

SML

U-2

① Simple Linear Regression

→ one dependent var y predicted from an indepnt var x .

→ R^2 : proportion of variation in dependent var y predictable from x .

→ reported as: $Y = a + bx$

→ slope gives the change in Y for a unit change in x

→ used in: economics, finance,

Multiple Linear Regression

→ one dependent var y predicted from a set of independent var (x_1, x_2, \dots, x_k)

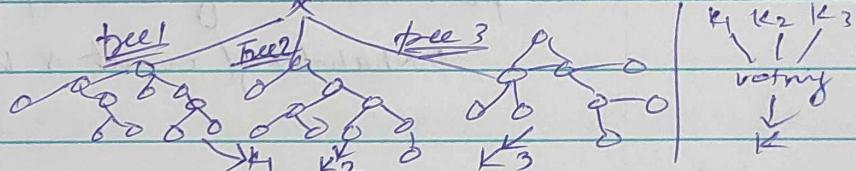
→ R^2 : proportion of variation in dependent var y predictable from set of indepnt var. (x 's)

→ reported as: $Y = a + b_1x_1 + b_2x_2 + \dots + b_kx_k$

→ slope gives the change in Y for a unit change in each var.

→ used in: marketing, health care

② Random Forest



→ Powerful technique used in data science field for solving various problems.

→ used for both classification & regression

→ Ensemble method: It is an ensemble method that combines multiple decision trees to create more accurate model

→ Decision Tree: Each tree in random forest is decision tree, which makes sequence of decisions to arrive at a result.

→ It uses feature randomness to construct tree. This means only random subset of features is selected

→ Bagging involves creating multiple random samples of dataset on each sample. final prediction is made by avg. of all predict. trees.

(B) Logistic regression & its advantages

- It's a statistical method used to model the relationship b/w binary outcome (0, 1, true, false) and predictor var.
- It uses logistic function to transform the output of linear regression model into probability values b/w 0 and 1
- Used in ml and data science for classification task like predicting whether email is spam or not -
- It is trained using maximum likelihood estimation
- Max Likelihood estimation is the process of estimating the parameters of model -
probability(event + non-event) = 1

Adv

- ① easily implemented & interpreted
- ② can easily extend to multiple classes.
- ③ very fast at classification
- ④ good accuracy
- ⑤ less over-fitting

Information Value (IV)

Akaike info. criteria (AIC)

Receiver operating char. (ROC)

Rank ordering

C-statistics

K-S Statistics

- ④ Random forest using German credit data
- ↳ popular ML algorithm for classification & regression
 - ↳ uses decision tree to improve accuracy.

For German credit data, the dataset contains:

- | | |
|----------------------------|----------------------|
| (1) status of existing acc | (6) property |
| (2) credit history | (7) installment-plan |
| (3) purpose | (8) hours |
| (9) savings acc. | (9) job |
| (5) employment duration | (10) Telephone |

(1) Let's import some packages for the German credit data
import pandas as pd

from sklearn.ensemble import RandomForestClassifier
credit_data = pd.read_csv("credit_data.csv")

(2) Create the dummy var.

```
>>> dummy_stat = pd.get_dummies(credit_data['status of existing data'],  
prefix = 'stat-of-exist-acc')
```

```
>>> dummy_ch = pd.get_dummies(credit_data["Credit history"],  
prefix = "cred-his")
```

— Similarly for all the datasets —

(3) Remove extra dummy var. of all categories

credit_data_new = pd.concat ([dummy_stat, dummy_ch, dummy_purpose, ...])

(4) Split the data into 70-30

x_train, x_test, y_train, y_test = train_test_split(credit_data_new, drop(['class'], axis=1), train_size=0.7, random_state=42)

(5) Assume hyperparameter values:

no. of trees is 1000

Max. decisions tree can grow is 100

Min obs. reg is 3

Min no. of obs in node is 2

and apply these into the model.

(6) Result:

		Random forest Train CM		Random forest Test CM	
		Bredicted 0	Bredicted 1	Bredicted 0	Bredicted 1
Actual	0	501	0	0	240
	1	18	185	1	43
Accuracy		0.979		0.855	

⑤ Compare Ridge and Lasso regression models.

Ridge Regression Model

- uses penalty regression in L2 norm
- shrinks the coefficient towards zero, but rarely exact zero
- It increases the bias of the model
- uses closed form algo to solve the problem
- If 2 or more predictors are highly correlated, Ridge Reg. will shrink their coefficients
- Handling dataset is difficult

Lasso Regression Model

- uses penalty regression in L1 norm
- shrinks the coefficient exactly to zero
- It decreases the bias of the model.
- uses iterative optimization algo to solve the problem.
- while Lasso can choose one of them & set remaining to zero
- Handling dataset is easy

⑥ Grid search on Random Forest

- While finding soln with Random forest, we get less accuracy. To increase the accuracy we will use Grid Search
- It is performed by changing hyperparameters for the model.
- These Hyperparameter models are:

No. of Trees (1000, 2000, 3000)

Max depth (100, 100, 300)

Min sample split (2, 3)

Min sample in leaf nodes (1, 2)

- It is used to get the best performance of the model.
- How to perform
 - (1) Define grid hyperparameters
 - (2) Train Random forest model.
 - (3) Evaluate performance of each model.
 - (a) Select the combi. of best hyperparameters on performance
 - (b) Test the final model.

⑦

$x_i \quad y_i$

1 1.2

2 1.8

3 2.6

4 3.2

5 3.8

To apply linear regression,
plot the data to find the
relationship b/w x and y.

$$\text{So, } y = mx + b$$

m is slope of line

b is y intercept.

$$m = (n \sum x \cdot y - \sum x \cdot \sum y) / (n \sum x^2 - (\sum x)^2)$$

$$b = (\sum y - m \sum x) / n$$

n is no. of data points.

$\sum xy$ sum of product of x and y

$\sum x$ sum of x values

$\sum y$ sum of y values

$\sum x^2$ sum of squared x values.

$$n = 5$$

$$\sum x = 15$$

$$\sum y = 12.6$$

$$\sum xy = 44.4$$

$$\sum x^2 = 55$$

$$m = (\sum xy - \bar{x} \sum y) / (\sum x^2 - (\bar{x})^2)$$

$$= (5 \times 44.4 - 15 \times 12.6) / (5 \times 55 - 225)$$

$$= 33 / 50 = 0.66$$

$$b = (\sum y - m \sum x) / n$$

$$= (12.6 - 0.66 \times 15) / 5$$

$$= 2.7 / 5 = 0.54$$

so, eqn becomes $y = mx + b$

$$y = 0.66x + 0.54$$

for 7th week sales, $x = 7$

$$y = 0.66 \times 7 + 0.54$$

$$= 5.16$$

for 9th week sales, $x = 9$

$$y = 0.66 \times 9 + 0.54$$

$$= 6.48$$

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} =$$

(8)	<u>outlook</u>	<u>Temp</u>	<u>Humidity</u>	<u>windy</u>	<u>Rainy</u>
1	sunny	Hot	High	weak	No
2	sunny	Hot	High	strong	No
3	overcast	Hot	High	weak	yes
4	ring	Mild	High	weak	yes
5	"	Cool	Normal	weak	yes
6	"	Cool	Normal	Strong	No
7	overcast	Cool	Normal	Strong	yes
8	Sunny	Mild	High	weak	No
9	Sunny	Cool	Normal	weak	yes
10	Ring	Mild	Normal	weak	yes
11	Sunny	Mild	Normal	Strong	yes
12	overcast	Mild	High	Strong	yes
13	overcast	Hot	Normal	weak	yes
14	Ring	Mild	High	Strong	No

Overall prob. of yes:

$$P(Y) = 9/14$$

Overall prob. of no

$$P(N) = 5/14$$

<u>Outlook:</u>	Y	N	P(Y)	P(N)
Sunny	2	3	2/9	3/5
overcast	4	0	4/9	0
Rainy	3	2	3/9	2/5

<u>Temperature:</u>	Y	N	P(Y)	P(N)
Hot	2	2	2/9	2/5
Mild	4	2	4/9	4/5
Cool	3	1	3/9	1/5

<u>Humidity:</u>	Y	N	P(Y)	P(N)
High	3	4	3/9	4/5
Nominal	6	0	6/9	1/5

<u>Windy:</u>	Y	N	P(Y)	P(N)
weak	6	2	6/9	2/5
strong	3	3	3/9	3/5

$$P(X/\text{play} = \text{yes}) \cdot P(\text{play} = \text{yes})$$

$$\Rightarrow \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}$$

$$\Rightarrow \frac{486}{91854} = 0.00529$$

$$P(X/\text{play} = \text{no}) \cdot P(\text{play} = \text{no})$$

$$\Rightarrow \frac{3}{5} \times \frac{1}{5} \times \frac{3}{5} \times \frac{4}{5} \times \frac{9}{14}$$

$$\Rightarrow \frac{180}{8750} = 0.0205$$

$$\begin{aligned}
 P(x) &= P(\text{yes}) + P(\text{no}) \\
 &= P(x | \text{play} = \text{yes}) + P(x | \text{play} = \text{no}) \\
 &= 0.00529 + 0.0205 \\
 &= 0.0253
 \end{aligned}$$

Now, $P(\text{play} = \text{yes} | x)$

$$\begin{aligned}
 P(\text{play} = \text{yes} | x) &= \frac{P(x | \text{play} = \text{yes}) \cdot P(\text{play} = \text{yes})}{P(x)} \\
 &= \frac{0.0053}{0.0253} = 0.209
 \end{aligned}$$

$$\begin{aligned}
 P(\text{play} = \text{No} | x) &= \frac{P(x | \text{play} = \text{no}) \cdot P(\text{play} = \text{no})}{P(x)} \\
 &= \frac{0.0205}{0.0253} = 0.810
 \end{aligned}$$

$P(\text{No})$ is high

\therefore The player cannot play.