE WANT AS MUCH DATA AS POSSIBLE.

16.00



JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May

22

04

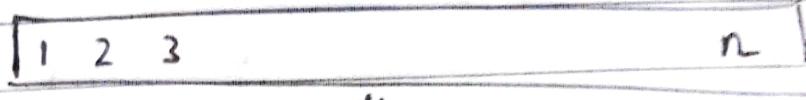
Wednesday

Week 10 Day (2016-2017)

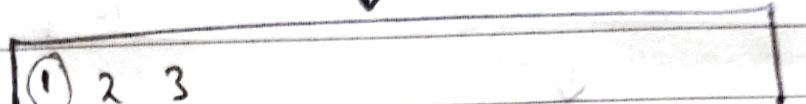
## LEAVE ONE OUT CROSS VALIDATION (LOOCV)

- ONLY ONE OBSERVATION IS USED FOR VALIDATION.

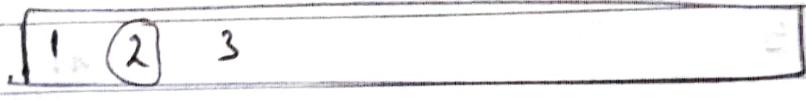
10.00



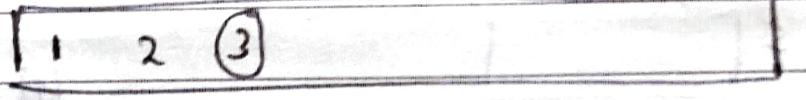
11.00



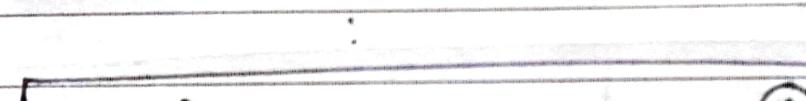
12.00



13.00



14.00



15.00



16.00

IN THIS METHOD, ONLY ONE DATA POINT IN CIRCLE IS USED FOR VALIDATION & REST IS USED FOR TRAINING.

WE REPEAT THIS PROCESS AS MANY TIMES AS WE HAVE DATA POINTS

18.00 - THE ERROR RATE IS APPROXIMATED AS THE AVERAGE OF ERRORS FOR EACH RUN

BENEFIT OF LEAVE ONE OUT CROSS VALIDATION (LOOCV)

- NO RANDOMNESS,

DRAWBACK

- NOT A VIABLE OPTION FOR LARGE DATASETS (IT WILL TAKE A LOT OF TIME)

05

'22

May

Week 19 Day (125-240)

Thursday

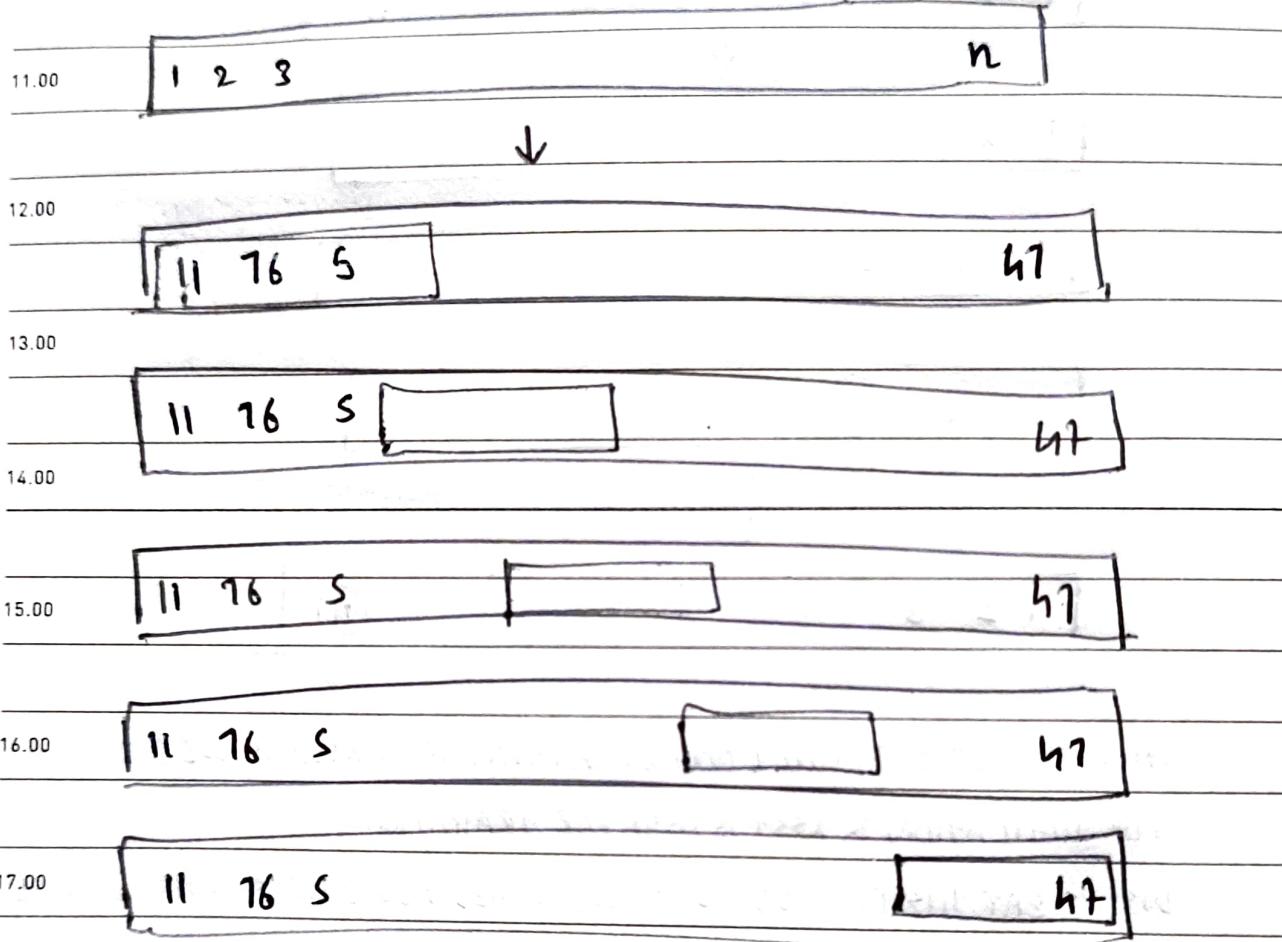
MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28
29	30	31											

## 09.00 K-FOLD CROSS VALIDATION

- RANDOMLY DIVIDE THE DATASET INTO K GROUPS OR FOLDS

10.00 OF EQUAL SIZE



18.00 WE USE THE COMPLETE BOX FOR TRAINING, AND USE THE SMALL BOX FOR VALIDATION. AND WE REPEAT THE PROCESS 'K' TIMES.

- LOOCV is k-fold with K=n
- Usually set K to 5 or 10
- BEST & MOST WIDELY USED METHOD.

JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May 06

Friday

Week 19 Day (126-239)

## 09.00 K-cross validation in PYTHON with LINEAR REGRESSION

10.00  $\text{lin\_reg} = \text{LinearRegression}()$ 11.00  $\text{MSEs} = \text{cross\_val\_score}(\text{lin\_reg},$  $X,$  $y,$  $\text{scoring} =$ 'neg-mean-squared  
- errors',  
 $cv=5$ )12.00  $\text{mean\_MSE} = \text{np.mean}(\text{MSEs})$ 13.00  $\text{print}(-\text{mean\_MSE})$ 

14.00

15.00

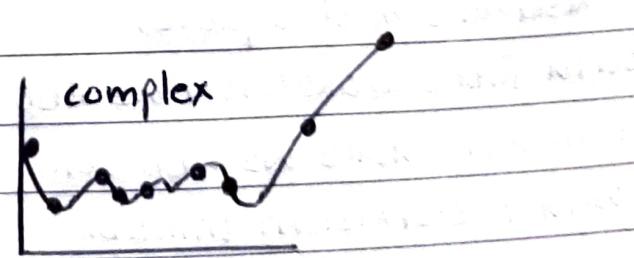
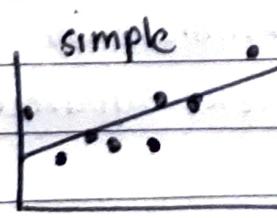
16.00

## REGULARIZATION

17.00 - MODELS CAN OVERFIT, MEANING THAT THEY WILL NOT GENERALIZE WELL

18.00 - THEY WILL NOT GENERALIZE WELL AND WILL FORM POORLY ON UNSEEN DATA.

→ BIAS VARIANCE TRADE OFF



07

22

May

Saturday

Week 19 Day (127-238)

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
					1	2	3	4	5	6	7	8	
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

09.00

WE HAVE TO FIND SOMETHING IN BETWEEN TO PREVENT MODELS FROM OVERRFITTING.

10.00

AND THAT'S WHEN WE USE REGULARIZATION.

'REGULARIZATION WILL DE-CCELERATE OUR MODELS TO PREVENT OVERRFITTING.'

12.00

REGULARIZATION METHODS -

- RIDGE REGRESSION

13.00

- LASSO

- ALSO CALLED - SHRINKAGE METHODS (RIDGE & LASSO)

14.00

RIDGE REGRESSION -

$$\text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$$

RSS - Residual Sum of Square.

Sunday 08

15.00

WE KNOW THAT TRADITIONAL LINEAR SETTING MINIMISES RSS (Residual Sum of Square)

WITH RIDGE REGRESSION WE ADD ANOTHER OPTIMIZATION FUNCTION, HERE WE ADD SUM OF PARAMETERS SQUARED WITH A COEFFICIENT LAMBDA.



MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May

22

Monday

Week 20 Day (131-234)

09

- 09.00 - LAMBDA IS CALLED TUNING PARAMETER.
- TO FIND BEST VALUE OF LAMBDA WE WILL USE CROSS VALIDATION WITH RANGE OF LAMBDA'S.
- \* BEST VALUE WILL BE THE ONE MINIMIZING THE TEST ERROR.
  - \* WITH THIS METHOD ALL PREDICTORS ARE KEPT.
  - \* ALSO CALLED AS L2 REGULARIZATION.
- 12.00
- RIDGE REGRESSION IN PYTHON -

13.00

ridge = Ridge()

14.00 parameters = {'alpha': [1e-15, 1e-10, 1e-8, 1e-4, 1e-3, 1e-2, 1, 5, 10, 20]}

15.00 ridge\_regressor = GridSearchCV(ridge, parameters, scoring='neg-mean-squared\_error', cv=5)

16.00

ridge\_regressor.fit(X, Y);

17.00

18.00

IN PYTHON TO USE IT WE FIRST INITIALIZE THE RIDGE MODEL THEN LAMBDA PARAMETER IS ACTUALLY ALPHA, WITH THE LIBRARY THAT WE USE IN PYTHON.

IN SNIPPET WE PROVIDE AN ARRAY WITH DIFFERENT VALUES OF ALPHA.

10

May

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

Week 20 Day (130-235)

Tuesday

09.00 THEN WE USE 5 FOLD CROSS VALIDATION TO MAKE USE OF BEST VALUE OF ALPHA FOR OUR MODEL

10.00

LASSO

- 11.00
- ADD A NEW TERM TO THE OPTIMIZATION FUNCTION
  - LASSO IS ALSO CALLED  $L_1$  REGULARIZATION.

12.00

$$RSS + \lambda \sum_{j=1}^p |\beta_j|$$

13.00

HERE WE ADD SUM OF ABSOLUTE VALUES OF COEFFICIENTS.

$\lambda$  - OUR TUNING PARAMETER

14.00

15.00

- 16.00
- IF LAMBDA IS LARGE ENOUGH, SOME BETA WILL GO TO ZERO  
SOME FEATURES WILL DISAPPEAR.

17.00 → THEN FEATURE SELECTION CAN BE DONE.

18.00 lasso = Lasso(tol=0.05)

parameters = {`alpha': [1e-15, 1e-10, 1e-8, 1e-4, 1e-3, 1e-2, 1, 5, 10, 20]}

lasso\_regressor = GridSearchCV(lasso, parameters, scoring='neg\_mean\_squared\_error', cv=5)

lasso\_regressor.fit(x, y)

of them.

JUNE 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12		
13	14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30										

May

'22

11

Wednesday

Week 20 Day (131-234)

## 09.00 RESAMPLING AND REGULARIZATION IN PYTHON -

10.00

import numpy as np

11.00

import pandas as pd

import matplotlib.pyplot as plt

12.00

%matplotlib inline

13.00

14.00

15.00

16.00

import numpy as np - library importng for analyses

17.00 import pandas as pd - "data analysis" library

import matplotlib.pyplot - pyplot module is collection of

18.00 functions that provide MATLAB like Interface for creating plots  
and charts.

%matplotlib inline - magic command that set up Jupyter  
notebook environment to display Matplotlib plot inline in  
notebook.



- Advertising.csv used



12

'22

May

Week 20 Day (132-233)

Thursday

09.00

`def scatter_plot(feature):`

10.00

`plt.figure(figsize = (10,5))`

11.00

`plt.scatter(data[feature], data['sales'],  
 c='black')`

12.00

`plt.xlabel(f"Money Spent on {features} (ds)")`

— 2

13.00

`plt.ylabel("Sales in $")`

— 3

14.00

`plt.show()`

— 4

15.00

`scatter_plot('TV')`

— 5

16.00

`scatter_plot('radio')`

— 6

17.00

`scatter_plot('newspaper')`

18.00

1) `def scatter_plot(feature):`: This defines a function called `scatter_plot` that takes a single argument called `feature`

2) This creates a new figure with size of 10 inches by 5 inches

using the `figure()` function from `pyplot` module.

3) X values are taken from `feature` column, Y values are taken from `sales` column, `c` specifies color of points

4) Add label to the X-axis

5) Add label to the Y-axis

6) plotting the Graph

MAY 2022

SUN	MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
							1	2	3	4	5	6	7	8
9	10	11	12	13	14	15	16	17	18	19	20	21	22	23



JUNE 2022

TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	
14	15	16	17	18	19	20	21	22	23	24	25	26
28	29	30										

May

22

13

Friday

Week 20 Day (133-232)

## Define Baseline Model & See How Regularization Improves it

```
from sklearn.model_selection  
import cross_val_score  
from sklearn.linear_model  
import LinearRegression
```

- 1

- 2

- 1) This function is used to perform cross-validation on model. WHICH HELP TO ASSESS HOW MODEL WILL GENERATE TO NEW DATA.
- 2) USED TO CREATE LINEAR REGRESSION MODEL, WHICH IS TYPE OF SUPERVISED LEARNING USED FOR PREDICTING CONTINUOUS VALUES.

```
X = data.drop(['sales'], axis=1)  
y = data['sales'].values.  
reshape(-1,1)
```

- 1

- 2

- 1) Creates a new column new data frame X by removing the column named 'sales' from original data set. This axis=1 argument specifies that the column should be removed. This line assumes the 'sales' column contains target variable or output that model is trying to predict, and remaining columns in data are features or inputs that model will use to make predictions.

14

22

May

Week 20 Day (134-231)

Saturday

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8						
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

- 09.00 2) The reshape (-1,1) method reshapes the array into a column vector with the same number of rows as original array. but with only one column.

11.00

`lin-reg = LinearRegression()`

— 1

12.00

`MSEs = cross_val_score(lin-reg, X, y,`

— 2

`scoring='neg_mean_squared_error',``(cv=5)`

— 3

`mean_MSE = np.mean(MSEs)`

— 4

`print(-mean_MSE)`

14.00

15.00

16.00

17.00

Sunday 15

- 1) Create an instance of Linear Regression class

- 18.00 2) uses the cross\_val\_score function from scikit-learn to perform 5-fold cross validation on LR Model.

- 3) calculates the mean of negative MSE Values across all folds of cross validation.

- 4) print negative values of MSE Values as output



Homework



Email



Question

2022

WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12
15	16	17	18	19	20	21	22	23	24	25	26
29	30										

May 16  
Monday  
Week 21 Day (136-229)

## REGULARIZATION

### RIDGE REGRESSION

```
from
import sklearn.model_selection
import GridSearchCV
from sklearn.linear_model
import Ridge
```

GridSearch is scikit-learn that performs exhaustive search over specified parameter grid to find set of hyperparameters.  
Ridge is class in scikit-learn that implements Ridge regression Model.

ridge = Ridge()

parameters = {'alpha': [1e-15, 1e-8, 1e-10,  
1e-4, 1e-3, 1e-2, 1, 5, 10, 20]}

ridge\_regressor = GridSearchCV(ridge,  
parameters, scoring='neg\_mean\_squared\_error', cv=5)

ridge\_regressor.fit(X, y);

17

'22

May

Week 21 Day (137-228)

Tuesday

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
					1	2	3	4	5	6	7	8	
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

- 09.00 Ridge() - instance of Ridge class
- Second line - defines dictionary of hyperparameters to be tuned using grid search. We are only tuning the Hyper parameters from  $1e-15$  to  $20$ , mix of small & Large values.
- 10.00 Create instance of GridSearchCV
- 11.00 Fits the ridge-regressor object to input data X and target-variable y.

## Lasso (Regularization)

```
14.00 from sklearn.model_selection import GridSearchCV
15.00 from sklearn.linear_model import Lasso
16.00
17.00
```

```
18.00 lasso = Lasso(tol=0.05)
parameters = {'alpha': [1e-15, 1e-10,
1e-8, 1e-4, 1e-3, 1e-2, 1, 5, 10, 20]}
lasso_regressor = GridSearchCV(lasso,
parameters, scoring='neg-mean_squared_error', cv=5)
lasso_regressor.fit(X, y)
print(lasso_regressor.best_params_)
print(-lasso_regressor.best_score_)
```