

Advanced Machine Learning

Advanced Machine Learning (guarantor Christina Bauer)

- Linear and polynomial regression: Explain how the squared error cost function is applied to linear/ polynomial regression and how the gradient descent algorithm is used to solve the optimization objective.
- Logistic regression: Explain the difference of logistic regression and linear regression in terms of the used hypothesis, cost function, and optimization objective.
- Neural networks: Describe the model structure of a neural network and how forward and backward propagation is used to solve the optimization objective.
- Support vector machines: Explain the difference of SVM to logistic regression and how kernels are used in this context.
- Evaluating learning algorithms: Describe the bias and variance problem and how regularization, cross-validation, and other methods can be used to solve for these problems

① Linear regression

- statistical method to model rel. b/w dependent variable and one or more independent variables by fitting a linear equation

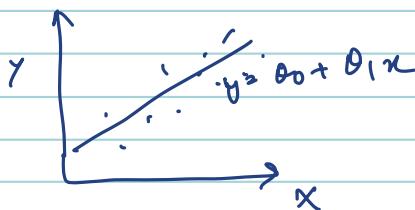
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

hypothesis function:

- suitable for linear trends

- few parameters simpler & less prone to overfitting

- not suitable for non-linear data



- e.g.: Predicting house prices
size ↑, price ↑

- Objective is same for both, to find optimal θ values (params) to minimize the error b/w predicted and target variable.

MSE → used in Regression to quantify error b/w predicted and target values (outputs).

hypothesis function: $h_{\theta}(x) = \theta_0 + \theta_1 x$

cost function : $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$

↑ prediction ↓ actual target
 error

Polynomial Regression

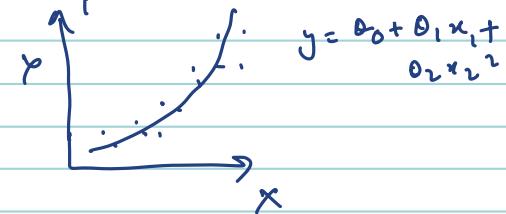
- extension of linear regression, when the relation b/w dependent and independent variable is modeled as polynomial function.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3 + \theta_4 x_1^4 + \theta_5 x_1 x_2 \dots$$

- suitable for non-linear trends

- more parameters, complex, more prone to overfitting

- suitable for non-linear complex data



- e.g.: Crop yield based on temperature, soil, rainfall
- growth of bacteria over time.

squared error ensures that both the error terms are treated equally.

Gradient Descent: Optimization algorithm to find the optimal value for θ (params) to minimize the cost function.

It iteratively updates the parameters in the direction of the negative gradient.

Cost function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$

Initialize with random params

GD update rule:

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

- ↓
- compute $h_\theta(x^i)$
- compute error
- calculate gradient
- update params
- repeat steps

Gradient for param θ_j

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) x_j^i$$

for polynomial includes x^j, x^3, x^4, \dots

In linear, cost function is convex.
∴ ensures convergence to global minima.

In polynomial, cost function is nonconvex,
can get stuck at local minima.



Logistic Regression: Statistical model used for binary classification problems.

Model estimates the probability that a given input belongs to certain class.

Hypothesis function: Sigmoid function (logistic) to map input features to probability value of 0 & 1.

g/p Independent features

$$\begin{bmatrix} x_{11} & x_{21} & \dots & x_{m1} \\ x_{12} & x_{22} & \dots & x_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1n} & x_{2n} & \dots & x_{mn} \end{bmatrix}$$

$y = \begin{cases} 0 & \text{if class deposit} \\ 1 & \text{if class feature} \end{cases}$



$$h(x) = \frac{1}{1 + e^{-x}}$$

$$0 < h(x) < 1$$

for linear regression: $h_0(x) = \beta_0 + \beta_1 x$

for logistic regression: $\sigma(z) = \sigma(\beta_0 + \beta_1 x)$

Decision boundary:

Setting up a threshold (common 0.5)

is set to classify

if $h_0(x) > 0.5$ classify 1
 $h_0(x) < 0.5$ classify 0

$$z = \beta_0 + \beta_1 x$$

$$h_0(x) = \text{sigmoid}(z)$$

$$h_0(x) = \frac{1}{1 + e^{-z}}$$

$$h_0(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$$h_0(x) = \sigma(w^T x + b)$$

$z = \text{linear combination of input features}$
 $z = \text{input to sigmoid function}$

The linear regression MSE cost function does not work in logistic regression, as it results in non convex graph

Cost function in logistic regression

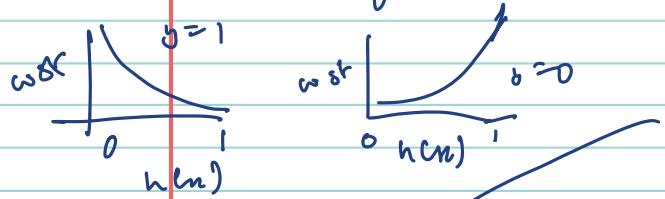


↳ Log loss function

(Binary cross entropy loss)

tailored

for binary classification



Cost function $(h_0(x), y)$	$\begin{cases} -\log(h_0(x)) \text{ if } y=1 \\ -\log(1-h_0(x)) \text{ if } y=0 \end{cases}$
--------------------------------	--

two equations can be compressed into 1

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left(y_i \log(h_0(x_i)) + (1-y_i) \log(1-h_0(x_i)) \right)$$

y_i = actual label

$h_0(x)$ = predicted prob. of i th class

Gradient Descent:

Reduce cost value

cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left(y_i \log(h_0(x_i)) + (1-y_i) \log(1-h_0(x_i)) \right)$$

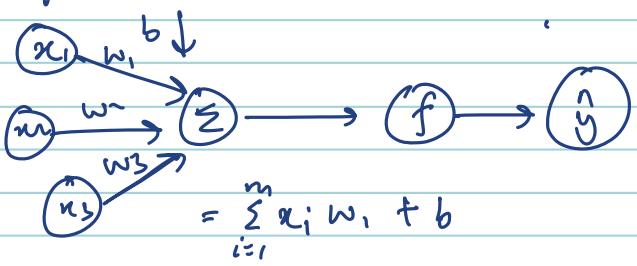
$$h_{\theta}(x_i) = \hat{y} = \sigma(w^T x + b) \quad \sigma(x) = \frac{1}{1+e^{-x}}$$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m (y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)) \right]$$

Aspect	Linear Regression	Logistic Regression
Hypothesis Function	$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$	$h_{\theta}(x) = \frac{1}{1+e^{-(\theta_0+\theta_1 x_1 + \dots + \theta_n x_n)}}$
Cost Function	Mean Squared Error (MSE): $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$	Log-Loss (Binary Cross-Entropy): $J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(h_{\theta}(x_i)) + (1-y_i) \log(1-h_{\theta}(x_i))]$
Gradient of Cost Function	$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_{ij}$	$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_{ij}$
Gradient Descent Update Rule	$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$ MSE min	$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$ log loss min
Prediction Output	Continuous value (e.g., price, height)	Probability (0 to 1)
Convergence Behavior	Generally straightforward due to convex MSE cost function	Convex Log-Loss cost function, but sigmoid function introduces non-linearity
Feature Scaling Importance	Helps speed up convergence, especially with diverse feature ranges	Crucial for efficient convergence due to sigmoid scaling issues

③ Neural Network

Class of ML models inspired from the structure & function of biological neurons in brain.



- Input layer
- Hidden layer
- Output layer.

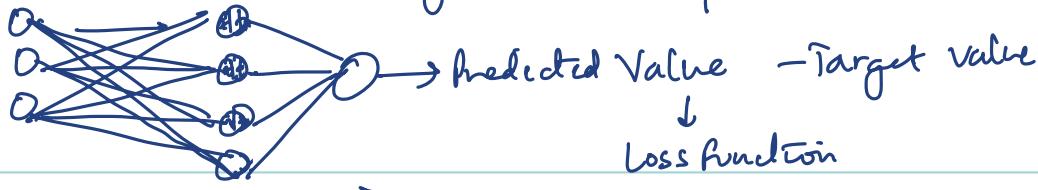
w = Represent strength of connection w/b neuron

b = shift activation function, improves model's flexibility.

→ forward propagation:

process by which input data passes through the network and generates output.

Input → hidden layer → Output



$O_i \rightarrow$ Input feature fed to network

Hidden \rightarrow weight \times Input + bias and ReLU is applied

$O_h \rightarrow$ computes output based on prob statement

Reg : continuous (logistic neuron)

Class : neuron = no. of classes

\rightarrow Example with Linear Regression

Forward

Computes output of
n/w from given
input

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Input: x - input
 y - true value

hidden layer: compute weight
 $z = \theta_0 + \theta_1 x$

Output: $h(z) = z$

calculate cost: $\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$

process iteratively continues
until model converges
to an optimal solution
(multiple epochs)

Backward

minimize the error
by updating params

[involves finding gradient ∇
of cost $J(\theta)$ w.r.t to params]

compute error:

$$\text{Error} = h_{\theta}(x) - y$$

$$\text{Error} = (\theta_0 + \theta_1 x) - y$$

compute gradient for each θ_0 & θ_1

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$

bias (intercept)

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)x_i$$

slope (w)

update weight & bias

$$\theta_0 := \theta_0 - \alpha \frac{\partial J(\theta)}{\partial \theta_0}$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial J(\theta)}{\partial \theta_1}$$

α = learning rate

SVM \rightarrow Support Vector Machines

Supervised Machine learning algo.
primarily for classification tasks.

Objective: finds an optimal hyperplane that best separates the data points of different classes in a high dimensional space.

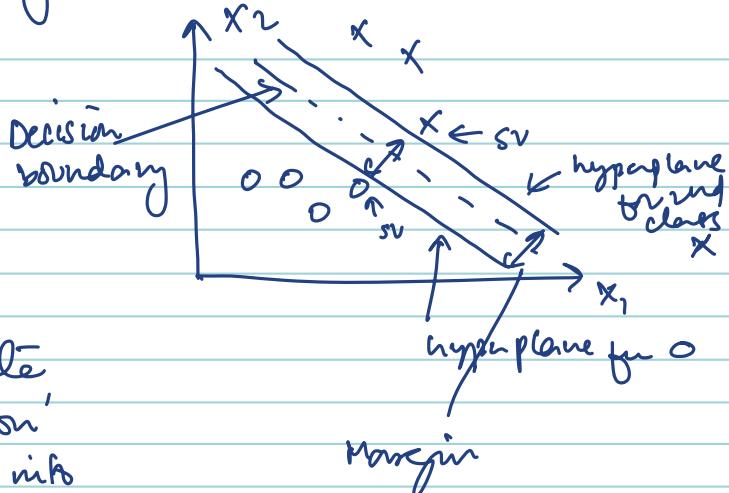
maximizes the margin b/w the classes

hyperplane - decision boundary

2D - line

3D - plane

high dim: hyperplane

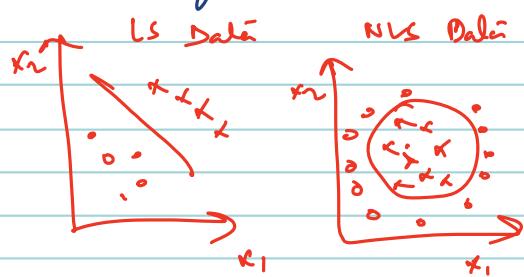
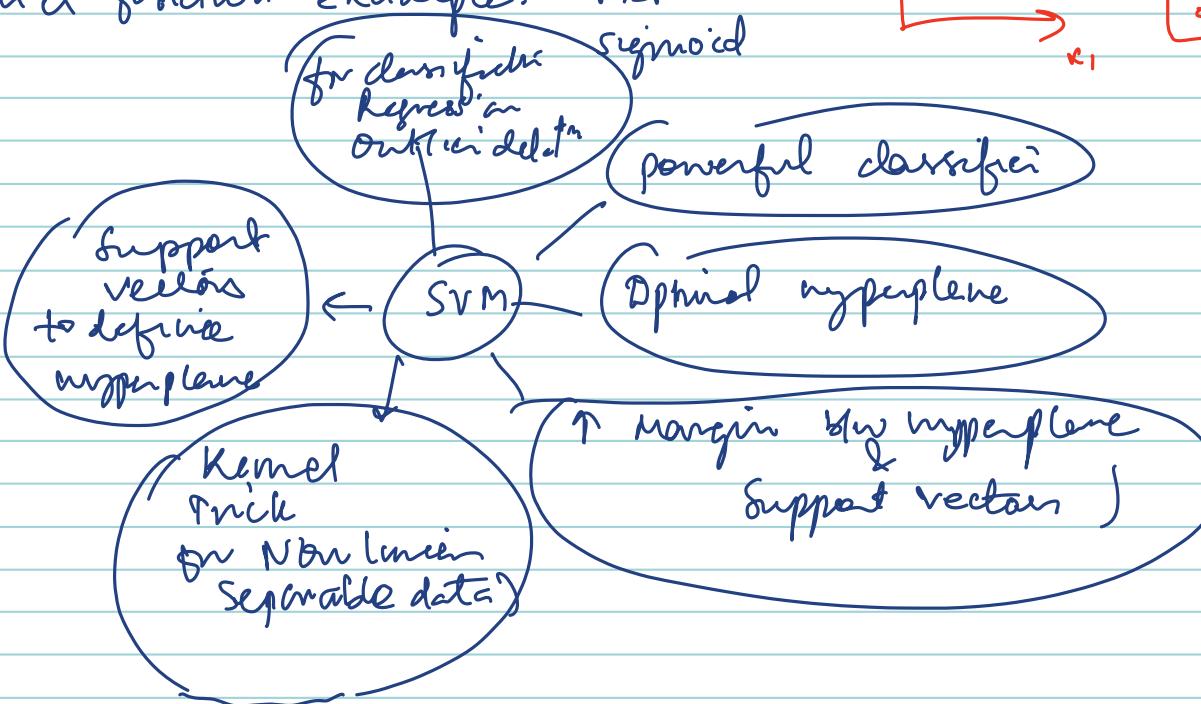


Kernel trick:

For non-linearly separable data, SVM uses Kernel function to map the input features into a higher dimensional space where,

Linear hyperplane can separate the data

Kernel function example: RBF



- ① SVM can handle non linear & linear decision boundary (finds hyperplane)
- ② More robust at handling outliers
- ③ Maximizes margin \rightarrow optimization objective
- ④ Computationally expensive

Logistic regression

- ① Linear decision boundary
- ② Sensitive to outliers
- ③ minimizes cost (log loss) using algo like GD.
- ④ Simple implementation & less expensive

(5) effective in ↑ dimensional space

(6) hypothesis function: Uses sigmoid to output probabilities

(6) bdim spaces (linear)

(6) uses hyperplane to separate class

(5) Evaluating Learning Algorithm:

Bias : Difference b/w predicted and actual values.

→ ↑ Bias model pays little attention to training data & oversimplifies the model

↑ Bias
Model too simple

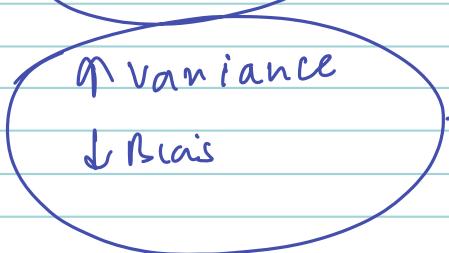
Variance: Model's sensitivity to fluctuations in training data
↑ Variance model → pays lot of attention to training data and learns noise & randomness in data leading to not able to generalize on test data (unseen)

↓ Variance
Model too complex



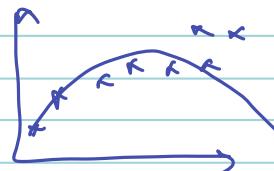
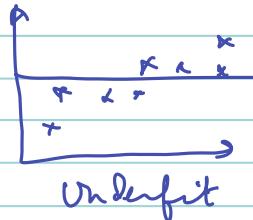
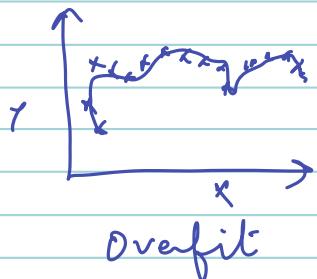
→ Underfitting

Performs bad on train & test data



→ Overfitting

Performs good on train
bad on test (unseen)



Bias variance Tradeoff:

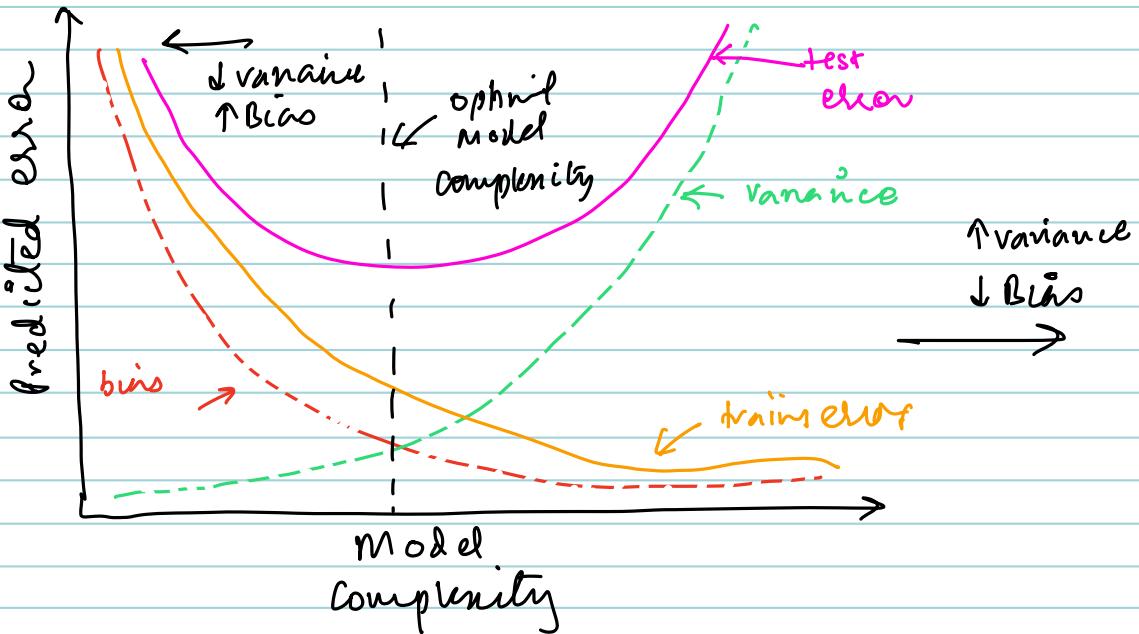
If model → simple and few params → ↑ Bias, ↓ Variance (underfit)

If model → complex and more params → ↓ Bias, ↑ Variance (overfit)

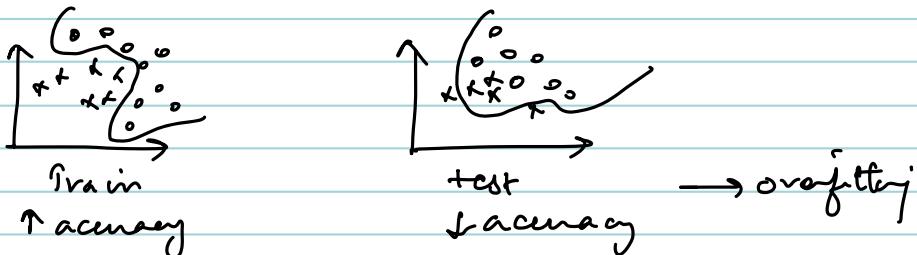
- Goal is to have low Bias & low variance, fitting

data well w/o being overly complex
for good model, good balance b/w bias & variance
such as it minimizes total error

$$\text{Total Error} = (\text{Bias})^2 + \text{Variance} + \text{Irreducible Error}$$



Regularisation:

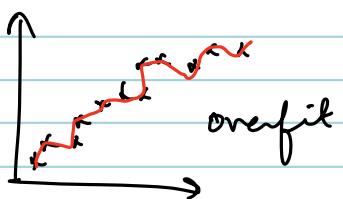


to avoid overfitting, we ideally want a smooth curve that can fit well in both training as well as test dataset.

↓ use
regularisation



e.g.: Regression



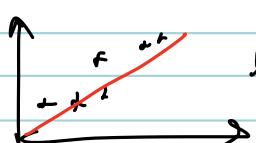
$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

500 params

technique → called regularisation
nullify effect of certain params/term to get smooth



$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots$$



$$\text{linear } z = \theta_0 + \theta_1 x$$

2 params → might be (ie, we want middle ground)

Regularisation. C_1 & C_2 \rightarrow [nullify the effect of certain params / neurons
↑ (linearity on model)]

Methods

- Dropout
- Early stopping

Ignore less important features \leftarrow Regularisation \rightarrow lower model complexity
 \rightarrow computational expense reduce.
by scaling down / shrink regression coefficients

$$y = 0.4 + 1.2x_1 + 2.5x_2 + 3.9x_3$$

↑ res. rimp ↑ imp

L_1 Regularisation LASSO

$$\text{Loss} = \text{Original loss} + \lambda \|\mathbf{w}\|$$

$\|\mathbf{w}\|$ absolute value of coefficient

can lead to some coeff value $\rightarrow 0$
for feature selection

L_2 Regularisation Ridge

$$\text{L2 or Ridge} = \text{Loss} + \lambda \|\mathbf{w}\|^2$$

(penalty)

$$\|\mathbf{w}\|^2 = w_1^2 + w_2^2 + w_3^2 + w_4^2 + \dots + w_n^2$$

- as $\lambda \uparrow$, coefficient mag \downarrow
- adds sum of squared coeff to loss
 - Reduces mag of coeff, but not $\neq 0$

Dropout Regularisation

Randomly dropping out (deactivating) a proportion of neuron during training,

helps in model to avoid relying too heavily on any single neuron and encourages robust feature engineering.

If dropout $\rightarrow 20\%$, then each neuron has 20% chance of being dropped at every training step.

Early Stopping to avoid overfitting : (Regularisation techn.)

Technique where, training process is halted once the model's performance on validation data starts degrading

CV (80.4, 78.5, 79.5, 82.4)

Robust

generalized
consistent

CV (68.2, 81.2, 72.1,
diverse 52.1)
not robust
overfitting

Cross Validation

To evaluate the model's performance to detect overfitting & underfitting

Dataset \div into 3 parts

Training Validate Test
66% 20% 20%

Methods : - K fold cross validation

All data



\rightarrow test
fold

K folds

Test

eg:
for
 $K=5$

split-1 1 2 3 4 5 \rightarrow accuracy of Itv 1

split-2 1 2 3 4 5 \rightarrow Acc Itv 2

split-3 1 2 3 4 5 \rightarrow Acc Itv 3

split-4 1 2 3 4 5 \rightarrow Acc Itv 4

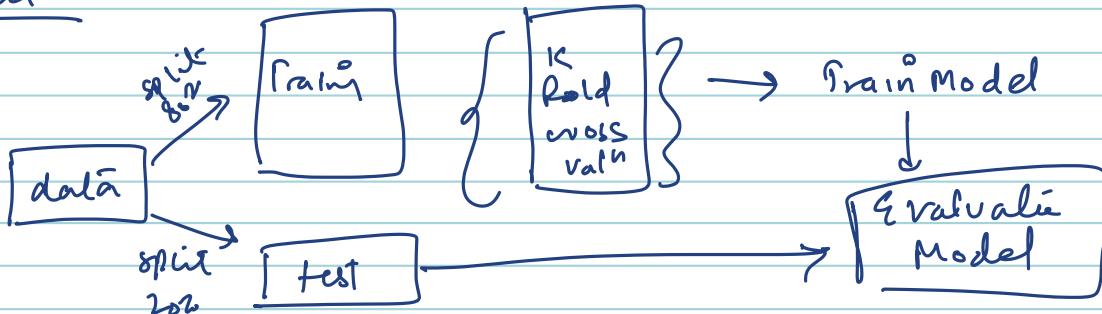
splits 1 2 3 4 5 \rightarrow Acc Itv 5

average accuracy score

low para acc

Test data \rightarrow accuracy

Flowchart:



for K , one fold is test & remaining $K-1 \rightarrow$ training

Can be used for model selection

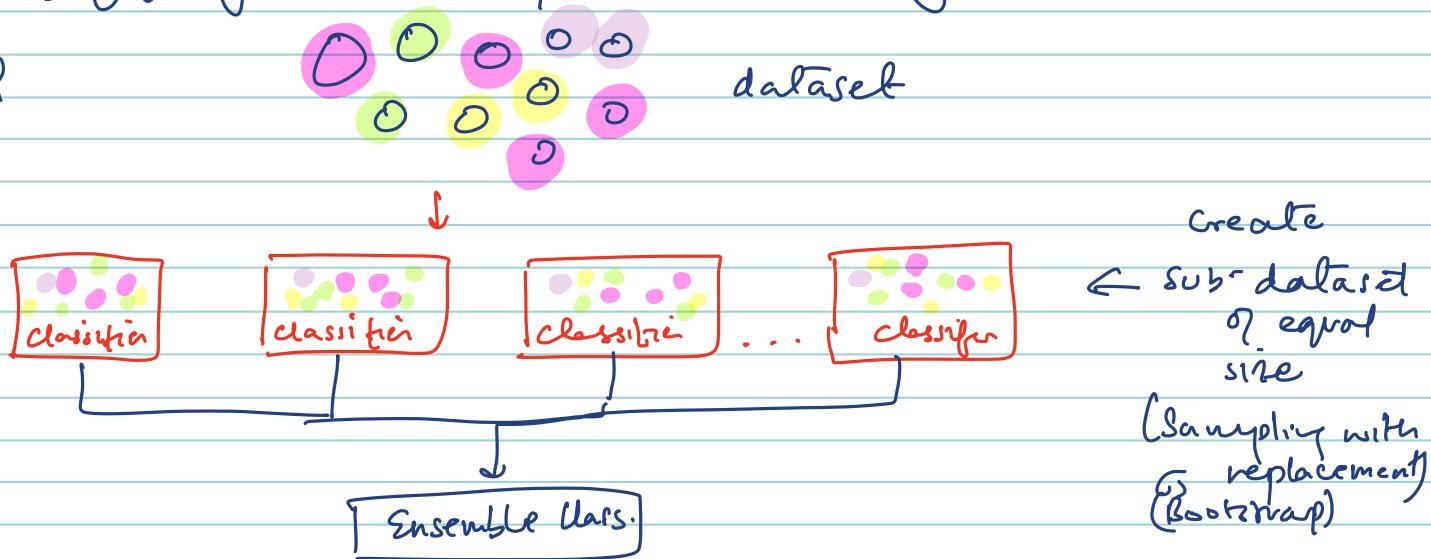
gives more robust estimate of model's performance on unseen data & hence prevents overfitting

data efficient: every data point is used for training validation, making most of available data.

Ensemble Methods: Address Bias-Variance tradeoff by combining multiple models to improve overall performance.
to reduce variance

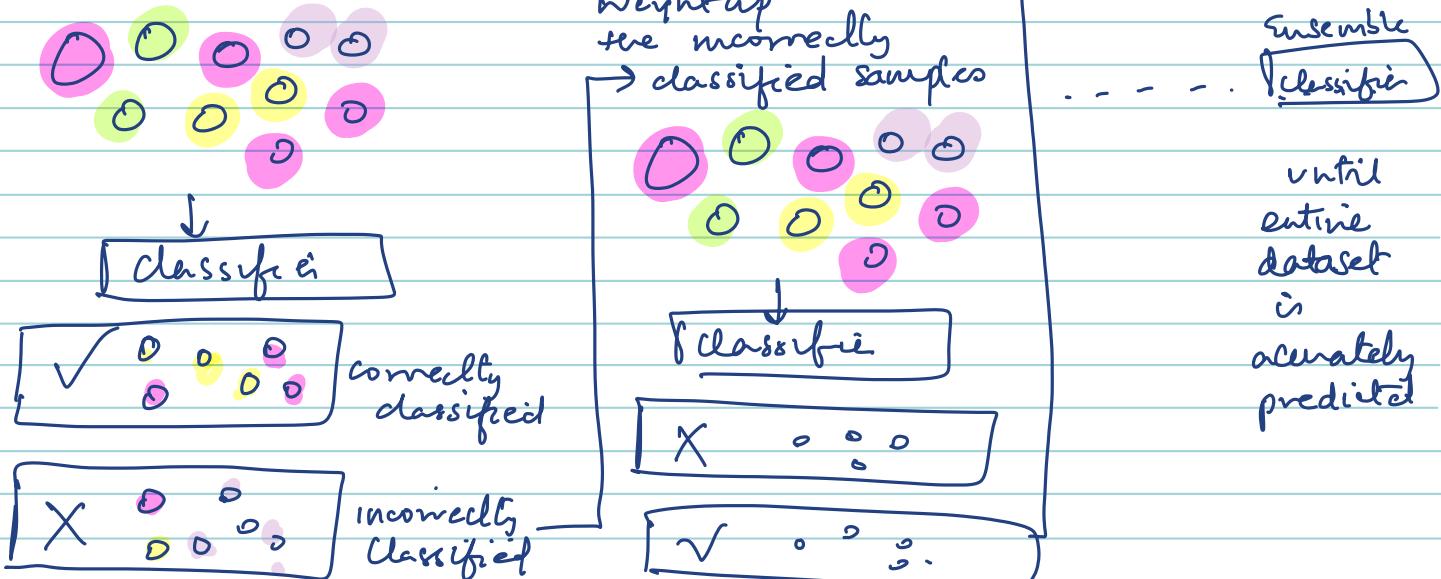
Bagging: Bootstrap + Averaging

parallel



Boosting

sequential



Both are good at reducing variance as they combine multiple models

But for bias, Bagging has close to zero bias as no weight of data is change.

Boosting: weighting added :: Reduces Bias.
but ↑ prone to overfitting