

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

Thursday

April

2011

14

Week 15 • 104-261

Logistic Regression.

Day-21.

Overview of linear regression → there is dependent & independent variable.
 Linear Reg. → dependence technique.
 Correlation → interdependence technique.
 there is no dependent & no independent var.

Statistical Variable.

(response var)

Dependn Independn.

O/P var.

\hat{y}_x , rainfall,

\hat{y}_x crop production

fertilizer usage

Purpose of LR.

it is used to
 find causal
 effect relationship
 b/w independent var
 on dependn vars.

(How much one var. is
 going to affect another
 var. → Regression)

independn affect dependn
 going to

When both vars are continuous:

depn & indepn → we go for Linear Regr.)

Regression

(one Indepn) ↓

Simple Regr.

↓
Linear

Non-linear

↓
(More than 1
indepn
var)
Multiple Regr.

↓
Linear

↓
Non-linear.

15

Friday

April

2011

$$(IV) \quad y = \beta_0 + \beta_1 x + e \rightarrow \text{random error.}$$

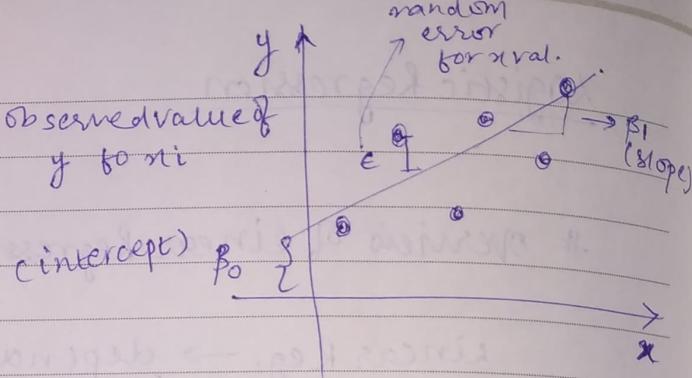
↓ ↓ ↓ ↓ ↓
 intercept wef indepn var. (IDV)

Week 15 • 105-260

	Mo	Tu	We	Th	Fr	Sa	Su	2011
Wk	13	14	15	16	17	18	19	
	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		

Conducting Regr. Analysis.

1. Plot Scatter Diagram
2. Formulate General Model
3. Estimate Parameters
4. Determine strength & significance of Association.
5. Test the for Significance.



axis.

$$\text{OR}^2 \rightarrow [58.08\%]$$

 $y \rightarrow \text{depn. var.}$ x → indep. var
° infer p-value ($p < 0.05$ ie it is affecting y)

° frame eqn. (coeff + intercept)

(formed this eqn by stats).

$$\text{Ex} \quad \text{houseprice} = 98.2483 + 0.10977 (\text{sq. feet})$$

const.

on an avg. value

(of house increases by $0.10977 (\text{sq. feet})$ for each sq. foot)

ie no one can construct a square

feet. min.

you need to pay

amount you have to pay.

for each one sq. foot

of size.

for an additional sq. feet, you have min 109.77 dollars.

$$y = \beta_0 + \beta_1 x + e$$

const avg. Sq. feet.

(coeff.)

C.I. Upper 95%

0.18580

X 1000

= \$185.80

max you have to pay this for house.

min you have to pay this for house.

= \$33.74

min you have to pay this for house.

= \$33.74

min you have to pay this for house.

° This is for simple linear regression

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

Saturday

April

2011

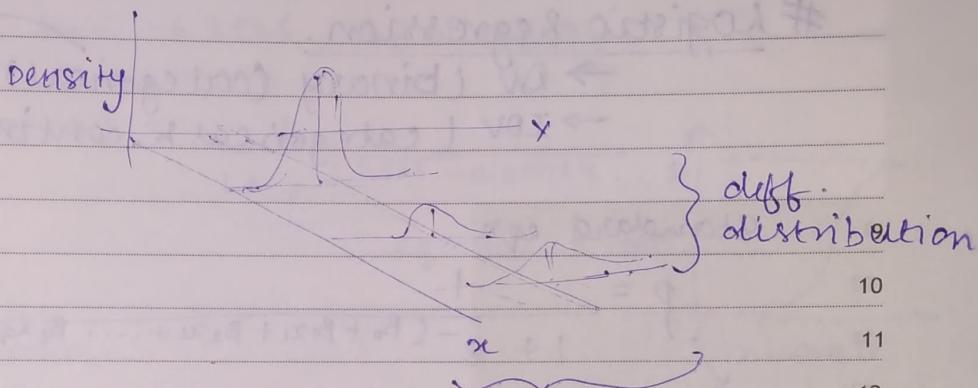
16

Week 15 • 106-259

• Multiple linear regression.

Regression Assumption.

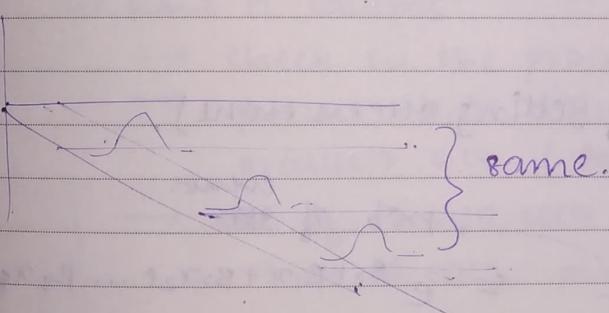
- ① Heteroskedasticity. → This test should pass.
 ↗ different spread.
 ↗ all the IDV should be distributed in different way.



homoskedasticity.
 same spread
 (equal)

This test should be passed.

(skewness,
 kurtosis
 etc.)



② Multicollinearity.

- ↘ IDV
 whatever you are taking & all the IDV are correlated with each other is called multicollinearity.

Sunday

17

107-258

When you are going for Reg. ~~test~~, multicollinearity test should fail.

(The IDV we are taking should not be correlated with each other)

- Heteroskedasticity → Pass
- Autocorreln → Pass
- Multicollinearity → Fail

18

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 16 • 108-257

③ Autocorrelation

↳ whatever DV we are taking in LR, The DV & IDV are correlated with each other.

Logistic Regression.

→ DV (binary (categorical)) ↗ only 2 (0 or 1, Yes/No...)

→ IDV (categorical & continuous.)

standard eqn:

$$0 \leq p \leq 1$$

$$p = \frac{1}{1 + e^{-(B_0 + B_1x_1 + B_2x_2 + \dots + B_qx_q)}}$$

prob. of belonging
to class 1.

Prob. of
getting outcome.

(P prob. of getting success event).

Odds ratio:

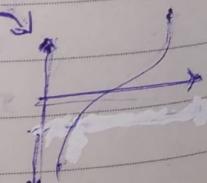
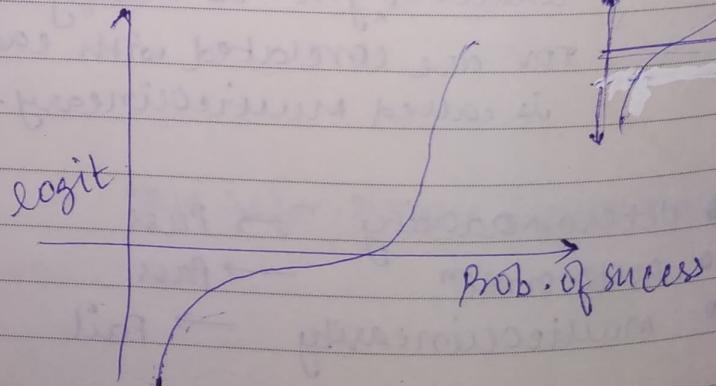
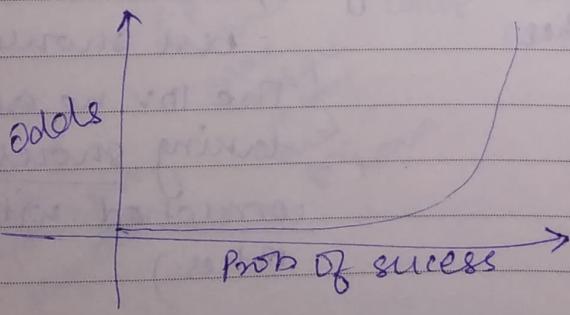
prob. of event.

$$\text{Odds} = \frac{p}{1-p} \Rightarrow e^{B_0 + B_1x_1 + B_2x_2 + \dots + B_qx_q}$$

log on both sides.

$$\log(\text{odds}) = B_0 + B_1x_1 + B_2x_2 + \dots + B_qx_q$$

$$\log(\text{odds}) = \text{logit}$$



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

Tuesday

April

2011

19

Week 16 • 109-256

single predictor Model \rightarrow 1 DV, 1 IDV

multiple predictor Model \rightarrow 1 DV, more than

$$P(\text{Loan} = \text{Yes} | \text{Income} = x) = \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$

after fitting, $b_0 = -6.3525$, $b_1 = 0.0392$

$$P(\text{Loan} = \text{Yes} | \text{Income} = x) = \frac{1}{1 + e^{6.3525 - 0.0392x}}$$

multiple predictor

for prediction of loan -

we have 12 columns:

→ Check all the pvalues of all coln.

→ pvalue < 0.05 (significant coln).

pvalue > 0.05 (less significant and may make no sense)

→ out of 12, 10 coln have pvalue < 0.05

→ in 10 coln, pvalue = 0. (most significant coln for predicting loan)

$P < 0.05 \times (\text{Prob of getting loan is less.})$

$P > 0.05 \checkmark (\text{Prob. of getting loan is high}).$

• When we go for logistic neg., multicollinearity should fail.

* Both linear & logistic regression, both are used to find causal relationship b/w variables.

• It can be used for explanatory task (profiling) or predictive tasks (classification).

20

Wednesday

April 2011

Week 16 • 110-255

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		

Practical Implementation

df = pd.read_excel('logisticRegression.xlsx', sheetname=0)
df.head()

y = df['Admit']

x = df[['GRE', 'GPA', 'Rank']]

import statsmodel.api as sm

x1 = sm.add_constant(x)

Logistic = sm.Logit(y, x1)

result = Logistic.fit()

result.summary()

→ check pvalue < 0.05 or not?

if $p < 0.05$ (var. are significantly imp.)

$p > 0.5$ (prob. of admission is high)

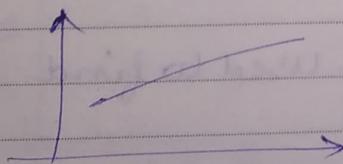
$p < 0.5$ (prob. of admission is less)

Linear Regression

① DV, IDV → continuous

② deciding factor R^2

③



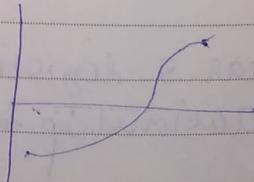
④ $y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$

Logistic Regression

① DV → 0/1, IDV → categorical / cont.

② deciding factor → pvalue

③



④ $p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}}$

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

Assignment → Bank loan modeling, Attrition Project

(Apply logistic
to both
datasets.)

Thursday

April

2011

21

Week 16 • 111-254

⑤ cond'n:

- Heterosk
- Multicollinearity (X) fail
- Autocorrelation

⑥ Predictive algorithm.

⑤ cond'n.

- multicollinearity
- variable selection.

⑥ Classification algorithm.

H.W.

Real Estate Analysis.

* Linear Regrⁿ → Build S model (inside LR)

- Price V/S sqft-living
- Price V/S bedrooms
- Price V/S bathrooms
- Price V/S floors
- Price V/S (all IDV). → multiple regression.

Day-22.

Correlation → To find reln. b/w 2 var. only.

Regression → To find causal effect reln. b/w IDV & DV

1. Linear: simple & multiple.

2. Logistic: single predictor model & multiple predictor model.

not:

info) Decision Tree & Random Forest.

Classification
Technique

Dependent
Var.

IDV

Purpose of Algorithm

1. Decision Tree

categorical

categ.+
continuous

It is used to classify
records in pictorial
format with help of
Gini index.

22

Friday

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 16 • 112-253

2. Random Forest.

categorical

cat.+
continuous

ensemble.

It is a sample DT.
algo, also used to
find imp. var. for DT.

Decision Tree Analysis:

1. Decision Tree Classification
2. Decision Tree Prediction
3. Variable Type & split Type
4. Measure of Node Impurity
5. Decision Tree Induction Algorithm - ID3, CART
6. Practical Issues of Decision Tree
7. Ensemble method - Random Forest.

only
algo → QT → used as classification as well as prediction.

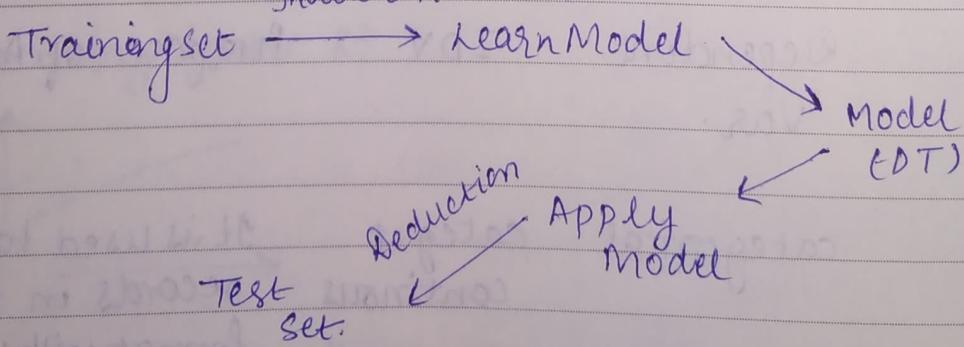
(target var → discrete (categorical))

5. We can generate more than 1 DT for a dataset. But which node is having lesser GINI Index, that node is taken as classification.

Based on GINI index - we classify DT.

#

Tree Induction
algorithm



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

Saturday

April 2011

23

→ How to specify Test cond'n?

Week 16 • 113-252

Depends on
attribute type

Depends on no. of
ways to split

Nominal

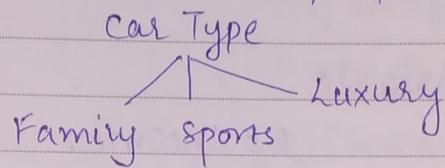
Ordinal

continuous

2 way split (Binary split)
multi-way split

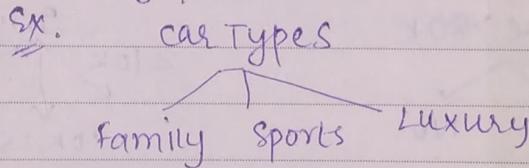
• Nominal Scale → used for Identification purpose.

Ex.
defining in
distinct (different)
way.



Multiway split → more than 2 way is called multiway split

→ use as many partitions as distinct values.



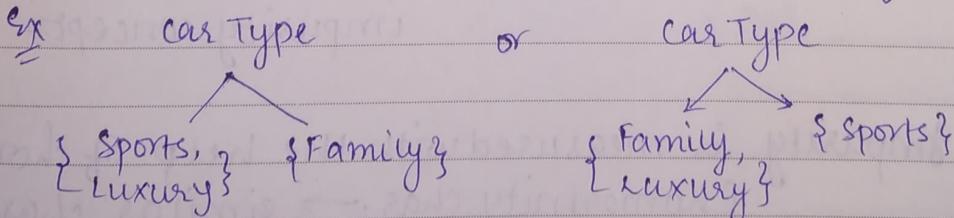
Binary Split → Only 2 categories are these.

(Splits values into 2 subsets,
Need to find optimal partitioning)

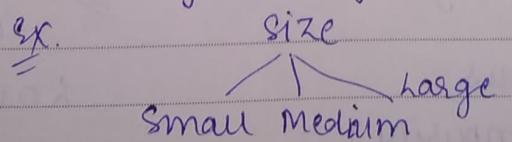
Sunday

114-251

24



• Ordinal Scale → ranking the object.



25

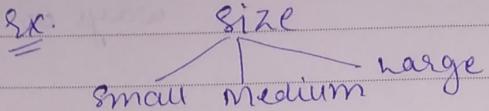
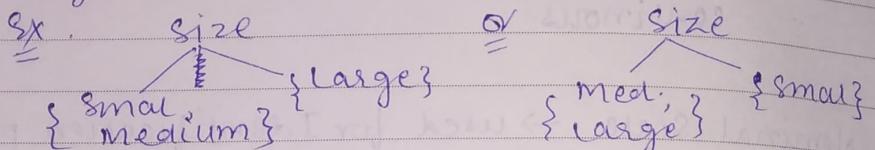
Monday

April

2011

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		

Week 17 • 115-250

Ordinal Multinway \rightarrow 2 or more partitions.Ordinal Binary split \rightarrow 2 categories / subset.

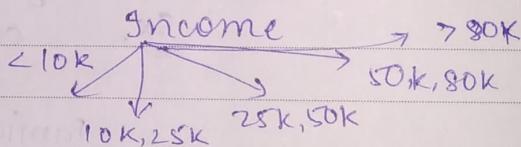
etc..

9 • continuous Variable

variables are
not defined
pre.

Ex Income.

10 It can be any value

Income $> 80K$ 

• Binary split

• Multinway split

which is best split?

we can't determine. \therefore uses node impurity concept.• Node Impurity. is measured with the help of homogeneous (homogeneity char. \rightarrow similar characteristics).Ex. $C_0 \rightarrow 5$
 $C_1 \rightarrow 5$ } heterogeneous.Non-homogeneous
(High degree of impurity)Ex. $C_0 : 9$
 $C_1 : 1$ } homogeneousLow degree of
impurity

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

Tuesday

April 2011

26

Week 17 • 116-249

We can measure node impurity →

- Gini Index
- Entropy.
- misclassification error.

GINI Index:

$$\downarrow \quad f(t) = 1 - \sum_j [p(j/t)]^2$$

ranges b/w
0 to 0.5

prob. of getting specific
event / total no. of obs.

- smaller gini index → more purity (similarity is more)

$$\text{Ex: } \begin{array}{l} C1: 0 \\ C2: 6 \end{array} \quad \begin{array}{l} C1: 1 \\ C2: 5 \end{array} \quad \begin{array}{l} C1: 2 \\ C2: 4 \end{array} \quad \begin{array}{l} C1: 3 \\ C2: 3 \end{array}$$

$$\text{Gini} = 0.00 \quad \text{Gini} = 0.278 \quad \text{Gini} = 0.444 \quad \text{Gini} = 0.500$$

(homogeneity
is high).

$$\downarrow \quad P(C1) = 0/6 = 0$$

$$P(C2) = 6/6 = 1$$

$$\begin{aligned} \text{Gini} &= 1 - P(C1)^2 - P(C2)^2 \\ &= 1 - 0 - 1 \\ &= 0 \end{aligned}$$

Ex Multiway split.

Family sport luxury

$$\begin{array}{cccc} C1 & 1 & 2 & 1 \\ C2 & 4 & 1 & 1 \end{array}$$

$$\text{Gini} = 0.393$$

use count
matrix to
make
decision

Entropy:

$$\downarrow \quad f(t) = - \sum_j p(j/t) \log p(j/t).$$

ranges from
0 to 1.

$$\text{Ex: } \begin{array}{l} C1: 1 \\ C2: 5 \end{array}$$

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$\begin{aligned} \text{Entropy} &= -(1/6) \log_2(1/6) - (5/6) \log_2(5/6) \\ &= 0.65 \end{aligned}$$

27

Wednesday

April

2011

Week 17 • 117-248

entropy $\rightarrow \downarrow$ (low)
 homogeneity is more.
 entropy \rightarrow high (\uparrow)
 homogeneity is ~~more~~ less

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

entropy is minimized when all values of target are same.

entropy is maximized when there is equal chance of all values for target attribute. (ie result is random).

→ purity.

0 → impurity is less (homogeneity is more)

high → impurity is high (purity is less).

Mis

Classification Error:

$$\checkmark \text{ Error } (t) = 1 - \max_i P(i/t).$$

ranges

from 0 to

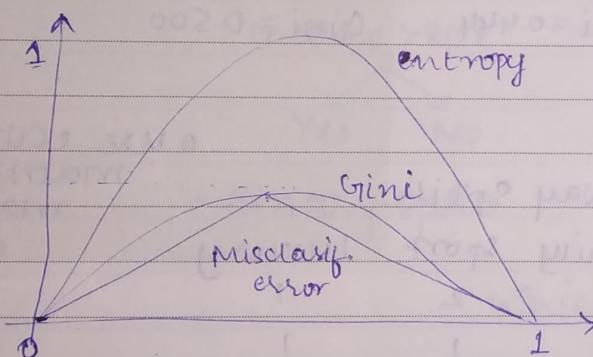
0.5

Ex

$$C_1 : 1 \quad P(C_1) = 1/6 \quad P(C_2) = 5/6$$

$$C_2 : 5 \quad \text{Error} = 1 - \max(1/6, 5/6)$$

$$= 1 - 5/6 = 1/6$$



out of these 3,

Gini index gives higher accuracy.

Stopping criteria for Tree Induction.

If node is having same characteristics, we can terminate node.

PPT.

- 1.) Stop expanding a node when all the records belongs to same class.
- 2.) Stop expanding a node when all records have similar attribute values.
- 3.) Early termination.

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	

- Binary split → only 2
 - Categorical split → More than 2 ie (Low, Med., High)
 - continuous split
- Income, Age

April

2011

28

Week 17 • 118-247

ID3 (Iterative Dichotomous)

two ie it will do binary split

CART (Classification & Regression Trees).

- Developed by:

Breiman, Friedman, Olshen,

Stone in early 80's

Binary → ID3.

cart { categorical
continuous

continuous target: → use sum of squared errors.
(Regression Trees)

categorical target → choice of entropy, Gini measure,
(classification trees) 'error rate', information gain,
'twining', splitting rule.

Difficulties in DT.

- Underfitting & overfitting.
- missing values
- costs of classification

Day - 23.

- Overfitting → If you are fitting more var. (have more IDV) it is difficult to classify records.

We go for random forest algo.

Q: How to avoid overfitting?



29

Friday

April

2011

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

Week 17 • 119-246

- 1.) Stop growing when data split not statistically significant
- 2.) Grow full tree then post-prune
- 3.) use ensemble of DT.

- Pre-pruning → Before generation of DT, you can resolve overfitting
- # Random Forest → ensemble sample DT.
→ identify imp. indep. var. for DT algo.

- Pruning (minimizing tree).

cut the
tree.

control
growing of
tree - either it is
depth wise or
width wise)

Two types of Pruning.

↓
Pre-Pruning
(forward pruning).

↓
Post-Pruning
(backward pruning)

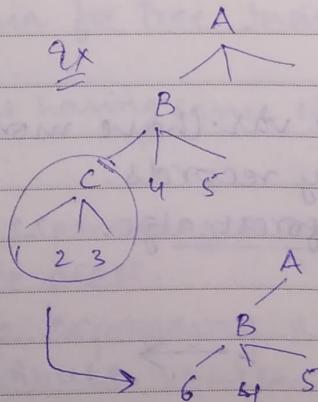
Before generation of DT,
if you are minimizing
the tree.

once tree is generated,
how to classify the
records.

- We go for Random Forest

it increases the
accuracy of tree).

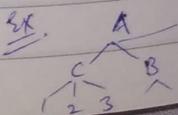
Subtree
replacement.



If tree is
growing
more,
identify which
node is growing
more & that
node is
summarized.

- Subtree Replacement
- Subtree Raising.

↓
If tree is raising
more, the
entire node
is attached to
main node.



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

Saturday
April / May 2011

30

Week 17 • 120-245

① Handling Missing Value.

- Missing value affect DT construction in 3 different ways.
 - Affect how impurity measures are computed
 - Affect how to distribute instance with missing value to child nodes.
 - Affect how a test instance with missing value is classified.

↓
missing
fb data
is less,

we can
substitute it

② If missing data
is more, we can
del entire col.

Ex 100 obsr.

① 15-20 missing → fill it

② 50-80 missing in a coln
(del. that coln).

how to fill missing value.

- var. is continuous → fill by mean (avg.)
- var. is categorical → fill by mode.

- Fast at classifying unknown record
- Easy to interpret

Strengths (Advantage of DT).

- use for generating rules.
- Perform classification without much computation.
- IDV can be categorical or continuous.
- Provide clear indication of which field are most imp. for prediction or classification.

Weaknesses (Disadvantage of DT)

- Prediction / Depn var. (DV) should be categorical.
(NOT suitable for continuous var.)
- Perform poorly with many class (categories), & small data.
-

Sunday

121-244

1

2

continuous var → changed to ranges.

DV → Sales / Income

Monday

May 2011

<10K, 10-15K, 15-20K...

converted to categories
then apply Regressor DT.

Week 18 • 122-243

Wk	May						
	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

Random Forest.

- which DV is having more depth in tree.
more easy to classify.
- find most imp. variable in DT.

Practical implementation.

Dataset used → Titanic.

train = pd.read_csv('train.csv')

9

10

if Age coln has missing value,

11

train['Age'].mean # 32.69985

12

1

new_age = np.where(train['Age'].isnull(), 32, train['Age'])

2

train['Age'] = new_age.

3

check
null
value

4

↓
Replace
by
32

5

text → category.

6

• Label Binarizer (only 2 categories).

7

• Label Encoder (2 categories, or more.)

8

from sklearn import tree

from sklearn import preprocessing.

lab = preprocessing.LabelEncoder()

sex = lab.fit_transform(train['sex'])

model = tree.DecisionTreeClassifier()

model.fit(x=pd.DataFrame(sex),

y=train['survived'])

DV is categorical,

if it is continuous

go for DTRegressor)

June	Mo	Tu	We	Th	Fr	Sa	Su
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			

DT always support dot file.

(user defined).
name of
file.

Tuesday

May

2011

3

Week 18 • 123-242

with open('DT1.dot','w') as f:

f = tree.export_graphviz(model, feature_names=['Sex'],
out_file=f);

copy the
code from
the file (opened in Notepad)

interface

I where to
store this
model.

webgraphviz.com (interface)

use for generating DT.

paste that code in this site & generate graph.

① model2 = pd.DataFrame([sex, train['Pclass']]).T

model.fit(X=model2, y=train['survived'])

with open('DT2.dot','w') as f:

f = tree.export_graphviz(model, feature_names=['Sex', 'Pclass'],
out_file=f);

Day-24.

Model3 = pd.DataFrame([sex, train['Pclass'], train['Age'],
train['Fare']]).T

model = tree.DecisionTreeClassifier(max_depth=8)

model.fit(X=Model3, y=train['survived'])

with open('DT3.dot','w') as f:

f = tree.export_graphviz(model, feature_names=['Sex', 'Pclass',
'Age', 'Fare'], out_file=f);

If you
add more
T.PV, size
gets
complicated.

4

Wednesday

May

2011

Week 18 • 124-241

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	16
20	16	17	18	19	20	21	22	23
21	23	24	25	26	27	28	29	

model.score(x= model3, y= train['survived'])

89% accurate

Q. Accuracy is 89%. but not more than 90%. Why?

To increase accuracy, we need to find imp. variable by Random Forest then again predict. The accuracy will increase.

- Till now we do classification, Now We will predict by DT.

test = pd.read_csv('test.csv')

new_age = np.where(test['Age'].isnull(), 28, test['Age'])
test['Age'] = new_age

sex = lab.fit_transform(test['sex'])

feature = pd.DataFrame([sex, test['Pclass'], test['Age'],
test['Fare']]).T

pred = model.predict(x= feature)

output = pd.DataFrame({'PassengerId': test['PassengerId'],
'survived': pred})

output.to_csv('output.csv', index=False);

DataFrame
converted to csv format.

- Random Forest.

from sklearn.ensemble import RandomForestClassifier

Titanic
train
dataset
(imported
already).

train.columns

labelencoder.

train['Sex'] = lab.fit_transform(train['Sex'])

June	Mo	Tu	We	Th	Fr	Sa	Su
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			

Thursday

May 2011

5

Week 18 • 125-240

train['Embarked'] = lab.fit_transform(train['Embarked'])

rf-model = RandomForestClassifier(n_estimators=1000,

max_features=2, → Binary split
oob_score=True)

features = ['Sex', 'Pclass', 'SibSp',

out-of-bag score,

'Embarked', 'Age', 'Fare']

Based on DT, each & every
node, we can find
accuracy.

rf-model.fit(X=train[features], y=train['Survived'])

print('OOB Accuracy:')

print(rf-model.oob_score_);

#→ 80%.

Q: Find imp. features:

for feature, imp in zip(features, rf-model.feature_importances_):

print(feature, imp);

Sex - 0.268 ✓

Pclass - 0.088

SibSp - 0.052

Embarked - 0.034

Age - 0.270 ✓

Fare - 0.285 ✓

accuracy score.

out-of-bag :

These 3 are imp.
features

HW:

With these 3
variables,
find
accuracy.
(Prediction)

writing the
rules
also.

Attrition → DT, RF

Bank loan → DT, RF

6

Friday

May

2011

Week 18 • 126-239

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	

Naive Bayes

classification technique

classify records with help of probability.

Day-25.

Classification Algorithm.

Model a classification rule directly.

Model the prob. of class membership given YP data

make a probabilistic model of data within each class.

Ex. K-NN, SVM,
DT, perceptronEx. Logistic Regressn,
perceptron with cross-entropy costEx. Naive Bayes,
Model Based classifier.

1, 2 are ex. of discriminative or conditional classification

2, 3 are both ex. of probabilistic classification

3 is an ex. of generative classification

• Naive Bayes.

→ Both DV, IDV are categorical

• Numerical value must be binned & converted to categorical.

→ can be used with very large datasets.

Probability Basics.

◦ Prior Prob. : $P(X)$ ◦ conditional Prob. : $P(X_1 | X_2), P(X_2 | X_1)$ ◦ Joint prob. : $X = (X_1, X_2)$

$$P(X) = P(X_1, X_2)$$

June 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			

Saturday

May

2011

7

Week 18 ■ 127-238

- Relationship: $P(X_1, X_2) = P(X_2|X_1).P(X_1) = P(X_1|X_2).P(X_2)$
- Independence: $P(X_2|X_1) = P(X_2)$, $P(X_1|X_2) = P(X_1)$
 $P(X_1, X_2) = P(X_1).P(X_2)$

Bayesian Rule.

$$P(C|X) = \frac{P(X|C).P(C)}{P(X)}$$

Posterior = Likelihood × Prior
Evidence.

Ex. Dataset:

Day	outlook	Temp.	Humidity	wind	Play Badminton
Day 1	sunny	Hot	High	Weak	No
Day 2	sunny	Hot	High	Strong	No
Day 3	overcast	Hot	High	Weak	Yes
:					
Day 14	Rain	Mild	High	Strong	No

Q. For day: (sunny, cool, high, strong), *
what's play prediction?

Sunday

128-237

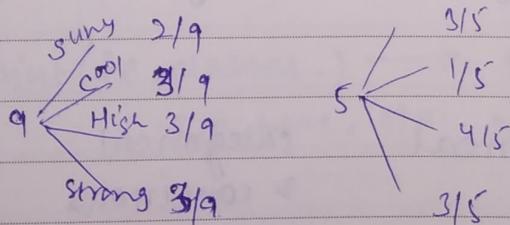
8

Total = 14

Yes = 9

No = 5

There is no such record in
dataset (exactly same
weather)



$$P(\text{Yes}/x') = [P(\text{sunny}/\text{Yes}) \cdot P(\text{cool}/\text{Yes}) \cdot P(\text{High}/\text{Yes}) \cdot P(\text{strong}/\text{Yes})] / P(\text{Yes})$$

$$= 0.0053$$

$$P(\text{No}/x') = 0.0206$$

Predicted \rightarrow No ($P(\text{No}) > P(\text{Yes})$)

9

Monday

May

2011

Week 19 • 129-236

- It doesn't involve statistical models

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

① Advantages of Naïve Bayes.

- Handles purely categorical data.
- Works well with very large datasets.
- Simple & computationally efficient. (easy & fast).

Disadvantages.

- Requires large no. of records
- Problematic when a predictor category is not present in training data.
(Assign 0 prob. of response, ignoring info. in other var.)

Applications.

- Spam classification.
- Medical diagnosis
- Weather Prediction
- Text classification
- Character recognition.

Classification Technique

DV

IDV

Purpose.

Naïve Bayes

categorical

categorical

Used to classify record with help of probability.

K-NN
(K-Nearest Neighbour)

categorical

categorical & continuous

Used to classify record with help of Euclidean Distance

SVM
(Support Vector Mc)

categorical

categorical & continuous

Used to classify records with help of hyperplane.

June 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				

Tuesday

May 2011

10

Week 19 ■ 130-235

- SVM is applicable for over-dimensional data.

Practical implementation.

Dataset → svmtrain (titanic dataset)

```
df = pd.read_csv('svmtrain.csv')
```

```
from sklearn import preprocessing
from sklearn.cross_validation import train_test_split
# If it is not working use this -
from sklearn.model_selection import train_test_split
```

```
from sklearn.naive_bayes import GaussianNB
```

Note {

- accuracy score → how much % the model is accurate.
- confusion matrix → how many records are classified more accurately

```
from sklearn.metrics import accuracy_score, confusion_matrix
```

→ convert text into numerical.,

```
le = preprocessing.LabelEncoder()
```

```
df['Sex'] = le.fit_transform(df['Sex'])
```

print(le.classes_) → gives you category name.

e.g here, o/p: male, females.

```
y = df['Survived']
```

```
x = df.drop(['Survived', 'PassengerId'], axis=1)
```

axis=0 (row.)

11

Wednesday

May

2011

Week 19 • 131-234

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

y.count() # 889

x-train, x-test, y-train, y-test = train-test-split(x, y, test_size=0.3, random_state=0)

from sklearn.naive-bayes import *

clf = BernoulliNB()

y-pred = clf.fit(x-train, y-train).predict(x-test).

accuracy-score(y-test, y-pred, normalize=True) # 76.48%

confusion-matrix(y-test, y-pred)

→
132 25
37 73

confusion matrix →

		Actual.		Training Data in coln	
		Class of Interest	All other classes	TP	FP
predicted	Class of Interest	True +ve	(True +ve)	(False +ve)	
	All other classes.	FN	(False -ve)	TN	(True -ve)

Ex

=

survived Training

	0	1	Total = 267
Survived	0	132	132
Test	1	37	37
	correct = 205	(diag.)	
	incorrect = 62	(non-diag.)	

$$\text{correct \%} = 205/267 = 0.76779$$

June 2011						
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		

Assignment - (Print confusion matrix also) 1

DV DV Accuracy. Thursday
 Pclass other ? 59%
 Gender all ? 78%
 sibsp " ? 67%
 Parch " ? 74%. use loop.
 Embarked " ? 69%.

Week 19 • 132-233

12

(watch it once)

Day-26

- All v Project Description (Day 24/21)

Day-27.

K-Nearest Neighbours.

- classification Technique
- classify records with help of euclidean distance.

find closeness of 2 points

Features.

- All instances correspond to pts. in n-dimensional Euclidean space
- classification is delayed till new instance arrives.
- Target fn may be discrete or real valued.
- classification done by comparing feature vectors of diff. points.

It is instance Based Learning

→ Based on characteristics of record we classify the o/p.

Euclidean distance →

Find closeness

$$\sqrt{(x_1 - u_1)^2 + (x_2 - u_2)^2 + \dots + (x_p - u_p)^2}$$

of 2 records.

$$D(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Ex:

$$\text{euclidean dist.} = \sqrt{\underbrace{(35-41)^2}_{\text{age}} + \underbrace{(95000-215000)^2}_{\text{salary}} - \underbrace{(3-2)^2}_{\text{CD}}}$$

1st response
2nd response

record 1

Ex Jay:
Age = 35
Salary = 95k
CD = 3

record 2.

Riya:
Age = 41
Salary = 215k
CD = 2