

20

Friday

May

2011

Week 20 • 140-225

May	Wk	Mo	Tu	We	Th	Fr	Sa	Su	2011
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	

Unsupervised Learning.

Day 28.

Association-Rule Technique.

↳ Help to build the recommended system.
recommendation

① Purpose:

To find the frequently occurring combination of itemsets.

when we are buying a product what is prob. of buying another similar product.

Ex. Amazon, Netflix, YouTube

9

10 ② Support → Prob. of occurring (buying). → \Rightarrow what is prob. of specific item.

buying Colgate Brush.

③ confidence

1 conditional prob. of occurring combination

2 of itemsets.

Ex. What is prob. of buying colgate brush

along with colgate paste.

5

Ques.

What is Association Rule?
Study of what goes with what customer who bought x also bought y.
what symptoms go with what diagnosis.

Methods of Association Rule Technique.

→ Apriori Algorithm

(Use to find the frequently occurring combination of items.)

→ MBA (Market Basket Analysis)

(combining fast moving combination with slow moving combination).
product product

④ Transaction Based or event based

(each every transaction, what is combination of products customer n).

Also called 'MBA' and 'affinity analysis'

purchased

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

Saturday

May

2011

21

Week 20 • 141-224

Association Rule Mining.

Finding frequent patterns, associations, correlations or causal structures among sets of items in transaction databases, relational DB.

Application: Basket data analysis, cross-marketing, clustering etc.

Ex: Rule form:

" Antecedent → consequence [support, confidence]"

Diapers → Beers [0.5%, 60%]

Frequent terms used in this Model.

- Items (I): ie product
- Transaction (t): ie set of items ($t \subseteq I$)
- Transaction Database (T): ie set of Transactions ,
 $T = \{t_1, t_2, \dots, t_n\}$

① Apriori principle:

↳ Any subset of frequent itemset must be frequent.

(Used to find frequently occurring combination of ~~itemsets~~ itemset).

Ex:

list of Items in store.

Sunday

22

142-223

1 → Milk

2 → Jam

3 → Bread

4 → Wheat Bread

5 → Butter

Database D.

TID	Items
100	1 3 4
200	2 3 5
300	1 2 3 5
400	2 5

• listing all occurrences

how many

times item

appear

in data

base

{1} → 2

{2} → 3

{3} → 3

Itemset sup

{4} → 1

{5} → 3

Itemset sup

Item has least no. of occurrence

(remove from list).

23

Monday

May

2011

Week 21 • 143-222

	Mo	Tu	We	Th	Fr	Sa	Su
May Wk 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

Apply combinations

updated list →	itemset	sup.	itemset.	sup.	how many
	{13}	2	{123}	1	customers purchased these combination
	{23}	3	{133}	2	
	{33}	3	{153}	1	
	{53}	3	{233}	2	
			{253}	3	
			{353}	2	

* put it here)

Updated:

	itemset	sup.
9	{1,33}	2
10	{2,33}	2
11	{2,53}	3
12	{3,53}	2

now removing least occurred.

Aim:

(Find the freq. occurring itemset)

combining
 rem one
 least occurred.

final combination (frequently occurred combination)
 $\{2,3,5\} \rightarrow 2$

* Association Rule is applicable to all the sectors.

6 Statistical significance of Rules,

- Test of proportions
- Formal adjustment of statistical significance.
(control prob. of Type I error).

Multilevel Association Rule.

If the customer who buy 1 product, what is prob. of buying another product. If he buys these 2 product, what is prob. of buying 3rd product

(Items at lower level are expected to have lower support).

June 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			

Tuesday

May

2011

24

Week 21 • 144-221

Market Basket Analysis

↳ combining fast moving product with slow moving products.
(making combo.)

Important aspect:

- Customer Demographics
- Order Characteristics
- Item Popularity.

(New Product Intro. Strategy).

But if we introduce new product separately, people may not buy it so we make & launch with combo.

→ Colgate Bassalt + Colgate Paste + added new product
if people like it, they will also buy separately (Vedshakhi) after using.

Variations/Extension in Association Rule.

- Disassociation Rules
- Sequential Analysis using Association Rules
- Association based Clustering.

Over a period of time, customer buying preferences changed. Based on this the recommendation also gets changed.

Practical Implementation.

import pandas as pd

```
df = [['Milk', 'Bread', 'Wheat Bread'], ['Jam', 'Bread', 'Butter'], ['Milk', 'Jam', 'Butter', 'Bread'], ['Jam', 'Bread']]
```

25

Wednesday

May

2011

Week 21 • 145-220

May	Wk	Mo	Tu	We	Th	Fr	Sa	Su
17/22		30	31					
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

```
from mixtend.preprocessing import TransactionEncoder  
te = TransactionEncoder()  
te_ary = te.fit(df).transform(df)  
df1 = pd.DataFrame(te_ary, columns=te.columns_)  
print(df1)
```

9 from mixtend.frequent-patterns import apriori
10 apriori(df1, min_support=0.1)
11 ↪ min. occurrence
12
freq = apriori(df1, min_support=0.1, use_colnames=True)
13 ↪ uses colnames
14 (items) instead of 0, 1, 2, 3...
15 # to find the occurrences length; ie how many times particular
16 item is occurred.

```
freq['len'] = freq['itemset'].apply(lambda x: len(x))  
freq  
    ↪ print along with length.
```

↳ length should be 2. (50% occurrence.) ↳ 2/4 total items

```
freq[(freq['len'] == 2) & (freq['support'] >= 0.5)]
```

Based on this you can frame diff. condn.

June						
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

May

2011

26

Week 21 • 146-219

Project Building Session. (Project 1.)

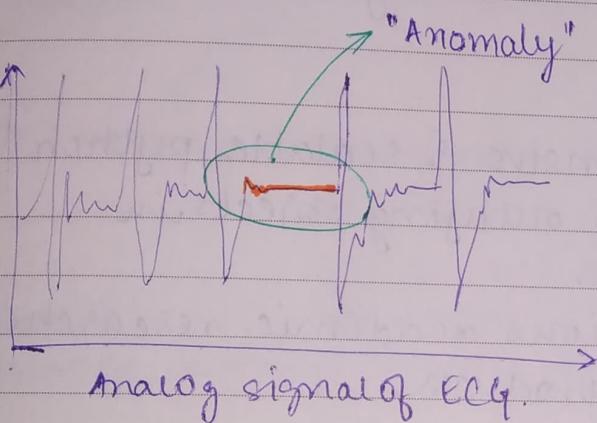
"it can be
said as
outliers."

Anomaly Detection in Machine Learning.

something
abnormal

(uncertainty & Noise / Disturbance).

Ex



Here is a pattern in all,
but the green circle
denotes that there is
some abnormal pattern
in it.

- Anomaly detection is a technique used to identify unusual patterns that do not conform to expected behavior, called outliers.
- It is also called Novelty Detection, Outlier Detection, Forgery Detection, & Out-of-distribution Detection.
- Anomalies are datapts. that are inconsistent with distribution of majority of data points.

Novelty Detection

↳ if you have any
new features
which is similar to
object features.

Anomaly Detection

↳ different features
which has been
included as
actual input.

27

Friday

May

2011

May	Mo	Tu	We	Th	Fr	Sa	Su	2011
Wk	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

Week 21 • 147-218

Types (Areal Domain) of anomaly Detection.

- Time series anomaly detection
- Video-level anomaly detection.
- Image-level
 - 1. Anomaly classification target
 - 2. out-of-distribution (OOD) Detection target
 - 3. Anomaly Segmentation target

PyOD

- PyOD is a comprehensive & scalable python toolkit for detecting outlying objects in multivariate data.
- PyOD is used in various academic researches & commercial products.
- optimized performance with JIT & parallelization when possible, using numba & joblib packages used to deploy any model

Practical Implementation,

- install library → pip install pyod.
or use
conda install -c conda-forge pyod

Introduction. ① Linear Models for outliers Detection

- PCA (Principal component Analysis)
 - use sum of weighted projected dist. to eigenvector hyperplane as outliers scores.

2011						
June	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

May 2011

28

Ex PCA
used to find significant features out of entire population.

identify contribution of every individual feature. Based on contribution can we remove any feature.
Week 21 • 148-217

A	B	C	D	E	F (target)
40	35	22	1	2	

Significant features. can be removed as they are not much affecting our target var.

• MCD (minimum covariance determinant).

When minimum covariance is there, then there is no outlier

• OCSVM (one class support vector machine)

(hyperplane concept of SVM).

We put all inline datapoint in one class & outliers in other class

② Proximity Based outlier detection.

• LOF (Local Outlier Factor)

If you have one set of data pts. which you have so close outliers, it is called LOF

• CBLOF (Clustering Based LOF)
if you have 2 classes, behind these you have any other category, it is outlier.

Sunday

Ex. 149-216

29

Proximity sensors used in cars / parking system if your car is going close to other car, alarm shouts!!!

30

Monday

May

2011

Week 22 • 150-215

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

- KNN (*Knight K-nearest neighbors).
↳(use distance of kth nearest neighbor as outlier).
- HBOS (Histogram Based outlier score)

③ Probabilistic Models for outlier Detection.

- ABOD (Angle Based Outlier detection)

④ Outlier ensembles & combination frameworks optimization.

- Isolation forest
- Feature Bagging.

Practical Implementation. → check file provided by LMS.

- Import Python Packages, Metrics Package, Pyod Package

```
import
from scipy.io import loadmat
```

matfiles are matlab files

- Define data file & read X & y.

Here X & y are
defined in
matfiles.

matfiles are
different from
csv, excel files.
They are in
dict. format.

```
= mat['y'].ravel()
contamination
20% ID.
```

Based on stat factors,
algo. is going to be
detect outlier.

clf.decision_function(x_test.)

similar to predictfn
in ml.

Here we are
using
pyod

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

Week 23 ■ 159-206

Day 29 → K means clustering.
(implementation from scratch.)
Build step by step every step
with Python.
- refer Notebook k.ipynb

Day 30, →

```
9     col_list = ['Age', 'Education', 'TotalWorkingYears',  
10    'YearsInCurrentRole']  
11    ↓
```

tatjine

only 4

coin
round

^{com}
randomly) df.head()

import kmeans algo.

```
from sklearn.cluster import KMeans
```

```
clust = kMeans(n_clusters=3)
```

clust.fit (df)

5)  randomly taken

as 3.
we don't know
how many
clusters
are there).

```
from sklearn.metrics
```

import silhouette-score

labels = clust.labels

silhouette-score(aff, labels)

$$\rightarrow 0.3532 \text{ (op.)}$$

If you have

O that
means

data points are overlapped.

\Leftarrow (Same datapts as
in cluster A & cluster
B also.)

To estimate clustering performance we have matrix →

method \downarrow compute score it lies in b/w -1 to +1. $\rightarrow 0.3532 \text{ (O/P).}$

worst
fit.
(poor
clustering)

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

Thursday

June 2011

9

Week 23 • 160-205

range

0.41 - 1.0

0.51 - 0.7

0.26 - 0.5

< 0.25

Interpretation

A strong structure has been found

A reasonable structure " "

The structure is weak & could be artificial

No substantial structure is found.

• Plotting elbow curve to get optimum k-value.

k = list(range(1, 11))

wcss = []

for i in k:

kmeans = KMeans(n_clusters=i, random_state=1, init='k-means++')

kmeans.fit(df)

labels = kmeans.labels_

wcss.append(kmeans.inertia_)

plt.plot(k, wcss, color='m')

plt.xlabel('Number of clusters')

plt.ylabel('within cluster sum of sq. errors')

plt.show()

default value
in kmeans.

→ By plot we got k=2,

* Plotting Bar Graph.

cluster = kmeans(n_clusters=2, random_state=1)

df['cluster'] = cluster.fit_predict(df)

df['cluster'].value_counts().sort_index().plot(kind='bar')

plt.xlabel('Cluster Number')

plt.ylabel('size')

plt.show()

10

Friday

June

wcss \rightarrow actual - predicted
(error).

2011

Week 23 • 161-204

June	Mo	Tu	We	Th	Fr	Sa	Su	2011
Wk	22	23	24	25	26	27	28	29
		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			

PLOTING silhouette score \rightarrow $K = \text{list}(\text{range}(2, 11))$ $ys = []$ for $i \in K :$ $kmeans = KMeans(n_clusters=i, random_state=1)$ $kmeans = fit(laf)$ $labels = kmeans.labels -$ $ys.append(silhouette_score(laf, labels))$

9 plt.plot(K, ys, color='m')

10 plt.xlabel('K')

11 plt.ylabel('silhouette score')

12

PCA V/S LDA



Principal component Analysis



Linear discriminant Analysis

• Feature Selection technique / Dimensionality Reduction.

Target Based Problem
(for supervised learning)↓
useLDA to select
significant
features.C LDA is specifically designed
for supervised learning)Target less.
(for unsupervised)

↓

use PCA to
get important
feature.(PCA can be used for
supervised or
unsupervised.)
But designed for unsupervised

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

- LDA, PCA makes model complex.
We don't use it much.
- (Try to avoid LDA, PCA in real time).
- If performance increases then use it

Saturday

June

2011

11

Week 23 • 162-203

LDA.

- Step 1. computing d-dimensional mean vectors.
- Step 2. computing scatter matrices.
 - within class scatter matrix
 - between class scatter matrix
- Step 3. find eigenvalues & eigenvectors by linear algebra packages.
- Step 4. calculate ratio of eigen values.
- Step 5. eliminate non-significant features.
& evaluate with significant ones.

from sklearn.discriminant_analysis import

linearDiscriminantAnalysis as LDA

lدا = LDA(n_components=2)

take top 2 eigen values.

(with highest %.)

lدا.fit_transform(x,y)

After this
give this transformed
value to my
algorithm, logistic,
linear etc.
~~or~~ linear etc.
(it will give
my 2 feature
to algo. → you
will select
new 2) :: training
time will be
less.

Sunday

163-202

12

PCA.

all the steps are same as LDA, only difference
is in step 2.

→ It will not compute 2 different scattermatrix.
as we don't have target defined in unsuperv.

from sklearn.decomposition import PCA

pca = PCA(n_components=2).

X_pca = pca.fit_transform(X)

13

Monday

June

2011

Week 24 • 164-201

June	Mo	Tu	We	Th	Fr	Sa	Su
Wk	22	23	24	25	26	27	28
	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				

Day-31.REVISION.

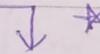
• Unsupervised Learning -

- It is a class of ml technique to find pattern in data.
- The data given is not labeled. (no target(y))
- Patterns in data are used to identify / group similar observation.

9 Clustering

- It discovers hidden structures of data.
- A way to decompose set into subset, with each subset representing group with similar characteristics.
- These groups $\xrightarrow{\text{are called}}$ clusters.

4 clustering



Top-down
approach

connectivity based
clustering

centroid based
clustering.

Ex: Hierarchical

Ex: K-means

K-Means

- ① Select random points as centroids.
- ② Calculate distance.

→ Distance being calculated by 3 methods

$$|(x_2 - x_1)| + |y_2 - y_1| \rightarrow \bullet \text{ manhattan}$$

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \rightarrow \bullet \text{ Euclidean (Preferred)}$$

$$\max(|x_2 - x_1|, |y_2 - y_1|) \rightarrow \bullet \text{ chebyshov}$$

July 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	31				

Tuesday

June

2011

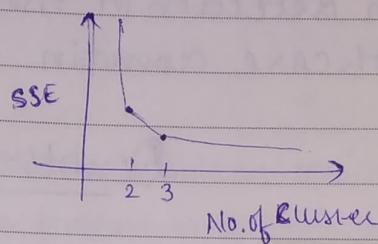
14

Week 24 • 165-200

- (3) update centroids by taking the average of its data points
- (4) repeat step ② & step ③ till distance of each data point is nearer to the own cluster mean than the other cluster.

Elbow Method.

To find optimum k-value.



termination condition.

Performance Estimation : (Silhouette-coefficient)

$$s(i) = \begin{cases} \frac{1 - a(i)}{b(i)} & \text{if } a(i) < b(i) \\ 0 & \text{if } a(i) = b(i) \\ \frac{b(i) - 1}{a(i)} & \text{if } a(i) > b(i) \end{cases}$$

How it is calculated

- It is a measure of how similar a data point is to its own cluster compared to that of other clusters.
- lies b/w -1 to +1

Here,

$a(i) \rightarrow$ dis. of individual pt.

$b(i) \rightarrow$ mean value of centroids

Pros:

- with large var., k-means is faster than other clustering algo.
- produces tighter clusters than hierarchical.
- less impacted by outliers

Cons:

- predicting no. of clusters is a tedious task
- If data has original clusters, the algo. doesn't work well

15

Wednesday

June

2011

Week 24 • 166-199

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

• img & video tracking → kmeans
is best.

④ Industry Applications →

- * ① Customer segmentation:
 - ② Anomaly detection
 - ③ Creating news feeds (cluster articles, based on similarity)
 - ④ Pattern detection in medical imaging for diagnostics.
- Health care domain.

Doubts-clearing Session

In Naive Bayes,

- Bernoulli → when DV is binary categorical
- multinomial → more than 2 categories in DV
- Gaussian → when numerical values converted into range.

Precision → zscore / ztest}

Fscore → ftest

Recall → Ttest

} as we did
in statistical
learning, these values are
same.

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

Thursday

June

2011

16

Week 24 • 167-198

Day - 32.

Ensemble Technique.

Agenda.

- Intro. to Ensemble Techniques.
- Bagging
- Boosting
- Stacking
- case study on classification problem optimization with ensemble
- case study on Regression problem optimization with ensemble
- How to take forward this production.

Ensemble Techniques /methods :

Aggregating

divide & aggregate
all the soln of task
then come with
a soln.)

- Ensemble is a meta-algorithms that combine several ML techniques into one predictive model in order to decrease ① variance (bagging), ② bias (Boosting), ③ improve prediction (stacking).

→ To improve performance of model, Ensemble is used.

* Bias should be minimized assumption will be close to outcome)

assumption (objective is to minimize variance).

actual - pred.

17

Friday

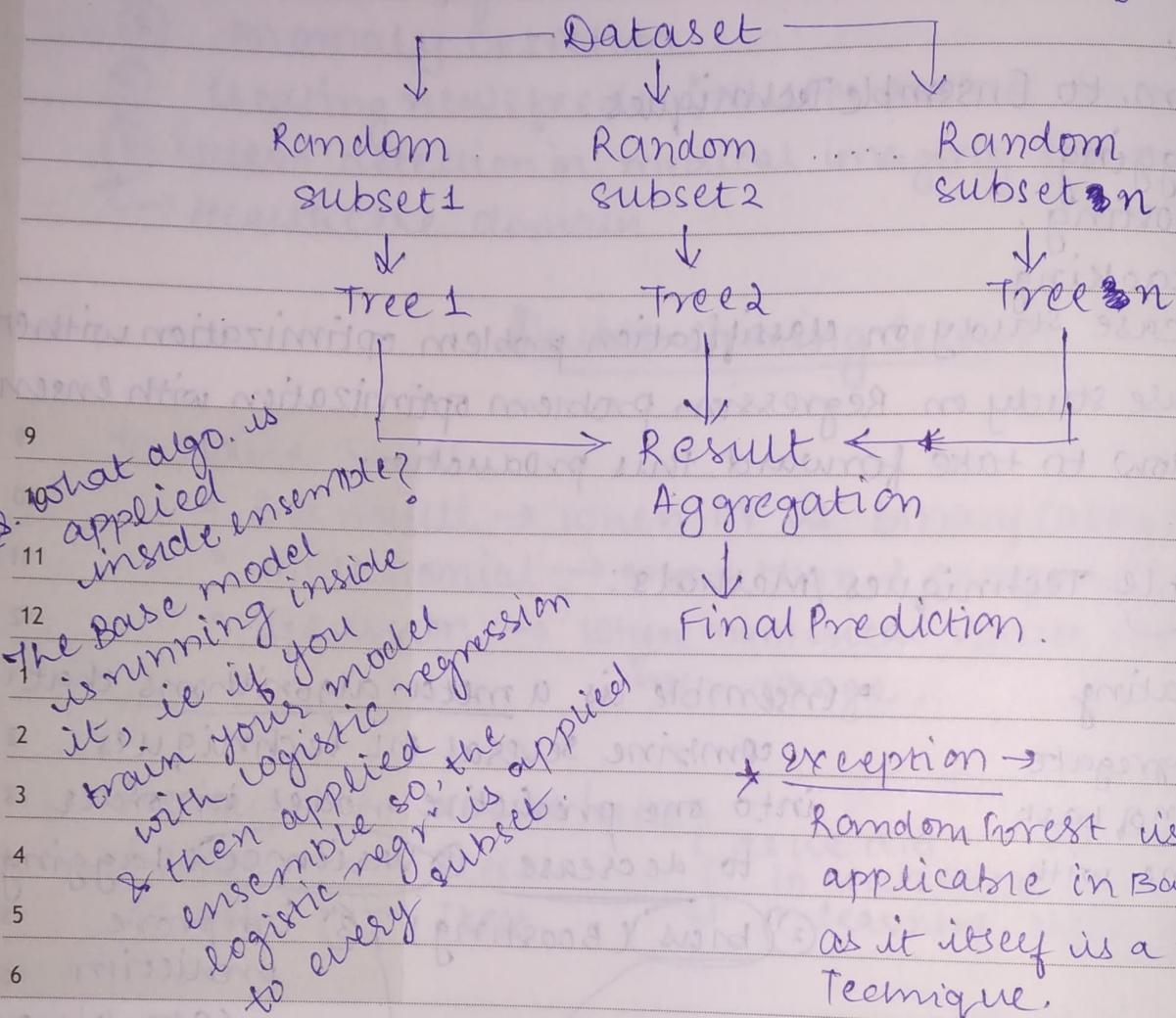
June

2011

Week 24 • 168-197

Bagging Ensemble.

it converts the dataset into randomly n subsets so that every subset include every type of data.

BaggingBoosting

Bootstrap aggregation
 which invokes automatically. (automatically, divides data)

If model is split into 10 subsample, for all 10 subsamples, Model is built/Trained at a time. (Parallel Model construction)
 ... Training time is very less.

It is not parallel execution, it is sequential execution.

Explanation.

July 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	31				

Saturday

June

2011

18

Week 24 • 169-196

Boosting

- ex. If the dataset is split into 10 records, then the 1st model is trained, once the model is built or trained, the result of 1st model goes to model 2. (If base algo. is logistic, then what hyperparameters are used by model 1, what was the error rate, why it is getting that error... are the questions model 2 takes.)
- * using the model 1 output as an inference, it trains the model 2.... this process continues till best model is trained, it goes sequentially, the final outcome is far better from model 1. (Best model trained).

∴ training time
is very high
as ~~we~~ till the time
1st model is
trained model 2
has to wait.
model has to
wait till completion
of n-1 model.

11
12
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

Stacking/Voting

- ↳ calling multiple algorithms for a same dataset & then find which algo. is giving best accuracy on dataset.

Sunday

19

- ex. If dataset is split into 10 subsamples, in first time, logistic is applied to all 10 samples, out of which 7 are saying/giving accuracy 1. ∴ for logistic model it will take the no. which is getting highest no. of votes (Here 1) then DT is performed... similar process goes on for all 10 samples. (It is concerned about majority of samples.)

20

Except Random Forest, Extra Tree all algo.
Monday can use these Ensemble techniques.

June 2011

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

Use Avg
all pred

Week 25 • 171-194

* use

voting
for
final
predn.

classification

- Bagging classifier()
- Random Forest classifier()
- Extra Trees classifier()

- Bagging → splits data into n subsamples

→ create model (& the samples are chosen with ~~with~~ replacement.)

the original
dataset remains
same.

without
replacement

- 11 original dataset
- 11 original dataset
- 11 it's not same
- 11 it's not model
- 1 choose 100 records from dataset.
- 2 Now dataset is left with 900 records.
- 3 Model 1 will choose samples from these 900 records.
- 4 Model 2 will choose samples from these 900 records.

Regression

- Bagging Regressor()
- Random Forest Regressor
- Extra Trees Regressor()

Random Forest

→ It follows Bagging technique but it selects a set of features which are used to decide best split at each node of DT.

→ Base estimator is DT.

→ You can use pruning (max depth, etc) can be used while classifying.

? Extra Trees (Extremely Randomized Trees)

→ You can't use pruning, this is unpruned algorithm.

→

(+99 more to follow) was written CONCERNED 10/13

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

XGBoost → (Extreme Boosting)
ie, it will

go till maximum.

Global minima.

Tuesday

June

2011

21

Week 25 • 172-193

Boosting Algorithms:

Classification

- AdaBoostClassifier()
- GradientBoostingClassifier()

Regression

- AdaBoostRegressor()
- GradientBoostingRegressor()

Adapting
Boosting.

(Based on weak learner
current model is updating
the weights ie it is
Adapted.)

Gradient Boost

(It is apt.
where we
encounter
cost fn)

In this Gradients
are trying to be
minimized
based on
previous
model.

Practical Implementation.

→ Dataset ('Pima Diabetes')

```
df = pd.read_csv('diabetes.csv')
df.head()
```

- split into x & y.

x = df.drop('Outcome', axis=1)

y = df.outcome.

- split into train & test data.

• Build logistic Model & find accuracy → ~78%.

• Build KNN model & find accuracy. → ~74%.

Ex. $K=3$

22

Wednesday

June

2011

Week 25 • 173-192

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

* Bagging.

```
from sklearn.ensemble import BaggingClassifier
bg = BaggingClassifier(base_estimator=log-reg,
                       n_estimators=12,
                       verbose=1)
```

bg.fit(xtrain, ytrain)

bgpred = bg.predict(xtest)

logistic
reg.
funcn
(defined
earlier)

Training score.

bg.score(xtrain, ytrain)

Testing score.

bg.score(xtest, ytest)

Accuracy.

accuracy_score(ytest, bgpred) ~ 78%.

(there was no change in logistic regression acc.,
so don't include bagging classifier for this
dataset).

Day-33.

* RandomForest.

```
from sklearn.ensemble import RandomForestClassifier
```

rf = RandomForestClassifier(n_estimators=10)

rf.fit(xtrain, ytrain)

rfpred = rf.predict(xtest)

accuracy_score(ytest, rfpred) → 71%.

	July						
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	31	

Thursday
June 2011

23

Week 25 • 174-191

* finetune the hyperparameters

`rfl = RandomForestClassifier(n_estimators=8, max_depth=5,
max_features=6, random_state=2,
verbose=2)`

`rfl.fit(xtrain, ytrain)`

`rfl_pred = rfl.predict(xtest)`

`accuracy_score(ytest, rfl_pred)`

{ accuracy-score(ytest, rfl_pred)
only.

from previous model the acc. increases from 71 to 77%.

* ExtraTree Ensemble

`from sklearn.ensemble import ExtraTreesClassifier`

`et = ExtraTreesClassifier(n_estimators=8)`

`et.fit(xtrain, ytrain)`

`et_pred = et.predict(xtest)`

`accuracy_score(ytest, et_pred) → 73%`

when you increase to 10, accuracy is 75%.

* Overfitting & Underfitting

Test ↑
Train ↓

Train ↑
Test ↓

`rf.score(xtrain, ytrain) → Training accuracy.`
`rf.score(xtest, ytest) → Test accuracy`

- ① High Variance
- ② more outliers

* since, the data is imbalanced, so there is high chance of overfitted da mode.

0 → 500

1 → 268

Not equal 0's & 1's ∵ imbalanced dataset.

24

Friday

June

2011

June						
Mo	Tu	We	Th	Fr	Sa	Su
22						
23	6	7	8	9	10	11
24	13	14	15	16	17	18
25	20	21	22	23	24	25
26	27	28	29	30		

Week 25 • 175-190

Q.

how to overcome this problem?

You may use balancing technique to balance dataset,
either 268 pushes to 500 or 500 pulls down to 268.
This Balancing might reduce/overcome overfitting problem.

Boosting.

ADA BOOST

```

from sklearn.ensemble import AdaBoostClassifier
ada = AdaBoostClassifier(n_estimators=10) → logistic
ada.fit(xtrain, ytrain)   regres.
ada_pred = ada.predict(xtest)
accuracy_score(ytest, ada_pred) → 67%.

```

Gradient Boost.

```

from sklearn.ensemble import GradientBoostingClassifier
gb = GradientBoostingClassifier(n_estimators=10)
gb.fit(xtrain, ytrain)
gb_pred = gb.predict(xtest)
accuracy_score(ytest, gb_pred) → 77%.

```