

14

Monday

Module 3.

March

2011

March							2011
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Week 11 • 73-292

(EDA). Exploratory Data Analysis

Day 14.

Numpy.

- Fast
- Support nd array.
- high level math fn.

Ques. Why do we need Numpy?

If we have 1000×1000 matrix, multiply.

- Python loop $> 10\text{min}$
- Numpy $\approx 0.03\text{s}$.

Overview.

- Arrays.
- Shaping & transposition
- Mathematical Opern
- Indexing & slicing
- Broadcasting

↓
data of
1 area
from 1D
array
to another
array.

Arrays.

- vectors
- Matrices
- Images
- Tensors } 2D part.
- conv Nets.

Numpy Arrays.

Ex. import numpy as np

b = [1, 2, 3, 4, 5]

list ↴

converted
list to
np.array ↴ arr = np.array(b)

arr

np.array ↴ array([1, 2, 3, 4, 5])

arr.shape
2D → (5, 1)

• arr.shape[1]
no. ↴ coln if array is in 2D form.

arr.shape[2]

no. ↴ channels (if array
is in 3D form).

(if it is
in 2D,

it represent
20 rows.)

arr.shape[0]
no. of
elem. in
2D ↴ #5

April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
13					1	2 3
14	4	5	6	7	8	9 10
15	11	12	13	14	15	16 17
16	18	19	20	21	22	23 24
17	25	26	27	28	29	30

Tuesday

March

2011

15

Week 11 • 74-291

2D Array Creation.

b1 = np.array ([[1,2,3], [4,5,6], [7,8,9]])

[[1,2,3]

2D.

[4,5,6]

[7,8,9]]

b[1][2] → 6

* array slicing can't be done from backward dirn., it can only be possible from left to right

numpy functions.

Data matrix → np.zeros((4,4)) → 4x4 matrix of 0.

matrix → np.ones((3,6)) → 3x6 matrix of 1

use it → np.eye(5,5) → diag. 1 & other 0. (identity matrix).

(S) = np.full((3,3), 89) → 3x3 matrix with all elem = 89.

↳ np.ones_like(S) → it captures the size of S ie. 3x3.
Creates 3x3 ones matrix. You don't need to specify the size in this

↳ np.zeros_like(S)

Create 3x3 zero matrix.

• a = np.array ([[1,2,3,4], [5,6,7,8], [9,10,11,12]]) → 3x4

b = a[1:3, :2]

↳

[5,6]

[9,10]]

↳ [[1,2,3,4]

[5,6,7,8]

[9,10,11,12]]

b[0][0] = 500.

16

Wednesday

March

2011

Week 11 • 75-290

March							2011
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

$b \rightarrow [[\textcircled{500}, 6], [9, 10]]$

change value
to 500 in b
& it is
impacted
on 'a'.

$a \rightarrow [[1, \textcircled{2}, 3, 4], [\textcircled{500}, \textcircled{6}, 7, 8], [9, 10, \textcircled{11}, 12]]$

it is
also
updated.

$b = np.array([1, 1, 2, 3])$

$[1, 1, 2, 3]$

$ss = (a[\text{np.arange}(3), b])$

$[1, 1, 2, 3]$

$[2, 6, 11]$

object,
array,
 $a[2, 2]$

$ss[0] = 200$

$ss \rightarrow [200, 6, 11]$

Q. will it make impact on a?

No impact on a. bcz,

• whenever you make subset with slicing, & then make any change in subset it will reflect in main array.

• when you make subset without slicing, it will not make any impact on original data.

place of 2, 6, 11 add +10

$a = a[10] \leftarrow$
at the
particular
place

$\underbrace{a[\text{np.arange}(3), b]}_{t=10} + 10$

OP: 1 (11) 34
500 (16) 78
9 10 (2) 12

April 2011						
Mo	Tu	We	Th	Fr	Sa	Su
				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	

Thursday

March 2011

17

Week 11 • 76-289

$a > 2 \rightarrow$ gives boolean op (T or F)

$a[a > 2] \rightarrow$ gives values with satisfy condn
if ~~also~~ a is 2D but this only gives 1D.

* Array Math.

default dtype → float.

$x = np.array([[1, 2], [3, 4]], dtype=np.float64)$

$y = np.array([[5, 6], [7, 8]], dtype=np.float64)$

$np.add(x, y)$ # $x+y$ 9

$np.subtract(x, y)$ # $x-y$ 10

$np.multiply(x, y)$ # $x*y$ 11

$np.divide(x, y)$ # x/y 12

$np.sqrt(x)$ 1

$v = np.array([9, 10])$ dot 3

$w = np.array([11, 12])$ $9x11 + 10x12$ 4

• $v.dot(w)$ # 219 5

• $np.dot(v, w)$ # 219 6

→ • $np.dot(x, y)$ $\rightarrow [[19 22]$ 7
 $43 50]]$

37 ↗ $x = np.array([[1, 2, 3, 4], [3, 4, 5, 6]])$

↳ $np.sum(x)$ # compute sum of all elements

↳ $np.sum(x, axis=0)$ # compute sum of each coln.

↳ $np.sum(x, axis=1)$ # compute sum of each row.

[17 20]

18

Friday

March 2011

	March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
9		1	2	3	4	5	6
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Week 11 • 77-288

Transpose.

$$x \rightarrow [[1\ 2 \\ 3\ 4]]$$

$$(x.T) \rightarrow [1\ 3 \\ 2\ 4]$$

interchange rows & coln.

Broadcasting

9 np.empty((5, 5)) → 5x5 random values.

10 x = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9], [10, 11, 12]])

11 y = np.empty_like(x)

12 empty matrix with same shape as x.
create empty matrix of 3x3.

1 x.shape → (4, 3)

(doesn't generate 3x3

3 v = np.array([1, 0, 1]) matrix with random values

4 for i in range(len(x)):

$$y[i, :] = x[i, :] + v[i]$$

len(x)=4
6 print(y)

$$= x[0, 0] + v[0] \\ y[0, 0] : x[0, 0] + v[1] \\ \vdots$$

$$\# [[2\ 2\ 4 \\ 5\ 5\ 7 \\ 8\ 8\ 10 \\ 11\ 11\ 13]]$$

y shape
(4, 3)

April 2011						
Su	Mo	Tu	We	Th	Fr	Sa
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	

Saturday

March 2011

19

Week 11 • 78-287

Day 15.

combining Arrays.

→ similar to np.array()

• from numpy import array.

from numpy import vstack, hstack
np.vstack()
np.hstack()

vertical stack
horizontal stack

if you want to
merge row
wise

↓ if you
want to
merge col
wise

a1 = array([1, 2, 3])

a2 = array([4, 5, 6])

a3 = vstack((a1, a2)) → [[1, 2, 3] merged column wise
4, 5, 6]]

, a4 = hstack((a1, a2)) → [1 2 3 4 5 6] merged row wise;
if P is
always in
1D.

In vstack,
Both the array should have same no. of elements.

It will show an error. a1 = array([1, 2, 3]) → 3 elem.

a2 = array([4, 5, 6, 7]) → 4 elemsunday
vstack((a1, a2))

20

79-286

for 2D: hstack.

a1 = array([[1, 2, 3], [4, 5, 6]])

a2 = array([[11, 22, 33], [44, 55, 66]])

hstack((a1, a2))

→ 1 2 3 11 22 33
if P. 4 5 6 44 55 66

✓ vstack((a1, a2))

executed

(no error).

21

Monday

March 2011

Week 12 • 80-285

Wk	Mo	Tu	We	Th	Fr	Sa	Su
9			1	2	3	4	5
10	7	8	9	10	11	12	13
11	14	15	16	17	18	19	20
12	21	22	23	24	25	26	27
13	28	29	30	31			

Note.

conversion is possible only from
lower to higher:

- 1D → 2D ✓
- 1D → 3D ✗ (already)

Array Reshaping.

Reshape 1D to 2D. & 2D to 3D:

shape.

(5,)

1D data

data = array([11, 22, 33, 44, 55])

O/P: data1 = data.reshape(data.shape[0], 1)

shape:
(5,1)
now { [11
22
33
44
55]]

data2 = data.reshape(

(data.shape[0],
data.shape[1], 1))convert next
dimension.

converted

2D to 3D now.

Pandas.

3 data structures.

→ Series (1D)

→ DataFrame. (2D)

→ Panel. (3D)

Series.

1D array
like structure.
with homogenous
data

- 1D labelled homogeneous array,

- size immutable.

- values of Data Mutable.

DataFrame.

- Heterogeneous data
- size mutable
- Data mutable.

Panel.

- Heterogeneous
- size mutable
- Data mutable.

	April							2011	
Wk	Mo	Tu	We	Th	Fr	Sa	Su		
13						1	2	3	
14	4	5	6	7	8	9	10		
15	11	12	13	14	15	16	17		
16	18	19	20	21	22	23	24		
17	25	26	27	28	29	30			

default dtype of all python
NOTE package is
= float

Tuesday

March 2011

22

Week 12 • 81-284

Implementation.

- pd.Series(data, index, dtype, copy)

Ex

$l = [1, 2, 3, 4]$

$s = pd.Series(l)$

s

$\begin{pmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{pmatrix}$

Op. of series
datastructure.

$s = pd.Series(data, index=[9, 10, 11, 12])$ (You can create
customized index).

$\begin{pmatrix} 9 & 1 \\ 10 & 2 \\ 11 & 3 \\ 12 & 4 \end{pmatrix}$ customized
Op.

$s[0] \rightarrow$ error

$s[11] \rightarrow 3$

$data = \{'a': 0, 'b': 1, 'c': 2\}$

$s = pd.Series(data)$

$\begin{matrix} a & 0 \\ b & 1 \\ c & 2 \end{matrix}$

$s['a'] = s[0] \rightarrow 0$

$s['b'] = s[1] \rightarrow 1$

A dictionary,
you can access
element by
default index or
user defined.

$s[::-1] \rightarrow$
 $\begin{matrix} c & 2 \\ b & 1 \\ a & 0 \end{matrix}$

Day-16.

DataFrame.

- explore or handle 2d data.

- both size & info(data) are mutable.

C:/users/.../diabetes-with-name.csv

$df = pd.read_csv('C:/users/.../diabetes.csv')$

don't use backward slash. (use '/' forward
slash.)

Pandas
don't support
Backward ('.')
slash.

23

Wednesday

March 2011

part of documentation

for reference.

	Mo	Tu	We	Th	Fr	Sa	Su	2011
9			1	2	3	4	5	6
10		7	8	9	10	11	12	13
11		14	15	16	17	18	19	20
12		21	22	23	24	25	26	27
13		28	29	30	31			

Week 12 • 82-283

- ML algo. supports only structured data.
unstructured data can't be used directly,
convert it to structured
by hadoop tool.
Ex. wordfile →

- Datafile or dataset used here:
 - Diabetes with names.csv.

- Null values affect the average. (Statistics gets differ.)

Ex You have 15 rows & 2 are Null (NaN).

sum/15

(taking null value) avg. by comp. for 15 rows = 34.25

Actual avg = sum of 13 rows / 13

removing null values,
then cal. avg.

every

df.isna().sum() → calculate missing values in a row &
df.isnull().sum() → prints sum
both are same. of each row
seperately.

df.isna().sum().sum() → calculate total no. of missing values.
in entire dataset.

df.fillna() → this fn fills every NaN value with that
value.

Ex df.fillna(df.mean())

df.info() → info. of dataframe. (get all file info.)

(gives datatypes, total values in a particular
row..)

April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
13					1	2 3
14	4	5	6	7	8	9 10
15	11	12	13	14	15	16 17
16	18	19	20	21	22	23 24
17	25	26	27	28	29	30

Thursday

March 2011

24

Week 12 • 83-282

`len(df)` → gives total no. of rows in dataset.

`df.size` → $row \times col$

Ex 8 rows, 5 columns: $8 \times 5 = 40$

`df.shape` → gives shape of dataset.

Ex $(8, 5)$

`df.describe()` → returns statistical summary of dataset.

Ques If NaN values are not there in dataset, instead of NaN any other variable like '?', 'Unknown' etc is there will `isnull().sum()` give same result?

Soln NO,

Find '?' in dataset by method, replace it with NaN & then try to fill it.

Alternative,

while loading dataset, define there which type of null column you have.

Ex `pd.read_csv('---.csv', na_values='?')`

df.isnull().sum()

We have ? as Null value in dataset.

→ You will get the count of null values.

soln1.py → Dataset (titles.csv) → movie title dataset.

`df1 = pd.read_csv('titles.csv', index_col=None)`

`df1.head()`

→ top 5 rows.

25

Friday

March

2011

March						
Wk	Mo	Tu	We	Th	Fr	Sa
9		1	2	3	4	5
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

Week 12 • 84-281

df3 = pd.read_csv('cast.csv', index_col=None)

Q1 earliest 2 film in title dataset

df2.sort_values('year').head(2)

always sort in ascending order.
(lower to higher).

Q2 How many movies have title 'Hamlet'?

len(df2[df2.title == 'Hamlet'])

Q3 when was first movie titled 'Hamlet' made?

df2[df2.title == 'Hamlet'].sort_values('year').head(1)

(Check Notes given)

Day 17.

Exploratory Data analysis.

df = pd.read_csv('_____ .csv')

df.head()

df.isna().sum()

Handling Missing Values.

- Should have domain knowledge.
- When we encounter huge no. of nan values, we use to fill nan by calling fillna().
- less no. of nan, if you wish you can drop.

	April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13					1	2	3
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

create dropped record
as subset.

Saturday

26

March

2011

how values should
be deleted.

Week 12 • 85-280

df.dropna (self, axis=0, how='any', thresh=None, subset=None,
drop
how
nam
or
'all'
you can
set the
limit
to del.
missing
value.
inplace=True)

NaN → NOT a timestamp.

(it is used where time/date is used).

df.dropna () → it drops rows where atleast 1 element is missing.

df.dropna (how='all') → it will drop when in a col/row
all the values are NaN.

df.dropna (thresh=2) → keep rows with atleast 2 non-NaN values.

• Filling missing Values.

df.fillna (df.mean())

df.fillna (method='pad').head()

(Previous value taken
is placed before
in place of NaN) & value &
place over it

Sunday

27

86-279

padding

↓ ↓

forward

(ffill)
(default)

backward

(bfill)
↓ syntax

fillna (method='bfill')

pad bfill

9x	6	6	6
NaN	6	8	
8	8	8	

28

Monday

March 2011

Week 13 • 87-278

March 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
9	1	2	3	4	5	6
10	7	8	9	10	11	12
11	14	15	16	17	18	19
12	21	22	23	24	25	26
13	28	29	30	31		

↑ ex
value
NAN → not fill in case of
bbill

To handle value of first last ~~err~~ (exception).

You can use both ffill & bfill (but not at a time, use simultaneously.)

dynamic
way of
filling
values.

df.interpolate()

it will fill missing
value itself.

→ Problem.

we can't judge which method
is used to fill values.

→ We should not give entire data

to M/C. bcz if model is not
have nice accuracy. You will not
be able to check why it
is not giving accuracy.

If you handle yourself
you can do increase or
decrease the accuracy.

interpolate feature:

- spline → use circle eqn
- polynomial → magnitude value in logic ..
- akima
- ..

EDA → know Nature of Data.

PLOT Graph → to get insights.

data.describe().T → Transposed
matrix.

data.hist(figsize=(12, 10))

Density
of record.

	April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13				1	2	3	
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Tuesday

March 2011

29

Week 13 • 88-277

Day-18

data.boxplot(figsize=(12, 6))

Try to explore IQR (25%, 50%, 75%)

check one by one:

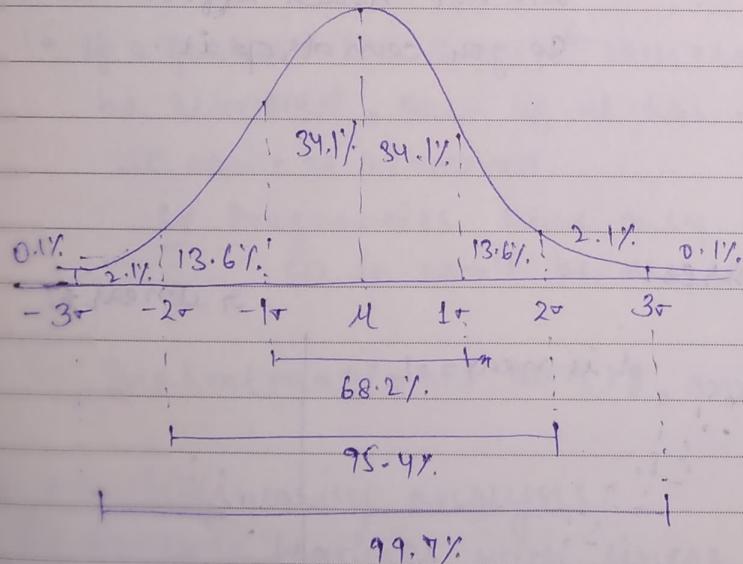
sns.boxplot(df['Glucose'])

sns.boxplot(df['Pregnancies'])

:

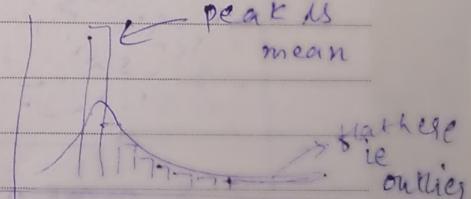
Inference from Boxplot.

Max outlier, Min outlier, Moderate outlier, No outlier



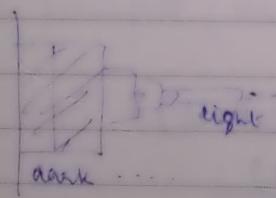
Distribution is used to give dataflow of feature

sns.distplot(df['coln'])
↳ gives distribution.



sns.boxenplot(df['coln'])

color variance outlier
→



data distribution can be identified by color of graph

1 has outliers
2 has moderately linearity.

30

Wednesday

March

2011

Week 13 • 89-276

Q: What is Max. Outlier?

Value above $4 \pm 2\sigma$.

Q: Should we remove or retain/as it is?

	Mo	Tu	We	Th	Fr	Sa	Su
Wk	9	1	2	3	4	5	6
	10	7	8	9	10	11	12
	11	14	15	16	17	18	19
	12	21	22	23	24	25	26
	13	28	29	30	31		

When you have max. outliers. → keep the moderate outliers which are near to max. pt.

→ drop extreme outlier?

- No, if we have not so large data.

Ex 10,000 insulin patients, 100 people have 800 as insulin level. (so it is advised to keep outlier)

(You should have Domain knowledge for outlier treatment).

Yes. → if we have very large amount of data.

(Rare) Too extreme - drop it

Ex 1 severe insulin patient, 1-2 people have 800 level of insulin. It will not much affect model, so you can drop it.

moderate outlier

minimum outlier

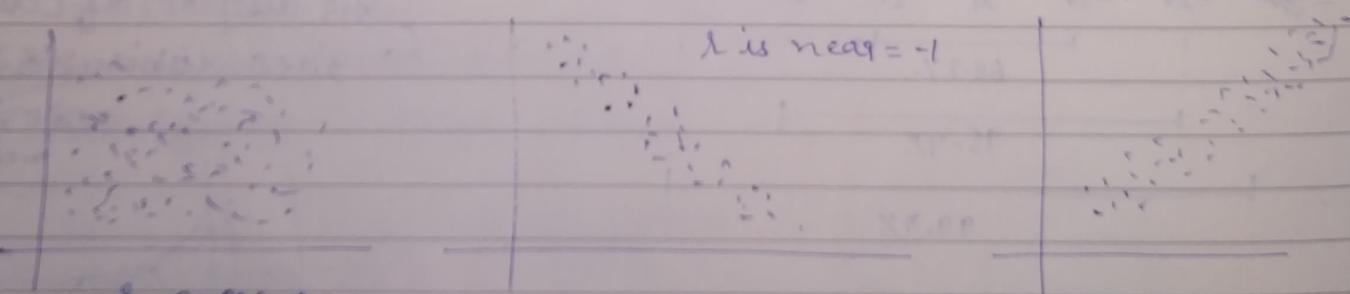
keep it/retain it

② Correlation.

↳ linear reln b/w features.

or linear +

r is near = -1



r near 0

No correln

Negative correln

Positive correln

coefficient of reln (Pearson's coefficient) :

$$r = f(x,y) = \frac{\text{cov}(x,y)}{\text{std. } x \cdot \text{std. } y}$$

	April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13					1	2	3
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Thursday

March 2011

31

Week 13 • 90-275

-1 → strong -ve corr.

+1 → strong +ve corr.

0 → no correln.

Tips for taking decision on corr. coeff.

-0.2 to +0.2 → consider as 0 (no corr.)

(-0.3 - -0.8) to (+0.3 - +0.6) → moderate reln (either +ve or -ve)

(-0.6 - -0.8) to (0.6 + +0.8) → good corr.

> 0.8 in both sides → strong corr.

Select imp. features

what decision to take →

→ no correln : drop these features → (feature selection technique)

- If a feature has 0 correln with every other feature then it can be removed, But if it has correln with atleast 1 feature it can't be removed.

Ex Pregnancies have only correln with age. (all other are zero)
so it can't be dropped.

Sns.heatmap(data.corr(), annot = True).

Multicollinearity analysis.

↳ multiple features with linear reln.

Sns.pairplot(data) → You can't determine reln.

* Sns.pairplot(data, hue = 'outcome')

19

Friday

April

2011

```
plt.xlabel('---')
plt.ylabel('---')
```

	April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13					1	2	3
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Week 13 • 91-274

$x =$
`sns.barplot(x=df['coln1'], y=df['coln2'])`

`sns.scatterplot(x=df['coln1'], y=df['coln2'])`

`sns.pointplot(x=df['coln1'], y=df['coln2'])`

gives trend.

`df.coln.value_counts()` ? to counts no. of values

$\cong df['coln'].value_counts()$ which belongs to some category.

`pd.crosstab(df.coln1, df.coln2)`

acts like pivot table.

counts value of each coln.

Σ

sal \rightarrow coln1, department \rightarrow coln2

Dep. IT HR CSE EC ME CIVIL

Sal.

	High	83	51	225	80	74	68
--	------	----	----	-----	----	----	----

low	609	368	364	180	402	451
-----	-----	-----	-----	-----	-----	-----

med.	535	180	370	225	376	363
------	-----	-----	-----	-----	-----	-----

Scatter plot with multivariate analysis,

`sns.lmplot(x='coln1', y='coln2', data=df, hue='coln', fitreg=False)`

`plt.show()`

It is like scatter plot with 'hue'.

• dropping coln.

`x=df.drop('coln1', axis=1)`

May	Mo	Tu	We	Th	Fr	Sa	Su	2011
17/22	30	31	1	2	3	4	5	6
18			7	8	9	10	11	12
19			13	14	15	16	17	18
20			19	20	21	22	23	24
21			25	26	27	28	29	

Saturday

April

2011

Week 13 • 92-273

Concepts of ML & Linear Regression.

Day-19.

Supervised

Regression classification

Unsupervised

clustering.

Ex: Stock Market Predn,

House Predn, Salary Predn

Target (y).

Predictive
Modelling

Regression \rightarrow x (Numerical/
categorical) \rightarrow y (numerical)

classification \rightarrow x (numerical) \rightarrow y (categorical)

Ex: predict
something

has High, Low,
Person, diabetes
or not

① Linear Regression.

wrt. if
something is
↑ or ↓ it is
linear.

predicting
real no.

$$y = mx + c \quad \begin{matrix} \text{slope} \\ \text{intercept} \end{matrix} \quad (\text{line eqn}).$$

$$m = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

Sunday

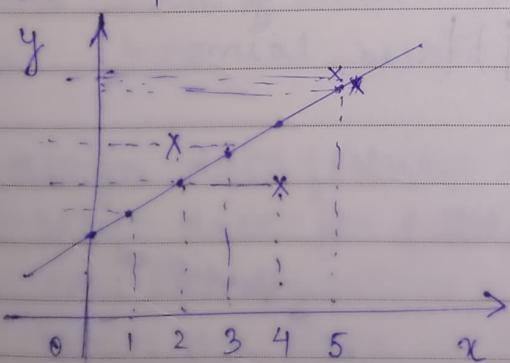
93-272

3

(D/P) \rightarrow depn var.

no. of
records

$$c = \frac{1}{n} (\sum y - m \sum x)$$



indep'n var.

Σx	x	y	Σy	Σx^2
0		2		
1		3		
2		5		
3		4		
4		6		
Σ	-	-	-	-

Put all the
values in
 $m \& c$,
& find y .

4

Monday

April 2011

Week 14 • 94-271

April						
Wk	Mo	Tu	We	Th	Fr	Sa
13					1	2
14	4	5	6	7	8	9
15	11	12	13	14	15	16
16	18	19	20	21	22	23
17	25	26	27	28	29	30

Ex. After calculation,

$$\hat{y} = 0.9x + 2.2$$

error \rightarrow predicted - actual.We can measure error in y .

x	y	\hat{y}	$y - \hat{y}$
0	2	2.2	2 - 2.2 = -0.2
1	3	3.1	
2	5	4	
3	4	4.9	
4	6	5.8	

Actual Predicted

Estimation.

Linear Regression Performance metrics / cost fn / Error

$$1 \quad \text{Error} = y - \hat{y}$$

$$2 \quad \text{Total error} = \sum (y - \hat{y})$$

$$5 \quad \text{Avg error / Mean error} = \frac{1}{n} * \sum (y - \hat{y})$$

$$6 \quad \text{Mean squared error} = \frac{1}{n} * (\sum (y - \hat{y}) * (y - \hat{y}))$$

$$7 \quad \text{Mean Absolute error} = \frac{1}{n} * |\sum (y - \hat{y})|$$

$$8 \quad \text{Root Mean Squared error} = \sqrt{\text{MSE}}$$

Predicted.

$$9 \quad \text{R-square Value: } R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

b/w \checkmark
lies 0 to 1

mean value.

May	2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
17/22	30	31				1	
18	2	3	4	5	6	7	8
19	9	10	11	12	13	14	15
20	16	17	18	19	20	21	22
21	23	24	25	26	27	28	29

Tuesday

April 2011

Implementation steps.

① Data gathering for problems
↳ collect the data

② Requirement Analysis to be done.
↳ what we want to explore.
(what is objective).

③ Data Exploration using pandas perform EDA.

↳ Data preprocessing / Wrangling...
cleaning data (find missing values...)

try to plot datapoints → come to know data patterns.
(We can't plot all data features, if 2 features → plot.
n features → you can't plot.)

④

Data preprocessing - check missing values.

- NAN & NAT value handling.

→ isna(), notna()

→ fillna() → use various methods to fill NAN

→ interpolation, bfill, ffill, pad..

⑤ Find nature of data distribution by various methods such as

↳ binomial, bernoulli, uniform ... soon

⑥ assign data features to x-input, y-output

⑦ split data for training & testing

⑧ call the algorithm from sklearn

⑨ train the data using algo. using fit method.

⑩ test the model using test data

⑪ call the performance metrics ex R², RMSE, accuracy score etc.

Adjusted-R²

when we are
not satisfied with R²
we go for adjusted R²
very precise = $1 - \frac{\sum (y_i - \hat{y}_i)^2}{n-p-1}$

$$\frac{\sum (y_i - \bar{y})^2}{n-1}$$

6

Wednesday

April 2011

Week 14 • 96-269

Wk	Mo	Tu	We	Th	Fr	Sa	Su	2011
13						1	2	3
14	4	5	6	7	8	9	10	
15	11	12	13	14	15	16	17	
16	18	19	20	21	22	23	24	
17	25	26	27	28	29	30		

Simple Linear Regression.Implementation.

Dataset used: Auto insurance payment.

```
df = pd.read_csv('auto-insurance-payment.csv')
df.head()
```

df.dtypes

df.info()

df.describe()

9

10

EDA.

11

df.hist()

12

df.corr()

1

```
sns.scatterplot(df['claims'], df['Payment'])
C check linearity.
```

2

3

sns.pairplot(df)

sns.boxplot() → check do we need to treat outliers?

x = df.iloc[:, :-1]

y = df.iloc[:, -1:]

Split records

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(x, y,
test_size=0.2, random_state=42)
20% to test size
80% to train set
```

	2011						
May	Mo	Tu	We	Th	Fr	Sa	Su
17/22		30	31		1		
18	2	3	4	5	6	7	8
19	9	10	11	12	13	14	15
20	16	17	18	19	20	21	22
21	23	24	25	26	27	28	29

Thursday

April

2011

Week 14 • 97-268

data.shape

x-train.shape

x-test.shape

Model Building.

```
from sklearn.linear_model import LinearRegression()
lin = LinearRegression()
```

Train Model.

```
lin.fit(x-train, y-train)
lin.coef_
lin.intercept_
```

Visualize training set result.

```
plt.scatter(x-train, y-train, color='red')
plt.plot(x-train, lin.predict(x-train), color='green')
plt.show()
```

Test Model.

y_pred = lin.predict(x-test)

Visualize test set result

```
plt.scatter(x-test, y-test, color='blue')
plt.plot(x-test, lin.predict(x-test), color='green')
plt.show()
```

include
put.title('')
put.xlabel('')
put.ylabel('')

8

Friday

April

2011

Week 14 • 98-267

	April							2011
Wk	Mo	Tu	We	Th	Fr	Sa	Su	
13						1	2	3
14	4	5	6	7	8	9	10	
15	11	12	13	14	15	16	17	
16	18	19	20	21	22	23	24	
17	25	26	27	28	29	30		

Estimate cost

from sklearn.metrics import mean_squared_error,
r2_score.

$$\text{RMSE} = \text{np.sqrt}(\text{mean-squared-error}(y_{\text{test}}, y_{\text{pred}}))$$

mp. * $R^2 = r2_score(y_{\text{test}}, y_{\text{pred}})$

Ques. How to predict for unseen value?

$$\text{unseen_pred} = \text{lin.predict(np.array([[108]]))}$$

→ give any value
it will predict
op.

Linear.

1 Multiple Regression. (Bivariate Analysis). Day-20.

Dataset → 'health insurance cost'

→ 19 coln only.

5 df = pd.read_csv('health-insurance-mising-data.csv')

6 df.head()

df.dtypes. → Algo. don't take object values. so
df.info() convert 'obj' to 'numeric'.

df.describe()

df.isna().sum()

drop irrelevant coln. (check ~~which~~ domain knowledge or
by stats. or

df = df.drop(['coln', 'coln'], axis=1) by plotting graphs..)
which coln is irrelevant
for predicting target.

	May 2011						
Wk	Mo	Tu	We	Th	Fr	Sa	Su
17/22	30	31				1	
18	2	3	4	5	6	7	8
19	9	10	11	12	13	14	15
20	16	17	18	19	20	21	22
21	23	24	25	26	27	28	29

Saturday

April 2011

9

$\text{df} = \text{df}.fillna(\text{method}='bfill')$

→ anything you can fill with Week 14 • 99-266
mean...etc

Q: When to do Normalization?

if instances of features differs numerically as larger value. ex one col has value in single digit, others has 5, 6, digits. Statistically higher digits features dominates lower dig. values ie,

which feature has higher value than std, mean etc can dominate all other features.

By Normalization → Performance of model ↑es.

① Standard Scalar. →

$$Z = \frac{x - \text{mean}}{\text{std}} = \frac{x - \bar{y}}{\sigma}$$

from sklearn.preprocessing import StandardScaler.

sc = StandardScaler()

df1 = sc.fit_transform(df)

Sunday

100-265

10

it transforms our output also.

Model1.

$x = df.iloc[:, 0:3]$

$y = df.iloc[:, -1:]$

} Here sex, smoker, region can be dropped.
& no standard scaling is applied.

11

Monday

April

2011

April 2011						
Wk	Mo	Tu	We	Th	Fr	Sa
13					1	2
14	4	5	6	7	8	9
15	11	12	13	14	15	16
16	18	19	20	21	22	23
17	25	26	27	28	29	30

Week 15 • 101-264

Split it in 75-25%.

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y,
                                                    test_size = 0.25, random_state
                                                    = 42)
```

Train Model.

```
from sklearn.linear_model import LinearRegression
multiple = LinearRegression()
```

we have parameter

normalize = False

default.

(we don't need to separately do normalization, when we make it True, algo will

automatically do normalization)

Test model.

y_pred = multiple.predict(xtest)

on unseen data.

unseen = multiple.predict(np.array([[19, 33, 0]]))

unseen

↳ -134.345

model is performing too bad.
(-ve prediction).Evaluation.

```
from sklearn.metrics import r2_score
r2_score(y-test, y-pred)
```

-2.1527

worst model.)

here R^2
is
negative
i.e. it
is too
bad
model.

R2 score $0 < R^2 < 1$

May 2011						
Wk	Mo	Tu	We	Th	Fr	Sa Su
17/22	30	31				1
18	2	3	4	5	6	7
19	9	10	11	12	13	14
20	16	17	18	19	20	21
21	23	24	25	26	27	28

Tuesday

April

2011

12

Week 15 • 102-263

~~df~~.corr()

→ no corr b/w features.

∴ we do EDA process in starting to know whether our data is compatible or we can go further with this or not.

df~~data~~.hist(figsize=(10,8))

sns.pairplot(df)

df2 = pd.read_csv('health-insurance-cost.csv')

df2.shape

→ (1338, 11)

converting categorical to Numerical →

from sklearn.preprocessing import LabelEncoder

lab = LabelEncoder()

df2['smoker'] = lab.fit_transform(df2['smoker'])

converting int/float... into category type:

• df['coln'] = df['coln'].astype('category')

• df['coln'] = ~~df~~ pd.Categorical(df['coln'])

category
types
can't be
in
statistical
summary.

Exploring categorical feature.

df['coln'].value_counts(normalize=False)

sns.countplot(x='Gender', data=df)

any
coln name.

13

Wednesday

April

2011

	April	2011					
Wk	Mo	Tu	We	Th	Fr	Sa	Su
13					1	2	3
14	4	5	6	7	8	9	10
15	11	12	13	14	15	16	17
16	18	19	20	21	22	23	24
17	25	26	27	28	29	30	

Week 15 • 103-262

Boxplot with 2 features -`sns.boxplot(x=df['col1'], y=df['col2'])`# Multivariate Analysis.

Ex
`g = sns.FacetGrid(df, col='Married', row='Gender')
g.map(sns.boxplot, 'Losses')`

Pivot table → summarization of more than 2 or 3 features

Ex
`pd.pivot_table(data=df, values='Losses', index='No. of vehicle', columns='Gender', aggfunc='mean')`

Convert categorical → dummy var.

`X = pd.get_dummies(X, drop_first=True)`# OLS Method`import statsmodels.api as sm``lm = sm.OLS(y_train, X_train)``lm2 = lm.fit()``lm2.summary()`→ check r^2 value.

OLS

ordinary least square

Linear Regression.`lm3 = LinearRegression()``lm3.fit(X_train, y_train)``y_pred = lm3.predict(X_test)`check r^2 then..