

Module - II

Regularization for Deep Learning

Data set augmentation (or) Data augmentation -

- Augmentation means adds more stuff or information to given data.
- Data augmentation is a technique used in machine learning to artificially expand the size and diversity of a data set by applying transformations or modifications to the existing data.
- An image can be transformed into number of images with few transformations.
- Data augmentation mainly used in Computer Vision, Speech Recognition and NLP.

Types of data augmentation -

1. Image data augmentation -

- It is used in Computer vision to increase the diversity.

Basic transformations -

- i) Flipping - Horizontal or vertical flip of images.
- ii) Rotation - Rotating the images by certain degree.
- iii) Cropping - Cropping the image by randomly.
- iv) Scaling - Resizing the image to different size.
- v) Brightness adjustment - Increase or decrease brightness.
- vi) Saturation - changing the colour intensity

2. Text data augmentation -

- it is used in NLP to generate various of text data.

Basic Techniques -

- i) Synonym replacement - Replace the words with synonyms.
- ii) Random insertion - adding random words from the vocabulary.
- iii) Random deletion - removing the arbitrary words from sentence.
- iv) Random swap - swapping the positions of the word.
- v) Back Translation - translating a sentence to another language and back.
- vi) Sentence paraphrasing - using AI models to rewrite sentences.

3. Audio augmentation -

- it is used in speech recognition.

Basic Transformations -

- i) Pitch shifting - changing the pitch of an audio sample.
- ii) Speed variation - slowing down or speed up speech.
- iii) Volume scaling - increase or decrease the loudness.
- iv) Time stretching - changing speed without affecting pitch.
- v) Background noise addition - adding real world noise like traffic, crowd.

Adv -

1. prevent overfitting
2. Improve model performance
3. Balances imbalanced data set

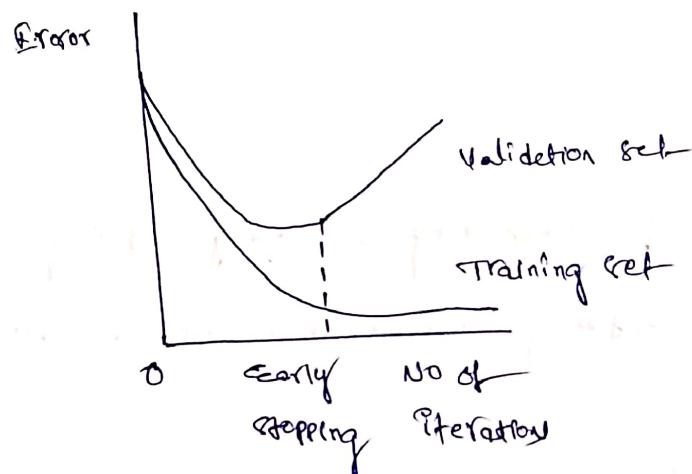
Dis Adv -

1. Computationally risk
2. not always effective
3. domain dependency

Early Stopping -

- Early stopping is a technique used in machine learning models to prevent the overfitting.
- When the model learns the training data too well and not well on new data or unseen data.
- Then it will cause best performance on training data and poor performance on validation data or testing data.

How early stopping works -



Step 1 - Training and validation sets -

- During the training, the data can be divided into two sets,
 - i) Training set - It is used to train the model,
 - ii) Validation set - It is used to evaluate the working performance of the model during training.

2. Monitoring the model performance - When the model is training

monitoring the model performance on both training set and validation set.

3. Identifying overfitting -

- initially the model performance was good in both sets. At a certain point the model may start to perform worse on the validation set and continue the performance well on training set. This difference indicates the overfitting.

4. Stopping Training -

- In the early stopping, the training is halted when the model performance reduced on the validation set. It means the model parameters captured where the model best performs on new data.

Advantages -

1. Prevention of overfitting - early stopping helps to prevent the overfitting and enhance its generalization to unseen data.
2. Computational Efficiency - early stopping reduces the computational resources like time and resources.
3. Automated Model Tuning - It provides an automated method to determine the optimal number of iterations and reducing the manual tuning.
4. Improved Generalization - early stopping enhances the model ability to generalize the unseen data or new data.
5. simple and easy implementation method.

Dropout —

- dropout is a technique used in neural networks used to prevent the overfitting.

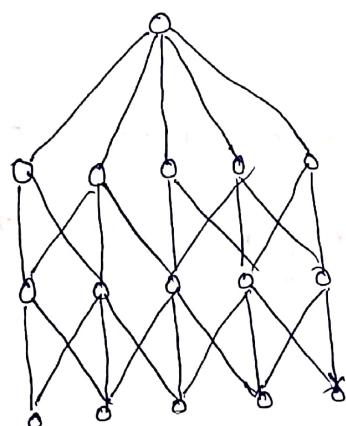
When overfit occur —

- When the model learns training data too well and not well on

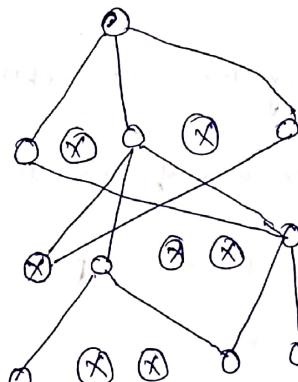
unseen data or new data.

- Then it will cause bad performance on training data and poor performance on validation data or testing data.

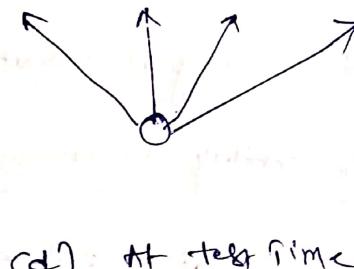
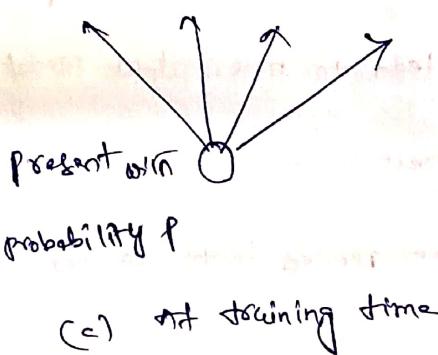
How Pt works —



(a) fully connected
neural network



(b) After applying the
dropout method



1. Randomly Dropping The Neurons -

During training, dropout method selects the fraction of neurons randomly in the network and temporarily drops them i.e., those neurons are ignored in training pass. This is done for each iteration.

2. Combining The Dropped Neurons -

- when a neuron is dropped then its all incoming and outgoing connections are also removed for that iteration.
- this will creates new and smaller network from the feed original one.

3. Combining The Sub-networks -

- At the end of the training, the model combine all these sub networks which will help to reduce the overfitting and improves the model ability.

Advantages -

1. preventing the overfitting

2. improves the model performance

3. Efficient model averaging - at the end of the training

model combining all sub networks, it will lead to more robust model.

4. Applicability to various network architecture -

- It's versatile architecture it can be applied in to various networks like CNN, RNN etc.

5. Reduced the inter-neuron dependency.

Parameter sharing and parameter Typing

- Generally parameters play an important role in machine learning models.
- When training the deep learning models managing the parameters is crucial step.
- Basically parameters are "weight" and "bias". The model learns these parameters from the data. These parameters are helpful to make the predictions.

Parameter sharing and parameter Typing are the two techniques that helps to reduce the no. of the parameters and making the models more efficient and faster.

$$y = w_1x + b$$

$y \rightarrow$ output

$x \rightarrow$ input

$b \rightarrow$ bias

$w_1 \rightarrow$ weight

Advantages -

1. Reduce the memory size
2. Improve the performance of a model
3. Speed up the training process
4. Maintain the Consistency in model behaviour.
5. Reducing the overfitting problems.

1. Parameter Sharing -

- It was same set of parameters are used in multiple parts of the model.
 - Parameter sharing mostly used in CNN.
- Ex - In a CNN, an image is broken into small no. of parts.
- A single filter i.e kernel is scans the all small no. of parts of the image like (top-left, centre, bottom-right etc)
 - Instead of using different filters for each part, the same filter is used for all regions or parts of the image.

2. Parameter Tying -

- Here multiple parts containing equal weight i.e, each part contains equal weight.
- Parameter tying is mostly used in RNN.
- In RNN, same set of weights are used for each step.

Importance of RNN -

- RNN processes time series data at each step, at time.
- At each time step, the same weight matrix is used for transforming the input.
- maintains consistency across different time steps.

```

EX - import torch
       import torch.nn as nn
       class Simple CNN(nn.Module):
           def __init__(self):
               super(Simple CNN, self).__init__()
               self.conv = nn.Conv2d(in_channels=1, out_channels=1,
                                   kernel_size=3)
           def forward(self, x):
               return self.conv(x)
       image = torch.rand(1, 1, 5, 5)
       model = Simple CNN()
       output = model(image)
       print("Input shape:", image.shape)
       print("Output shape:", output.shape)

```

output -

Input shape : torch.size([1, 1, 5, 5])

Output shape : torch.size([1, 1, 3, 3])

Adversarial Training -

- Adversarial Training is a technique used in machine learning models to make stronger against the tricky attacks.
- These attacks use invisible changes in the initial input data to fool the model into making the wrong predictions.

Ex - Imagine a model that recognizes animals. If we change a few pixels in a cat image (so small that human can't notice) The model might mistakenly file a dog.

How Adversarial works -

1. We create these tricky examples like Adversarial Example.
2. We train the model using both normal and tricky examples.
3. The model learns to ignore the tiny changes and make correct predictions.

Types of Adversarial Attacks -

i) FGSM (Fast gradient sign method) -

It's a quick way to add tiny changes to fool the model.

ii) PGD (Projected Gradient descent) -

It's very strong multi attack to add small or minor changes.

iii) CEW (Carlini & Wagner) -

It's very advanced trick to fool the machine learning models.

Adversarial Training Techniques -

To protect models we use some defensive techniques.

i) Basic Adversarial Training -

This technique train the model on both normal and tricky examples.

ii) PGD Based Training -

It's advanced method to train the model using stronger attacks.

iii) Trades - it is a method that balances accuracy and robustness.

challenge in Adversarial training -

1. Take more time
2. Accuracy drops
3. Attackers keep improving

importance of Adversarial training -

1. Cyber security - It protects AI systems from attackers.
2. Self-driving Cars - Stop small changes in traffic signs from fooling the car.
3. Facial Recognition - It prevents the attackers from tricking security systems.

Parameter Initialization strategies -

- Parameter initialization means weight initialization i.e., adding some weight to the each neuron in deep learning.
- Proper weight initialization will help to
- i) Avoiding vanishing gradients i.e., when gradients become too small then stop the learning process.
- ii) Avoiding exploding gradients i.e., when gradients grow too large training process make unstable.
- iii) Breaking symmetry i.e., it ensures that different neurons learns different features.
- iv) Speed up Convergence i.e., The model learn faster.

Types of Initialization —

1. Zero Initialization — Formula

$$W = 0$$

problem — If all weights are initialized as zero, then every ^{neuron} in a layer will make same output and gradient during back propagation, i.e. all neurons learn the same thing, making the network useless.

When it will use —

- Never used for weight initialization, but in some cases used for bias initialization because bias start with '0',

2. Random Initialization —

$$W = \text{random small values}$$

- It avoids the problem of symmetry.
- if the values are too large then the activations explode.
- if the values are too small then learning slows down.

When it use —

- Random values can efficiently work in simple networks only, but not work properly in large networks.

3. Xavier Initialization —

$$W = N(0, \frac{1}{n^{in}})$$

where n^{in} is the no. of input neurons

- It keeps the variance of activations consistent in layers.
- It prevents the activations from becoming too small or too large.

When it use -

- It is used for sigmoid functions not for ReLU functions.

4. He Initialization -

$$W = N(0, \frac{2}{n^{in}})$$

- It can be helps to remove the negative values.
- It prevents the ReLU problem i.e. The neurons are stops learning.

When it use -

- It is used for ReLU and Leaky ReLU functions.
- It is suitable for all deep learning networks.

5. LeCun Initialization -

$$W = N(0, \frac{1}{n^{in}})$$

- It is designed for specifically for SELU functions.
- It ensures the activations remain self-normalizing.

When it use -

- It is used for SELU functions.
- It supports all deep neural networks very well.

Noise Robustness

- Noisy data means it contains irrelevant information, missing features, errors, outliers etc. These information will lead to noisy problem.
- Noise Robustness is machine learning model ability that handle the noisy data without reducing the performance of a model.

Types of Noises -

1. Feature Noise (Input Noise) -

- It will occurs when input variables are incorrect or missing.

Ex - In a house prediction problem, if the no. of rooms is added incorrectly the model learns wrong patterns.

Causes -

1. measurement errors.
2. data corruption
3. Human mistakes.

2. Label Noise (Output Noise) -

- It will occurs when output variables are incorrect or missing.

Ex - In an image classification task, an image of cat is mistagged as dog.

Causes -

1. Human Errors.
2. sensor malfunction
3. misclassification in data.

3. Adversarial Noise -

- It will occurs when intentionally added some noise to the original data it leads to mislead the model.

Ex - changing a few pixels in image it will lead to mislead.

Causes -

1. Malicious Attacks
2. Security Threat Attacks.

4. Environmental noise - External factors that effect the data collection.

Ex - Background noise in speech Recognition.

Drawbacks of Noise -

1. Increasing the training time.
2. Reduces accuracy.
3. Leads to overfitting problems.
4. Bias-variance trade-off.

Methods to Improve Noise Robustness -

1. Data Preprocessing Techniques -

- i) outlier detection & Removal - It detect and remove the extreme values like Z-score method and IQR method.
- ii) data augmentation - It will add some variations to training data like rotation, flipping, colour changes etc).
- iii) missing data values - It will replace the missing values by adding mean or mode.

2. Robust model architecture -

- i) Regularization Technique - It will helps to preventing the overfitting problems by L_1 and L_2 Regularization techniques.
- ii) dropout method - It randomly drops the neurons during the training period to prevent overfitting problems.
- iii) ensemble learning - It will combine all multiple weak models to make stronger and improves the stability.

Ex - Random Forest, Bagging

3. Noise Robust Training

- i) Early stopping — It will stop the learning process when the loss of functions becomes higher.
- ii) Adversarial training — Here train the model with intentionally added noise, to improve the robustness.
Ex — Add small changes in an image for better appearance.
- iii) Data weighting method — By assigning weights to each sample based on the importance of the sample.

Applications —

1. Speech Recognition
2. medical diagnosis
3. finance
4. Autonomous Vehicles

Multitasking learning —

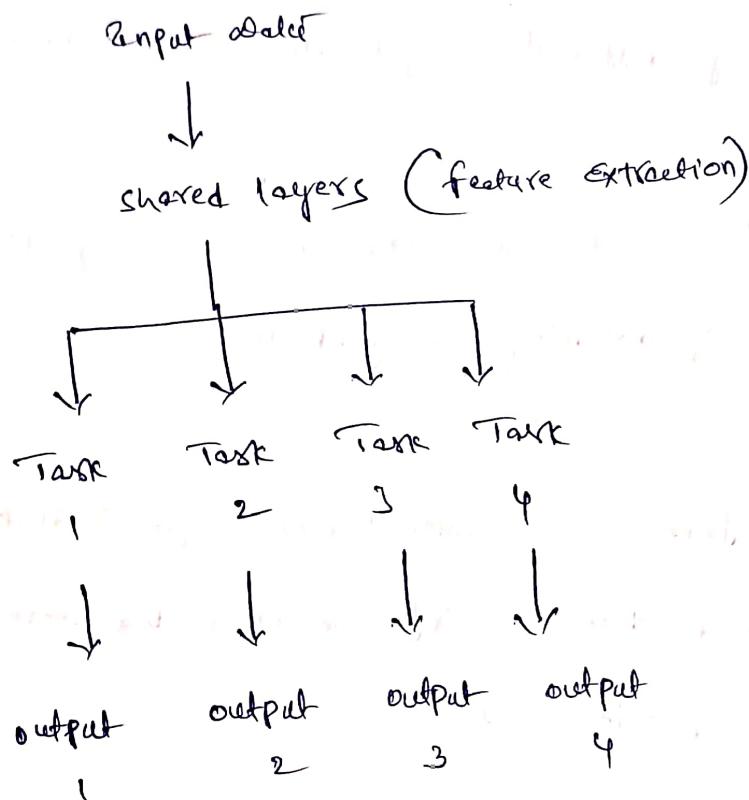
- multitasking learning is a machine learning technique where a model is trained to perform multiple relative ~~task~~ simultaneously.
- Generally for each task different model will be trained, but in multitask learning for ~~each~~ tasks we use one better trained model.

Advantages —

1. Improves generalization
2. Reduces overfitting
3. Efficient Learning
4. Handles data scarcity

How multitasking works - Basically it works based on 2 types.

1. Hard parameter sharing - Here we are using a shared network with separate output layers for each task.
 - It will help to reduce the overfitting.



2. Soft parameter sharing -

- Here for each task it has its own model but parameters are regularized to be similar.

Applications of Multitasking learning -

1. Computer Vision - for detecting the objects
2. NLP - for sentiment analysis, entity Recognition
3. Health Care - diagnosing the multiple diseases
4. Autonomous Vehicles - for object detection, predictions.

Multitasking algorithms —

1. multitask-neural network (MTNN) —

- it uses shared layers for feature extraction and separate layers for each task.

2. uncertainty based weighting —

- giving different weight to the each task based on the importance.

3. Grad Norm — it will balances the training process across all tasks by adjusting the gradients.

challenges —

1. Task interference — some tasks might conflict.

2. Unbalanced data — if one task has more data, the model can prioritize it.

Sparse Representations —

- Sparse Representation is the way of encoding the data where the most elements are zero or not active. This is the opposite of dense representation.

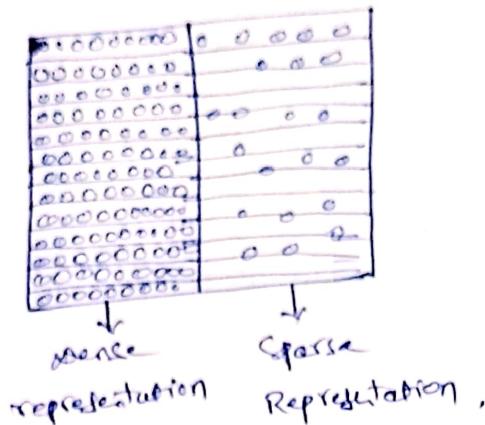
- Instead of representing a signal with all basis vectors,

- sparse representation uses only a few significant ones.

Why sparsity is need —

1. To improve the efficiency
2. To improve Robustness
3. To maintain interpretability

Dense vs Sparse representation



- ### Dense Representation —
1. The grid is mostly filled with black cells.
 2. most of the elements are non-zero values.
 3. It uses all features.

- ### Sparse Representation —
1. The grid is filled with few black cells.
 2. most of the elements are zero values.
 3. It uses a few features.

How to achieve Sparse Representations —

1. L_1 Regularization — It encourages the sparsity in models.
2. Dictionary learning — It learns a set of functions where data can be represented sparsely.
3. Auto encoders — It encourages the hidden layers to have sparse activations.
4. Thresholding — Set small values to zero.

- ## Applications —
1. Computer Vision
 2. NLP
 3. Neuro Science

Tangent distance, Tangent prop and manifold Tangent classifier

1. Tangent distance — It is a method used in image recognition to compare two images in a more flexible way than just pixel-to-pixel comparison.
 - Instead of directly comparing two images we check how much we need to transform one image to match another image.

Why it's useful —

- Small changes in position, rotation can make two images look different in pixels but they still represent same object.
- Tangent distance allows us to compare two images based on their shape not in pixel values.

How it works —

1. find Tangent vectors — Identify small transformation like shifting, rotation and scaling that don't change the image meaning.
 2. measure distance — Compute the shortest distance between the transformed versions of the two images.
 3. Better matching — If the two images belong to the same object, their transformed versions should be close.
2. Tangent propagation —
- It's a technique used in neural networks to make models to learn invariant features i.e., the model recognize if few transformations are applied in an image.

- Instead of just training the neural network on normal data, we also train it to recognize slightly transformed versions of the data.
- why it is useful -
 - it helps the model ignore small changes like rotation or brightness.
 - it improves the generalization i.e. the model works well on new unseen data.

How it works -

1. Compute Tangent vectors - find small changes like rotation, shifting
2. modify Training - Instead of just learning from raw inputs, the model also learns from these transformed versions.
3. Better prediction - the network becomes more robust and makes better predictions.
3. Tangent manifold classifiers -
 - it is a method used in deep learning to make neural networks learn better representations of data by considering the underlying structure of the data.
 - These data points can be lie on high-dimensional surface in high-dimensional space. We want our model to learn this structure so it can generalize better.

why it is useful -

- data exists in high dimensions, but
- it helps the model focus on meaningful variations and ignore unnecessary noise.

How it will work -

1. Find the manifold - Identify small changes that do not change the data meaning.
2. Train with Tangent Constraints - make the neural network understand and respect these small changes.
3. Classify based on the manifold - Instead of treating each point separately the model classifies data based on its location on the manifold.

Learning vs pure optimization -

Learning is when a model improves its performance by understanding its patterns from the data. It generalizes to new situations instead of just memorizing.

Ex - Training data \rightarrow Learning algorithm \rightarrow model generalized
(cats) (Neural Network) (recognize new cats)

features -

1. Uses past experience with training data.
2. Generalizes to unseen data.
3. Tries to minimize the errors while still understanding patterns.

When to use learning -

1. You need to generalize.
2. The data is complex and has patterns.
3. You want to predict future outcomes.

Pure optimization — Pure optimization is a method to finding the best possible solution to the given problem without worrying about the learning pattern.

Ex — Solving a math equation find x that makes

$$3x + 5 = 25$$

objective function \rightarrow optimization algorithm \rightarrow best solution

(minimize cost)

(gradient descent)

(lowest error)

features —

1. No generalization if focus only the specific problem.
2. It will focus on minimize or maximize the function.
3. Does not improve over time it only finds best possible solution.

differences —

feature

Learning

Pure optimization

1. Goal

understand the pattern

find the best solution for

and generalities

specific problem.

2. Memory of

Yes, remember it

No, only solves once

data

Improves over time

3. Example

Recognize cat in images

finding the shortest

path in a maze

4. Generalization

Yes, works on new data

No, only for the given problem

5. Method

using training data and adjust parameters

Uses formulas or algorithms to min/max the function.

challenges in neural network optimization

- Neural network optimization is a process of improving the performance of a neural network by adjusting the parameters like weight and biases.
- The main goal of optimization is to reduce the loss function and improves the accuracy.
- However optimizing neural network is difficult due to several challenges.

1. Vanishing & exploding gradients

what happens -

- When training deep neural network the gradients (updated to weights) becomes too small or too large.
- This may learning slow or cause the model fail.

Solution -

- Use ReLU function instead of sigmoid activation function.
- Use batch normalization to stabilize the updates.
- Use gradient clipping to prevent gradient exploding.

2. Getting stuck in local minima

what happens -

- during training The algorithm tries to find best weights by minimize the loss function.

- Sometimes it will get stuck in local minimum i.e small dip instead of find global minimum (deepest point)

Solution -

- use optimizers like Adam, RMSprop to escape local minima.
- use learning rate scheduling to make learning process dexter.
- Try Random weight initialization to avoid bad starting points.

3. overfitting (Poor Generalization) -

What happens -

- The model performs well on training data but fails on unseen data. It memorizes the training example instead of learning pattern.

Solution -

- use Regularization techniques L_1 , L_2 , dropout to prevent overfitting.
- increase training data with data augmentation.
- use early stopping to reduce overfitting.

4. choose right learning rate -

What happens -

- if the learning rate is too high the model jumps around and never converges.
- if the learning rate is too low learning is slow and may get stuck.

Solution -

- use learning rate scheduling to reduce learning rate over time.
- use Adaptive learning rate Adam, RMSprop to adjust the learning rate dynamically.

5. Hyper Parameter Tuning -

What happens -

- choosing the right batch size, learning rate, number of layers, number of neurons is very tricky.
- A wrong combination can make the model performance poor.

Solution -

- use Grid Search to test different values.
- Use Bayesian optimization for automated tuning.

Semi Supervised Learning -

- Semi Supervised learning is type of machine learning algorithm where we use both labeled and unlabeled data to train a model.
- It is in middle of the both supervised and unsupervised learning algorithms.

Why we need semi supervised learning -

- Labeled data is expensive and time consuming to collect.
- Unlabeled data is too large and harder to use.
- Semi supervised learning helps by using small amount of labeled data and large amount of unlabeled data to improve accuracy.

How it will work -

- first, the model learns from the labeled data.
- Then it uses these patterns to make predictions on the unlabeled data.
- finally it learns from both to improve overall performance.

Techniques used in semi-supervised -

1. Self-Training - The model predicts labels for unlabeled data, then reuse those predictions for training.
2. Consistency regularization - The model makes small changes to the data like rotation, flipping and expects the prediction to stay the same.
3. Graph Based learning -
 - data points can be treated as nodes in a graph
 - labeled points help spread knowledge to unlabeled points that are similar.
4. Pseudo-labeling -
 - The model assigns fake labels to unlabeled data and learns from them.

Advantages -

- uses less labeled data and saves the money.
- performs better than unsupervised learning method.
- it can work with large data sets.

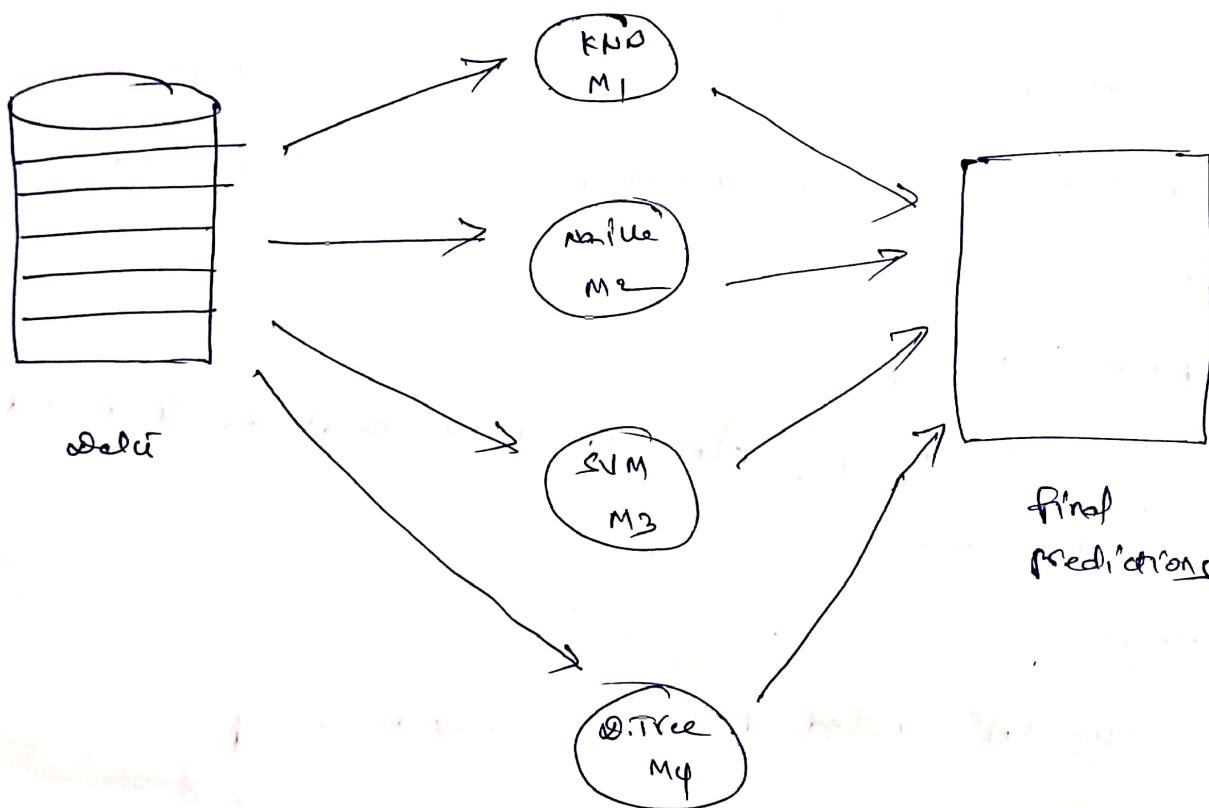
Disadvantages -

- wrong labels can reduce the accuracy.
- more complex than supervised learning.
- not always reliable, it depends on the how unlabeled data matches with labeled data.

Ensemble learning methods -

Bagging and Boosting -

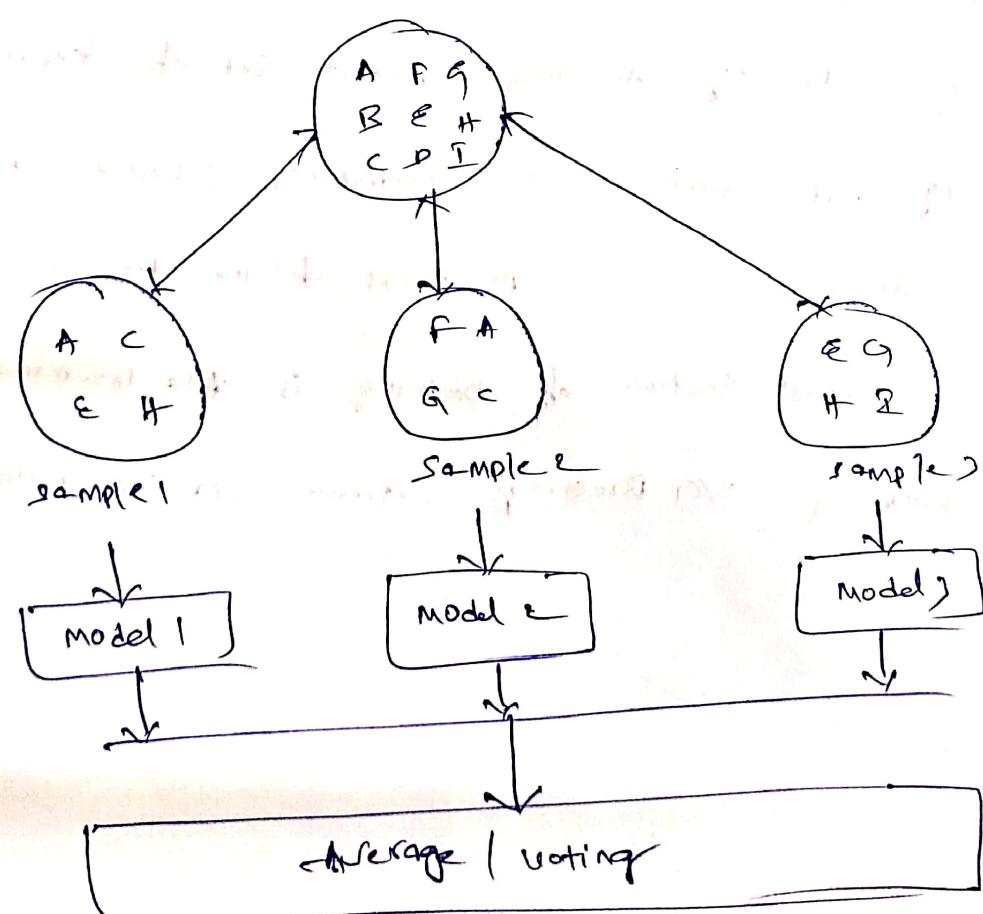
- Ensemble method is a technique that combines all the predictions from the multiple learning algorithms to make more accurate predictions than individual model.
- A model is comprised by many model is called an ensemble model.



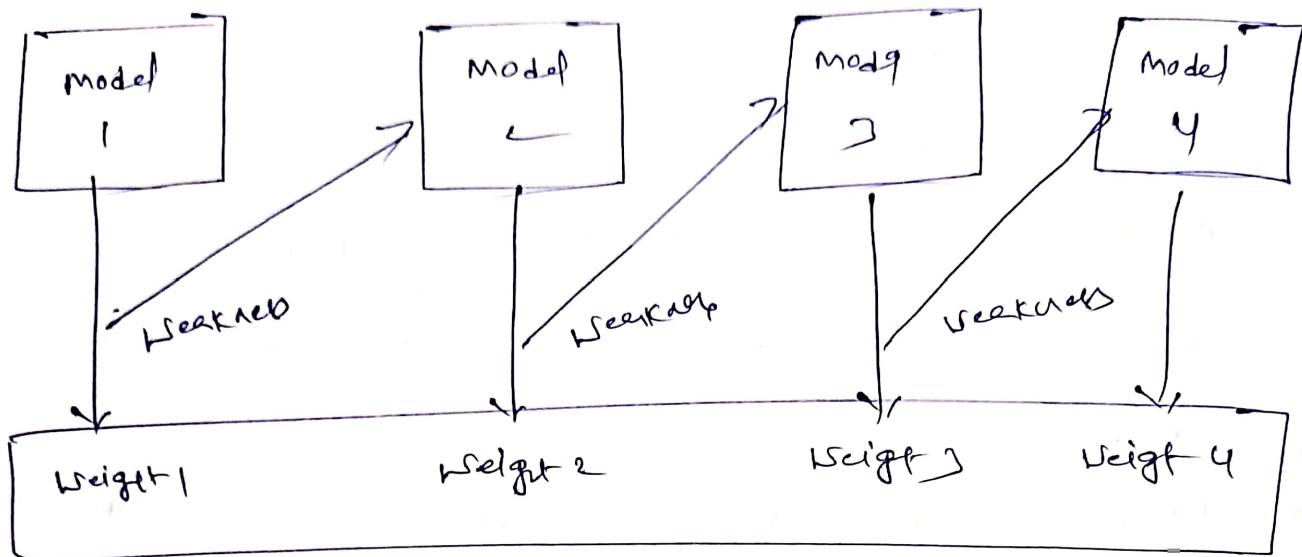
- Many ensemble methods contains same type of learning algorithms is called Homogeneous ensemble and some methods contains different types of learning algorithms is called Heterogeneous Ensemble learning methods.

1. Bagging - (Bootstrap Aggregation)

- Bagging is a procedure it is used to reduce the high variance of an algorithm i.e., reduce the overfitting problems.
- Bagging makes each model to run independently and aggregates the output at the end without preference to any model.
- Bagging is performed by two factors i.e., Bootstrap and Aggregation.
- Bootstrap means suppose a data set contains "n" samples we take "m" samples and trained with an model.
- Aggregation means after training data combines the predictions and take majority or mode or median.



2. Boosting -



Ensemble with all predictions.

- Boosting is an algorithm that converts weak learners to strong learners. Here base learners are trained sequentially. Then the next learner is trying to reduce the errors and updating the parameters and compared to previous learners.
- Boosting works by training different set of learners sequentially and combining the predictions, where later learners focus more on the mistakes of the previous learner.
- Most important feature of boosting is the Gradient boosting, XG Boosting, extreme gradient boosting.

GP