

Machine learning:-

A computer program P is said to learn with experience E on class of tasks T , if the performance on tasks T , measured by P improves with experience E .

- Recognizing spoken words
- Driving autonomous cars
- Identifying newly found astronomical structures
- Learning to play word-class backgammon

Designing a learning system:-

1) Choosing training experience:-

- Type of feedback ^{direct} on the each move or the performance ^{indirect}
- Degree of control on sequence of training examples
- Choosing a representation for training examples ^{distributing}

2) Choosing target function:-

$$\cdot T: B \rightarrow M \quad (\text{chooseMove})$$

$$\cdot V: B \rightarrow R^{\text{real no.}}$$

• Learning the target function is called function approximation.

$$V(b) = 100 \quad V(b) = -100 \quad V(b) = 0 \quad b = \text{final state}$$

win loose draw

$$\cdot V(b) = V'(b)$$

3) Choosing a representation for target function:

$$\hat{V}(b) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_6 x_6$$

4) Choosing a function approximation algorithm

$$\langle b, V_{\text{train}}(b) \rangle \quad b \in B$$

ordered pair of each move & score associated with it is given.

estimating training values :-

Assign the training value of $v_{train}(b)$ for any intermediate board state b to be $\hat{v}(\text{successor}(b))$, \hat{v} is current approximation & \hat{v} successor is next board state.

$$v_{train}(b) \leftarrow \hat{v}(\text{successor}(b))$$

Adjusting the weights :-

We need define best hypothesis, or set of weights which minimizes the squared error E between the training values & the values predicted by hypothesis \hat{v} .

$$E = \sum (v_{train}(b) - \hat{v}(b))^2$$

LMS weight update rule:

$$w_i \leftarrow w_i + \eta (v_{train}(b) - \hat{v}(b)) x_i$$

Issues in ML:-

- How much data is sufficient
- Which algorithm perform best
- When and how can prior knowledge be helpful
- Strategy for choosing training experience
- How can learner automatically alter its representation

Chapter 2:-

Concept Learning :-

A task of finding a potential hypothesis that best fits the training examples.

Inductive learning hypotheses :- Any hypothesis found to approximate the target function over a large set of training examples will also approximate target function over other unobserved examples.

attributes/instances & target function

Instance space = The target hypothesis that we want to learn is represented as a conjunction of all attributes

Possible distinct instances = multiplication of possible values for each attribute.

Hypothesis space:
Set of all possible hypotheses

If sky has 3 possible values and other attributes have 2 values

$$\text{instances space} = 3 \times 2 \times 2 \times 2 \times 2 \times 2$$

$$\text{Hypothesis space} = (3+2) * (2+2) * (2+2) * (2+2) * (2+2) * (3+2) + 2 \text{ because one for } \phi \text{ & one for ?}$$

Concept learning space: - Task of searching through a large space of hypotheses to find the hypothesis that best fits the training example.

$$\text{General hypothesis} = h = \langle ?, ?, ?, ?, ?, ? \rangle$$

$$\text{specific hypothesis} = h = \langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$$

Given hypotheses h_j and h_k , h_j is said to be more general or equal to h_k , if any instance that satisfies h_k will also satisfies h_j .

$$h_j \geq h_k \equiv \forall x (h_k(x) = 1 \rightarrow h_j(x) = 1)$$

Find-S algorithm:-

Finding the most specific hypothesis.

- 1) Initialize h to most specific hypothesis ($\langle \phi, \phi, \phi, \phi, \phi, \phi \rangle$)
- 2) For every positive training example x ,
 - For every attribute constraint a_i in h
 - if the constraint is satisfied by x , then do nothing
 - else replace a_i in h by the next more general constraint that satisfies x .
- 3) Output h

* Ignores the negative example

* Hypothesis is consistent with positive training examples

Issues in find-s algorithm

- Has the learner converged to the target concept? - Is that is only consistent hypothesis present?
- Why most specific hypothesis?
- Are training examples consistent? What if noise/errors are present and its ignoring the negative examples as well.
- What if there are several maximally specific consistent hypotheses? Learner should be able to backtrack to get the more generalized one.

Consistent hypothesis:-

Hypothesis h is consistent over training example D with target function c , iff $h(x) = c(x)$ for each $\langle x, c(x) \rangle$ in D .

$$\text{Consistent}(h, D) \equiv \forall \langle x, c(x) \rangle \in D, h(x) = c(x)$$

Version space :- Set of all consistent hypotheses

$$VS = \{ h \in H \mid \text{consistent}(h, D) \}$$

List-then-eliminate algorithm:-

- 1) List all possible hypothesis in H in version space
- 2) for each hypothesis h in H
 - remove any from version space if $h(x) \neq c(x)$
- 3) output version space

* Works if finite hypotheses space is present

Candidate-elimination algorithm :-

- Works on List-then-eliminate principle
- More compact representation of version space
- Gives more specific & general hypothesis
- Gives the boundaries that delimits the version space

Algorithm :-

Initialize more general hypothesis G

Initialize more specific hypothesis S

for each training example d , do

if it is a +ve example:

for each s in S

- Remove s if its not consistent with d .

- Add ~~next~~ next to general constraint in s which satisfies d .
- Remove g from G if any are more specific than s in d .

if it is a -ve example:

for each g in G

- Remove g if not consistent.

- Add next to specific constraint in g which satisfies d .

- Remove g if any hypothesis is less general than another hypothesis in G .

Remarks :- There should not be errors in training examples

→ There should be a hypothesis that correctly describes target function

Biased hypothesis :- Only conjunction of all attributes.

$$bf(x_i, D_c) \rightarrow L(x_i, (D_c \wedge x_i)) \rightarrow L(x_i, D_c)$$

Inductive bias :- Inductive bias of L is the set of assumptions

B such that for new instances x_i ,

$$(B \wedge D_c \wedge x_i) \vdash L(x_i, D_c)$$

Chapter 3 :-

Decision trees :- Disjunction of conjunction of constraints of attribute values

When to use decision trees:-

→ Attributes have multiple values

→ Distinct output/class value

→ Disjunctive principle needed

→ If training data consists missing values or errors.

ID3 - Iterative Dichotomiser

Repeatedly divides the attributes into 2 or more subgroups.

→ Top-down approach

→ Never back-track

→ Maximum information gain or minimum entropy

→ Entropy is the degree of uncertainty corresponding with a random attribute.

$$\text{Entropy}(S) \equiv - (P_1 \log_2 P_1 + P_0 \log_2 P_0)$$

→ Information Gain is the expected reduction in entropy caused by partitioning the examples according to the attributes.

$$\text{Gain}(S, A) \equiv \text{Entropy}(S) - \sum_{i \in S} |S_i| \text{Entropy}(S_i)$$

Advantages of decision tree :-

→ Searches complete hypothesis space.

→ Uses all training examples

→ Less sensitive to errors

Disadvantages :-

→ Only one current hypothesis. No capability of representing all consistent hypothesis.

→ No backtracking

Training set :- * Data used to train the model.

* For every epoch, model learn over same data to learn the features.

Validation set :- * Validating the model

* How well the model is learning

* Give information on how to adjust hyperparameters (learning rate)

- * How to generalize the training examples
- * Ensure the model is not overfit

Test set:-

- * Data set is a set of data used to test the trained model
- * Unlabelled data is given

Cross-validation:-

- Statistical model to check ability of training model
- Select & compare the model.
- Divide the dataset into k -equal size groups
- Shuffle the dataset & then divide.
- Take 1 group for testing & $k-1$ for training
- $k-1$ used to fit the model
- 1 group to evaluate the model
- Get the evaluation score
- Discard the model
- Repeat this for k times

Types

- Representative $k=10$ $|S|=n$
 high variance & high bias (over-estimation)
 (change of data)
- Representative — too large amount of dataset
- $k=10$ — safer side. Low bias + modest variance
- $k=n$ — leave out one cross validation (LOOCV)
 - each training example considered for train & testing

K-fold APT:-

- sci-kit library
- Kfold (split no., whether to split or not, seed (no. used prior to shuffle))

Issues in decision tree:

- 1) How to avoid how deeply to grow the decision tree:
 - Over fitting
 - Training - fitting the model over a training data-set to get reliable prediction on unseen data
 - If there is a hypothesis h_1 , which lead to overfit such that h_1 has less error than h_2 on training examples but h_2 has less error than h_1 over complete distribution of instances.
 - Reason: Error in training example
 - Small number of examples associated with leaf nodes
 - More complex \rightarrow Decision Tree depth more
 - Train performance \uparrow but \downarrow testing performance
- 2) Handling continuous value attributes
 - changing the continuous value to discrete
 - Adding a new discrete valued attribute that partitions the continuous values into discrete value intervals
- 3) Choosing an appropriate attribute measure
 - Alternate way of choosing an attribute
 - If we take 'day' as root it classifies correctly but it has broader distribution.
 - It is a poor classifier because, its actually not a useful attribute
 - Gain ratio. - uses split information to know how broadly & uniformly splits the data
 - $$\text{split}(S, A) = - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$
 - $$\text{Gain Ratio}(S, A) = \frac{\text{Gain}(S, A)}{\text{split}(S, A)}$$

4) Handling missing values
 → most common data

5) Handling attributes with different cost
 → Gain(S, A)
 cost(A)

- If possible attributes with low cost taken
- for reliable classification - high cost attributes.

Gini-index or impurity :-

- Degree probability of a variable being wrongly-classified when randomly chosen.
- If lower gini-index preferred
- $0 +$ Between 0 and 1
- 0 : only one class or all elements in one class
- 1 : elements randomly distributed over all different classes
- 0.5 : elements ^{equally} randomly distributed over different classes.

$$\rightarrow 1 - \sum_{i=1}^C (P_i)^2$$

Avoid over-fitting :-

How to know the correct length of tree.

- 1) Direct approach - 2/3 training & 1/3 validation
- 2) Create a decision tree & then check if extending node helps in improving the accuracy
- 3) Minimum Length Description principle - encoding the training example & hypothesis - Halting the growth when encoding size minimized

Pruning:-

- 1) Reduced error pruning

→ Pruning of decision node

→ Remove sub-tree

→ Consider it as a leaf node & add most commonly classified value

→ Do it only if the resultant works fine
 → This will remove of any leaf nodes added because of irregularities

Rule post-pruning:-

- Over-fit the model
- Convert the tree as rules
- Generalize the rules, removing preconditions
- Arrange it based on accuracy

Internals :-

Computational learning theory:-

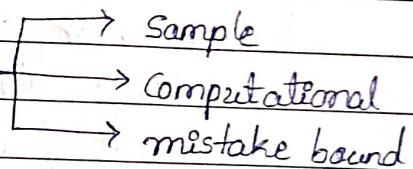
Sample complexity - Number of training examples

Computational complexity - Computational resources

Mistake bound - Incremental learning

- No. of misclassified training example

Learning problems → models/algorithm



Let us take the tree error for a hypothesis (h) w.r.t target concept C and instance distribution D . It is the probability that h will misclassify an instance drawn at random according to D .

$$\text{error}(h) = \Pr_{x \in D} (h(x) \neq c(x))$$

Approximating True error for hypotheses

Selecting h with true error 0:

→ Select h with true training error 0

→ But there may be many h , we don't know which to select

$\text{error}(h)$ Let true error $\leq \epsilon$, (small)Not all training instances succeed, only with $1-\delta$
probability will

(small)

 ϵ -exhaustedThe version space VSHD is said to be exhausted
wrt target concept C & D (instance distribution)
if for all hypothesis (h) in VSHD has true error
less than ϵ .

$$(\forall h \in VSHD) \text{error}(h) \leq \epsilon$$

How many training examples will ϵ -exhaust the VS?
Haussler's theoremFor a finite hypotheses space H, D is sequence of
random independent examples for target concept C,
then for any $0 \leq \epsilon \leq 1$, the probability of VS wrt H and
D for C is not ϵ -exhausted is less than $|H|e^{-\epsilon m}$

Proof:-

→ Let $h_1, h_2, h_3, \dots, h_k$ be hypotheses space H, which have
true error $\geq \epsilon$ for D and C

$$H_{bad} = \{h_1, h_2, h_3, h_4, \dots, h_k\}$$

→ We will fail to exhaust the VS if one of these hypotheses
happens to be consistent with all training examples in

→ The probability of a hypothesis having $\text{error}(h) \geq \epsilon$
being consistent with one of the training examples is
at most $(1-\epsilon)^m$

→ This hypothesis getting consistent with all m training
examples at most $(1-\epsilon)^m$

→ If we have k hypothesis
 $k(1-\epsilon)^m$

→ But $k \leq |H| \Rightarrow |H|(1-\epsilon)^m$

→ General inequality eq $a \leq \epsilon \leq 1 \Rightarrow (1-\epsilon) \leq e^{-\epsilon}$

$$\Rightarrow k(1-\epsilon)^m \leq |H|(1-\epsilon)^m \leq |H|e^{-\epsilon m}$$

→ Let us derive training examples no. wrt error rate $\leq \delta$
 $|H|e^{-\epsilon m} \leq \delta$

$$m \geq \frac{1}{\epsilon} (\ln(|H|) + \ln(\frac{1}{\delta}))$$

→ Number of training examples m is sufficient to tell the model is probably approximately correct hypothesis

PAC Learning:

Let us consider a concept C over n instances X by length n , learner L uses hypotheses space H . The C is PAC learnable by L if $\forall c \in C$, distribution D over X , such that $0 < \epsilon < \frac{1}{2}$, $0 < \delta < \frac{1}{2}$, then the L will with prob atleast $(1 - \delta)$ output a hypothesis $h \in H$ $\text{error}_{(D)}(h) \leq \epsilon$, with time polynomial to $1/\epsilon$, $1/\delta$, n , $\text{VC}(C)$

Agnostic learning:- The learner will not make any assumption that the target concept c is representable by H , it will simply just produce hypothesis with $\text{training error} \leftarrow \epsilon$ ~~if~~ don't assume $c \in H$

$$\Pr[\text{error}_D(h) \rightarrow \text{error}_D(h) + \epsilon] \leq e^{-2\epsilon^2 m}$$

$$\Pr[\text{error}_D(h) > \text{error}_D(h) + \epsilon] \leq e^{-2m\epsilon^2}$$

$$\Pr[(\exists h \in H)[\text{error}_D(h) > \text{error}_D(h) + \epsilon]] \leq |H| e^{-2m\epsilon^2}$$

$$m \geq \frac{1}{2\epsilon^2} (\ln(|H|) + \ln(\frac{1}{\delta}))$$

Shattering a set of instances :-

Let the such the instance set S is shattered by H if $h \in H$ consistent with dichotomy.

$2^{|\mathcal{S}|}$ dichotomies through H

Vapnik-Chervonenkis dimension :-

$\text{VC}(H)$ is defined over H and X is the size of the largest finite subset of X shattered by H

$$m \geq \frac{1}{\epsilon} (4 \log_2(\frac{2}{\delta}) + 8 \text{VC}(H) \log_2(\frac{13}{\epsilon}))$$

Conditional probability = Probability of an event happening given that it has relationship with one or more other events.

Bayes theorem - formulates the test. Takes the test result probability and calculates the real probability that the

Bayes theorem - Takes the test result and gives the real probability that the test identified the event.

$$P(A|B) = P(B|A) \times P(A)$$

$P(B)$

$$P(h|D) = \frac{P(D|h) * P(h)}{P(D)}$$

→ prior prob of h (prior knowledge of chance that h is correct)

The prior prob of D , prob that D is observed

The prob of D observed given that h is correct

Posterior prob. h holds after D has been observed

MAP hypothesis

Maximally probable hypothesis - Maximum Posterior hypothesis

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$$

$$= \operatorname{argmax}_{h \in H} \frac{P(D|h) * P(h)}{P(D)}$$

$$= \operatorname{argmax}_{h \in H} P(D|h) * P(h)$$

$P(D|h)$ → likelihood of D given h

Any hypothesis that maximizes $P(D|h)$ - Maximum likelihood

$$h_{ML} = \operatorname{argmax}_{h \in H} P(D|h)$$

Brute-force

→ here the learner considers a finite hypotheses space H over training instances X , where the task to learn target concept $c: X \rightarrow \{0, 1\}$

$$(x_1, d_1) (x_2, d_2) \dots (x_m, d_m)$$

$x_i^o \rightarrow$ training instance

$d_i^o \rightarrow$ target value $d_i^o = c(x_i^o)$

$$\forall h \in H \quad P(h|D) = \frac{P(D|h_i) * P(h_i)}{P(D)}$$

Then HMAP

Assumptions

- 1) D is noise free $d_i = c(x_i)$
- 2) c is contained in H. $(\exists h \in H)(\forall x \in X)(c(x) = h(x))$
- 3) No need to consider any hypothesis has different prior probability

$$p(h) = \frac{1}{|H|} \quad \forall h \in H$$

$$P(D|h) = \begin{cases} 0, & \text{if } h(x_i) \neq c(x_i) \quad \forall x_i \in X \\ 1, & \text{if } h(x_i) = c(x_i) \quad \forall x_i \in X \end{cases}$$

$(d_1, d_2, d_3, \dots, d_m)$ sequence of $(x_1, x_2, x_3, \dots, x_m)$ & h holds

$$\begin{aligned} P(h|D) &= \frac{P(D|h) * P(h)}{P(D)} \rightarrow \text{Inconsistent hypotheses} \\ &= \frac{0 * P(h)}{P(D)} = 0. \end{aligned}$$

$$\begin{aligned} P(h|D) &= \frac{P(h|D) * P(D|h) * P(h)}{P(D)} \rightarrow \text{Consistent} \\ &= \frac{1 * P(h) / |H|}{|VSH,D|} = \frac{1}{|VSH,D|} \quad \text{hypotheses} \end{aligned}$$

Bayes optimal classifier

h_1, h_2, h_3 . posterior prob = 0.4, 0.3, 0.3

h_i is MAP if new instance x is classified as '+'ve by h_i , it is -ve.

let us take h_1 says -ve, $h_2 \& h_3$ says '+ve'.

Here most probable classification is different than HMAP.

Most probable classification is given as by combining prediction of all hypotheses weighted posterior hypothesis

$$p(v_j|D) = \sum_{h_i \in H} p(v_j|h_i) * P(h_i|D)$$

$$\operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} p(v_j|h_i) * P(h_i|D)$$

Naïve Bayes classifier:-

→ x is conjunction of attributes

→ $v_j \in V$

$$V_{MAP} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j | a_1, a_2, a_3, \dots, a_n)$$

$$= \underset{v_j \in V}{\operatorname{argmax}} \frac{P(a_1, a_2, a_3, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

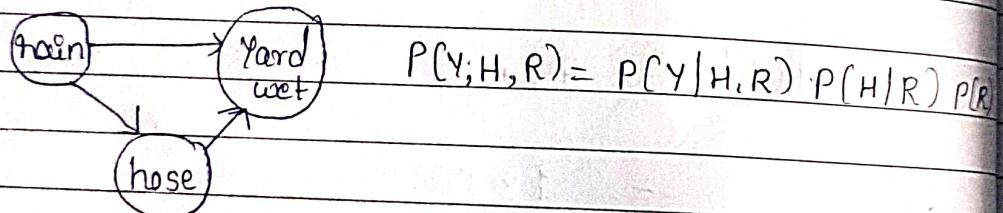
$$= \underset{v_j \in V}{\operatorname{argmax}} P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

$$V_{NB} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) \prod_p P(a_p | v_j)$$

Naïve Bayes network

→ Compact, flexible, interpretable representation of joint probability distribution

→ Acyclic directed graph - relation b/w variables

Instance based learning:-

KNN - Select 1s value

- Select 1s nearest neighbors for new datapoint

- Count no. of datapoints in each category

- Assign it to the category where got most no. neighbors

- Logy - Non parametric (functional form)

Care Based Reasoning:-

→ Intelligent system

→ experiences stored as care

→ Recall similar experience from memory

→ Reuse it for new instance

→ Save this new care

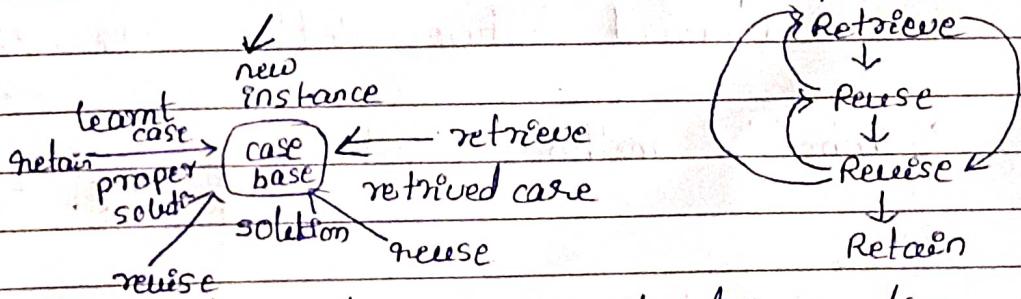
→ Similar problem has similar solution. Reuse & Retain
problem-context

- solution - pointer to other relevant cases
- solution quality assessment
 - steps to solution

Ex:- Technical diagnosis of a car fault

R4 based

- Retrieve - Check for similar case
 - similarity index, explanation method, case based.
- Repair - Transform/adapt
 - interactive, description, rule based, model, constraint
- Repair - simulation, real world
- Revise - Verify, Improve
- Retain - Add it to case base
 - similarity, case-base, solution adoption



Applications :- Help-desk, customer support, disease diagnosis, product recommendation, Text classification.

Confusion matrix :-

- $N \times N$ matrix used to evaluate the performance of the classification model, where 'N' is the no. of target classes
- compares actual & predicted target values

$$\text{Actual values} \quad \text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

		T	F	
Predicted Values	T	TP	FP	→ Type-1.
	F	FN	TN	
				↓ Type-2

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

specificity (GNR) True negative rate.

$$\text{TN} \quad \text{F1-score} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}$$

$$\text{harmonic mean} \quad \frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}$$

Accuracy - balanced data (TP & TN)

F1-score - Imbalanced data (FP & FN)

Precision - FP

Recall - FN

Feature selection:

- Selecting few features from set of features pool.
- Avoid data redundancy/overfitting / irrelevant features
- Data reduction, less storage, less complexity, less time
reduced dimension, increased accuracy.

→ Filter based method: chi-square, correlation, information gain

→ Wrapper method: features → subset → Train → Performance

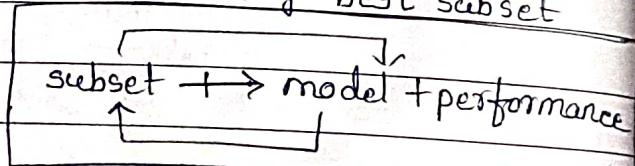
→ Expensive

→ Forward feature select^n

→ Backward feature eliminat^n

→ Recursive feature eliminat^n
selecting best subset

→ Embedded method: features →



Scikit-learn:

→ Dropping feature with 0 variance or low variance
varianceThreshold (0)

if we want 80% of 0's

$$0.8 * (1 - 0.8) = 0.16$$

varianceThreshold (threshold = 0.16)

→ Univariate feature selection - Selecting features based on
statistical ranking: chi-square test, correlation

Feature scoring, ranking, selection

→ Model based selection - ML models to select

→ Pipeline - Data preprocessing → mutual information → cross validation
→ hyper parameter tuning

Approximating the target function
 $\rightarrow y \in \mathcal{Y} = f(x)$

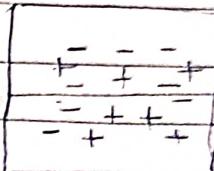
- Generalization means how well the concept learned by model can apply on new unseen data
- generalizing new data
- $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Types of error - ① Noise error

② Bias - Make assumptions

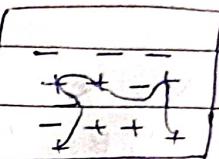
Simplest model with low variance

Underfitting

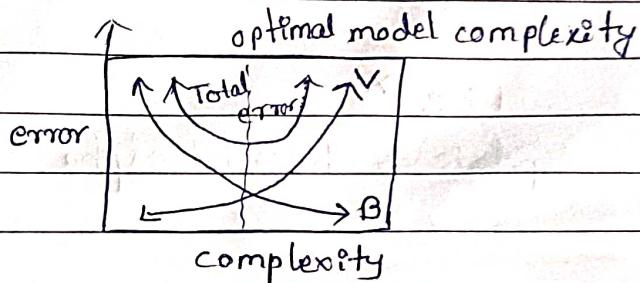
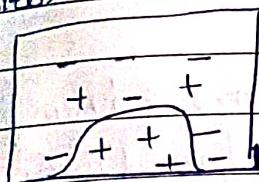


③ Variance - Fluctuations in data

- Learns too much
- Overfitting



Bias-variance tradeoff



Overfitting - Learns too much on the fluctuations, pattern, noise that hurts the model performance

- Interprets noise as patterns
- Good training performance but poor test data generalization
- Complex features

Underfitting - Neither learns nor generalize

- Performance not good
- No complexity

K-fold cross validation & hyperparameter tuning

Reducing no. of features

$$\text{Regularization} - J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Logistic regression:
 → Binary classification

Linear regression \rightarrow $y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$

$$y = b_0 + b_1x \quad P = \frac{1}{1 + e^{-y}}$$

$$\ln\left(\frac{P}{1-P}\right) = b_0 + b_1x$$

$$\text{transformed} = \frac{1}{1 + e^{-x}}$$

~~$X = b_0 + b_1x_1 + b_2x_2$~~

~~Let $b_0 = b_1 = b_2 = 0$~~

~~$\frac{1}{1 + e^{-x}}$~~

$$\text{update} \Rightarrow b = b + \text{alpha} * (X - \text{transformed}) * \text{transformed} * \frac{(1 - \text{transformed}) * X_i}{}$$

Ensemble method:-

- Learning algorithm
- Two or more algorithms are combined of similar or dissimilar types called bare learners, so that it incorporates the prediction values from all the bare learners.
- Averaging :- Taking average of the predicted values in regression models or while calculating the probability in classification models
- Majority vote :- Considering the prediction with majority vote.
- Weighted average :- Providing weight for every model and then calculating the average value. Giving high or low preference to specific models.

Bagging:-

- Error due to bias
- Error due to variance

→ Err

Low bias → ↑ flexible model

High bias → ↓ flexible model

Variance → Fluctuation in the data

↓ bias → ↑ Variance

↑ bias → ↓ Variance

We need ↓ bias & ↓ variance.

→ Averaging slightly different version of same model → ↑ performance

→ Building models in parallel on various sample data and voting on the prediction from all models.

→ Aggr. Bootstrap aggregation

→ It is a sampling technique

→ Selecting 'm' sample rows from the 'n' original dataset with replacement

→ All records are equally likely to be a part of bootstrap

→ Multiple bootstraps are created and provided as input to generate decision trees

→ Reduces the variance and overfitting

→ Increases stability & accuracy

Ex:- Random forest.

→ Probability of instance present in the bootstrap sample $(1 - \frac{1}{n})^n$.

Boosting:-

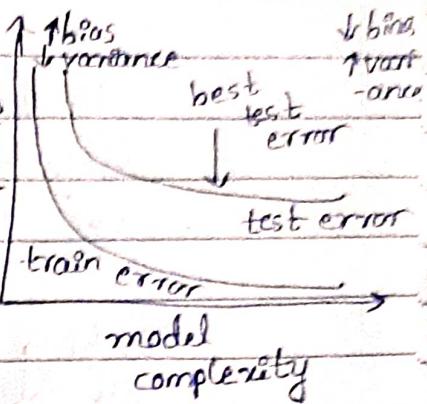
→ Models built in series

→ Weights are adjusted on what it learnt in previous model.

→ weak ones become strong ones.

→ First algorithm trained on entire dataset

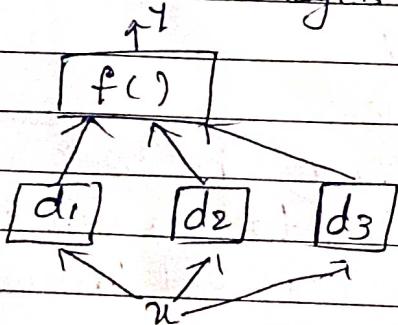
→ Weights are adjusted or increased for the poorly predicted dataset in previous model.



- It is done because one training model may not be good for entire dataset but can be good for some data. Therefore it increases the accuracy
- It boosts the overall accuracy
- It reduces the bias
- There can be the chance of model being overfitted. Therefore hyper parameter tuning is needed.
- Ex:- XGBoost, GBM, ADABoost.

Stacking :-

- Multiple layers of machine learning models placed one over another.
- Top layer model takes decisions based on output of models in layers below it.



→ fold predictions

→ Any no. of layers can be used

→ Averaging, majority vote, weighted average can be used in top layer

Advantages of ensemble method :-

- Accuracy ↑
- Stability ↑
- Robustness ↑

Disadvantages of ensemble method :-

- Time consuming
- Selection of models for creating an ensemble is difficult.