

28-02-'19

LEARNING FROM EXAMPLES

N-I-V: LEARNING PROBABILISTIC MODELS

I-T = LEARNING WITH HIDDEN VARIABLES

→ Learning From Examples

→ Forms of learning

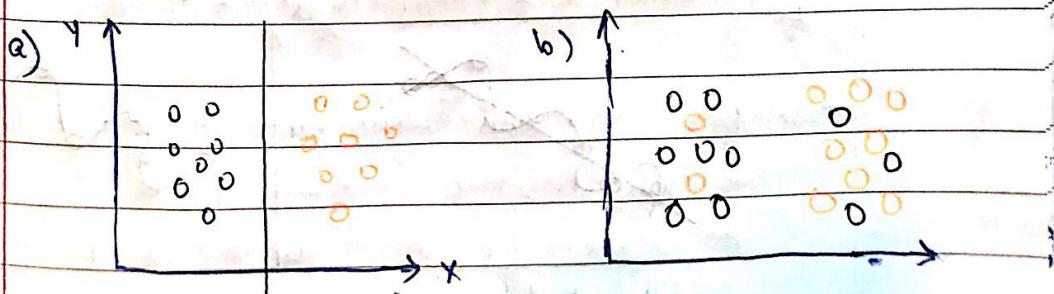
- Any agent can improve its learning by observing the data.
- We have to consider four factors for an agent to learn i.e., to improve the agent:
 - (i) Which component is to be improved in the agent?
 - (ii) What prior knowledge does the agent already have?
 - (iii) What representation is used for data and components in the agent?
 - (iv) What feedback is available to learn from for the agent?
- The agent has to learn six components:
 - c₁ (i) A direct mapping from conditions on current state to actions.
 - c₂ (ii) A means to infer relevant properties of the world from percept sequence.
 - c₃ (iii) Information about the way the world evolves and about the results of possible actions the agent can take.
 - c₄ (iv) Utility information indicating desirability of world states.
 - c₅ (v) Action value information indicating the desirability of actions.
 - c₆ (vi) Goals that describe classes of states whose achievement maximizes the agent's utility

- Representation and Prior Knowledge

(i) Factored Representation

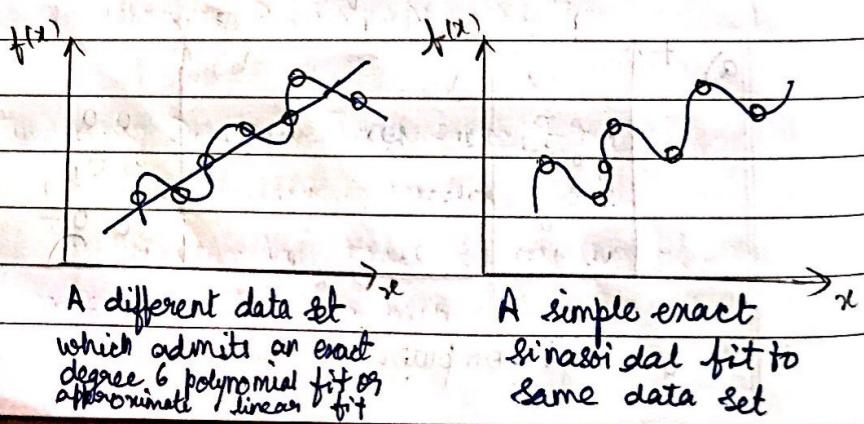
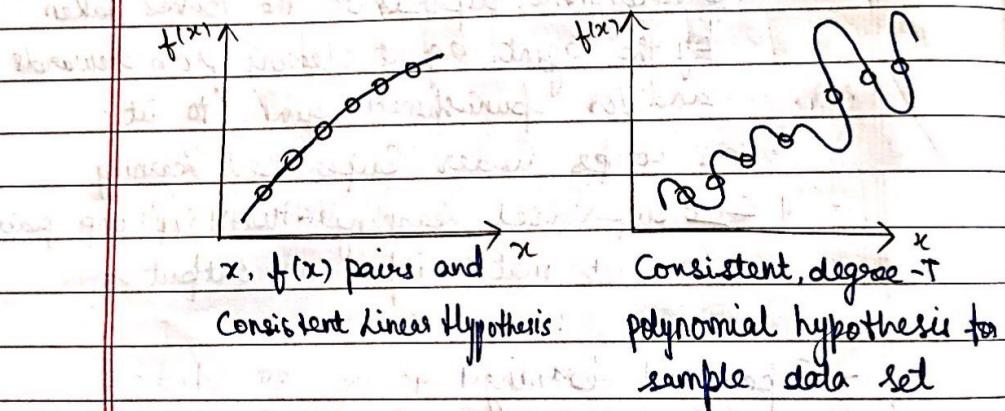
- We have a vector of attribute values which has both input values and output values.
- The values may be discrete or continuous.
- Analytical / Deductive Learning: From known general rule to new rule.
- Inductive Learning: Learning general rules with specific and input and output we get next state.
- Unsupervised Learning - Learn without any feedback
- ^{semi}Supervised Learning - Learn from labels and label unlabelled instances.
- Reinforcement Learning - Learn from the environment, depends on the moves taken by the agent, agent learns from rewards and loss punishments given to it.
- c1 comes under Supervised Learning
- ^{semi}Supervised Learning - With i/p, o/p pairs, we must match input to output.

- Supervised learning



0 → buys computer 0 → doesn't buy computer

- (a) requires a polynomial function with one degree \rightarrow simple.
- (b) requires a polynomial function with six or seven degrees \rightarrow complex.
- Supervised Learning:
N example i.i.p - DIP pairs
 $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
 $y = f(x)$
- h-hypothesis
- learning - Agent searches from the search space such that it performs well.
- We have to consider a simple hypothesis so that the result is accurate.
- Types of hypothesis



- We have to check how probable the hypothesis is to that data.

$$h^* = \operatorname{argmax}_{h \in H} P(h \mid \text{data})$$

By Baye's Rule

$$h^* = \operatorname{argmax}_{h \in H} P(\text{data} \mid h) P(h)$$

→ When the degree is low, the prior probability is high.

→ For a high degree polynomial function, the prior probability is low (less priority).

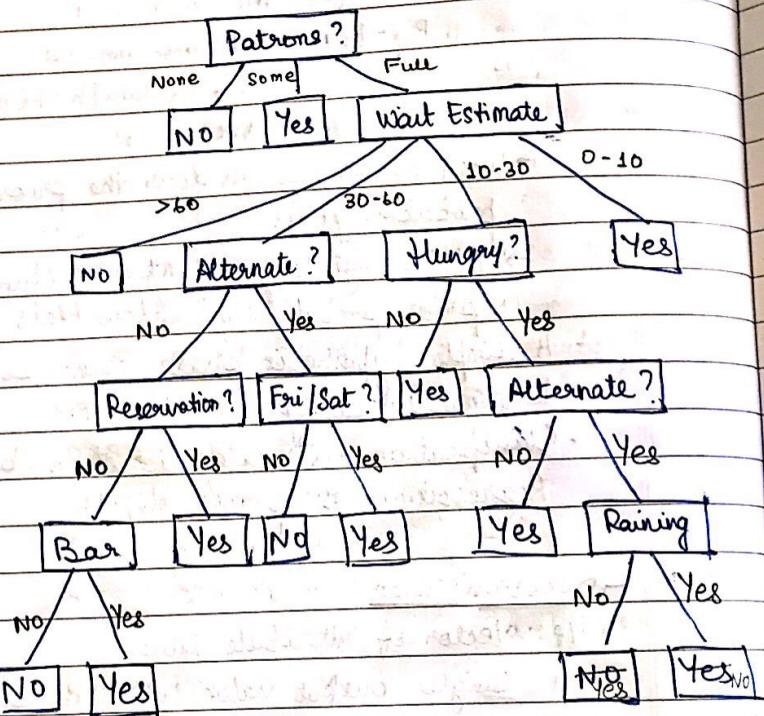
- A simple hypothesis phase can be terminated easily.
- Classification is categorization, but Regression is prediction.

→ - Decision Tree

- i/p: Vector of Attribute Values
- o/p: Single output value with one decision
- All attribute values are given as input; it classifies and gives a decision as an output.
- Boolean Classification is done here.
- Eg: Build a decision tree to decide whether to wait for a table at a restaurant

→ Attributes: Alternate (other restaurants nearby), Bar (available or not), Friday / Saturday, Hungry, Patrons (ppl in restaurant), Price (of restaurant), Raining (outside), Reservation (done or not), Type (Type (French, Italian, etc.)), Wait Time (00-10, 10-20, 20-30, 30-40 min)

- Ignore 'price' and 'type' attribute while constructing the decision tree.
- Decision Tree (Boolean):



01-03-'19

- Expressiveness of Decision Trees

$$\text{Goal} \iff (\text{Path 1} \vee \text{Path 2} \vee \dots)$$

- Inducing Decision Trees from Examples

- We have positive examples and negative examples.
- The decision tree must be as simple as possible and we must get a path.

- So, we go for a decision tree learning algorithm such that it finds a decision tree which is as simple as possible, so that the depth is minimum.
- It considers the attributes that are of first priority. It eliminates / does not consider useless attributes, and ignores them so that we can construct a simple decision tree.

- If all the examples that we have are positive (e.g. an unfair coin), we get a decision and we don't need to construct a tree.

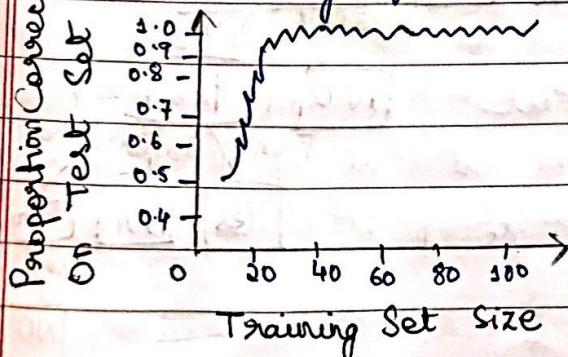
- If all the examples are empty, then return plurality value; i.e., most common output values among set of values (tie-breaker).

- If all examples have same classification, then return that type of classification.

- If the examples are half positive and half negative, we must consider the best attribute that can be considered.

- If there are no attributes in our decision tree, it means that there is some error or noise in the data, so we cannot construct the decision tree.

- Accuracy of Decision Tree Learning Algorithm



* Accuracy is high since each and everything is done 20 times.

-Decision Tree Learning Algorithm

function Decision-Tree-Learning (examples, attributes,

parent-examples) returns a tree

if examples is empty return PLURALITY-VALUE
(parent-examples)

else if all examples have same classification then
return the classification

else

$A \leftarrow \text{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$

tree \leftarrow a new decision tree with root test A

for each value V_k of A do

exs $\leftarrow \{e : e \in \text{examples} \text{ and } e_A = V_k\}$

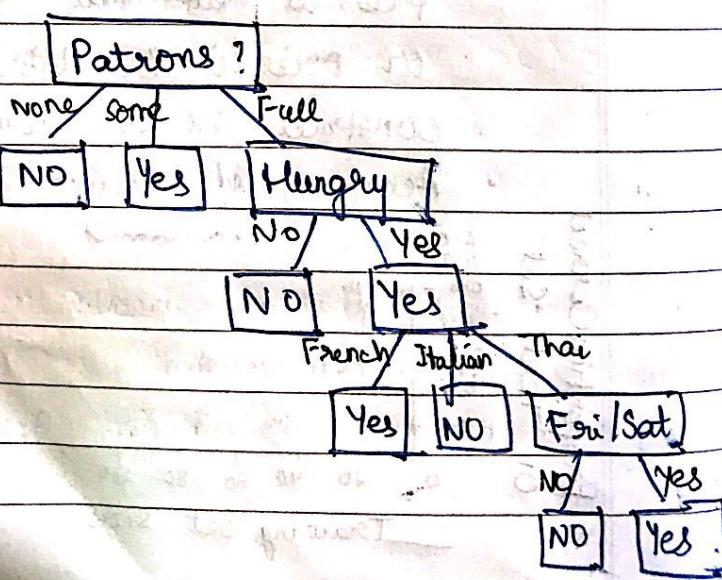
subtree \leftarrow Decision-Tree-Learning

(exs, attributes - A, examples)

add a branch to tree with label ($A = V_k$)

and subtree

return tree



- Choosing Attribute Tests

- We must consider only important attributes in the decision tree such that it becomes simple.
- We do this using entropy and calculating information gain.
- Entropy is randomness.
- For a fair coin, entropy = 1.
- For an unfair coin, entropy = 0.
- Entropy is the uncertainty of a random variable.
- Entropy : $H(V) = - \sum_k P(V_k) \log_2 \frac{1}{P(V_k)} = - \sum_k P(V_k) \log_2 P(V_k)$

$$H(Fair) = -(0.5 \log_2 0.5 + 0.5 \log_2 0.5) = 1$$

$$H(\text{loaded}) = -(0.99 \log_2 0.99 + 0.01 \log_2 0.01) \approx 0.08 \text{ bits}$$

$\approx 0.08 \text{ bits}$

- c) Entropy of a Boolean Random Variable

$$B(q) = -[q \log_2 q + (1-q) \log_2 (1-q)]$$

Goal Attributes

$$H(\text{Goal}) = B\left(\frac{P}{P+n}\right)$$

To know imp
attributes to
be chosen
for decision
tree and
remove
unnecessary ones

→ For 'k' distinct values

$$B\left(\frac{P_k}{P_k + n_k}\right)$$

- Randomly chosen example from training set

$$(P_k + n_k) / (P + n)$$

$$\bullet \text{Gain}(A) = B\left(\frac{P}{P+n}\right) - \text{Remainder}(A) \quad [\text{Information Gain}]$$

$$\text{where } \text{Remainder}(A) = \sum_{k=1}^d \frac{P_k + n_k}{P + n} B\left(\frac{P_k}{P_k + n_k}\right)$$

- Generalization and Overfitting

- When there is no target function, we get all the attributes in the decision tree and it will be overfit.
- For an overfit decision tree, we have to cut off unnecessary branches called as pruning.
- We check from the leaf and prune all unnecessary test nodes, until we find a decision tree that is relevant and simple.
- How large should Information Gain be?
 - Start with NULL Hypothesis.
 - If the degree of deviation is less than 5% of probability, then NULL hypothesis is accepted / accurate.
 - Else, it is rejected.
- The degree of deviation can be calculated as:

$$\hat{P}_K = P \times \frac{P_K + \bar{P}_K}{P + \bar{P}} \quad \hat{n}_K = n \times \frac{P_K + \bar{P}_K}{P + \bar{P}}$$

$$[P = +ve, \bar{P} = -ve]$$

Total Deviation:

$$\Delta = \sum_{k=1}^d \frac{(P_k - \hat{P}_k)^2}{\hat{P}_k} + \frac{(\bar{P}_k - \hat{\bar{P}}_k)^2}{\hat{\bar{P}}_k}$$

- Early Stopping: Stop generating that particular node if it is not required.
- If an attribute is not good enough, we should not split it again and early stop generating that particular node.

- Broadening the Applicability of Decision Trees

• Issues in Decision Trees:

(i) Missing Data

(ii) Multivalued Attributes

(iii) Continuous and integer valued input attributes

(iv) Continuous valued output attributes

• Missing Data

→ Examples may be unrecorded/expensive to record.

→ How to classify examples?

→ How to calculate information gain for missing data?

• Multivalued Attributes

→ Many possible values for an attribute.

→ Inappropriate values for information gain.

→ Check each and every value, gain ratio.

• Continuous and Integer valued input attributes

→ Split the attributes for continuous attributes.

→ But it is costly.

• Continuous valued output attributes

→ When we want to get an integer output such as predicting the value of a flat, we go for regression rather than classification.

- Evaluating and Choosing the Best Hypothesis

- It connects from the past to the future.
- It may take the past evidences, but a hypothesis must be selected in such a way that the future is best fit.
- The examples that satisfy these assumptions are called independent identically distributed (i.i.d), which connect the past to the future.
- $P(E_j | E_{j-1}, E_{j-2}, \dots) = P(E_j)$
- $P(E_j) = P(E_{j-1}) = P(E_{j-2}) = \dots$
- If we give already seen evidences in test set, then the accuracy is less. We need give unseen evidences. \Rightarrow Best Fit.
- Hold-out cross validation: We give the test set from the training set, 100% accuracy.
- K-fold cross validation: Divide into 'k' subsets, at each i-th iteration give $\frac{1}{k}$ as test set, do k rounds such that the accuracy is more.

Popular values for $k = 5, 10$

- Leave One Out Cross Validation to get a more accurate point.
- Peeking: Using test set performance to both choose a hypothesis and evaluate it. We wait till the time the lock is released to validate the data.

- Model Selection Complexity vs Goodness of Fit
 - We cannot estimate the accuracy of the model based on the degree of the polynomial.
 - Model Selection and Optimization is finding the best hypothesis within the given space.
 - We have a wrapper. It takes a decision tree learning algorithm as a parameter. The wrapper enumerates models depending on the parameters and size. It uses cross validation and it checks the average error rate on training and test set.
 - The training set error decreases when the validation decreases.
 - The model can be fit when the error rate is less.
 - We decrease the size of the dataset.

- From error rates to loss

→ Generalization loss for hypothesis

$$\text{GenLoss}_L(h) = \sum_{(x,y) \in E} L(y, h(x)) P(x, y)$$

→ Minimum Expected Generalization Loss

$$h^* = \underset{h \in H}{\operatorname{argmin}} \text{GenLoss}_L(h)$$

→ Empirical Loss

$$\text{EmpLoss}_{L,E}(h) = \frac{1}{N} \sum_{(x,y) \in E} L(y, h(x))$$

$$\hat{h}^* = \underset{h \in H}{\operatorname{argmin}} \text{EmpLoss}_{L,E}(h)$$

- Regularization

$$\text{Cost}(h) = \text{EmpLoss}(h) + \lambda \text{Complexity}(h)$$

$$\hat{h}^* = \underset{h \in H}{\operatorname{argmin}} \text{Cost}(h)$$

07-03-19

→ Theory Of learning - Learning with an agent

- How many examples do we need to get a good hypothesis?
- What hypothesis should we use?
- How complex should the hypothesis space be?
- How to avoid overfitting?
- How many examples do we need for learning?
- If the hypothesis is wrong, we get wrong result.
- If our hypothesis is consistent with a large number of examples, then it is right.
- The Probably Approximately Correct (PAC) learning algorithm: If our hypothesis is consistent with a large number of examples, then the error rate is low.
- Error Rate

$$\text{error}(h) = \text{GenLoss}_0(h) = \sum_{x,y} L_{0,1}(y, h(x)) P(x, y)$$

$\text{error}(h) \leq \epsilon \Rightarrow$ Hypothesis is Approximately Correct

$\epsilon \Rightarrow$ small constant

$\text{error}(h_b) > \epsilon \Rightarrow$ Hypothesis not approximately correct

- $P(h_0 \text{ agrees with } N \text{ example}) \leq (1-\epsilon)^N$

$$P(H_{\text{bad}} \text{ contains a consistent hypothesis}) \leq |H| \text{ bad} / (1-\epsilon)^N \\ \leq |H| (1-\epsilon)^N$$

Reduce Probability

$$|H| (1-\epsilon)^N \leq d$$

$$1-\epsilon \leq e^{-t}$$

$$N \geq \frac{1}{\epsilon} \left(\ln \frac{1}{d} + \ln |H| \right) \rightarrow \text{Algorithm is consistent, Hypothesis is correct}$$

N: No. of examples

Equation should see more than

H: Hypothesis

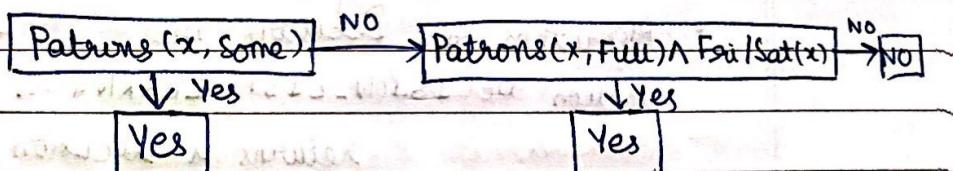
'n' examples for our algorithm

ϵ : Small constant

to be consistent

? - PAC learning examples: learning decision lists

- Applying PAC algorithm for new hypothesis space.
- Decision tree is like a decision tree but less complex. It branches only in one direction.
- It is a series of tests. If it passes the tests, it gives the result, i.e., one branch.
- Decision list for Restaurant Problem:



Will Wait $\Leftrightarrow (\text{Patrons} = \text{Some}) \vee (\text{Patrons} = \text{Full} \wedge \text{Fri/Sat})$

- We have a K-dimensional decision test.
- Our decision list above is 2-dimensional, since it has only 2 literals.
- A K-dimensional decision list or KDL has at most K literals and then it can generalize successfully.

08-03-14

- We have two main tasks in K-DL:
 - (i) We have to show that KDL is learnable.
 - (ii) We have to find an efficient algorithm that returns a consistent decision list.

- Number of conjunctions:

$$|\text{Conj}(m, k)| = \sum_{i=0}^{k+1} \binom{2^m}{i} = O(n^k)$$

$$|K\text{-DL}(n)| = 2^0 (n^k \log_2(n^k))$$

$$n \geq \frac{1}{e} \left(\ln \frac{1}{\epsilon} + O(n^k \log_2(n^k)) \right)$$

↳ we need 'n' no. of examples atleast for the PAC algorithm to learn from the examples and become accurate.

- To find an efficient algorithm, we use greedy algorithm such that it is accurate. It considers only the examples that are relevant and removes or does not consider the other examples.

- Algorithm for Decision Tree:

function DECISION-LIST-LEARNING (examples)

 returns a decision, or false

 if examples is empty then return the decision list

 t ← a test that matches a non empty subset examples + of examples

 such that the members of examples + are all positive or all negative

 if there is no such t then return false

 if the examples in examples + are positive then else if YES
 return a decision list with initial test t and outcome 0
 and remaining test given by DECISION-LIST-LEARNING
 (examples - examples +)

Regression and Classification with Linear Models

- We use regression to predict the future with the help of the past -

- Univariate Linear Regression:

i/P: x o/P: y

$$y = w_1 x + w_0 \quad \text{Here, } w_1 \text{ and } w_0 = \text{weights}$$

$$\text{Vector: } w = (w_0, w_1) \quad h_w(x) = w_1 x + w_0$$

$h \rightarrow \text{hypothesis}$

$$\text{Loss}(h_w) = \sum_{j=1}^N L_2(y_j, h_w(x_j)) = \sum_{j=1}^N (y_j - h_w(x_j))^2$$

$$= \sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2$$

$$w^* = \operatorname{argmin}_w \text{Loss}(h_w)$$

$$\sum_{j=1}^N (y_j - (w_1 x_j + w_0))^2 \text{ is minimized when } w_0 \text{ and } w_1 = 0$$

$$w_1 = 0.232 \quad w_0 = 2.46 \rightarrow \text{For given big}$$

$$w_0 = (\sum y_j - w_1 (\sum x_j)) / N$$

$$w_1 = N (\sum x_j y_j) - (\sum x_j) (\sum y_j)$$

$$N (\sum x_j^2) - (\sum x_j)^2$$

We have to use gradient descent to minimize loss.

We have to make the weights of w_0 and w_1 near to 0 to minimize the loss

Gradient Descent Algorithm:

$w \leftarrow$ any point in parameter space

loop until convergence do

for each w_i in w do

$$w_i \leftarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Loss}(w)$$

We try to minimize

the loss at each iteration

α : Learning Rate

ernable.
then that

for the
amples

greedy

It

consider

amples)

false

ecision list

subset

all

or YES
then else or no
d outcome
LEARNINGS
rules)

- In univariate, the loss function is a quadratic function: (Partial Derivation for the linear fn., for one training eq. with x_j) (x, y)

$$\begin{aligned}\frac{\partial}{\partial w_i} \text{Loss}(w) &= \frac{\partial}{\partial w_i} (y - h_w(x))^2 \\ &= 2(y - h_w(x)) \times \frac{\partial}{\partial w_i} (y - h_w(x)) \\ &= 2(y - h_w(x)) \times \frac{\partial}{\partial w_i} (y - (w_0 + w_1 x)) \\ \frac{\partial}{\partial w_0} \text{Loss}(w) &= -2(y - h_w(x)) \\ \frac{\partial}{\partial w_1} \text{Loss}(w) &= -2(y - h_w(x)) \times x\end{aligned}$$

Unspecify Learning Rate α

$$w_0 \leftarrow w_0 + \alpha(y - h_w(x))$$

$$w_1 \leftarrow w_1 + \alpha(y - h_w(x)) \times x$$

If $x = +ve \rightarrow$ Reduce w_1

If $x = -ve \rightarrow$ Increase w_1

- Multivariate Linear Regression

$$\begin{aligned}h_w(x_j) &= w_0 + w_1 x_{j,1} + \dots + w_n x_{j,n} \\ &= w_0 + \sum_i w_i x_{j,i}\end{aligned}$$

- Augmented Vectors: Add a feature to each x by taking on $a: x_j^1 : x_j, a = 1$

$$\text{Then: } h_w(x_j) = w \cdot x_j = w^T x_j = \sum_i w_i x_{j,i}$$

- And batch gradient descent update becomes:

$$w_i \leftarrow w_i + \alpha \sum_j (y_j - h_w(x_j)) x_{j,i}$$

- Regularization - Used to remove overfitting.

$$\text{Cost}(h) = \text{Emp Loss}(h) + \alpha \text{Complexity}(h)$$

$$\text{Complexity}(h, w) = L_q(w) = \sum |w_i|^q$$

It will not intersect on axes.

In L1, weights will become 0, but not in L2.

L1: Intersection done with diamond box - intersects on axes - weights will be zero - minimum loss.

L2: Intersection done with circle - intersection not on axes - weights are not zero - no minimum loss

- Linear Classification - Hard Thresholds

- Decision Boundary
- Linearly Separable
- Classification Hypothesis

- Perceptron Learning Rule

- Perceptron Convergence Theorem
- For a single example (x, y)

$$w_1 \leftarrow w_1 + \alpha (y - h_0(x)) x_j$$

- Linear Classification with Logistic Regression

$$h_w(x) = \text{Logistic}(w \cdot x) = \frac{1}{1 + e^{-w \cdot x}}$$

For a single sample (x, y) and L2 loss function

$$w_1 \leftarrow w_1 + \alpha (y - h_0(x)) h_0(x) (1 - h_0(x)) x_j$$

14-03-'19



Artificial Neural Networks

- Machine learning approach where a machine works like a ^{human} brain.
- It has neurons which are connected to each other.
- Thus, ANN is a network connected with a large no. of neurons, associated with some weights.
- It is called as connectionist model.
- The neuron is the main component.

→ Neuron Model

→ Architecture

→ Learning Algorithms

- Neuron model is the processing unit.
- Architecture has neurons connected to each other with weights.
- Learning algorithm makes the neurons learn.

→ Neuron Model

* Inputs

* Weights

* Adder

* Activation Function

$$\star U = \sum_{i=1}^m w_i * x_i$$

$$\text{Op: } Y = \psi(U + b) = \psi(U)$$

Here:

U = linear combiner = sum of wts x inputs

m = no. of inputs

ψ = Activation fn. to limit outputs

$$V = U + b$$

Activation fn: $v = \sum_{i=0}^m w_i * x_i$ where $w_0 = b$ and $x_0 = 1$

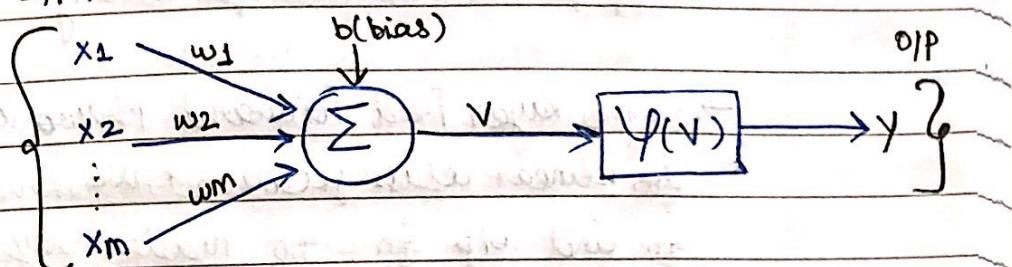
Here: $b = \text{bias}$

$w = \text{weights}$

$x = \text{inputs}$

→ Neuron

I/P values



→ Activation Functions

- Step Function

$$\Phi(v) = \begin{cases} a, & \text{if } v < c \\ b, & \text{if } v \geq c \end{cases}$$

- Ramp Function

$$\Phi(v) = \begin{cases} a, & \text{if } v < c \\ b, & \text{if } v > d \\ a + [(v - c) * (b - a) / (d - c)], & \text{otherwise} \end{cases}$$

- Sigmoid Function

$$\Phi(v) = \frac{1}{1 + \exp(-x * v + y)}$$

- Gaussian Function

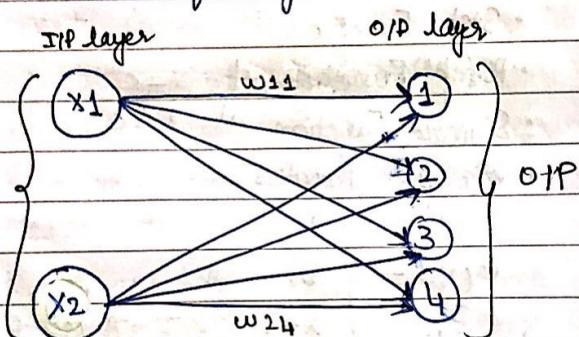
$$\Phi(v) = \exp\left(-\frac{(v - c)^2}{\sigma^2}\right)$$

→ Neural Network Architecture

- Single Layer Feed Forward Network
- Multi Layer Feed Forward Network
- Recurrent Network

→ Step function is used for Binary Classification.

→ Single Layer Feed Forward Network is used for linear classification. We have only IIP for and OIP fn., no hidden layers.
• IIP is connected to links, where we have the sum of weights.



Threshold Value.

OIP Value > Threshold value, neuron takes activated value.

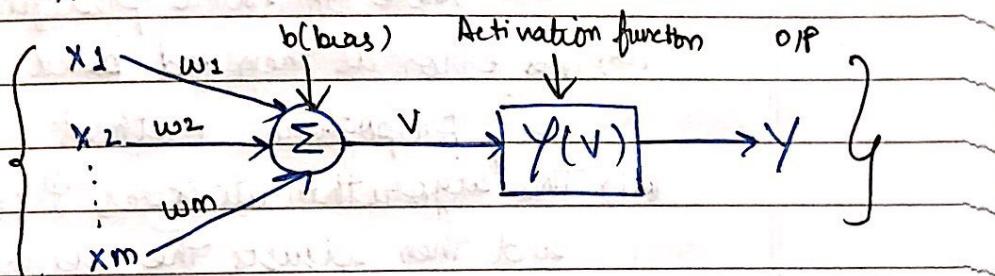
OIP Value < Threshold value, neuron takes deactivated value.



→ Perception Neuron Model

- Special model of single layer feed forward network.
- Used for binary classification.
- Classifies input into two categories and gives one O/P.
- Step fn. is activation fn.
- Before Activation fn. is used, it must be trained.
- If it classifies unknown instances correctly, output accuracy is high.
- +ve O/P : C1 -ve O/P : C2

I/P value



$$\Phi(V) = \begin{cases} -1, & \text{if } V < \theta \\ +1, & \text{if } V \geq \theta \end{cases}$$

- OR function graph → linearly separates T and F

I/P		O/P
x_1	x_2	$x_1 \text{ OR } x_2$
0	0	0
0	1	1
1	0	1
1	1	1

Graph of the OR function:

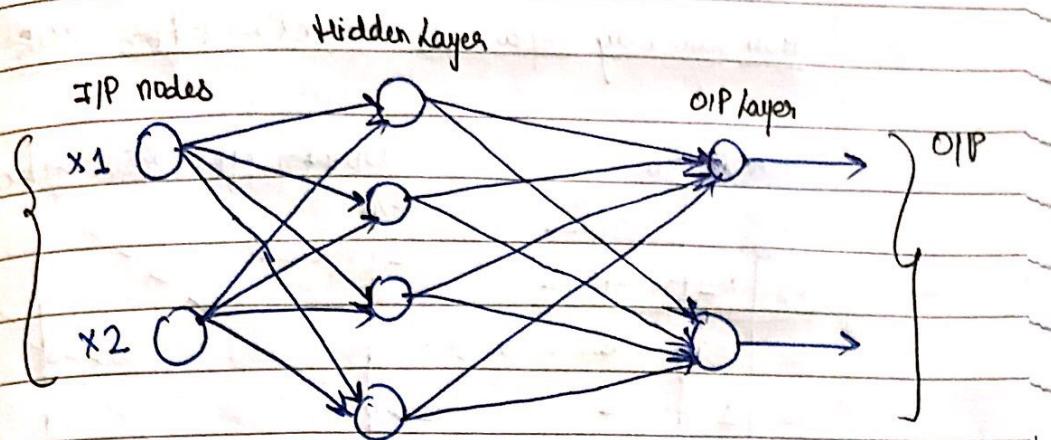
The graph shows the OR function plotted against x_1 and x_2 . The axes are labeled x_1 and x_2 . Four points are marked: $(0,0)$ labeled "false(0,0)", $(0,1)$ labeled "true(0,1)", $(1,0)$ labeled "true(1,0)", and $(1,1)$ labeled "true(1,1)". A diagonal line separates the "false" region (bottom-left) from the "true" region (top-right).

- Limitations: Can model only linearly separable functions, cannot go to more extent, only two classes used, can model AND by not XOR function.

→ Multi Layer Feed Forward Network

- We have a hidden layer here.
- Used for data which is not linearly separable.
- It will evolve with continuously learning examples.
- We have a backpropagation method
- Learning Technique: Backpropagation Method
 - (i) The output values obtained for a known set of input values is compared with the correct answer in order to determine the value of some predefined error function
 - (ii) This error is then fed back to the network using propagation method.
 - (iii) The algorithm utilizes this information and then adjusts the weights of all connections so that the corresponding error function is reduced to some extent
 - (iv) This procedure is repeated for a large number of training cycles. So the network obtains a state where the error becomes negligible.
 - (v) You can say that the network has learnt target function.

P.T.O.
→



→ Learning Algorithm for Perceptron

- Initial random weights as +0.5 and -0.5.
- Training data is presented to perceptron and the output is observed.
- We check the output, if it is incorrect, the weights are adjusted
 - $\Delta = \text{learning rate}$
 - O/P correct $\Rightarrow \Delta = 0$
 - O/P too low $\Rightarrow \Delta = +\text{ve no}$
 - O/P too high $\Rightarrow \Delta = -\text{ve no}$
- $w_i = w_i + (\Delta * x_i + E)$

$E = \text{Error}$ $\Delta (0 < \Delta < 1) \rightarrow \text{learning rate}$

↳ Adjustment of Weights

- Again, training data is checked and if the output is wrong, weights are adjusted again.
- Training continues till all weights are correct and all errors are zero.
- Each iteration in this process is known as an 'epoch'.

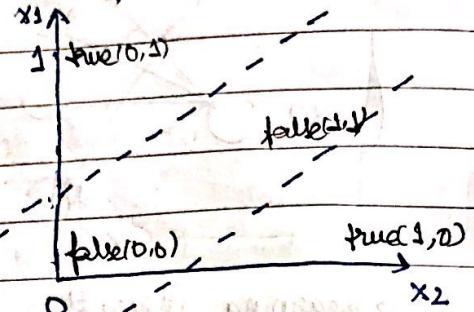
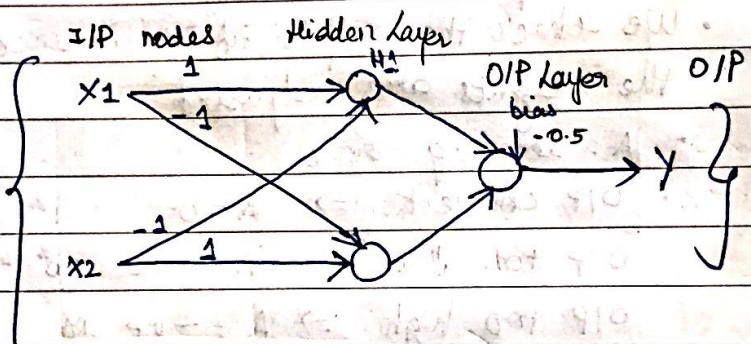
15-03-'19

Non linearly Separable Function for XOR function

XOR Fn

I/P	O/P	
x_1	x_2	$x_3 = x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0

Graph for XOR function

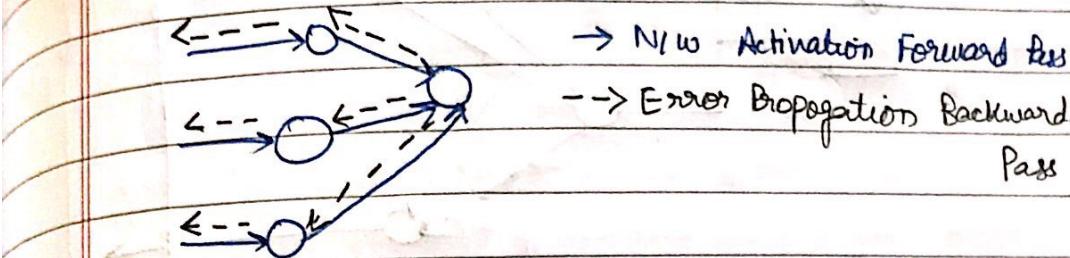
Network for XOR

I/P's	O/P's of Hidden Layer	O/P Node	$x_1 \text{ XOR } x_2$
$x_1 \ x_2$	$H_1 \ H_2$		
0 0	1 0	$-0.5 \rightarrow 0$	0
0 1	-1 → 0 1	$0.5 \rightarrow 1$	1
1 0	1 -1 → 0	$0.5 \rightarrow 1$	1
1 1	0 0	$-0.5 \rightarrow 0$	0

Back Propagation Training Algorithm for FFNN Network

for XOR

→ Refine Weight Values



weight Update Rule

$$w_{ji} = w_{ji} + \Delta w_{ji}$$

$$\Delta w_{ji} = n * [\delta E / \delta w_{ji}]$$

n is learning in [0.1, 0.9]

Stopping Criterion

→ Total mean squared error change

→ Generalization based criterion

Delta Rule for Error Minimization.

$$E_{av} = \sum_{i=1}^N \sum_{j=1}^M (y_{ij} - \hat{y}_{ij})^2$$

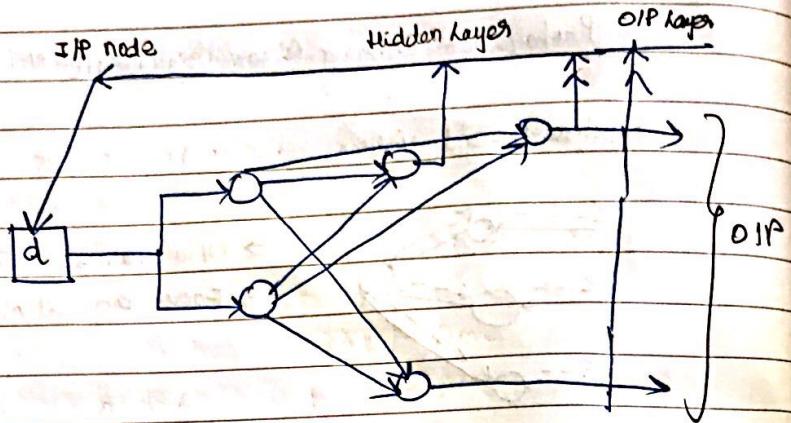
→ Recurrent Network

- Non Cycle Graph: Backpropagation

- Recurrent: Form a cycle

- In recurrent network, output will be given as i/p

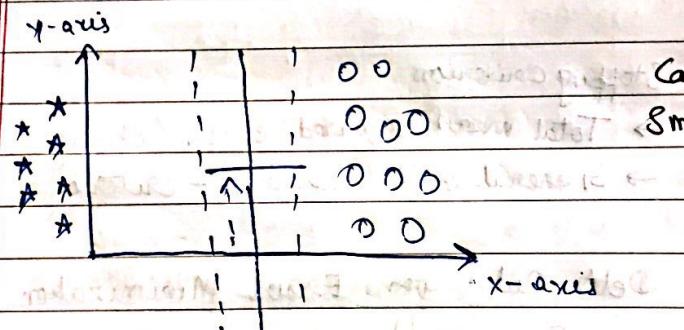
- In backpropagation, error value will be backpropagated



→ Support Vector Machine

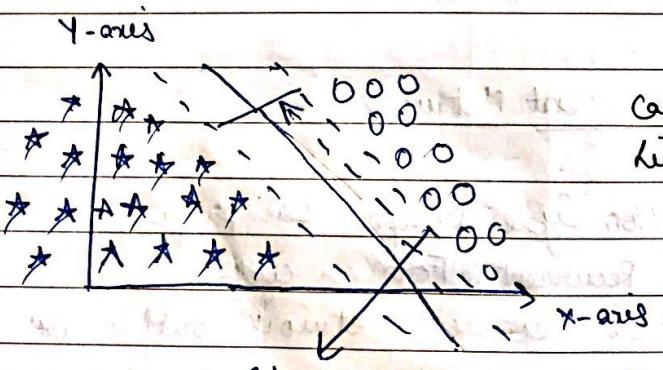
- linearly Separable and Non-linearly Separable Functions.

(ii)



Case 1:
Small Margin

(iii)



Case 2:
Zero Margin

Support Vector

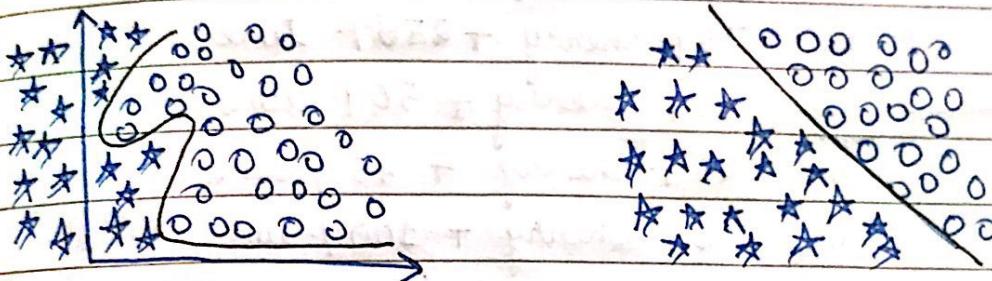
→ In non-linearly separable, we use Kernel functions to separate data.

Transformation into higher dimensional space

- we use kernel function to linearly separate classes
- we use higher degree polynomial

Input Space

Feature space



18-03-'19

→ Learning Probabilistic Models

- Two flavours are present for a chocolate, cherry and lime, which have the same wrapper.

- Hypotheses:

$$h_1: 100\% \text{ cherry} + 0\% \text{ lime}$$

$$h_2: 75\% \text{ cherry} + 25\% \text{ lime}$$

$$h_3: 50\% \text{ cherry} + 50\% \text{ lime}$$

$$h_4: 25\% \text{ cherry} + 75\% \text{ lime}$$

$$h_5: 0\% \text{ cherry} + 100\% \text{ lime}$$

Probability of each hypothesis

$$P(h_i | d) = \alpha P(d | h_i) P(h_i)$$

To Predict Lime / Cherry:

$$\begin{aligned} P(X | d) &= \sum P(X | d, h_i) P(h_i | d) \\ &= \sum P(X | h_i) P(h_i | d) \end{aligned}$$

→ Learning with Complete Data

- Agent tries to know by using the given complete data.

- We do not have any hypothesis here.

- Likelihood of particular dataset:

$$P(d | h_\theta) = \prod_{j=1}^N P(d_j | h_\theta) = \theta^c (1-\theta)^l$$

- Maximizing log likelihood

$$L(d | h_\theta) = \log \prod P(d_j | h_\theta) = \sum_{j=1}^N \log (d_j | h_\theta) = c \log \theta + l \log (1-\theta)$$

- Set resulting expression to '0'

$$\frac{dI(d(\theta))}{d\theta} = \frac{c}{\theta} - \frac{l}{1-\theta} = 0$$

$$\Rightarrow \theta = \frac{c}{c+l} = \frac{c}{N}$$

- Naive Bayes Model

$$P(C|x_1, \dots, x_n) = \propto P(c) \prod P(x_i|c)$$

→ Learning with hidden Variables

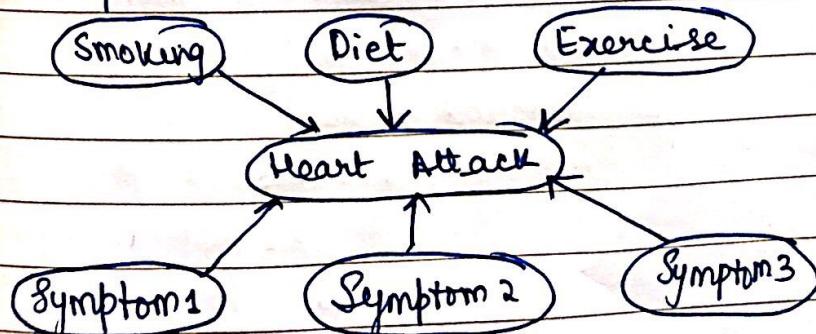
- Missing hidden variables arises complications
- Links will become more.
- Eg - Heart Attack is caused due to smoking, drinking, etc.

Remove hidden variable ~~smokes~~ heart attack.

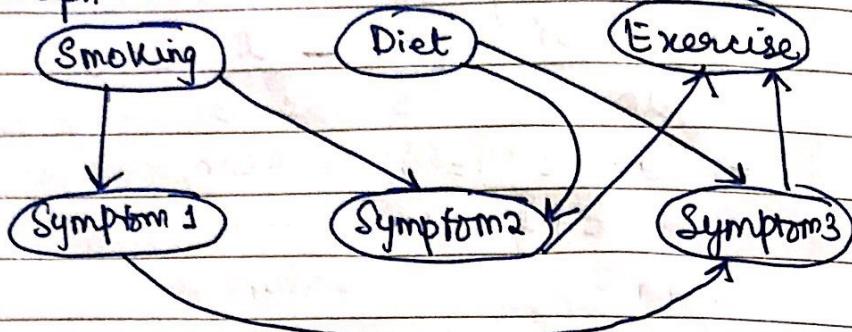
Links will be more, smoking and drinking

may cause other things:

- Graph will become more complex if we remove the hidden variables.
- Graph with hidden variable heart attack:



- Graph Without hidden variable Heart Attack.



→ EM Algorithm

Here:

x = observed value

z = hidden variables

θ = all the parameters

E-Step : Computation of Summation

M-Step : Maximization of log likelihood

$$\theta^{(i+1)} = \underset{\theta}{\operatorname{argmax}} \sum_z P(z=z|x, \theta^{(i)}) L(x, z=z|\theta)$$



NATURAL LANGUAGE PROCESSING

N-GRAM: NATURAL LANGUAGE FOR COMMUNICATION

Natural Language Processing

- How the machine/agent learns English.
- Language is a set of strings.
- $P(S) = \text{words}$ $S = \text{sentence}$
- N-gram character model:
 - Model of probability distribution of n-letter sequence.
 - Defined as Markov Chain, depends on immediate previous state.
 - Trigram Model (three states):
$$P(c_i | c_{i-1}, c_{i-2}) = P(c_i | c_{i-2}, c_{i-1})$$
 - Probability of Sequence of Characters
$$P(c_1, c_2, \dots, c_N) = \prod_{i=1}^N P(c_i | c_{i-1}, c_{i-2}) = \prod_{i=1}^N P(c_i | c_{i-2}, c_{i-1})$$
 - N-gram model is widely used in the language identification.
 - Initially, we must build trigram character model for identification of language:
$$P(c_i | c_{i-2}, c_{i-1}, \emptyset)$$

- Corpus: Body of text
- After building trigram model, we divide it into training model and evaluation model, and give it a dataset that it has not seen.
- After building trigram model, we get the probability of text from language and we get the most probable language by using the Baye's rule and also Markov Assumption.

$P(\text{text} | \text{language})$

$$\lambda^* = \underset{\lambda}{\operatorname{argmax}} P(\lambda | c_{1:N})$$

$$= \underset{\lambda}{\operatorname{argmax}} P(\lambda) P(c_{1:N} | \lambda)$$

$$= \underset{\lambda}{\operatorname{argmax}} P(\lambda) \prod_{i=1}^N P(c_i | c_{i-1:i-1}, \lambda)$$

- Smoothing N-gram model :

- Word starting with 't' is frequent

'th' → frequent

'ht' → unknown

- Even if we do not train that dataset, it should understand and give the correct output.

- Process of adjusting the probability of low frequency count is smoothing.

- Laplace Smoothing : Poor performance, but it goes in the right direction.

- Back-Off - Model

→ Better Method

→ Starts counting / estimating 'n' gram counts.

→ Linear Interpolation Smoothing is a backoff model that combines unigram, bigram and trigram.

$$\hat{P}(c_1 | c_{i-2:i-1}) = \lambda_3 P(c_i | c_{i-2:i-1}) + \lambda_2 P(c_i | c_{i-1}) + \lambda_1 P(c_i)$$

$$\lambda_3 + \lambda_2 + \lambda_1 = 1$$

→ λ value can be fixed or varied

If unigram / bigram / trigram value is high we increase λ at that gram and decrease at the other

- Model Evaluation

- How a model is built and how are we evaluating it.
- Perplexity describes the probability of a sequence.
- We must build a training and testing model, and give unknown dataset to estimate accuracy.
- Perplexity ($c_1:N$) = $P(c_1:N)^{-1/N}$
- Perplexity is calculated for:
 - unigram → more mistakes
 - bigram → less grammar mistakes
 - trigram → perfect sentence

- N-gram word models

- For an unknown word in the vocabulary, we add it as: <UNK>