

01.22.24

"Intro to Statistical Learning"

Statistical learning :
Theoretical foundation of ML :

used for 2 things :

- Find a "predict" → supervised learning
- Grouping → unsupervised "

Supervised vs Unsupervised :

Supervised : Build a Model to predict some output from past data
eg: From wt to height.
Past data $\rightarrow y \rightarrow$ Labels.

Unsupervised : Inputs but no output (y) labels

Notations : y = output, response, dependent variable

x_i^o = input, independent variable, predictors, features

n = number of observations (samples)

p = number of features / predictors

\in = "is an element of"

$a \in B$
a is an element of set B

(RT. 0.)

example:

n	TV	Radio	Newspapers	Sales
1	1	1	1	1

$y = \text{Sales}$

$x_i = \text{TV, Radio, Newspapers}$

$P = 3$ (feature گز نیویت)

Data in matrix form گز نیویت :

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ x_{31} & x_{32} & \dots & x_{3P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nP} \end{bmatrix} \quad \begin{array}{l} n = \text{number of samples} \\ \downarrow \\ \text{or observations} \end{array}$$

Overview of notations with an example:

$$x = \begin{bmatrix} \text{TV} \\ x_1 \end{bmatrix}, \begin{bmatrix} \text{Radio} \\ x_2 \end{bmatrix}, \begin{bmatrix} \text{Newspaper} \\ x_3 \end{bmatrix}$$

y Predicting \rightarrow Given x
Output Input

Predicting / estimating the function:

We know x_i but no y

$$\hat{y} = \hat{f}(x)$$

y vs \hat{y}

True Label

estimated / predicted
label

y vs $\hat{y} \rightarrow$ functⁿ need to be better
so that error can be reduced



- i) Reducible error
- ii) Irreducible error

$$[f(x) - \hat{f}(x)]^2 \rightarrow$$
 squaring for

→ getting + values

→ Amplifying the error

so that model can perform better.

$$E[(y - \hat{y})^2] = [f(x) - \hat{f}(x)]^2 + \text{var}(\epsilon)$$

reducible error

irreducible error

estimated functⁿ → so that error can be reduced (minimized during training)

(P.T.O.)

After training:

↓
Testing inference

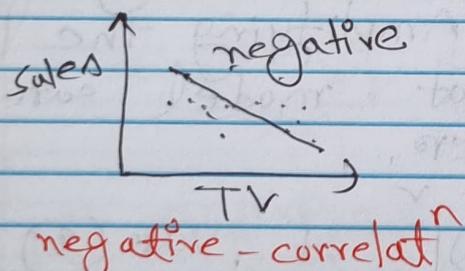
① Which predictors are the most important?

The model should predict

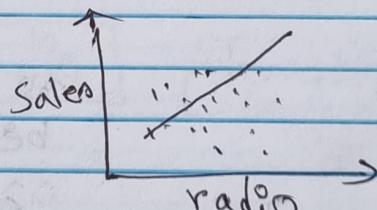
correctly invest 70%,
TV / news

② Is there a relationship between X and Y?

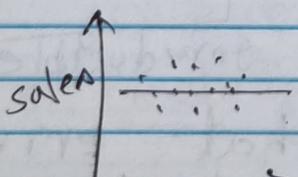
③ What kind of relationship exists between X and Y?



negative-correlat



positive correlation between sales and radio



newspaper

→ no relationship / doesn't affect the sale anyway

Take some data
train a function
Test a "

How do we estimate f ?
we wanna find f such that
 $y \approx f(x)$

2 ways to estimate the function:
→ parametric method
→ Non- " "

parametric Method:

2 step:

⑨ Assume the function be of a certain form

e.g.: Linear fn

$$f(x) = \beta_0 + \beta_1 x_1$$

x	y

⑩ Apply a procedure to fit / train the model.

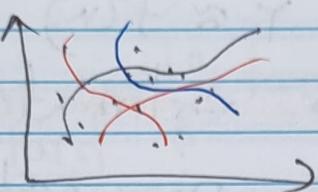
Non-parametric method:

We don't make any explicit assumption

We try to find the funct' itself from the data.

which uses what
without a prior
is fast

- No assumptⁿ of the fun form
- Estimate the fun from the dataset.



It requires lots of data

e.g.: Deep-learning methods

01.24.24
Lec: 03

$$y = f(x)$$

parametric

$$\hat{y} = \beta_1 x_1 + \beta_2 x_2$$

$$\beta_1 x_1 + \beta_2 x_2 \rightarrow \begin{array}{l} \uparrow \\ \text{net} \\ \downarrow \end{array}$$

Non-parametric → Data

which functⁿ can fit the data

uncountable model can fit # the data

XAI →

Explainable AI → computer vision model
for classification

e.g.: Pic of cat, which parts
are model taking

How do we measure the Accuracy of
a Model:

qualitative: Labels, classes

e.g.: Default / Not-default

quantitative: Deals with numbers

Predict → GPA (numbers, sales, age.)

Methods for qualitative → classification
Method

Quantitative → Regression methods

Regression methods :

$$y = f(x) \Rightarrow \text{True function}$$

after training, $\hat{y} = \hat{f}(x) \Rightarrow \text{estimated } y$

	x_1	x_2	y
Tr	radio		sales
-	-	-	-
-	-	-	-
-	-	-	-

\hat{y} will give the estimate of sales

$$y = 40 \quad \hat{y} = 30 \quad y - \hat{y} = 40 - 30 = 10$$

overall the test data

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

Mean Squared
Error

accumulate every error

↓
error will grow up

↓
so we take avg [$\frac{+++}{\text{total}}$]

$$\left[\frac{1}{n} \sum_{i=1}^n [y - \hat{f}(w)]^2 \right] \rightarrow \hat{y} = \hat{f}(w)$$

↓
MSE → for Regression Model (accurate)

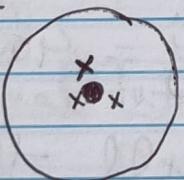
Bias-variance Trade-off:

Simple prob → complicated tool
compli → simple

Bias → How good/accurate a model is
Accuracy, good

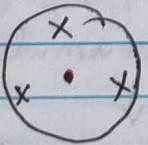
Variance → How different a model behaves with different data

Model 1
model hits →



→ doesn't hit but close enough
bias → ↑↑
variance → Low

Model 2



bias → ↓
variance → ↑

Bias + variance] balance needed

↑ Bias → good but ~~easy~~ change ~~but~~ different thing show ~~to~~ model.

$$E[(y_o - \hat{f}(x_o))] = \text{Expected test MSE}$$

↓
Accuracy of model

↓ we expect that

$$\text{Var}[\hat{f}(x_o)] + [\text{Bias}(\hat{f}(x_o))] + \text{Var}(\epsilon)$$

reducible error

irreducible error

Classification:

$$\hat{y}_o = \hat{y}_o$$

If prediction is = to true value,

$$\hat{y}_o \neq \hat{y}_o \text{ (error)}$$

Human → eye
but in comp → indicator function

$$\frac{1}{n} \sum_{i=1}^n I(\hat{y}_o \neq \hat{y}_o)$$

either true/false

x_1	x_2	Ban
		yes
		no
		yes
		no

or $\text{Ave} \left[I(y_o \neq \hat{y}_o) \right]$

finding Indicator function
 ↓
 variety of methods

(35 page)
 Basic:
 TSLP
 LIP

Base - Bayes Classifiers: probability

$y_o \neq \hat{y}_o \rightarrow$ we want to minimize
 this

this classifier minimize this based
 on Probability:

2 - blue markers

5 - red

blue marker probability $\rightarrow \frac{2}{7}$

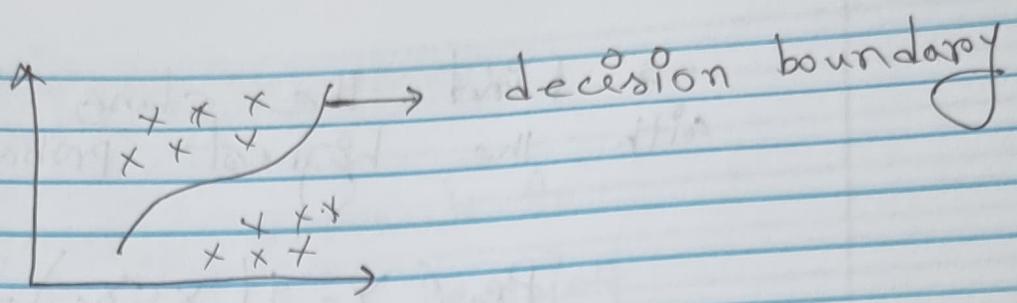
conditional probability
 $\Pr(Y=j | x=x_o)$ that largest

for 2 classes \rightarrow

$$\Pr(Y=j | x=x_o) > \frac{1}{2} \quad \text{or } \frac{1}{3}$$

if class $1/2$ or $2/5$ or $3/7$ then

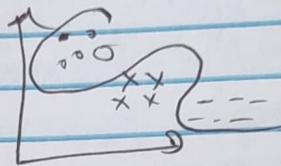
yes
 no
 maybe



Using decision boundary \rightarrow try to minimize Ave I

by separating the classes

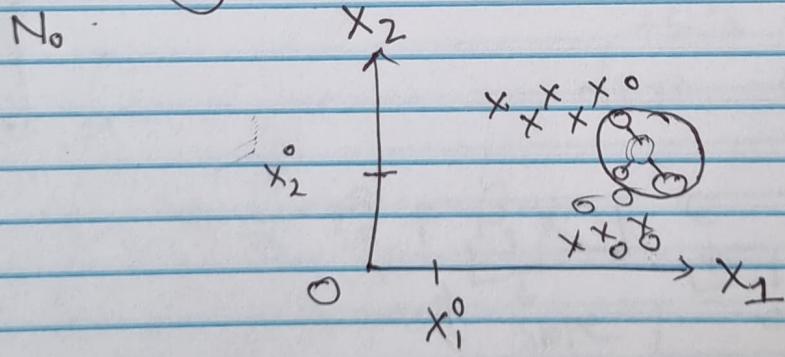
3 - classes :



K - Nearest Neighbours (KNN) :

Given k, x_0 , test data

\rightarrow we wanna identify k -points in the training data closest to x_0 , denoted as



x
o

→ find the class in N_0 (neighbors) with the highest probability.

$$Pr(Y=j \mid x \in N_0) = \frac{1}{k} \sum_{i=1}^k I(y_i=j)$$

For $x_1 = 0$

for $y_i = 1$
class

($\frac{1}{3} \times 3$)

$$\boxed{\frac{1}{3} \times 2} = \boxed{0.666}$$

$$\frac{1}{3} \times 1 = \frac{1}{3}$$

$$\boxed{\frac{1}{3} + \frac{2}{3}} \\ = \frac{3}{3} = 1$$

$$\boxed{k=3 \dots 7}$$