

# Machine Learning.

Date / /

Page No.

# gradient descent (Batch)

↳ optimisation algorithm

↳ find values of parameters (coeff) of function that minimises cost

create a plot with error (cost) loss func vs parameter.



↳ aims to reduce

error and find

optimal value of parameter.

In straight line there are 2 parameters and need to optimise both.

↳ we take small step sizes and reduce the error.

delta = derivative (Cost function)

Step size = (Learning rate) \* (delta)

Learning rate controls how much the coefficients can change on each update.

When to Stop?

↳ Step size close to zero ( $< 0.001$ )

↳ total steps (epoch), usually  $> 1000$

# In gradient descent, for linear regression,  
~~the~~ the cost / loss function is  
the residual error

$$= (Y - Y_{pred})^2$$

i.e. ~~the sum~~ <sup>mean</sup> of all squared errors.

the derivative is taken of this  
function wrt the parameters.

# Stochastic Gradient descent

↳ For big dataset the regular  
gradient descent is slow as there  
is a lot of calculation in each  
iteration, to reduce this

↳ Select a point randomly  
and apply algo instead whole  
data set

↳ usually used when redundant  
data set ~~exists~~

# Schedule:

↳ The way learning rate changes  
from relatively high to relatively  
low.

# mini batch

↳ selects some points (small subset)  
randomly in each iteration.

Note: If there is a new point we can just apply the algo from the point we left (i.e. using the previously obtained  $m, c$ ) instead of ~~resetting~~ resetting the values of parameters.

### Pros of Batch Gradient

- ↳ More stable convergence and error gradient than stochastic
- ↳ direct path is taken towards min
- ↳

### (Cons)

- ↳ Can converge at local minima
- ↳ slow learning rate since an update is performed only after going thru all observation.

### Pros of mini

- ↳ more stable convergence than stochastic
- ↳ computation efficient
- ↳ fast learning as more freq updates

### (Cons)

- ↳ need to figure mini batch size hyperparameters.

## Pros and of stochastic

- ↳ easy to fit in memory
- ↳ likely to reach near min faster than Batch in case of large dataset
- ↳ freq updates
  - ↳ " creates plenty of oscillation which can be helpful for getting out of local mins.

## Cons

- ↳ go in wrong directions
- ↳ frequent updates are computationally expensive due to using all resources for processing one training sample at a time.

## Formulas / steps

$$f_n = \frac{1}{n} \sum_{i=1}^n (Y - Y_{pred})^2$$

$$\frac{de}{dm} = \frac{d}{dm} \left( \frac{1}{n} \sum_{i=1}^n (Y - (mx + c))^2 \right)$$

$$= d - \frac{2}{n} \sum_{i=1}^n (Y - (mx + c))x$$

$$\frac{de}{dc} = -\frac{2}{n} \sum_{i=1}^n (Y - Y_{pred})$$

Note: Capital  $Y, Y_{pred}, x$  means they are series of data.

$$m = m - L \frac{de}{dm}$$

↳  $\frac{de}{dm}$  may be written as  
 $dm \rightarrow \delta(m)$

## LOGISTIC REGRESSION

↳ used to predict binary outputs

↳ True or False

↳ for classification

↳ instead of straight line a Sigmoid function is drawn

$$\text{sig}(t) = \frac{1}{1+e^{-x}} [0, 1]$$

Here:

Data is fit into linear regression model, which then be acted upon a logistic function predicting the target categorical variable.

Type

↳ Binary  $\rightarrow$  2 outcomes (T or F)

↳ Multinomial  $\rightarrow$  3 or more without ordering

↳ Ordinal  $\rightarrow$  3 or more with ordering

## # Decision boundary :

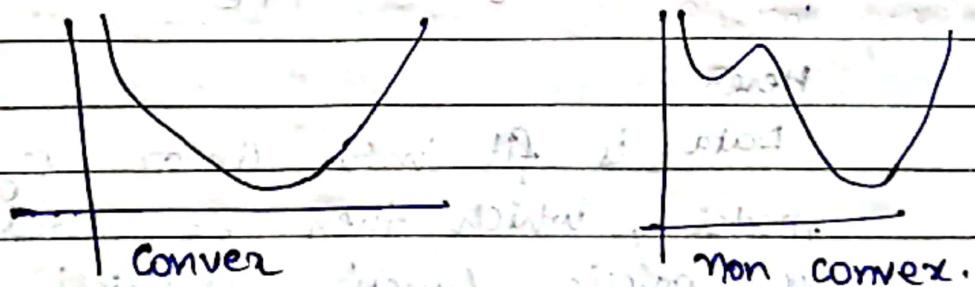
- ↳ To predict which class a data belongs a threshold can be set
- ↳ can be linear or non linear

mean

In logistic, Squared mean ~~are~~ error

Cannot be taken as cost function because then it will be a non convex function of parameters.

Gradient descent will converge into global minimum only if function is convex.



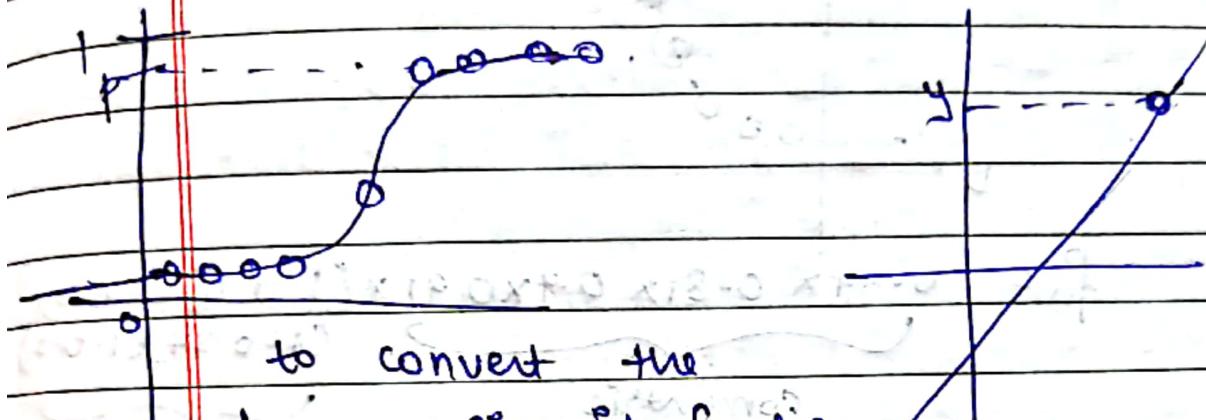
## # Wald's test

- ↳ To know if variable helps in prediction ~~or~~ or not.

# Totless ~~coress~~ "Totally useless" → not helping

- # There loss function is  $\rightarrow$  max likelihood function.

# logistic Regression is a part of Generalised Linear models.



to convert the  
pt on sigmoid function  
i.e.  $p$  on a straight

line, i.e.  $y$

$$y = \ln(p/(1-p))$$

} Transformation.

$$\text{at } p=1 \Rightarrow y = \infty$$

$$p=0 \Rightarrow y = -\infty$$

$$p=0.5 \Rightarrow y = 0$$

Reversely

$$p = \frac{1}{1+e^{-y}}$$

Coefficient are for intercepts and slope  
(Linear line)

Coeff

estimated  
standard

Intercept

$c$

slope

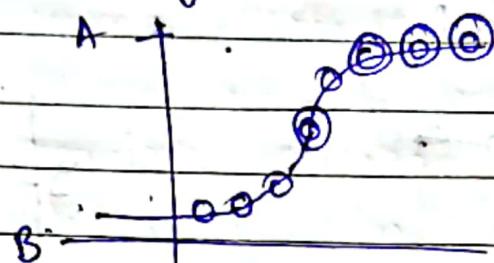
$m$

error

value,  $z = \frac{\text{est}}{\text{se}}$

} Wald's  
test

# calculating values of likelihood function



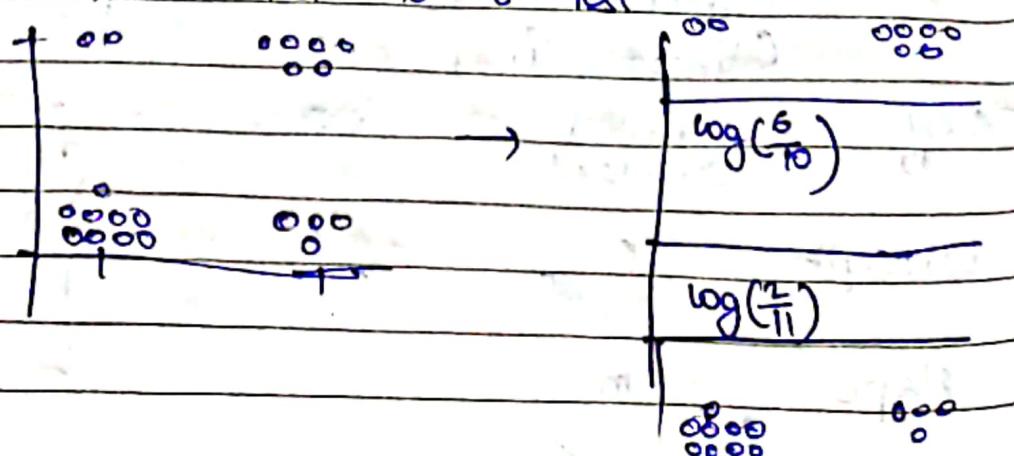
$$f_n = 0.49 \times 0.81 \times 0.9 \times 0.9 \times (1-0.6)(1-0.3) \times (1-0.27)(1-0.2)$$

concentric  
circle  
probability of  
getting A

single circle  
This done to  
scale it in the  
same term  
i.e. probability  
of not getting  
B.

- we can also take log of whole thing to analyse
- aims to get maximum likelihood.

# logistic regression in terms of discrete  
↳ similar to t-test



$$\text{size} = \log(\text{odds normal})B_1 + B_2 (\frac{\log(\text{odd mutant})}{\text{odd normal}})$$

Note:

- ↳ logistic function = sigmoid function
- ↳ features can be categorical or continuous but outcome is discrete
- ↳ linear decision boundary = hyperplane
- ↳ in equation

$$z = \frac{1}{1+e^{-y}} \quad \left. \begin{array}{l} y = mx + c \\ \Rightarrow \end{array} \right.$$

$m$  controls slope of rise

$c$  controls location of midpoint

pros

- ↳ Makes no assumption about distribution of classes in feature space
  - ↳ easily extensible to multiple classes (multinomial Regression)
  - ↳ Natural probabilistic view of class prediction
  - ↳ quick to train
  - ↳ very fast at classifying unknown records
  - ↳ Good accuracy
  - ↳ Resistant to overfitting
  - ↳ can find useful variables (important coefficients "wald's test")
- dis
- ↳ linear decision boundary

- ↳ assumes linear relationship among dependent and independent
- ↳ dependent is bound to be discrete.
- ↳ requires average or no multicollinearity b/w independent vars.

## # DECISION TREE (DT)

- ↳ used to visually and explicitly represent decision and decision making.
- ↳ predictive modelling, based on branch classification trees → classify regression trees → continuous values

- ↳ referred to as CART (Class<sup>n</sup> and Reg<sup>n</sup> trees)

## # Recursive Binary splitting

- ↳ All features are considered and different split points are tried and tested using a cost function.

- ↳ best cost one is selected.

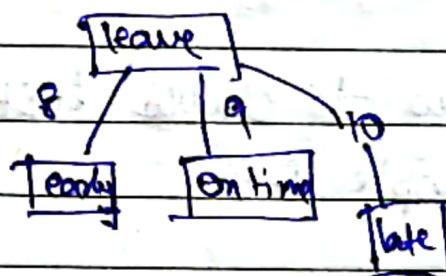
DT uses "inductive learning task"

- ↳ we particular facts to make more generalised conclusions.

We can also represent as a set of rules.

DT

rules



if leave == 8

status = early

else if leave == 9

status = ontime

else if leave == 10

status = late.

# ID3 → Iterative Dichotomiser 3

↳ classification algo

↳ greedy

↳ building DT

↳ Selects best attribute i.e. yields maximum information gain ( $G_{TC}$ ) or minimum entropy ( $H$ )

common terms in DT

Node → feature / attribute

branch → decision / rule

leaf → outcome.

DT is constructed in top-down recursive divide and conquer manner in ID3 (no backtracking)

pseudo code

↳ ppt "Lec 9 Decision Tree", page 14

JD3's bias (How to choose best DT)

↳ prefers simple DT

Entropy (H)

↳ measures the disorderliness/impurity

↳ minimum (0) when all records belongs to one class - i.e. most information

↳ maximum  $\log(C)$  [in case of Binary] when records are equally distributed among all classes

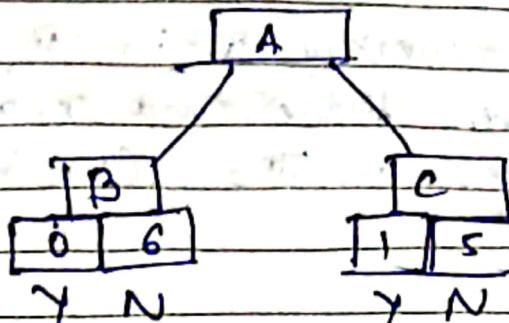
↳ small entropy → more pure

\* pure → when records ~~not~~ belong to same class.

$$\text{Entropy}(t) = H(t)$$

$$= - \sum_j p(j|t) \log_2 p(j|t)$$

ex:



$$H(B) = -\frac{0}{6} \log\left(\frac{0}{6}\right) - \frac{6}{6} \log\left(\frac{6}{6}\right) = 0$$

$$H(C) = -\frac{1}{6} \log\left(\frac{1}{6}\right) - \frac{5}{6} \log\left(\frac{5}{6}\right) = 0.65$$

$p(j|t) \rightarrow$  relative frequency of class  $j$  at node  $t$ .

### Information Gain (IG)

~~$$IG = H(p) - \left( \sum_i \frac{n_i}{n} H(i) \right)$$~~

$$IG = H(p) - \left( \sum_i \frac{n_i}{n} H(i) \right)$$

$$n = \sum n_i$$

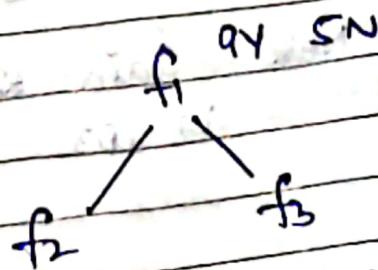
i.e., to split  $p$  into  $K$  partitions and  $n_i$  is the no. of records in  $i$ th partition.

↳ aim to max this

↳ measures reduction in  $H$

↳ con: tends to prefer splits that results in large no. of partition each small.

example:



$$IG = H(f_1) - \frac{8}{14} H(f_2) - \frac{6}{14} H(f_3)$$

$$\begin{aligned} &= -\frac{9}{14} \log\left(\frac{9}{14}\right) - \frac{5}{14} \log\left(\frac{5}{14}\right) - \left[ \frac{8}{14} \left( \right. \right. \\ &\quad \left. \left. - \frac{6}{8} \log\left(\frac{6}{8}\right) - \frac{2}{8} \log\left(\frac{2}{8}\right) \right) - \frac{6}{14} \left( -\frac{3}{6} \log\left(\frac{3}{6}\right) \right. \right. \\ &\quad \left. \left. - \frac{3}{6} \log\left(\frac{3}{6}\right) \right) \right] \end{aligned}$$

$$IG = 0.049$$

To overcome the ion of IG

↳ Gain Ratio

↳  $\frac{IG}{\text{Split Info}}$

$$\text{Split Info} = 1 - \sum_i \frac{n_i}{n} \log\left(\frac{n_i}{n}\right)$$

[terms same as IG]

↳ this adjusts GAIN (IG) by the entropy of partitioning

↳ higher entropy partitioning is penalised

↳ used by ID3 (extension of ID3)

Note: Gain Ratio is not used by ID3

↳ ID3 only used IG.

Another ~~is~~ Criteria to measure impurity

↳ GINI index for node t

$$\text{GINI}(t) = 1 - \sum_j p(j|t)^2$$

{ same  $p(j|t)$  as H }

$$\text{GINI split} = \sum_{i=1}^k \frac{n_i}{n} \text{GINI}(i)$$

↳ aims to minimise GINI split value

↳ min 0 → best

↳ max  $1 - \frac{1}{n}$  → worst

↳ used by CART, SLIQ, SPRINT

#

Specify test conditions

Name attribute type

no of ways to split

→ nominal

→ Binary

→ ordinal

→ multinomial way

→ continuous

## overfitting and underfitting

- ↳ similar to regression models

overfitting results in DT that are more complex  $\rightarrow$  than necessary

- ↳ Tree growth went too far
- ↳ Number of instances gets smaller
  - (more leaves match single example)

↳ Training errors  $\rightarrow$  not good estimate

To avoid overfitting

↳ pre-prune

↳ post-prune

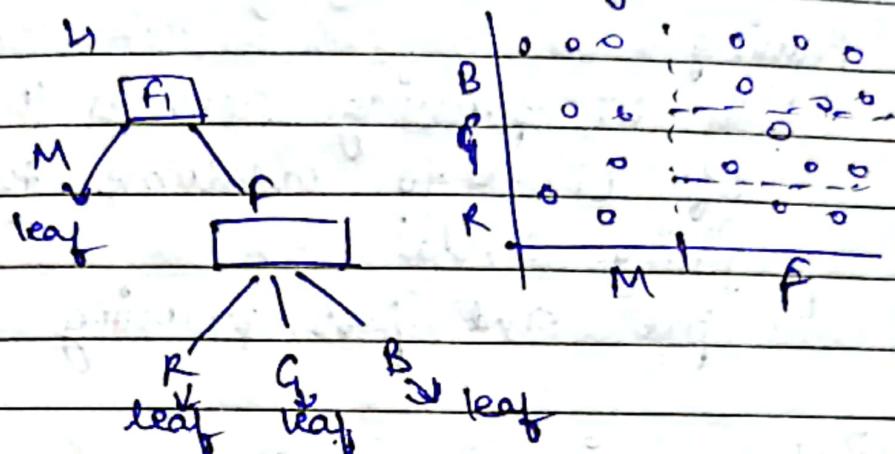
## DT pros

- ↳ Inexpensive to construct
- ↳ Extremely fast in classification
- ↳ easy to interpret (for small)
- ↳ Good accuracy

## cons

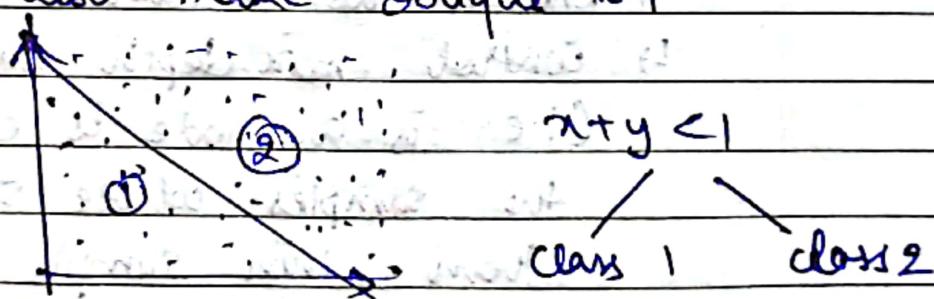
- ↳ Axis parallel decision boundary
- ↳ Redundancy
- ↳ need data to fit in memory
- ↳ need to re-train new data

### axis parallel boundary



Decision surfaces are axis aligned.

Can also make oblique DT



- ↳ Test cond<sup>n</sup> → more attributes
- ↳ expensive to find optimal test cond<sup>n</sup>
- ↳ more expressive representation.

Note:

Subtree replication

- ↳ Same subtrees appear in multiple branches.

## # Pruning

- ↳ In DT pruning, removes the branches of DT to overcome the overfitting condn
- ↳ pre and post pruning

## Post pruning

- ↳ use after construction of DT
- ↳ use when very large depth
- ↳ known as backward pruning
- ↳ control max-depth and min samples  
(i.e. if a node is classifying the samples whose size is less than min sample condn)

## Pre pruning

- ↳ used before construction of DT
- ↳ done using hyperparameter tuning
- ↳ overcome the overfitting issue

## Hyper parameter tuning

- ↳ In begining we don't know which criterion (GINI or ENTROPY), what min depth, what min sample size should be used.

- ↳ use range of values for each parameter and find the best value.

pros of pre and post

- ↳ creates simple and interpretable trees
- ↳ reduces overfitting
- ↳ since pruning selects best cross validated subtrees, pruned trees tends to fit the data well.

cons of both

- ↳ may eliminate some interesting information

Pre vs post

Pre

Post

- | Pre   | Post  |
|---|---|
| ↳ faster  | ↳ slower  |
| ↳ may not give good results as the approach is greedy and it may ignore splits that have subsequent imp. splits | ↳ results in better pruned tree (cross validation is performed at every step) |
| ↳ forward pruning   | ↳ backward pruning  |

# Measures for performance

Date / /

Page No.

# Bias: The inability of ML methods like linear regression to capture the true relationship (i.e. curved line)

# variance: The difference in fits b/w data sets is called variance. (b/w training and testing)

Note: 3 commonly used methods for finding sweet spot b/w simple and complicated model.

↳ regularization

↳ boosting

↳ bagging

High bias → underfitting of training  
↳ over generalising

High variance → overfitting

$$\text{Total error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible error}$$

We need to trade off bias and variance such that model doesn't suffer from over or under fitting.

L2

## Ridge Regression [~~L2~~ Regularisation]

- If the model is going in overfitting we introduce some amount of bias to help reducing the variance.

Linear regression minimises sum of squared errors by finding mean ~~the~~ values of parameters

by Ridge regression minimizes mean squared error +  $\lambda(\text{slope})^2$

By doing this we find better fit.  
to find the optimal value of  $\lambda$  we use cross-validation (usually 10 fold)

$\lambda(\text{slope})^2 \rightarrow$  Ridge regression ~~penalty~~ penalty.

- Ridge Regression works for continuous and discrete.

L1

## ~~L1~~: Lasso Regression

- ↳ Lasso Regression penalty  
 $= \lambda \times |\text{slope}|$
- ↳ Can be applied to same context as Ridge Regression

Ridge Regression can shrink the slope asymptotically close to 0 while Lasso regression can shrink the slope to 0.

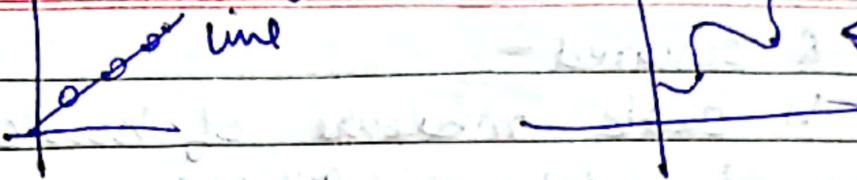
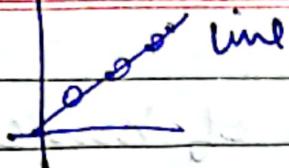
\* Lasso Regression can exclude useless variable from equations, it is little better than Ridge Regression at reducing the variance in model that contains a lot of useless variables.

but ridge better when there are more useful.

Note: As we increase  $\lambda$ , slope  $\downarrow$  but not to 0 ~~in L2~~ in L1

but there is a kink at 0 in L1 as we increase  $\lambda$ .

↳ linear  
↳ straight



## # Regression analysis

- ↳ form of predictive modelling tech which investigates relationship b/w dependent and independent variables

### Types

- ↳ linear
- ↳ logistic
- ↳ polynomial.

### Linear regression selection criteria

- ↳ classification and regression capability
- ↳ data quality
- ↳ computational complexity

### Formula

$$y = mx + c$$

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

find  $c$  using  $m, \bar{y}, \bar{x}$  in eqn.

## # R squared

↳ static measure of how close the data are fitted to regression model.

↳ coefficient of determination

↳ " multiple "

↳ cousin of R (correlation), just square of it.

if  $R^2$  of a = 0.2 and  $R^2(b) = 0.4$

then that means b is 2x better than a (but can't say with R)

$R \in [-1, 1]$   $-1, 1 \rightarrow \text{best}$

$0 \rightarrow \text{worst}$

$R^2 \in [0, 1]$   $0 \rightarrow \text{worst}$

$1 \rightarrow \text{best}$

formula

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$$

$$\text{RSS} = \sum (y_p - y)^2$$

$$\text{TSS} = \sum (y - \bar{y})^2$$

$y \rightarrow \text{actual}$

$\bar{y} \rightarrow \text{mean}$

$y_p = \text{predicted}$

low  $R^2$  is not always bad because in some cases prediction is tough and a 0.5  $R^2$  is also acceptably good.

## # confusion Matrix.

|        |  | Predict |    |
|--------|--|---------|----|
| Actual |  | P       | A  |
| A      |  | TP      | FN |
| B      |  | FP      | TN |

TP → True positive

FN → false negative

FP → false positive

TN → True negative

How to remember?

draw fig, mark A → +, B → -

~~+ tell~~

|      | A(+)  | B(-)  |
|------|-------|-------|
| A(+) | (+) + | (+) - |
| B(-) | (-) + | (-) - |
|      | - +   | + -   |

F P              T N

I term ↗ ↑ ↘ II term

I term → false/true positive/Negative  
easy.

Date / /

Page No.

$$TP \text{ rate} = \frac{TP}{TP + FN}$$

$$FP \text{ rate} = \frac{FP}{FP + TN}$$

$$\text{Success rate (Accuracy)} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Error rate} = 1 - \text{Success rate}$$

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

↑ measure  
detect true class

$$\text{Specificity} = \frac{TN}{TN + FP}$$

↑ not detecting  
too many  
false positive

Specificity ↴ used to confirm  
results of sensitivity.