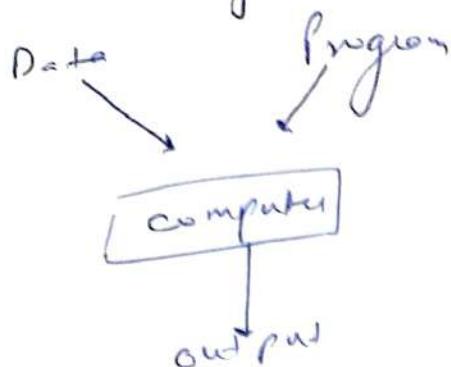
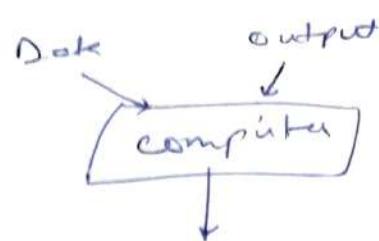


Introduction Machine Learning

Algoorth



ML



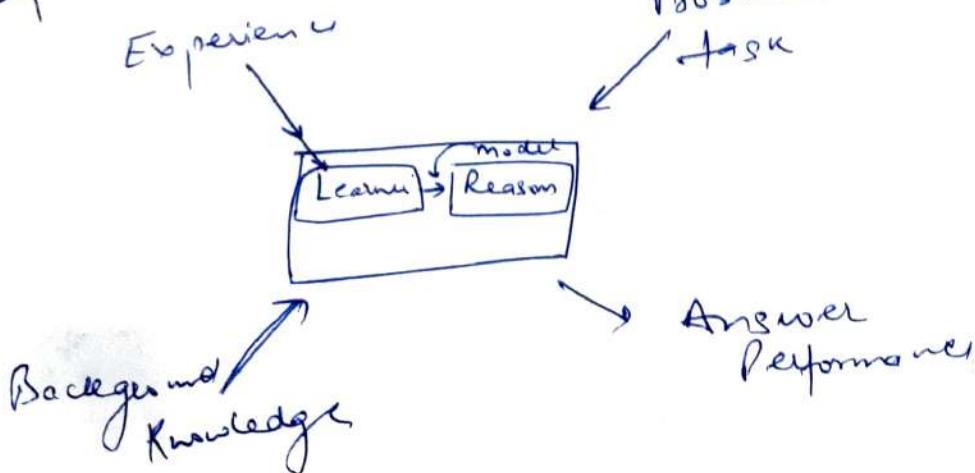
Program

Learning: is the ability to improve our behaviour with experience

ML explores algorithms

- learn / build model for data
- model used for prediction, decision making or solving tasks

Mitchell: A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance on tasks in T as measured by P improves with experience E



Many domain and application

Medicines

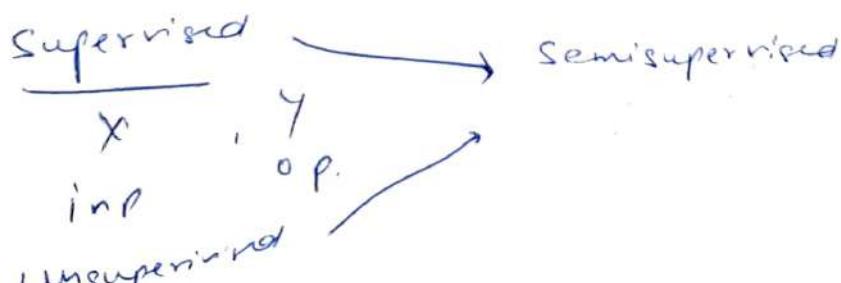
- Diagnosis: a disease
- Data mine historical medical record

to learn which future patients will respond best to which treatment

Vision, Robot control, NLP, speech recognition, machine translation, financial

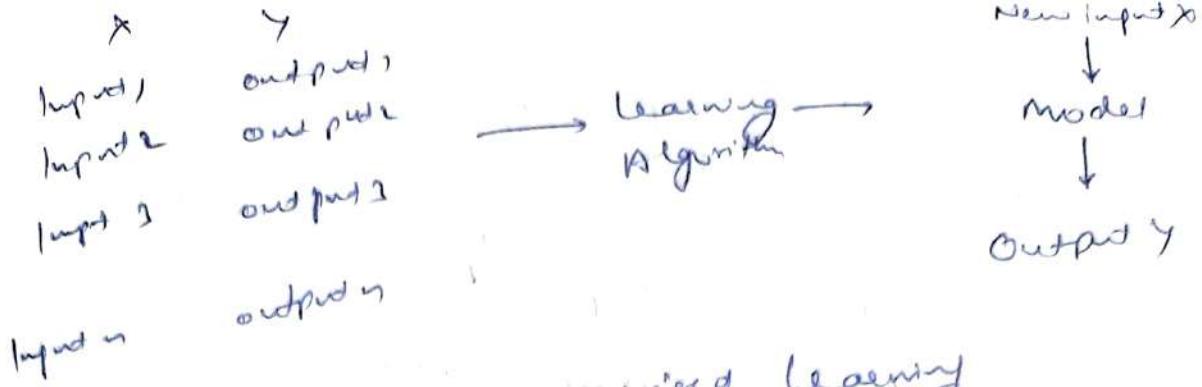
- ① Choose the training experience features
- ② Choose the target function
(that to be learned)
- ③ Choose how to represent the target function
- ④ Choose a learning algorithm to infer the target function
 - rich representation more difficult to learn
 - hypothesis language

Chapter 3 Different type of Learning

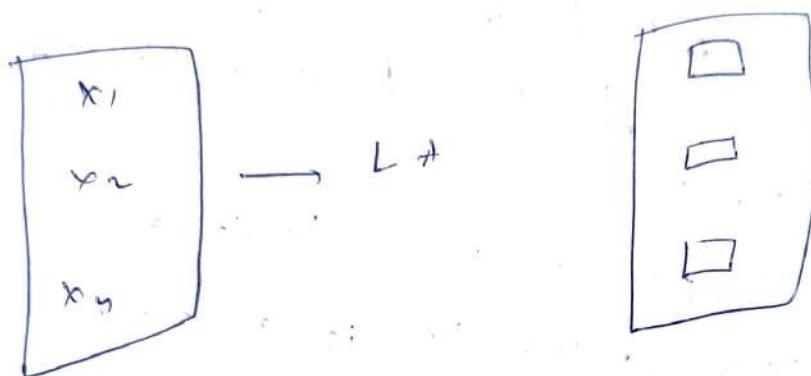


Reinforcement Learning

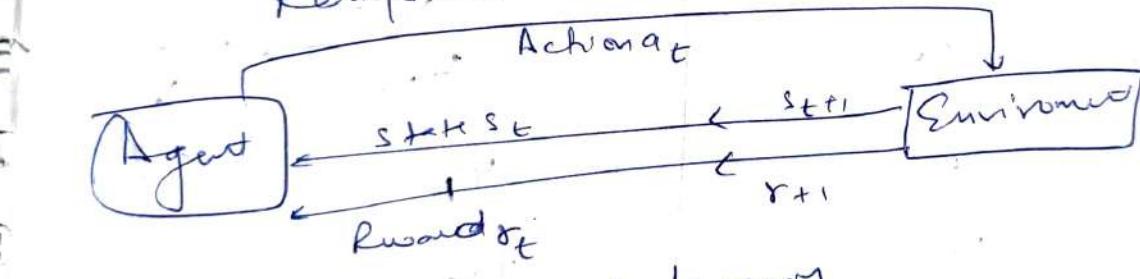
Supervised Learning



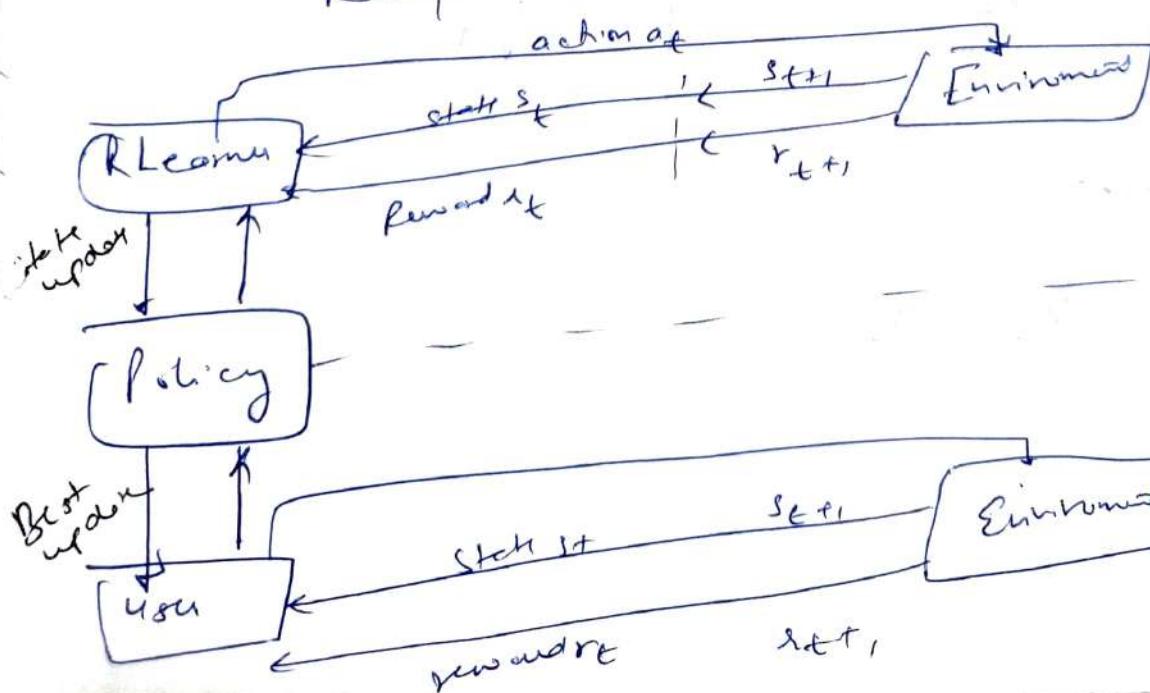
Unsupervised Learning



Reinforcement Learning



Reinforcement Learning



Supervised learning

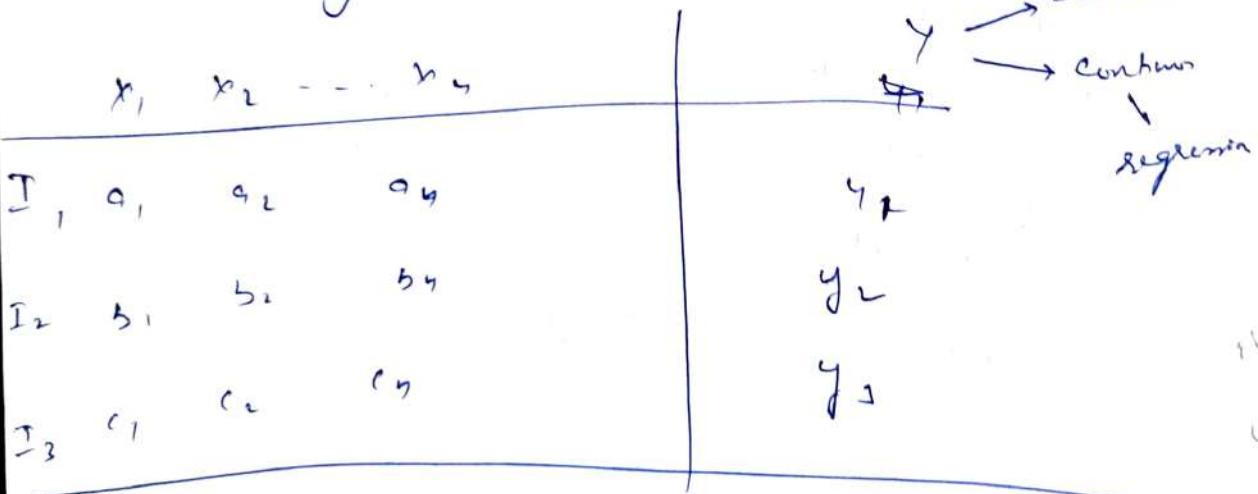
Given

- a set of input features x_1, \dots, x_n
- a target feature y
- a set of training examples where the values for the input features and the target features are given for each example
- a new example, where only the values for the input features are given

Predict the values for the target feature for

the new example

- classification when y is discrete
- regression when y is continuous



Test instance

$$z_1, z_2, z_3, \dots, z_n$$

Classification

Example: Credit scoring

Differentiating between low-risk and high-risk customers from other income and saving

- Discretization: if income > 0, and saving > 0, then low-risk else high risk

Regression

Example: Price of a used car

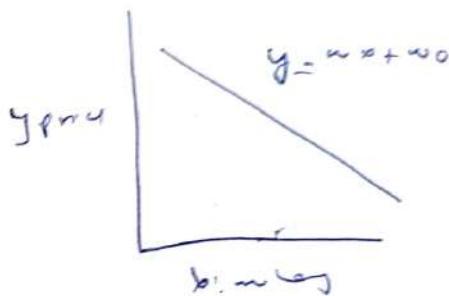
x : car attributes

y : Price

$$y = g(x, \theta)$$

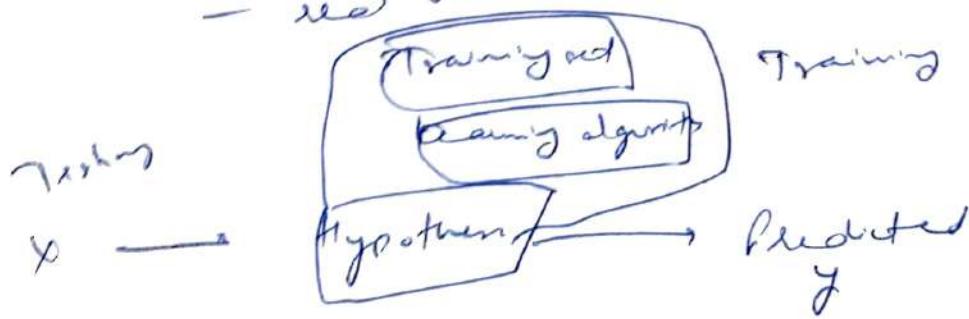
$g(\cdot)$ = model

θ parameters



features

- categorical
- ordinal
- integer valued
- real valued



Classification Learning

- Task 7

Representation

- ① Decision tree
- ② Linear function
- ③ multivariate linear function
- ④ single layer perceptron
- ⑤ multilayer neural network

Hypothesis space:

one way to think about a supervised learning machine is as a device that explores "hypothesis space"

- Each setting of the parameters in the machine is a different hypothesis about the function that maps input vectors to output vector.

Chapter 4 Hypothesis space
 and Inductive Bias

Inductive Learning or Prediction

Given example (\hat{x}, y)

$(\hat{x}, f(\hat{x}))$

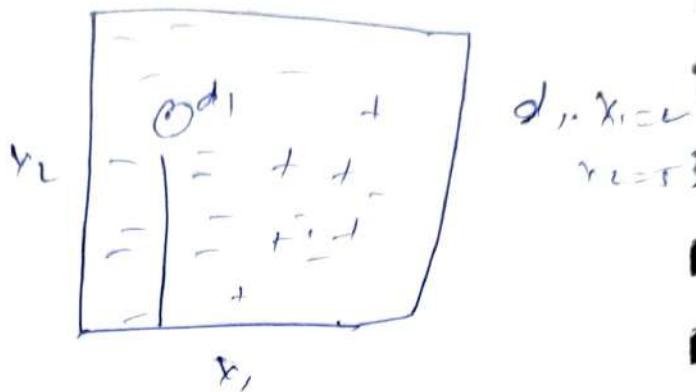
Classification $f(\hat{x})$ is discrete

Regression $f(\hat{x})$ is continuous

Probability estimation = probability of \hat{x}

Feature properties that describe each instances
feature vector

Language



hypothesis space

H $h \in H$

output of a learning

Examples (x, y)

Training Data set of examples

Instance Space

Concept C $C \subseteq X$

Target function $f: X \rightarrow Y$

find
hypothesis $h \in H$ that approximets f

H restricted hypothesis

by defining bias

- constraints
- preference

Input $s \subseteq x$
Output $h \in H$

x_1, x_2, x_3, x_4 | Possible instances
 $2^4 / 2^4$

Homology Book of
possible.

$$2^{16} / 2 =$$

- Hypothesis language reflect on inductive bias
of the learner
Hypothesis language

Occam's Razor

A classical example of Inductive Bias,
A simplest consisted hypothesis about
the target function is actually the best
Underfitting and overfitting.

Underfitting: model is too "simple" to represent all the relevant class characteristics.

- high bias and low variance

- high bias and low variance
- high training error and high test error
the "complex" and fit

overfitting: model is too "complex" (with many) interactions

irrelevant characteristics (noise)
low bias and high variance

- low bias and high variance
- low training error and high test error
-

Chapter 4 Evaluation and cross-validation

- Evaluating the performance of learning is important because
 - learning systems are usually designed to predict the class of "future" unlabeled data points.
- Typical choices for performance evaluation
 - Error
 - Accuracy
 - Precision/Recall
- Typical choices for Sampling Methods
 - Train / Test sets
 - k-fold cross validation

H₁: S

Experimental Evaluation

Error
Accuracy
Precision/Recall

Evaluation Sample

Training set
Test set
cross validation

Evaluating Prediction

- Suppose we want to make a prediction of value for a target feature on example x
 - y is the observed value of target feature on example x
 - \hat{y} is the predicted value of target feature on example x .
 - How is the error measured?
 - $\hat{y} = h(x)$
 - ① Absolute error $\frac{1}{n} \sum_{i=1}^n |h(x_i) - y_i|$
 - ② sum of sq. error $\frac{1}{n} \sum_{i=1}^n (h(x_i) - y_i)^2$
- ③ Number of misclassification

		True class		$i=1$
Hyp. class	Pos	Pos	Neg	
		FP	FN	
Pos				
Neg				
	P		H	

$$\text{Accuracy} = \frac{T_P + T_H}{P + H}$$

$$\text{Precision} = \frac{T_P}{T_P + FP}$$

$$\text{Recall} = \frac{T_P}{P}$$

Sample Error and True Error

- The sample error of hypothesis f with respect to target function c and data sample S is:
$$\text{error}_s(f) = \frac{1}{n} \sum_{w \in S} \delta(f(w), c(w))$$
- The true error ($\text{error}_D(f)$) of hypothesis f with respect to target function c and distribution D , is the probability that f will misclassify an instance drawn at random according to D
$$\text{error}_D(f) = \Pr_{w \sim D} [f(w) \neq c(w)]$$

Difficulties in evaluating hypotheses with limited space

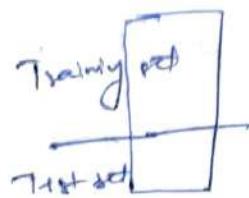
- Bias in the estimate: - The sample error is a poor estimator of true error
 - \Rightarrow test the hypothesis on an independent test set
- We divide the example into
 - Training examples that are used to train the learner
 - Test examples that are used to evaluate the learner

Variance in the estimate - the smaller the test set, the greater the expected variance.

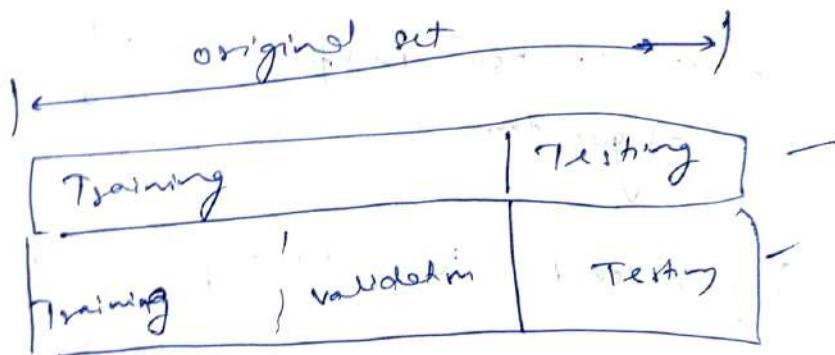
split the examples

Training set - to train the learner

Test set - to test the learner

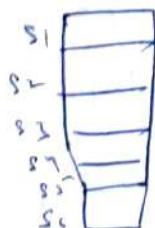


Validation set



Validation fails to use all the available data

k fold cross validation



Round i: use s_i for testing /
 $s - s_i$ for training

Trade off

In ML, there is always a trade off b/w
- complex hypothesis that fit the training
data well

- Chapter hypothesis about generalization error
- As the amount of training data increases, the generalization error decreases

Chapter 6 Tutorial 1

- Supervised Vs Unsupervised learning
- Categorical Vs Continuous features
- Regression Vs Classification
- Bias Vs Variance
- Generalization Performance of a learning algorithm.

Supervised Vs Unsupervised learning

unsupervised

Supervised

Training sample

→ (Input, Target)

The descriptions
of presentation

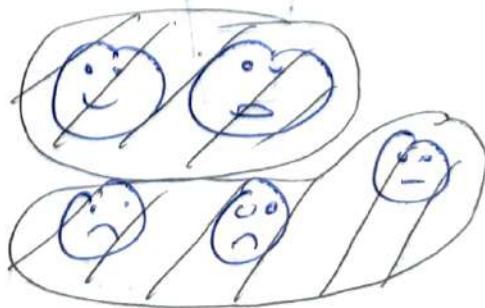


✓ Happy/sad

- You don't have target values

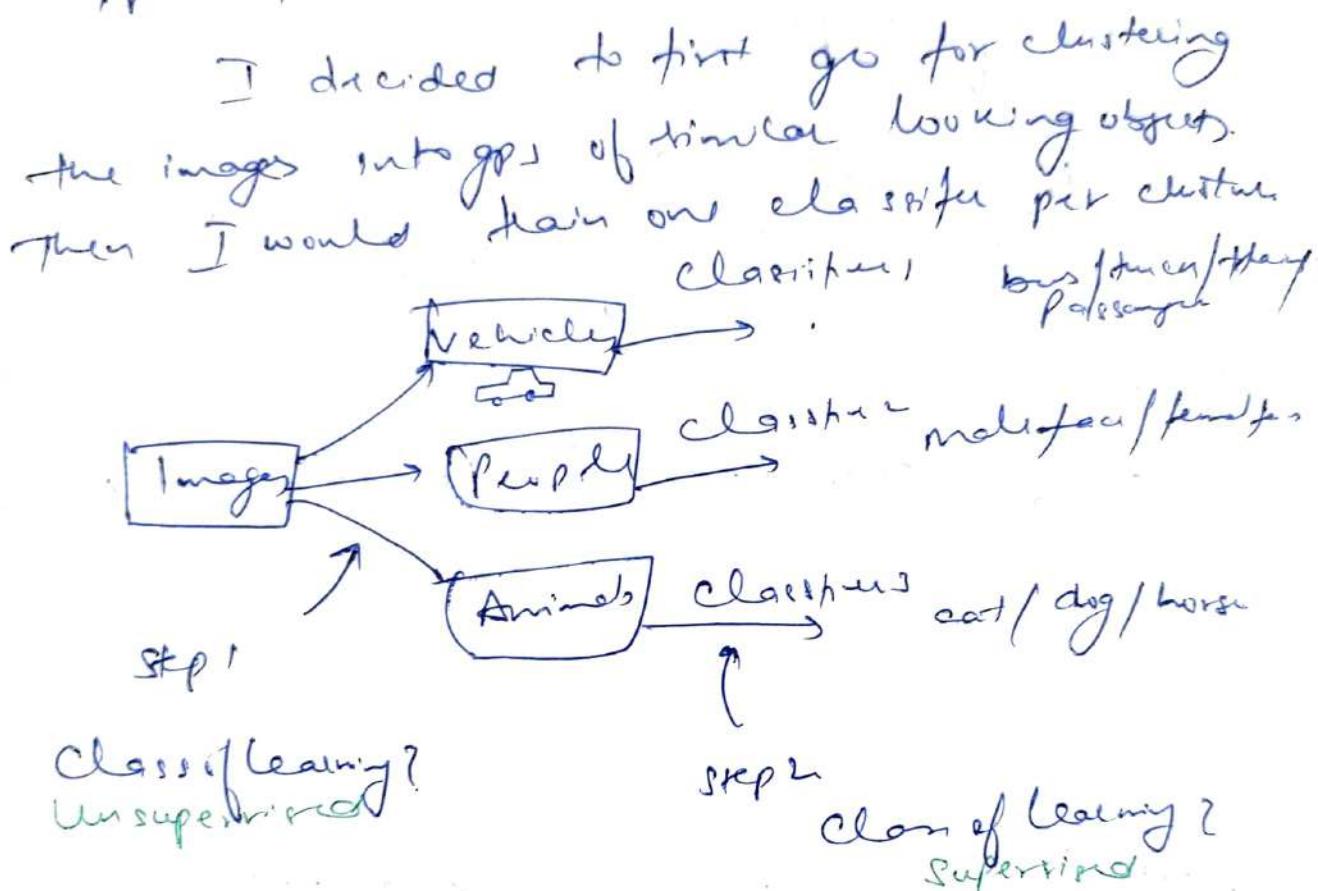
Training example

(Input)



Supervised Vs Unsupervised Learning

Q1: Suppose I have been given 1 million images with their ground truth labels. I looked at the dataset and found that the images are from a wide range of categories. There are vehicles, plants, animals, human faces etc.



kind of features

categorical Vs continuous

categorical features

↳ Finite no. of categories or classes

e.g.: Gender, age-group

continuous features

Input

no. of values

e.g. age, height, weight, price, etc

Types of supervised learning

Regression & classification

Type of output variable

Classification

Output / Target variable

is categorical / discrete

e.g. Given images of 5 animals → species of animal
Predict the animal

Given image of a tumor

Predict cancerous / not cancerous

Regression

Output / Target variable is continuous

e.g. Given certain features of a car

Predict

Price of car

Predict age / weight / height of person

continuous

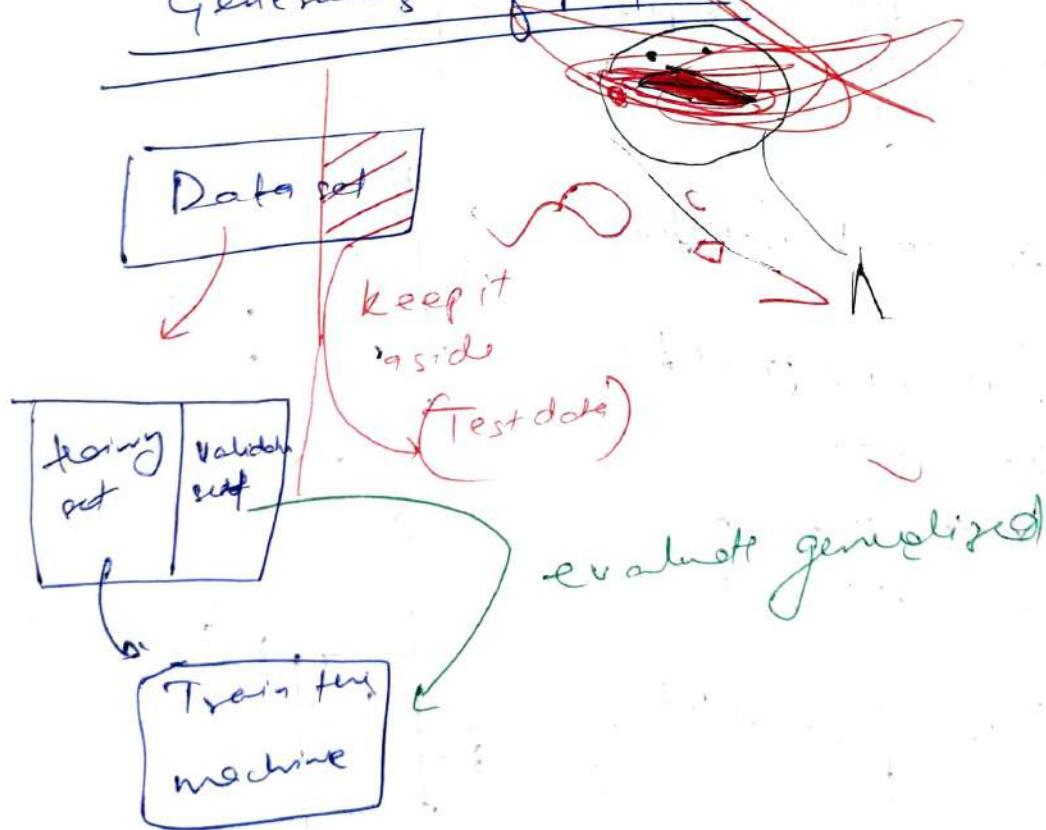
Bias vs Variance

Bias: erroneous Assumption in the learning algorithm

= Variance sensitivity of model towards noise rather than important features of the relevant relationship. B/w input and output

	# of features ↓ N	# of parameters ↓ N _{par}	# of training examples Remains the same
Bias	↑ error	↑ error	↓ error
Variance	↓ error	↑ error	↓ error

Generalized of Performance



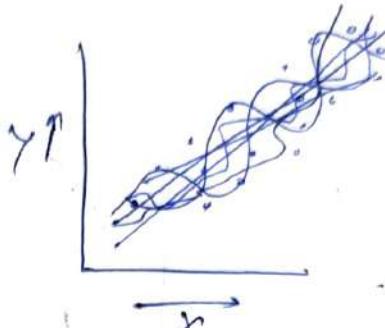
Chapter 5

Linear Regression

Regression is a supervised learning problem

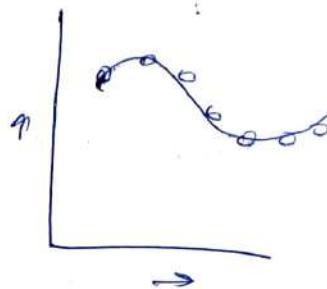
x, y
 $n \rightarrow \mathcal{D} \rightarrow \text{continuous}$

Simple regression



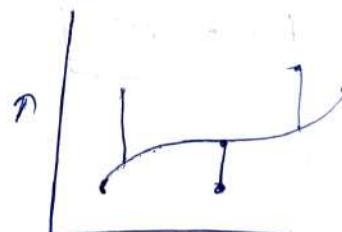
Linear regression

Given an input x compute an output y



Eg:

- Predict height from age
- Predict house price from house size
- Heart disease from wall fat sensors.

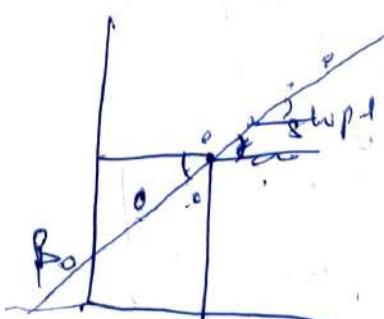


$$y = \beta_0 + \beta_1 x + \epsilon$$

y -intercept

Random error

slope



$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon \quad (\text{probable error})$$

$$E(Y|x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \quad \text{least sq.}$$

Least sq. line

$$\boxed{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p}$$

minimize SSE

β_0 is the intercept, β_j is the slope for j^{th} variable x_j .

Assumption

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

$$E(\varepsilon_i) = 0$$

$$d_1: y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1$$

$$\sigma(\varepsilon_i) = \sigma_\varepsilon$$

$$d_2: y_2 = \beta_0 + \beta_1 x_2 + \varepsilon_2$$

Errors are independent

$$d_n$$

ε_i are normally distributed

$$\text{Loss function} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

LSC part

The regression line or the least sq.

regression line is the unique line such that
the sum of the squared vertical (d) distance
b/w the data pt and the line is the
smallest possible

$$d_i = (x_i, y_i) \sum_{j \in D} (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$$

How to learn β_0, β_1 (parameters)

from 2d problem

$$Y = \beta_0 + \beta_1 x$$

- To find the values for the coefficients which minimize the objective function we take the partial derivatives of the objective function (losses) with respect to the coefficient. Set these to 0, and solve

$$\beta_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$\beta_0 = \frac{\sum y - \beta_1 \sum x}{n}$$

Multiple Linear Regression

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$\text{hence } = \sum_{i=0}^n \beta_i x_i$$

- there is a closed form which requires matrix inversion etc
there are iterative techniques to find weights

- delta rule (also called LMS method) which will update towards the objective of minimizing the SSE

Linear Regression

$$h(\mathbf{x}) = \sum_{i=0}^n \beta_i x_i$$

To learn the parameters $\Theta(\beta_i)$?

- make $h(\mathbf{x})$ closer to y , for the available training example
- Define a cost function $J(\Theta)$

$$J(\Theta) = \frac{1}{2} \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

Find Θ that minimizes $J(\Theta)$

LMS Algorithm

- Start a search algorithm (e.g. gradient descent algorithm) with initial guess of Θ
- Repeatedly update Θ to make $J(\Theta)$ small, until it converges to minimum

$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\Theta)$$

- J is a convex quadratic function, has a single global minima. gradient descent eventually at the global minimum

- At each iteration of the algorithm, it does
 - * step in the direction of steepest descent
(i.e. direction of gradient)

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2$$

LMS update Rule

- if you have only one training example (x, y)

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2$$

$$= 2 \cdot \frac{1}{2} (h(x) - y) \frac{\partial}{\partial \theta_j} (h(x) - y)$$

$$= (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right)$$

$$= (h(x) - y) x_j$$

- for a single training eg., this gives the update rule

$$\beta_j = \beta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)}$$

$$h(x) = \sum_{i=0}^n \beta_i x^i$$

To learn the parameters $\theta(\beta_i)$?

make $h(x)$ close to y , for the available training example

- Define a cost $J(\theta)$,

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

- find θ that minimizes $J(\theta)$,

LMS update rule

- If you have only one training example x, y)

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2$$

$$= x \cdot \frac{1}{2} (h(x) - y) \frac{\partial}{\partial \theta_j} (h(x) - y)$$

$$= (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x^i - y \right)$$

$$= (h(x) - y) x_j$$

- for a single training eg, this gives the update rule

$$\beta_j' = \beta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)}$$

in training examples

Repeat until convergence {

$$\theta_j = \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x_j^{(i)}$$

Batch Gradient Descent: Looks at every example on each step

Stochastic gradient descent

- repeatedly run through the training set
- whenever a training pt is encountered update the parameters according to the gradient of the error with respect to that training example only

Repeat {

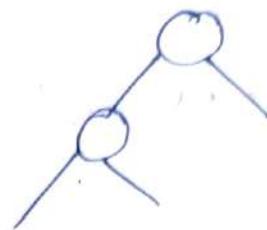
for $I = 1$ to m do

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)}$$

for every J , end for until converge

Chap 6 Induction of Decision Tree

Decision Tree classifier

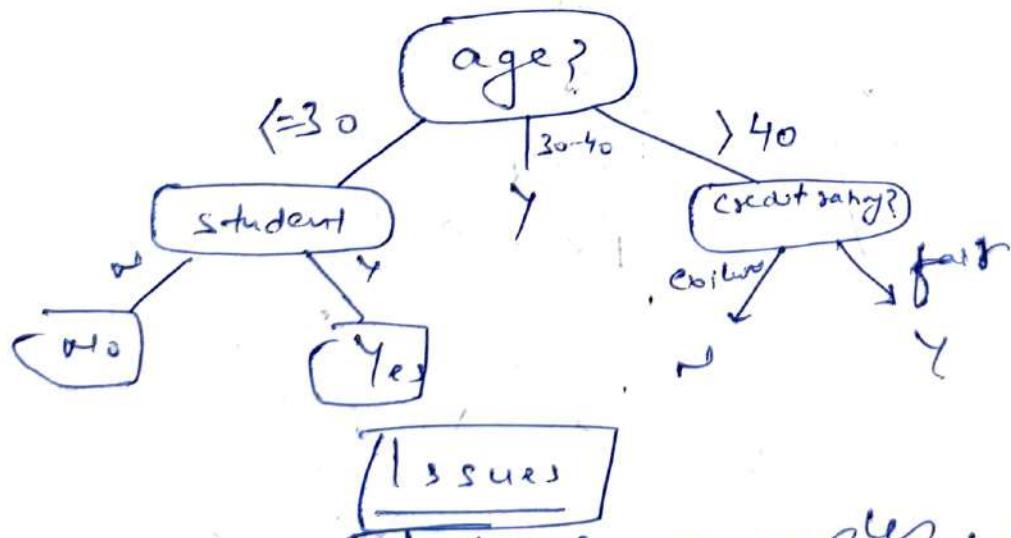
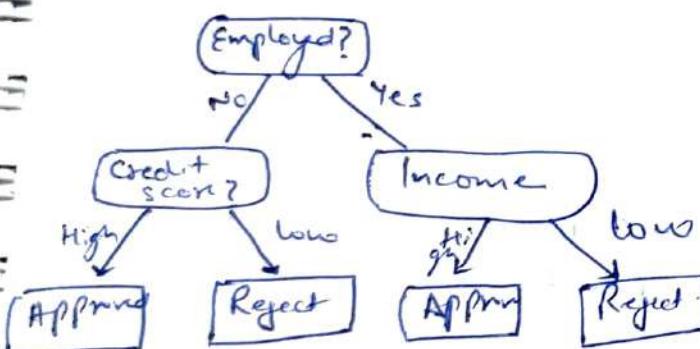


tree structure

- Decision Nodes →
- Leaf Nodes



classification value



- Given some training examples, what decision tree should be generated
- One proposal: prefer the smallest tree that consistent with the data (B-93)

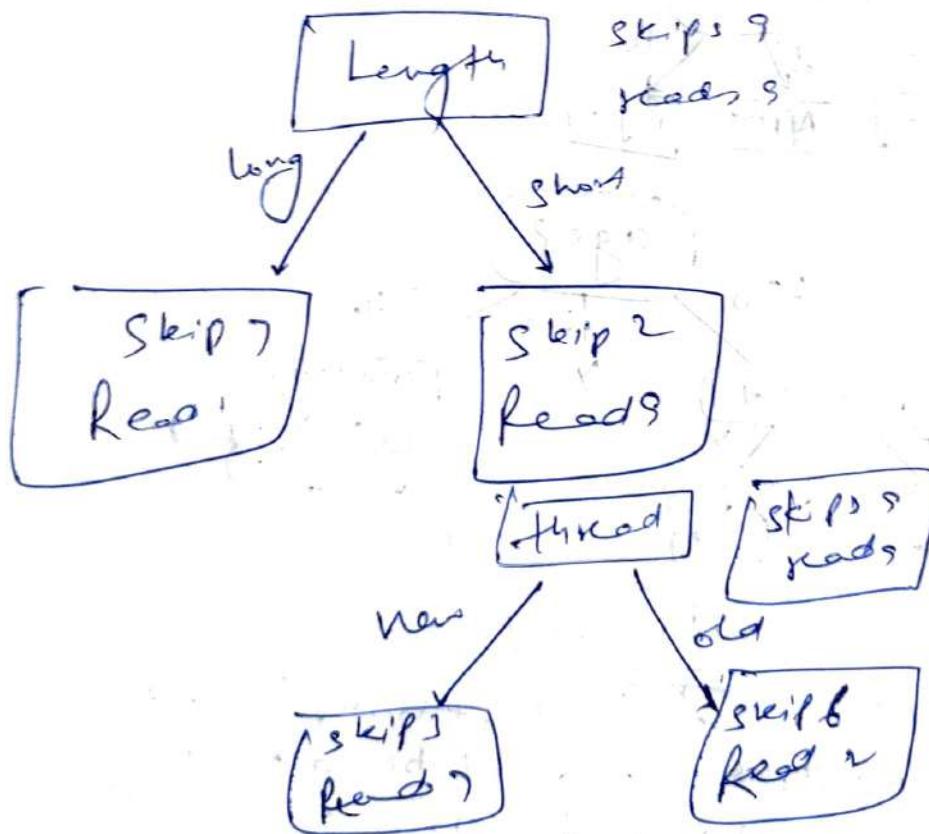
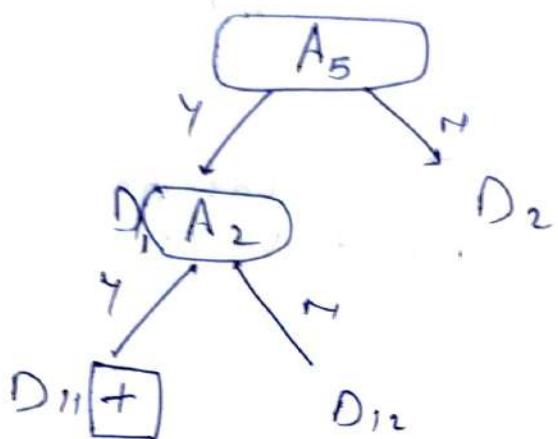
— Possible method.

— Search the space of decision trees for the smallest decision tree that fits the data.

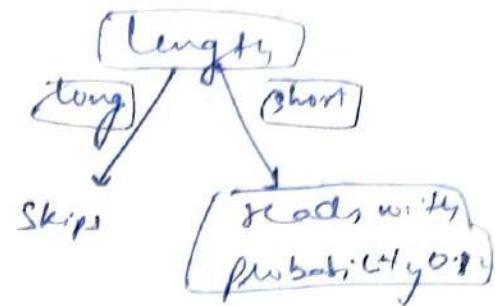
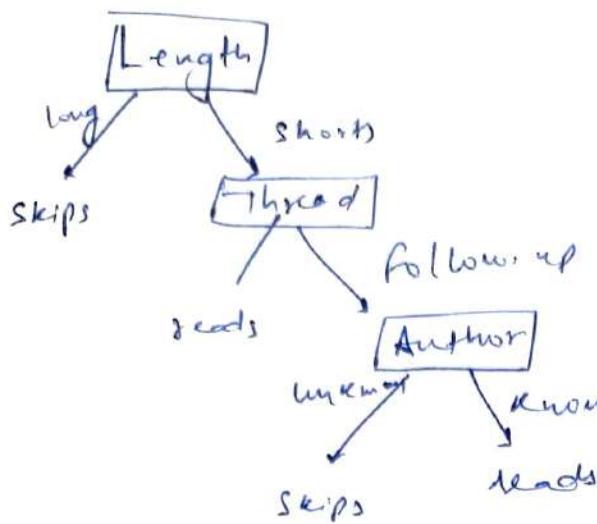
Prefer smaller tree
low depth

— small no. of nodes

Possible split



Two examples DTs



Decision Tree for PlayTennis

• Attributes and their values

- outlook: sunny, overcast, rain
- humidity: high, normal
- wind: strong, weak
- Temperature: hot, mild, cool

• Target concept - Play tennis: Yes, No.

Searching for a good tree

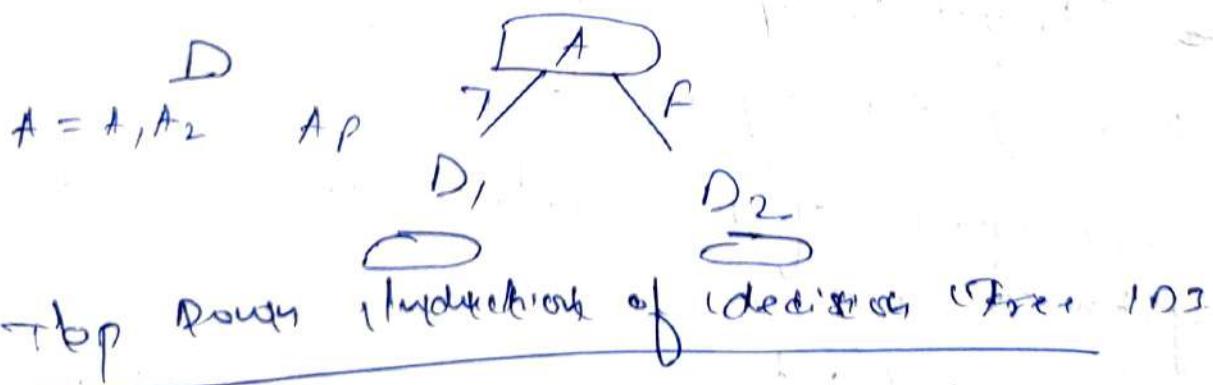
- the space of decision tree is too big for systematic search
- stop and return.
 - return the value for the target feature or
 - a distribution over target feature value
- choose a test (e.g. an input feature) to split on
 - for each value of the test, build a subtree for those examples with this value for the test

Top-Down Induction of Decision Tree ID3

1. $A \leftarrow$ the "best" decision attribute for next node
2. Assign A as decision attribute for node
3. for each value of A create new descendant
4. sort training examples to leaf node according to the attribute values of the branch
5. If all training examples are perfectly classified (same value of target attribute) stop, else iterate over leaf nodes

- ① which node to proceed with?
- ② when to stop.

Chapter 7: Learning Decision tree



Choices

- when to stop
 - no more input features
 - all examples are classified the same
 - too few examples to make any informative split.

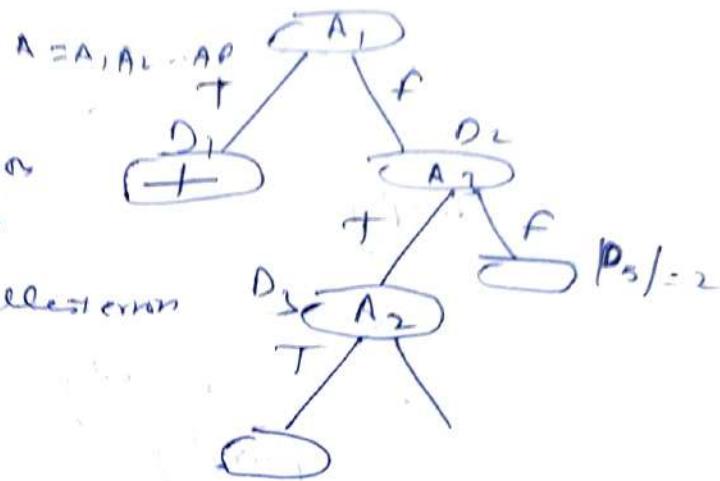
— which attribute to split on

- split gives smallest error
- with multi-valued feature
 - split on all values or
 - split values into half.

① when to stop

② which attribute to split on

1. split gives smallest error



which Attribute is best?

$$A_1 = ? \quad 28+, 35-$$

$$A_2 = 8+, 29+, 35-$$

T

F

$$24+, 5-$$

$$8+, 30-,$$

$$18+, 33-$$

$$(11+, 12-)$$

Principled criterion

- selection of an attribute to test at each node choosing the most useful attribute for classifying examples
- Information gain

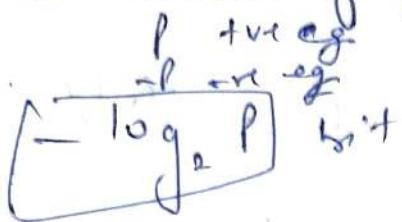
- measures how well a given attribute separates the training examples according to their target classification
- This measure is used to select among the candidate attributes at each step while growing the tree.
- Gain is measure of how much we can reduce uncertainty (Value lies b/w 0%)

Entropy

- A measure for
 - uncertainty
 - purity
 - information content
- Information theory, optimal length code assigns $(-\log_2 p)$ bits to message having probability p .
- S is a sample of training example
- p_i is the proportion of the examples in S
- p_r is the proportion of r examples in S
- Entropy of S : average of optimal no. of bits to encode information about certainty / uncertainty about S

$$\text{Entropy } (S) = P_+ (-\log_2 P_+) + P_- (-\log_2 P_-) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

Entropy \rightarrow measure of purity



$$P_+ (-\log_2 P_+) + P_- (-\log_2 P_-) \\ = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

- S is a sample of training examples
- P_+ is the proportion of +ve examples
- P_- is the proportion of -ve examples
- Entropy measures the impurity of S

$$\text{Entropy } (S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

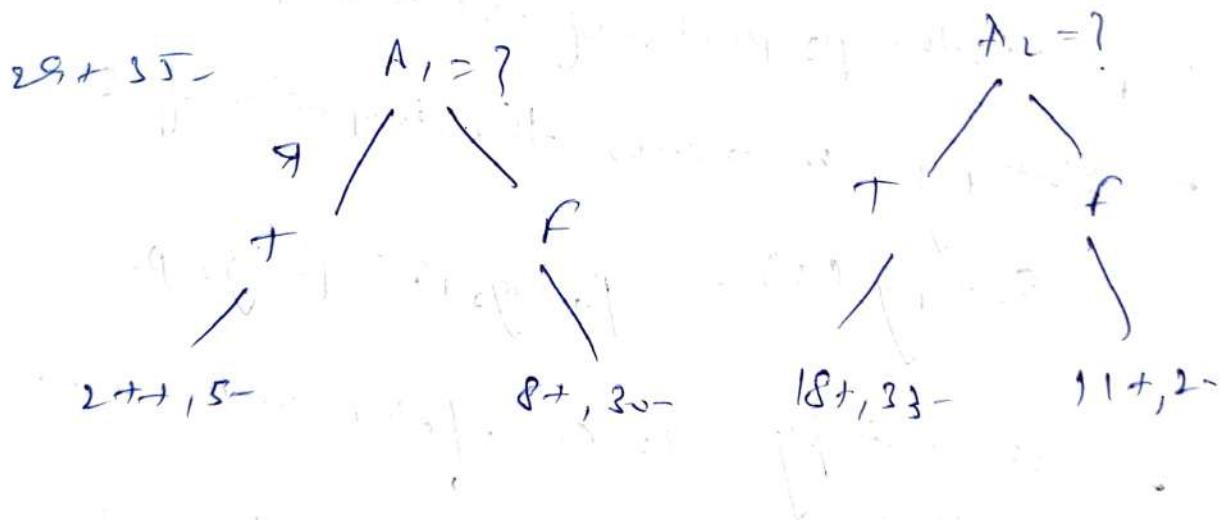
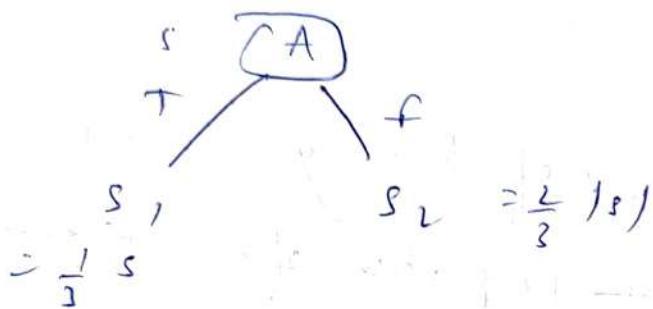
- The entropy is 0 if the outcome is certain
- The entropy is maximum if we have no knowledge of the system (or any ~~outcomes~~ outcome is equally possible)

Information Gain

$\text{Gain}(S, A)$: expected reduction in entropy due to partitioning S on attribute A

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{\text{values}(A)} \frac{|\text{level}|}{|\text{S}|} \text{Entropy}(\text{level})$$

$$\begin{aligned} \text{Entropy}([29+, 35-]) &= -29/64 \log_2 29/64 - 35/64 \log_2 35/64 \\ &= 0.99 \end{aligned}$$



$$\text{Entropy}([21+, 5-]) = 0.71$$

$$\text{Entropy}([8+, 30-]) = 0.75$$

$$\text{Gain}(S, A_1) = \text{Entropy}(S)$$

$$-26/64 * \text{Entropy}([21+, 5-])$$

$$-38/64 * \text{Entropy}([8+, 30-])$$

$$= 0.27$$

$$\text{Entropy } ([18+, 33-]) = 0.74$$

$$\text{Entropy } ([8+, 30-]) = 0.62$$

$$\text{Gain } (S, A_2) = \text{Entropy}(S)$$

$$- 51/64 * \text{entropy } ([10+, 33-])$$

$$- 13/64 * \text{entropy } ([11+, 2-])$$

$$= 0.12$$

ID3 Algorithm

splitting rule: GINI Index

- GINI index

— measure of node impurity

$$GINI_{node} (Node) = 1 - \sum_{c \in \text{classes}} [P(c)]^2$$

$$GINI_{split} (A) = \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} GINI(S_v)$$

- splitting Based on continuous Attribute

— Continuous Attribute → Binary split

— For continuous attribute

— Partition the continuous value of attribute A into a discrete set of intervals

— Create a new boolean attribute Ac,

looking for a threshold c,

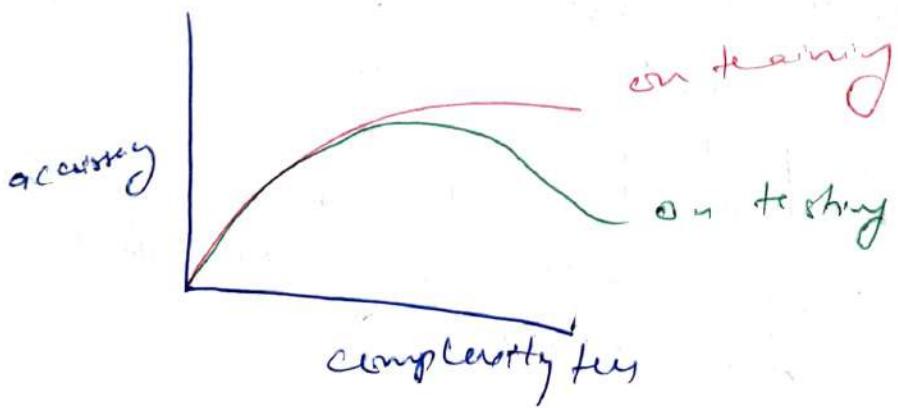
$$A_c = \begin{cases} \text{true} & \text{if } A_c < c \\ \text{false} & \text{otherwise} \end{cases}$$

- How to choose c ?
- consider all possible splits and find the best cut

Lecture 8 overfitting

A hypothesis h is said to overfit the training data if there is another hyp. h' on such that h' has more error than h on training data, h' has less error than h on test data

- Learning a tree that classifies the training data perfectly may not lead to the best generalization performance tree with the there may be noise in the training data
- may be based on insufficient data



overfitting due to insufficient examples

Notes on Overfitting

- overfitting results in decision trees that are more complex than necessary Dj pop
- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records.

Avoid Overfitting

- How can we avoid overfitting a decision tree?
 - Pruning: stop growing when data split not statistically significant
 - Postpruning: grow full tree then remove nodes
- methods for evaluating subtrees to prune
 - Minimum description length (MDL)
 - minimize: size(tree) + size(misclassification(tree))
 - cross validation

Pr-Pruning (Early stopping)

- Evaluate splits before installing them:
 - Don't install splits just because look worthwhile
 - when no worthwhile split to install,
 - don't

- Typical stopping conditions for nodes:
 - stop if all instances belong to the same class
 - stop if all the attribute values are the same
- more restrictive conditions
 - stop if no. of instances is less than some user specified threshold
 - stop if class distribution of instances are independent of the available features (e.g. using χ^2 test)
 - stop if expanding the current node does not improve impurity measures (e.g. Gini, or information gain)

Reduced error Pruning

- A post-pruning, cross validation approach
 - partitions training data into grow set and validation set. Use tree for the grow data
 - build a comp.
 - until accuracy on validation set \downarrow less

do:

for each non leaf node in tree
 Temporarily pruned tree below; replace it by majority vote

Test the accuracy of the hypothesis on the validation set.

Permanently prune the nodes with the greatest error in accuracy.

Problem: Uses less data to construct the tree sometimes done at the rules level.

General strategy: overfit and simplify

Triple Trade off

There is a trade off b/w three factors.

- complexity of $H, c(H)$

- Training set size, n

- Generalization error, E on new data

As $n \rightarrow \infty$, $E \downarrow$

As $c(H) \rightarrow \infty$, first $E \downarrow$ and then $E \uparrow$

As $c(H) \rightarrow \infty$, the training error \downarrow for constant H and they stay constant (frequently at 0)

Dealing with overfitting

- use more data

- use a tuning set

- regularization

- Be a Bayesian

Regularization in a linear regression model
overfitting is characterized by large weights.

Penalize large weights in linear regression

- Introduce a penalty term in the loss function.

$$E(\vec{w}) = \frac{1}{2} \sum_{n=0}^{N-1} \{(t_n - y(x_n, \vec{w}))\}^2$$

Regularized regression

1. (L_2 -regularization or ridge regression)

$$E(\vec{w}) = \frac{1}{2} \sum_{n=0}^{N-1} (t_n - y(x_n, \vec{w}))^2 + \frac{\lambda}{2} \|\vec{w}\|^2$$

1. L_1 Regularization

$$E(\vec{w}) = \frac{1}{2} \sum_{n=0}^{N-1} (t_n - y(x_n, \vec{w}))^2 + \lambda |\vec{w}|,$$

Chapter 9 Python Exercise 03
Decision tree and Linear Regression

scikit

Linear Regression

Step 1: visualization of the dataset

This is a synthetic dataset. The function used is $y = \frac{x}{2} + \text{R}(0, 1) + N(0, 1)$

$$y = \frac{x}{2} + \text{R}(0, 1) + N(0, 1)$$

Step 2: split the data
we split the data into
- training 70%
- validation 15%
- Test 15%
etc.

Step 3: fit a linear model
we fit a line that minimizes the
squared error on the training set

— Visualization of the model
The best line is the one that fits the
model has learnt to fit to the data. It
minimizes the squared error on the training set

Step 4: Evaluation of Mean Sq. error
The validation and test error give an
estimate of how well the model would
generalize to unseen data! the lower the better

Logistic regression

Classification by a linear decision boundary

Step 1: Visualization of dataset

- we will use a diminished version of iris dataset in this example
- we will consider only two input features
 - sepal length (in cm)
 - sepal width (in cm)
- And only two output classes
- Iris setosa
- Iris versicolor

Step 2 move during validation and

test split

Training - 70%, Validation - 15%, Test - 15%

Step 3: fit the logistic regression model

we fit a decision boundary that minimizes
the classification error

Step 4 Visualize the decision boundary

we show the idea of the feature
space assigned to the two different classes
blue belongs to class 0 and orange
class 1

Step 5 Evaluate the model

10% validation and test error was
only possible because the classes were
linearly separable

Decision tree regression

are used to learn piecewise linear models to fit the data

Step 1: Data preparation

we use the same synthetic dataset and make the same kind of train, validation and test splits as in the linear regression example.

Step 2: Fit the model

we study models of different maximum tree depth

Step 3: Evaluation and visualization

with very depth, the model fits tighter to the data.

Variation of error with max. depths of tree

we see that overfitting sets in for max_depth >?

Decision tree classification

learn piece wise linear decision boundaries

Step 1: Prepare dataset

Fit model

Step 3: Fit visualization the model

Evaluate model.

Tutorial 11

- Linear regression
- decision tree

Linear regression

Q:- Prepare you have given a set of training examples

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

Find the eqn. of the line that best fit the data in the sense that it minimize the sum of squares

Ans 1: Estimation of sq. error

$$J = \sum_{i=1}^n [y_i - f(x_i)]^2$$

where $f(x)$ is your linear regression hypothesis

Here $f(x) = mx + c$

$$f: X \rightarrow Y$$

$$y = mx + c = f(x)$$

$$J = \sum_{i=1}^n [y_i - m_i - c]^2$$

minimizing
w.r.t m and c
parameters of
straight line

$$1. \frac{\partial J}{\partial w} = 0 \Rightarrow 2 \sum_{i=1}^n [y_i - w x_i - c] (-x_i) = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^n [y_i - w x_i - c] x_i = 0} \quad \textcircled{C}$$

$$2. \frac{\partial J}{\partial c} = 0 \Rightarrow 2 \sum_{i=1}^n [y_i - w x_i - c] [-1] = 0.$$

$$\Rightarrow \boxed{\sum_{i=1}^n [y_i - w x_i - c] = 0} \quad \textcircled{D}$$

Decision tree

$$\left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] \quad \begin{array}{l} \text{Data point} \\ \in \mathbb{R}^n \end{array}$$

n : dimensionality of
the feature space

You have n features
of input data

1. A feature f_i

2. A pt. of split on that feature
axis f_i

such that the end of data pt. going into
the child nodes are homogeneous for most

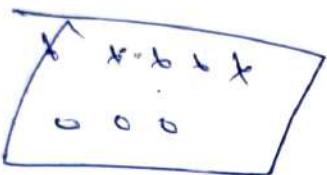
Entropy of set say we have pt. from n diff

classes is n , say the probability / fraction
of pt. of the different classes are

$$\{f_1, f_2, \dots, f_n\} \text{ sum to } f_1 + f_2 + \dots + f_n = 1$$

Entropy of this set is given by

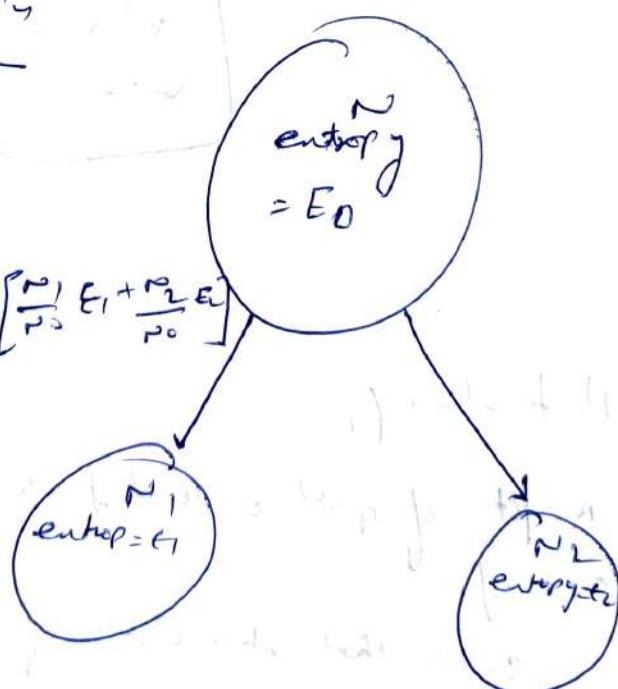
$$E = \sum_{i=1}^n f_i \log_2 f_i \text{ bits}$$

Eg  $f_{x=1} = \frac{5}{8}, f_{x=0} = \frac{3}{8}$

$$\text{Entropy} = -\frac{5}{8} \log \frac{5}{8} - \frac{3}{8} \log \frac{3}{8}$$

Information Gain

$$\text{Information Gain} = E_0 - \left[\frac{N_1}{N} E_1 + \frac{N_2}{N} E_2 \right]$$

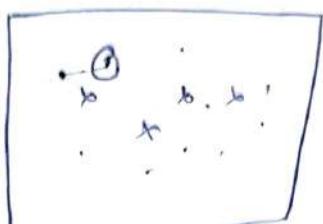


Chapter 10 K-Nearest Neighbour

(x_i, y_i)

$(x_i, f(x_i))$

Instance based learning
Lazy algorithm.



Find similar instances

Training phase :- save the example

Prediction time :- Get the test instance x_t
find the training ex. $(x_1, y_1), (x_2, y_2)$
that is closest to x_t

Predict y , as the output y_t

Classification predict the majority class from $\{y_1, y_2, \dots, y_n\}$

Regression predict the average of $\{y_1, y_2, \dots, y_n\}$

Voronoi diagram

Basic K-nearest neighbour classification

- Training method
 - save the training examples
- At prediction time!
 - find the k training examples $(x_1, y_1), \dots$
 - find the k training examples $(x_1, y_1), \dots, (x_k, y_k)$ that are closest to the test example (x_t, y_t)
 - classification: predict the average of among y_1 's

- Improvements:

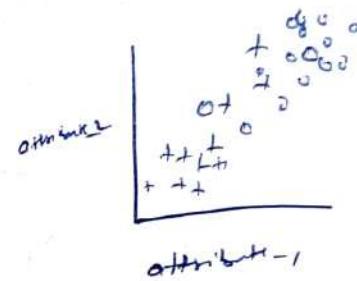
- weighting examples from the neighbourhood
- measuring "closeness"
- finding "close" examples in large training set quickly

K-Nearest Neighbors

$$\text{Dist}(c_1, c_2) = \sqrt{\sum_{i=1}^n (\text{attr}_i(c_1) - \text{attr}_i(c_2))^2}$$

Prediction $\hat{y}_{\text{test}} = ?$

- Average of k pt. more reliable when:
- noise in attributes
 - noise in class labels
 - classes partially overlap



$$x_i = (x_{i1}, x_{i2}, \dots, x_{in})$$

$$m_j = (m_{j1}, m_{j2}, \dots, m_{jn})$$

Dist. Euclidean (x_i, m_j)

$$= \sqrt{\sum_{k=1}^n (x_{ik} - m_{jk})^2}$$

- noise in attributes
- noise in class label
- classes may partially overlap

Equal weights to all attributes?

- only if the scale of the attributes are similar and the different attributes have equal range and equal variance

- classes are spherical in shape

- use larger k

- use weighted distance function

Small k captures fine structure of problem space
better may be necessary for small training set

Large k less sensitive noise (particularly class noise)
better prob. estimates for discrete classes

- larger training set allows use of large k

weighted Euclidean Distance

$$D(c_1, c_2) = \sqrt{\sum_{i=1}^n w_i \cdot (\text{attr}_i(c_1) - \text{attr}_i(c_2))^2}$$

- large weights \Rightarrow attribute is most important
- small weights \Rightarrow attribute is less important
- zero weight \Rightarrow attribute doesn't matter
- weights allow KNN to be effective with axis-parallel elliptical classes
- where do weights come from?

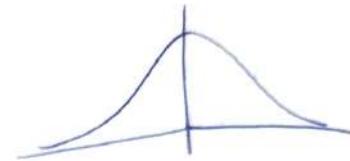
Distance - weighted KNN

- tradeoff between small and large k
can be difficult
- use large k , but more emphasis on nearer

neighbors?

$$\text{prediction}_{\text{test}} = \frac{\sum_{i=1}^k w_i * \text{class}_i}{\sum_{i=1}^k w_i} \quad (\text{or} \quad \frac{\sum_{i=1}^k w_i * \text{value}_i}{\sum_{i=1}^k w_i})$$

$$w_k = \frac{1}{\text{Dist}(c_k, c_{\text{test}})}$$



Locally Weighted Averaging

- Let k = no. of training points.
- Let weights fall-off rapidly with distance

$$\text{Prediction}_{\text{test}} = \frac{\sum_{i=1}^k w_i * \text{class}_i}{\sum_{i=1}^k w_i} \quad (\text{or} \quad \frac{\sum_{i=1}^k w_i * \text{value}_i}{\sum_{i=1}^k w_i})$$

$$w_k = \frac{1}{\text{kernel width Dist}(c_k, c_{\text{test}})}$$

- Kernel width controls size of neighborhood that has large effect on value (analogous to k)

Chapter 11 feature selection

Feature Reduction

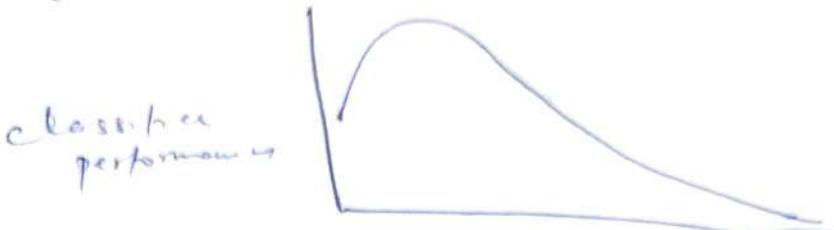
Features contain information about target

More features \Rightarrow More information

Better discriminative power

Curse of Dimensionality

- no. of training egs is fixed
 \Rightarrow the classifier's performance usually will degrade for a large no. of features!



Irrelevant features

\rightarrow Limited range of

Redundant features

\rightarrow Limited Computational

Feature reduction

resources

Feature selection

$$F = \{x_1, x_2, \dots, x_n\}$$

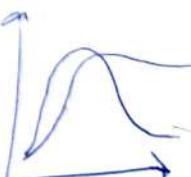
$$F' \subset F = \{x'_1, x'_2, \dots, x'_{n'}\}$$

- Feature Extraction

Project to $m < n$ dimensions

optimizes \rightarrow either Improve or maintain accuracy

\searrow Simplify classifier complexity



2nd Possible subsets

- optim methods

- Heuristic methods

- Randomized

Evaluation

- unsupervised - clustering
- supervised - wrappes

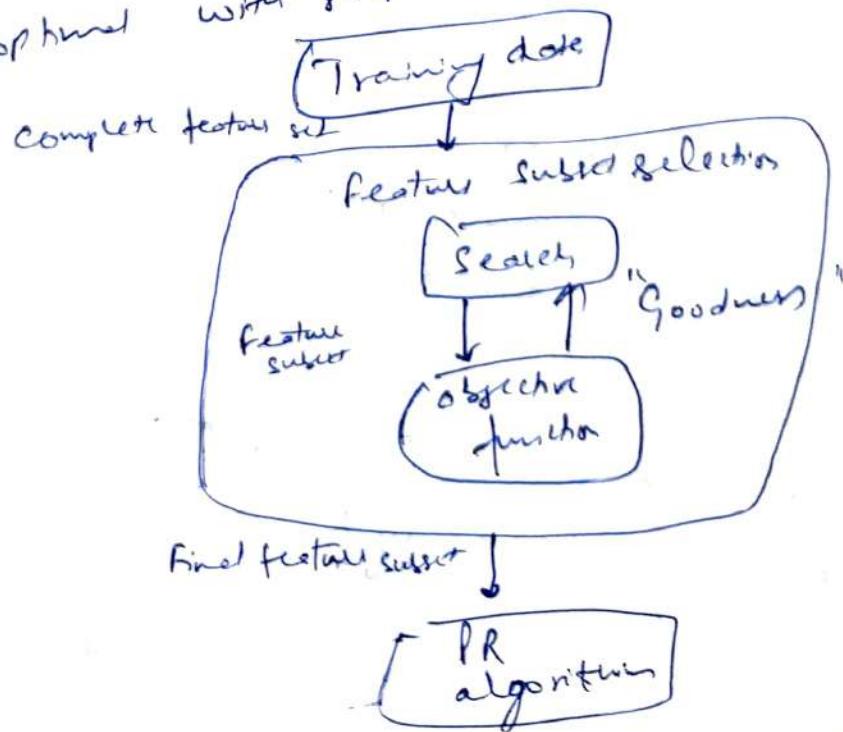
Subset selection

- d initial features
- there are 2^d possible subsets
- criteria to decide which subset is the best
 - classifier based on these m features has the lowest probability of error of all such classifiers.
- can't go over all 2^d possibilities
need some heuristics.

Feature selection steps

Feature selection is an optimization problem

- step 1: search the space of possible feature subsets
- step 2: pick the subset that is optimal or near optimal with respect to some objective function



Forward selection

- start with empty feature set
- Try each ~~removing~~ remaining feature
- Estimate classification/seg. error for adding each feature
- select feature that gives max. improvement
- stop when there is no significant ~~stop~~ improvement

Backward search

- start with full feature set
- Try removing features
- Drop feature with smallest impact on error

feature selector

univariate (look at each feature independently
of others)

- Pearson correlation coefficient

- f-score

- chi-square

- signal to noise ratio

- mutual information

- etc

Rank features by importance
fanning cut-off is determined by user

Univariate methods
measure some type
of correlations b/w
two random variables
- the label (y_i)
and fixed feature
(x_{ij} for $j=1 \dots d$)

Pearson correlation coefficient

- Measure the correlation b/w two variables
- formula for Pearson correlation:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- The correlation r is b/w +1 and -1
 - +1 means perfect +ve correlation
 - -1 ~~means~~ in other direction.

Signal to noise ratio

- Difference in means divided by difference in standard deviation b/w the two classes

$$S_{2n}(x, y) = (\mu_x - \mu_y) / (\sigma_x - \sigma_y)$$

- Large values indicate a strong correlation

Multivariate feature selection

- Multivariate (considers all features simultaneously)
- consider the vector w for any linear classifier
- classifier classification of x opt w is given by $w^T x + b$
- small entries of w will have little effect
- on the dot product and therefore those features are less relevant
- for example if $w = (10, 0.01, -9)$ then feature 0 and 2 are contributing more to the dot product than feature 1
- A ranking of features given by their w_1, w_2, \dots

- The w's can be obtained by any of linear classifiers
- A variant of this approach is called margin feature elimination
 - compute w on all features
 - Remove feature with smallest w_j
 - Recompute w on reduced data
 - If stopping criterion not met then go to step 2

Chapter 12 feature extraction

$$\begin{matrix} \mathbf{x} \\ \left[\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \right] \end{matrix} \longrightarrow \begin{matrix} \mathbf{z} \\ \left[\begin{matrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{matrix} \right] \end{matrix} = f \left[\begin{matrix} x_1 \\ \vdots \\ x_n \end{matrix} \right]$$

find a projection matrix \mathbf{w}

~~**~~

$$\tilde{\mathbf{z}} = \mathbf{w}^T \mathbf{x}$$

- Uncorrelated, cannot be reduced further
- Large variance

Geometrical picture of Principal components

Algebraic definition of PCs

Given a sample of p observations in vector of M variables

$$\{x_1, x_2, \dots, x_p\} \in \mathbb{R}^n$$

define the principal component of the sample by the

linear transformation

$$z_i = w_i^T x_j = \sum_{i=1}^n w_{ij} x_{ij}, \quad j = 1, 2, \dots, p$$

where the vector $w_i = (w_{i1}, w_{i2}, \dots, w_{ip})$

$$x_j = (x_{1j}, x_{2j}, \dots, x_{nj})$$

is chosen such that $\text{Var}[z_i]$ is maximum

$$x_1 = \{ \quad \}$$

$$x_2 = \{ \quad \}$$

$$x_3 = \{ \quad \}$$

$$\vdots$$

$$x_p = \{ \quad \}$$

$\text{Var}[z_i]$ is max.

PCA 2

feature 2

PCA 1

feature 1

(orthogonal)

①

maximize total variance

②

minimize correlation

choose $m \leq n$ orthogonal directions
which maximize total variance

Train data $\text{Cov}(x) = \Sigma$
select m largest eigenvector

- n -dimensional feature space
- $n \times n$ symmetric covariance matrix estimated from sample $\text{Cov}(x) = \Sigma$
- select m largest eigenvalues of the covariance matrix and associated eigenvectors
- first eigenvector will be direction with largest variance

Is PCA a good criterion for classification?

- Data variation determined the projection direction.

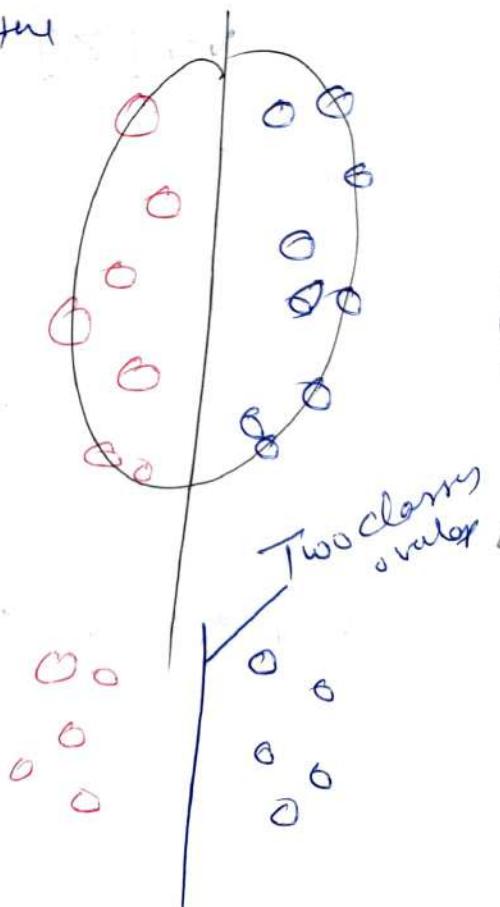
what's missing

- class information

what's a good projection?

- similar, what is a good criterion?

- separating different classes



What class information may be useful?

- B/w - class distance
 - Distance b/w the centroids of different classes

within-class distance

- Accumulated distance of an instance to the centroid of its class
- Linear discriminant analysis (LDA) find most discriminant projection by
 - maximizing b/w-class distances
 - and minimizing within-class distance

Linear Discriminant Analysis

- find a low-dimensional space such that when x is projected, classes are well separated

Good projections

- means are as far away as possible
- scatter is small as possible

Fischer Linear Discriminant

$$J(\omega) = \frac{(m_1 - m_2)^2}{S_1^2 + S_2^2}$$

Chapter 13 collaborative filtering

Recommended system

① Item recommended

② Rating prediction system

$U = \text{set of users}$ $p: U \times I \rightarrow R$

$I = \text{set of items}$ - Learn p from data

- user p to predict the utility
value of each item to each user

① Content Based

- Based on Content Similarity

② - collaborative filter system

- similar user

	s_1	s_2	s_3	s_4	s_5	s_m
u_1	5	-	-	2	3	
u_2	2	2	-	(2)	4	
u_3	2	2	4	-	-	
u_4						
u_5						
u_n						

User based nearest method

Item based method

collaborative filtering for Rating Prediction

- User-based Nearest Neighbour

- Neighbour = similar user

- Generates a prediction for an item
 i by analyzing rating for i from user in
 u 's neighbourhood

Neighborhood formation phase

- Let the record of the target user be u , and the record of another user be v (vet)
- The similarity b/w the target user, u , and neighbor, v , can be calculated using Pearson's correlation coefficient

$$\text{sim}(u, v) = \frac{\sum_{i \in C} (r_{u,i} - \bar{r}_u)(r_{v,i} - \bar{r}_v)}{\sqrt{\sum_{i \in C} (r_{u,i} - \bar{r}_u)^2} \sqrt{\sum_{i \in C} (r_{v,i} - \bar{r}_v)^2}}$$

Recommendation Phase

- Use the following formula to compute the rating prediction of item i for target user u

$$p(u, i) = \bar{r}_u + \frac{\sum_{v \in N} \text{sim}(u, v) \times (r_{v,i} - \bar{r}_v)}{\sum_{v \in N} |\text{sim}(u, v)|}$$

where N is the set of k similar user, $r_{v,i}$ is the rating of user v given to item i ,

Issue with the user-based KNN CF

- The problem with the user-based formulation of collaborative filtering is the issue of scalability:

- it requires the real time comparison of the target user to all user record in order to

generate predictions

- A variation of this approach that remedies this problem is called item based CF

- Loss of scalability

- item based CF

by comparing items based on their pattern of rating across user. The similarity of item i and j is computed as follows

$$\text{sim}(i, j) = \frac{\sum_{u \in U} (r_{u,i} - \bar{r}_u)(r_{u,j} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,i} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,j} - \bar{r}_u)^2}}$$

Chapter 19 Python exercise PCA and KNN

Chpt 15

Chpt 16

Bayesian learning

Probability : modeling concepts

Bayesian Probability : Partial beliefs

Bayesian Estimation, calculate validity of a proposition

- Prior estimate
- New relevant evidence

Bayes Theorem

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{P(D)} \rightarrow \text{Prior}$$

$$P(h|D) = P(h) \cdot P(D|h)$$

$$P(D|h) = P(D) \cdot P(h|D)$$

Bayes theorem:

Goal: to determine the most probable hypothesis given the data D plus any initial knowledge about the prior probability of various hypotheses with

An example Does patient have cancer or not?
A patient takes a test and the result comes back positive. The test returns a correct result in only 98% of the cases in which the disease is actually present, and a correct - result in only 97% of the cases in which the disease is not present. Furthermore, 0.002% of the entire population have this cancer.

$$P(\text{cancer}) = 0.002, P(\text{-cancer}) = 0.998$$

$$P(+|\text{cancer}) = 0.98, P(-|\text{cancer}) = 0.02$$

$$P(+|\text{-cancer}) = 0.02, P(-|\text{-cancer}) = 0.98$$

$$P(\text{cancer} | +) = \frac{P(+ | \text{cancer}) P(\text{cancer})}{P(+)}$$

$$P(-, \text{cancer} | +) = \frac{P(+ | \text{-cancer}) P(\text{-cancer})}{P(+)}$$

Maximum A Posterior (MAP) hypothesis

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

The goal of Bayesian learning: the most probable hypothesis given the training data (Maximum A posterior hypothesis)

$$h_{\text{MAP}} = \arg \max_{h \in H} P(h | D)$$

$$= \arg \max_{h \in H} \frac{P(D|h) P(h)}{P(D)}$$

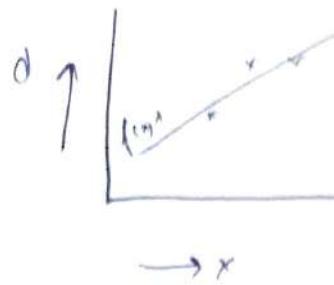
$$= \arg \max_{h \in H} P(D|h) P(h)$$

Learn a real valued function

f. Target function

$$(x_i, d_i) | d_i = f(x_i) + \epsilon_i$$

$$d_i \sim N(f(x_i), \sigma^2)$$



$$h_{ML} = \arg \max_h P(D|h)$$

$$= \arg \max_h \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{1}{2} \left(\frac{d_i - h(x_i)}{\sigma} \right)^2}$$

$$= \arg \max_h \sum_{i=1}^m -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2} \left(\frac{d_i - h(x_i)}{\sigma} \right)^2$$

$$\arg \max_h \sum_{i=1}^m (d_i - h(x_i))^2$$

Bayes Optimal Classifier

Training data h_{MAP}

$$h_1, h_2, h_3 \quad P(h_1|D) = .4$$

$$P(h_2|D) = .3$$

$$P(h_3|D) = .3$$

x $h_1(x) = +$
 $=$ $h_2(x) = -$
 $h_3(x) = -$

$$\text{argmax}_{\{h_i\}} \sum_{v_j \in v} P(v_j | h_i), P(h_i | D)$$

Bayes Optimal Classifier

- Q: Given new instance x , what is the most probable classification
- $h_{MAP}(x)$ is most probable classification

Gibbs Sampling

Choose a hypothesis randomly acc to $P(h|D)$

Use it to classify the new instance

Chpt. 17 Naive Bayes

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

$$P(Y|x) \propto P(\bar{x}|Y) \cdot P(Y)$$

$$P(x_1, x_2, \dots, x_n | Y) \cdot P(Y)$$

Naive Bayes

$$\text{Assumption} \quad = P(x_1 | Y) \cdot P(x_2 | x_1, Y) \cdot P(x_n | x_1, \dots, x_{n-1}, Y) \cdot P(Y)$$

$$= P(x_1 | Y) \cdot P(x_2 | Y) \cdot P(x_3 | Y) \cdot P(Y)$$

Bayes rule:

$$P(Y=y_k | x_1, \dots, x_n) = \frac{P(Y=y_k) P(x_1, \dots, x_n | Y=y_k)}{\sum_j P(Y=y_j) P(x_1, \dots, x_n | Y=y_j)}$$

Assuming conditional independence among x_i :

$$P(Y=y_k | x_1, \dots, x_n) = \frac{P(Y=y_k) \prod_i P(x_i | Y=y_k)}{\sum_j P(Y=y_j) \prod_i P(x_i | Y=y_j)}$$

So classification rule for $X^{\text{new}} = \langle x_1, \dots, x_n \rangle$ is:

$$y^{\text{new}} \leftarrow \arg \max_{y_k} P(Y=y_k) \prod_i P(x_i^{\text{new}} | Y=y_k)$$

Main Bayes Algorithm - discrete x_i

- Train Main Bayes Examples

for each value y_k

$$\text{estimate } \pi_k \equiv P(Y=y_k)$$

for each value x_{ij} of each attribute x_i

$$\text{estimate } \theta_{ijk} \equiv P(x_i = x_{ij} | Y=y_k)$$

- classify (x^{new})

$$y^{\text{new}} \leftarrow \arg \max_{y_k} P(Y=y_k) \prod_i P(x_i^{\text{new}} | Y=y_k)$$

$$y^{\text{new}} \leftarrow \arg \max_{y_k} \pi_k \prod_i \theta_{ijk}$$

- Probability must sum to 1, so need to estimate only $m-1$ parameters

Estimating Parameters: Y, X_i ; discrete-valued

Maximum likelihood estimate (MLE):

$$\hat{\pi}_k = \hat{P}(Y=y_k) = \frac{\#\{D \mid Y=y_k\}}{|D|}$$

~~$\hat{\theta}_{ijk} = \hat{P}(X_i=x_{ij} \mid Y=y_k)$~~

$$\hat{\theta}_{ijk} = \hat{P}(X_i=x_{ij} \mid Y=y_k) = \frac{\#\{D \mid X_i=x_{ij} \wedge Y=y_k\}}{\#\{D \mid Y=y_k\}}$$

no. of items in set D for which
 $Y=y_k$

Estimating Parameters: Y, X_i ; discrete-valued

If unlucky, our MLE estimating for $P(X_i \mid Y)$
may be zero

$$\hat{\pi}_k = \hat{P}(Y=y_k) = \frac{\#\{D \mid Y=y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i=x_{ij} \mid Y=y_k) = \frac{\#\{D \mid X_i=x_{ij} \wedge Y=y_k\}}{\#\{D \mid Y=y_k\}}$$

MAP estimates

$$\hat{\pi}_k = \hat{P}(Y=y_k) = \frac{\#\{D \mid Y=y_k\} + 1}{|D| + 1} \xrightarrow{\text{only diff.,}} \text{"imaginey example,"}$$

$$\hat{\theta}_{ijk} = \hat{P}(x_i=x_{ij} \mid Y=y_k) = \frac{\#\{D \mid x_i=x_{ij} \wedge Y=y_k\} + 1}{\#\{D \mid Y=y_k\} + 1}$$

Assumption of Conditional Independence

often the x_i are not really conditionally independent

- we can use Naive Bayes in many cases anyway
 - often the right classification, even when not the right probability

Gaussian Naive Bayes (Continuous x)

- Algorithm: continuous valued features
- conditional probability often modeled with the normal distribution

$$P(x_i=x \mid Y=y_k) = \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}}$$

Sometimes assume variance

- is independent of Y (i.e., σ_i)
- or independent of x_i (i.e. σ_k)
- or both (i.e., σ)

Bairn-Moor Bayes (example)

for each value y_k

$$\text{estimate } \pi_k = P(Y=y_k)$$

for each attribute X_i ; estimate

class conditional mean μ_{ik} values σ_{ik}

- classify (x^{new})

$$y^{\text{new}} \leftarrow \arg \max_{y_k} P(Y=y_k) \prod_i P(x_i^{\text{new}} | Y=y_k)$$

$$y^{\text{new}} \leftarrow \arg \max_{y_k} \pi_k \prod_i \text{Normal}(x_i^{\text{new}}, \mu_{ik}, \sigma_{ik})$$

Estimating Parameters: Y discrete, X ; continuous

maximum likelihood estimate!

by fitting curve

$$\hat{\mu}_{ik} = \frac{\sum_j x_i^j \delta(Y^j = y_k)}{\sum_j \delta(Y^j = y_k)}$$

(i^{th} feature)
 k^{th} class

$$\delta(2) = 18^2 \text{ true}$$

else 0

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (x_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

Unit 18 Bayesian Networks

- is a type of graphical model
- ✗ - $\{X_1, X_2, X_n\} \cup Y$

Nodes

Arcs Causality

conditional independence

①

Kat^t
wants to

Relationships

Accident

②

Cat for
work

Raining
Day

③

Traffic jam

④

Cat for
meaty

⑤

Waking up
postponed

DAG

27-1

Represents efficiently the joint probability distribution
of the variables

Arcs representation

- Arcs represent probabilistic dependencies among variables
- DAG
- A node has local probability distribution (CPT)

Conditional probability table associated with each node specifies the conditional distribution for the variable given its immediate Parent in the graph.

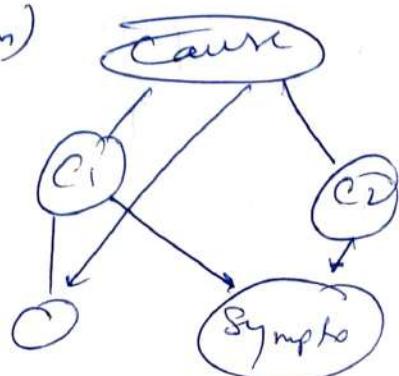
Each node is asserted to be conditionally independent of its non-descendant i , given its immediate parents.

Inference captures posterior probabilities given evidence about some nodes

- Exact Inference is tractable for arbitrary BN
- Approximate techniques eg Monte Carlo methods can also be used
- Efficient algorithms that leverage the structure of graphs

Application of Bayesian Networks

- Diagnosis: $P(\text{cause} \mid \text{symptom}) = ?$
- Prediction: $P(\text{symptom} \mid \text{cause}) = ?$
- Classification: $P(\text{class} \mid \text{data})$
- Decision-making
(given a cost function)



- structure of the graph \leftrightarrow conditional independencies
 - dense relations
- In general

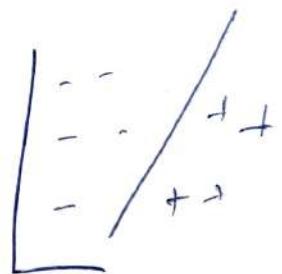
$$P(X_1, X_2, \dots, X_n) = \prod P(X_i | \text{parent}(X_i))$$

The full joint distribution

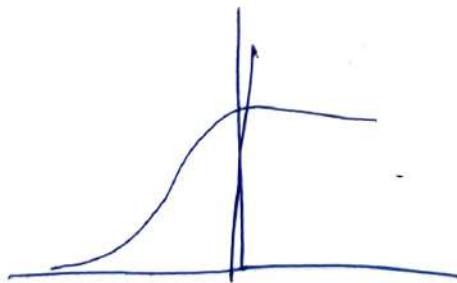
The graph-structure approximation

- requires that graph is acyclic
 - 2 components to a Bayesian network
 - the graph structures
 - The unnormalized probabilities
- Chpt Linear logistic regression

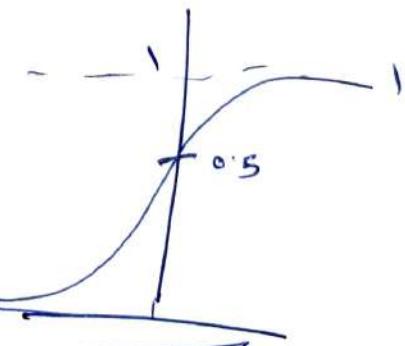
linear regression $h(x) = \sum_{i=0}^m \beta_i x_i$



Logistic Regression



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



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

$$h(x) = g(\beta^T x)$$

$g(z)$ sigmoid function

$$= \frac{1}{1 + e^{-\beta^T x}}$$

$$g^{(2)} \rightarrow 1 \text{ as } z \rightarrow \infty$$

$$g^{(2)} \rightarrow 0 \text{ as } z \rightarrow -\infty$$

$$g'(z) = \frac{1}{(1 + e^{-z})^2} * e^{-z}$$

$$= \frac{1}{1 + e^z} - \left(1 - \frac{1}{1 + e^{-z}}\right)$$

$$= g(z) \cdot (1 - g(z))$$

$$\boxed{p(Y|X)} = h(x)^y (1 - h(x))^{1-y}$$

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

$p_y(\cdot | \beta)$ is our estimate of $p(Y|X)$

Learn optimal values of β

$$L(\beta) = p(Y|X; \beta)$$

$$= \prod_{i=1}^m p(y_i|x_i, \beta)$$

$$= \prod_{i=1}^m \left(h(x_i)^{y_i} (1 - h(x_i))^{1-y_i} \right)$$

$$l(\beta) = \sum_{i=1}^m y_i \log(h(\mathbf{x}_i)) + (1-y_i) \log(1-h(\mathbf{x}_i))$$

$$\beta = \beta + \alpha \nabla_{\beta} l(\beta)$$

$\boxed{x_i, y}$

$$\frac{\partial}{\partial \beta} l(\beta) = \left(y \frac{1}{g(\beta^T x)} \right) - (1-y) \frac{1}{1-g(\beta^T x)} \left(\frac{\partial}{\partial \beta_j} g(\beta^T x) \right)$$

$$= y \frac{1}{g(\beta^T x)} - (1-y) \frac{1}{1-g(\beta^T x)} (1-g(\beta^T x)) \frac{\partial}{\partial \beta_j} g(\beta^T x)$$

$$= (y(1-g(\beta^T x)) - (1-y)g(\beta^T x)) \alpha_j$$

$$= (y - h_{\beta}(x)) \alpha_j$$

$$= (y - h_{\beta}(x)) \alpha_j$$

$$\beta_j = \beta_j + \alpha (y - h_{\beta}(x)) \alpha_j$$

Put above values
in eqn

$$\boxed{\beta = \beta + \alpha \nabla_{\beta} l(\beta)}$$

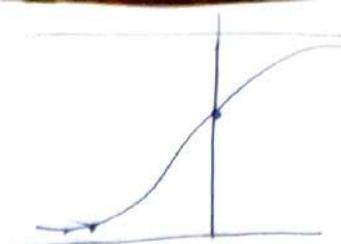
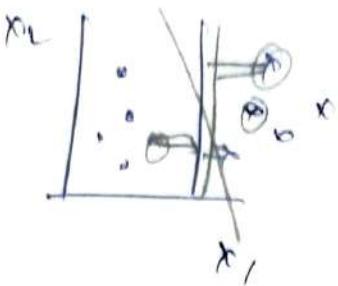
Lecture 20 Introduction to support vector machine

Logistic regression

$$P(y=1|x) = h(x) = \sigma(\beta^T x)$$

Predict 1 when $\bar{h}(x) > 0.5$

① when $h(x) < 0.5$



functional margin distance of (x_i, y_i) from (w, b)

$$c(x_i, y_i) \text{ wrt } (w, b)$$

$$\gamma^i = y_i (w^T x_i + b)$$

$$2w_1 + 3w_2 + 1 = 0$$

$$4w_1 + 6w_2 + 1 = 0$$

$$2w_1 + 3w_2 + 1 = 0$$

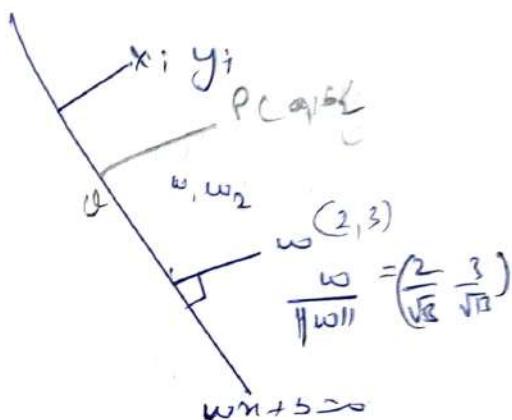
$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

$$\gamma = \min_{i=1}^m \gamma^i$$

Geometrical Margin

$$2w_1 + 3w_2 + 1 = 0$$

$$\rho = \frac{\gamma - b}{\|w\|}$$



$$(a_1, a_2) = (b_1, b_2) + \frac{w}{\|w\|}$$

$$w^T(a_1, a_2) - \frac{\|w\|}{\|w\|} + b = 0$$

$$b = \frac{w^T(a_1, a_2) + b}{\|w\|}$$

$$= \frac{w^T(a_1, a_2)}{\|w\|} + \frac{b}{\|w\|}$$

$$y = y \cdot \frac{w^T(a_1, a_2)}{\|w\|} + \frac{b}{\|w\|}$$

$$y = my$$

Maximizing margin width

- maximizing $\frac{y}{\|w\|}$ subject to

- $y_i(w^T x_i + b) \geq 1$ for $i = 1, 2, \dots, n$

- scale so that $y = 1$

- maximizing $\frac{1}{\|w\|}$ is the same as minimizing $\|w\|^2$

- minimizing $\|w\|^2$ subject to the constraint

- for all (x_i, y_i) , $i = 1, \dots, n$:
 $w^T x_i + b \geq 1$ if $y_i = 1$

$$\omega^T x_i + b \leq -1 \text{ if } y_i = -1$$

Large Margin Linear Classifiers

Formulation

$$\text{minimize}_{\omega} \frac{1}{2} \|\omega\|^2$$

such that

$$y_i(\omega^T x_i + b) \geq 1$$

Solving the optimization problem

$$\text{minimize}_{\omega} \frac{1}{2} \|\omega\|^2$$

$$\text{s.t. } y_i(\omega^T x_i + b) \geq 1$$

Chpt 2 | The Dual formulation

$$\text{minimize}_{\omega} \frac{1}{2} \|\omega\|^2$$

$$\text{s.t. } y_i(\omega^T x_i + b) \geq 1$$

- optimization problem with convex quadratic objective and linear constraint
- can be solved using QP

Lagrange duality to get the optimization problem

- dual form, allows us to use kernel to get optimal margins
- allows classifiers to work efficiently for very high dimensional space
- allows us to derive an efficient algorithm for solving the above optimization problem that will

typically do much better than generic QP software

The primal Problem

$$\min_w f(w)$$

$$\text{s.t. } g_i(w) \leq 0, i=1, \dots, k$$

$$h_i(w) = 0, i=1, \dots, l.$$

the generalized Lagrangian

$$L(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

this α 's ($\alpha_i \geq 0$) and β_i 's are called Lagrange multipliers!

$$\max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta) = \begin{cases} f(w), & \text{if } w \text{ satisfies} \\ & \text{primal constraints} \\ \infty, & \text{otherwise} \end{cases}$$

A re-written primal

$$\min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$$

Lagrangian Duality! cont.

$$\text{The primal Problem } P^* = \min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$$

$$\text{The dual Problem! } d^* = \max_{\alpha, \beta, \alpha_i \geq 0} \min_w L(w, \alpha, \beta)$$

Theorem (weak duality)

$$d^* = \max_{\alpha, \beta, \alpha_i \geq 0} \min_w L(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta, \alpha_i \geq 0} L(w, \alpha, \beta)$$

Theorem (strong duality)

$$L(w, \alpha, \beta) = P^*$$

If there exist a saddle pt. of $L(w, \alpha, \beta)$

$$\text{we have } d^* = P^*$$

The KKT conditions

If there exists an saddle pt of L
then if satisfies the following "Karush-Kuhn-Tucker" (KKT) conditions:

$$\frac{\partial}{\partial \omega_i} L(\omega, \alpha, \beta) = 0, \quad i = 1, \dots, n$$

$$\frac{\partial}{\partial \beta_i} L(\omega, \alpha, \beta) = 0, \quad i = 1, \dots, l$$

$$\alpha_i g_i(\omega) = 0, \quad i = 1, \dots, m$$

$$g_i(\omega) \leq 0, \quad i = 1, \dots, m$$

$$\alpha_i \geq 0, \quad i$$

Theorem: If ω^*, α^* , and β^* satisfy the KKT condition,
then it is also a soln. to the primal and the
dual problem.

Support vector

- Only a few α_i 's can be nonzero
- Call the training data pts whose

α_i 's are nonzero

$$\alpha_i g_i(\omega) = 0, \quad i = 1, \dots, m$$

if $\alpha_i > 0$ then $g_i(\omega) = 0$

Solving the optimization problem,

Quadratic
Programming
with linear
constraints

$$\begin{cases} \text{minimize } \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 \end{cases}$$

Lagrangian function

$$\text{minimize } L_p(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (y_i(w^T x_i + b))$$

$$\text{s.t. } \alpha_i \geq 0$$

minimize
over w and b
for fixed α

$$\frac{\partial L_p}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\frac{\partial L_p}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$L_p(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (w_i^T x_j) - b \sum_{i=1}^m \alpha_i y_i$$

$$L_p(w, b, \alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

The Dual Problem

Now we have the following dual opt problem

$$\text{max}_{\alpha} J(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

$$\text{s.t. } \alpha_i \geq 0, i = 1, \dots, n$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

Primal

This is a quadratic programming problem

- A global maximum of α_i can always be found

- Once we have Lagrange multipliers $\{\alpha_i\}$ we can reconstruct the parameter vector $w = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$ as a weighted combination of training examples.

$$w = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$w = \sum_{i=1}^{S_v} \alpha_i y_i \mathbf{x}_i$$

- for testing with a new data \mathbf{z}

- compute

$$w^T \mathbf{z} + b = \sum_{i \in S_v} \alpha_i y_i (\mathbf{x}_i^T \mathbf{z}) + b$$

and classify \mathbf{z} as a class 1 if the sum is

+ve, and class 2 otherwise
Note: w need not be formed explicitly
(and b)

Solving the optimization problem

- the discriminant function is

$$g(\mathbf{x}) = w^T \mathbf{x} + b = \sum_{i \in S_v} \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

- If relies on a dot product between the test point \mathbf{x} and the support vector \mathbf{x}_i

- solving the optimization problem involved computing the dot product $\mathbf{x}_i^T \mathbf{x}_j$ b/w all pairs of having pts.
- The optimal w is a linear combination of small no. of data pts.

Chapter 2.2 Maximum margin classifier

Linear SVM formulation.

Find w and b such that

$$\frac{2}{\|w\|} \text{ is maximized}$$

Each having pt.

$$y (w \cdot \mathbf{x}_i + b) \geq 1$$

$$\frac{1}{2} \sqrt{w \cdot w} \text{ is minimized}$$

$$\|w\|^2 = w \cdot w \text{ minimized}$$

Limitation of previous SVM formulation

- what if the data is not linearly separable
- one noisy data pts. 2.

Extend the definition of maximum margin to allow non-separating plans.

Bell W. Golombok

$w \cdot w + C \epsilon$ having error.

Objective to be minimized

- Minimize

$w \cdot w + C \epsilon$ (distance of error pts. to the constraint)

- Add slack variable ξ_i

$$M = \frac{C}{\|w\|}$$

15456

Minimize

$$w \cdot w + C \sum_{k=1}^m \xi_k$$

M constraints

$$\left\{ \begin{array}{l} w \cdot x_k + b \geq 1 - \xi_k \text{ if } y_k = 1 \\ w \cdot x_k + b \leq -1 + \xi_k \text{ if } y_k = -1 \end{array} \right\}$$

\equiv

$$y_k (w \cdot x_k + b) \geq 1 - \xi_k, \quad k = 1, \dots, m$$

$$\xi_k \geq 0, \quad k = 1, \dots, m$$

Lagrangian

$$L(w, b, \xi, \alpha, \beta)$$

$$= \frac{1}{2} w \cdot w + c \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i [y_i (w \cdot x_i + b) - 1 + \xi_i] \\ - \sum_{i=1}^m \beta_i \xi_i$$

α_i 's and β_i 's are Lagrange multipliers (≥ 0)

$$\xi_i \geq 0 \text{ for } i = 1, \dots, m$$

Dual formulation

$$\text{max}_{\alpha} J(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j) \quad \text{s.t.}$$

$$\text{min}_{\alpha} J(\alpha) = \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j) - \sum_{i=1}^m \alpha_i \quad \text{s.t.}$$

Linear SVM

$$\text{s.t. } \alpha_i \geq 0 \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y_i = 1$$

$$\text{s.t. } 0 \leq \alpha_i \leq C \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y_i = 0.$$

Solution to soft margin classification

- x_i with non-zero α_i will be support vector
- Solution to the dual problem is:

$$w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$b = y_K (1 - \xi_K) - \sum_{i=1}^m \alpha_i y_i x_i^T x_K$$

for any κ s.t. $a_\kappa > 0$

for classification

$$f(x) = \sum_{i=1}^n \alpha_i y_i x_i + b$$

(we need to compare w explicitly)

Chapter 23

Nonlinear SVM and Kernel

← function

$$x \rightarrow \phi(x)$$

Map to high dim feature space
computable
costs?

$$\sum \alpha_i y_i y_j x_i^T x_j \frac{\text{---}}{O(D^2)} O(D^2)$$

$$D \gg d.$$

Kernel functions

$$x_a \quad x_b$$

$$\phi(x_a) \quad \phi(x_b)$$

$$x_a \cdot x_b \\ \phi(x_a), \phi(x_b)$$

$$k(x_a, x_b) = \phi(x_a) \cdot \phi(x_b)$$



the kernel trick

$$g(x) = w^T \phi(x) + b = \sum_{i \in SV} \alpha_i [\phi(x_i)^T \phi(x)] + b$$

- we only use the dot product of feature vector in both (the training and test)
- A kernel function is defined as a function that corresponds to a dot product of two feature vectors in some expanded feature space

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

$$x = [x_1, x_2]$$

$$K(x_i, x_j) = (1 + x_i \cdot x_j)^2$$

$$= 1 + x_i^2 x_j^2 + 2 x_i \cdot x_j x_i^2 x_j^2 + x_i^2 x_j^2$$

$$+ 2 x_i \cdot x_j + 2 x_i^2 x_j^2$$

$$= [1 \ x_i^2 \sqrt{2} x_i \ x_i^2 \ x_i^2 \sqrt{2} x_i \ \sqrt{2} x_i^2].$$

$$[x_j^2 \sqrt{2} x_j \ x_j^2 \ x_j^2 \sqrt{2} x_j \ \sqrt{2} x_j]$$

~~$\phi(x_i) \cdot \phi(x_j)$~~

commonly used kernel functions

Linear kernel: $K(x_i, x_j) = x_i \cdot x_j$

Polynomial of power P:

$$K(x_i, x_j) = (1 + x_i \cdot x_j)^P$$

Gaussian (radial basis function)

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

Sigmoid

$$k(x_i, x_j) = \tanh(\beta_0 x_i \cdot x_j + \beta_1)$$

In general functions that satisfy Mercer's condition can be kernel function.

- Kernel function can be thought of similarity measure b/w the input objects
- Not all similarity measure can be used as kernel function.
- Mercer's condition states that any +ve semi-definite kernel $k(x, y)$, i.e

$$\sum_{i,j} k(x_i, x_j) c_i c_j \geq 0$$

- Can be expressed as a dot product in a high dimensional space.

Nonlinear SVM: optimization

- Formulation (Lagrangian Dual Problem)

$$\text{maximize}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

such that $0 \leq \alpha_i \leq C$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

The sum of the discriminant function is

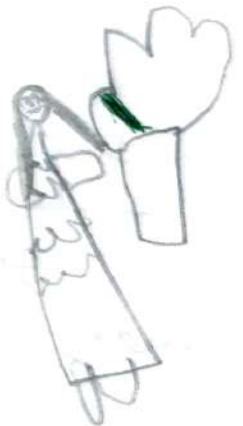
$$g_{\text{ent}} = \sum_{i \in \text{sv}} \alpha_i k(x_i, x_j) + b$$

Performance

- SVM works very well in practice
 - the user must choose the kernel function and its parameters
- They can be expensive in time and space for big data set
 - The computation of the maximum-mARGIN hyper-plane depends on the square of the number of training ~~the~~ cases
 - we need to store all the support vectors
- The kernel trick can also be used to do PCA in a much higher-dimensional space, thus giving a non-linear version of PCA in the original space

Multiclass classification

- SVMs can only handle two class outputs
- Learn \Rightarrow SVMs
 - SVM₁ learns class 1 vs REST
 - SVM₂ learns class 2 vs REST
 - SVM_n learns class n vs REST
- Then to predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the right class.



Sigmoid

$$k(x_i, x_j) = \tanh(\beta_0 x_i \cdot x_j + \beta_1)$$

In general, functions that satisfy Mercer's condition can be kernel function.

- kernel function can be thought of similarity measure b/w two input objects
- Not all similarity measure can be used as kernel function.
- Mercer's condition states that any we semi-definite kernel $k(x, y)$, i.e.

$$\sum_{i,j} k(x_i, x_j) c_i c_j \geq 0$$

- Can be expressed as a dot product in a high dimensional space.

Nonlinear SVM: optimization

- Formulation (Lagrangian Dual Problem)

$$\text{maximize}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

such that $0 \leq \alpha_i \leq C$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

The sum of the discriminant function is

$$g_{\text{ent}} = \sum_{i \in S} \alpha_i k(x_i, m_j) + b$$

Performance

- SVM work very well in practice
 - the user must choose the kernel function and its parameters
- They can be expensive in time and space for big data set
 - the computation of the maximum-mARGIN hyperplane depends on the square of the number of training ~~the~~ cases
 - we need to store all the support vectors
- the kernel trick can also be used to do PCA in a much higher-dimensional space, thus giving a non-linear version of PCA in the original space

Multiclass classification

- SVM can only handle two class output
- Learn \approx SVMs
 - SVM₁ learns class₁ vs REST
 - SVM₂ learns class₂ vs REST
 - SVM₃ learns class₃ vs REST
- Then to predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the region



Chpt 2.1 Solution of the Dual Problem

SMO (sequential optimization)

Coordinate Ascent

$$\max_{\alpha} w(\alpha_1, \alpha_2, \dots, \alpha_n) \quad \text{constraint}$$

Loop until convergence $\{\}$

for $i = 1 \text{ to } n \quad \{\}$

$$\alpha_i = \arg \max_{\alpha_i} w(\alpha_1, \alpha_2, \dots, \alpha_n)$$

Coordinate ascent

- ellipses are the contours of the function
- At each step, the path is \parallel to one of the axes.

Sequential minimal optimization

constrained optimization:

$$\max_{\alpha} I(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \quad (\alpha_i \geq 0)$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

- choose a set of α_i 's satisfying the constraints
- α_i is exactly determined by the other α_j 's
- we have to update at least two of them simultaneously to keep satisfying the constraints.

The SMO algorithm

Repeat till convergence [

1. select some pair α_i and α_j to update next
 (using a heuristic that tries to pick the two that will allow us to move the biggest progress toward the global maximum.)

2. Re-optimize $w(x)$ with respect to α_i and α_j
 while holding all the other α_k 's ($k \neq i, j$) fixed }

- . the update to α_i and α_j can be computed very efficiently

Chap 21 Python Exercise on SVM

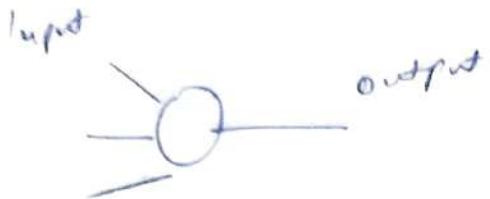
Chap 27 Neural Network Introduction

10¹² billion of neurons

1. massive parallelism

2. Connectionism

3. Distributed Association memory

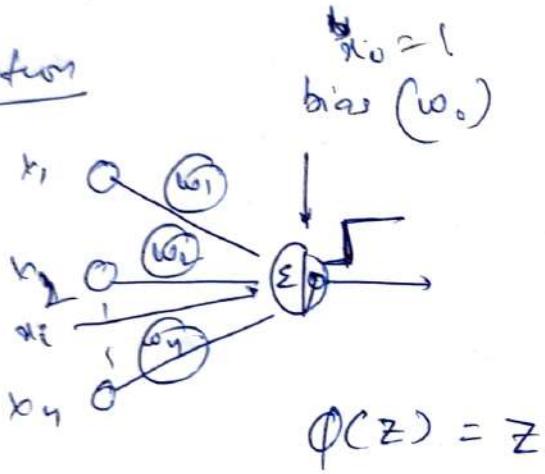


An ANN incorporates the two fundamental components of biological neural net

① Nodes - Neurons

② weights - Synapses

Perception



$$\phi(z) = \sigma$$

$$y = \sum_{i=0}^n w_i x_i + b$$

$$\phi_2(z) = \sigma$$

$$\frac{D}{X_1 Y_1}$$

$$x_1 \quad y_1$$

$$x_m \quad y_m$$

Perception learning rule

$$w_i = w_i + \Delta w$$

$$\Delta w_i = \eta(y - \hat{y})^n$$

- Perception learning convergence if D is linearly separable..

Gradient descent can be used in the most general phase

$$E = \frac{1}{2} \sum_{d \in D} (y_d - \hat{y}_d)^2$$



$$\phi_1(z) = z$$

Stochastic gradient
descent



$$E = \frac{1}{2} (y - \hat{y})^2$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = x_i(y - \hat{y}_i)$$

$$\varphi_3(z) = \sigma(z) = \frac{1}{1+e^{-z}}$$

Sigmoid function.

$$\varphi_3(z) = \varphi_1(z) \cdot (1 - \varphi_1(z))$$

$$E = \frac{1}{2} \sum_{d \in D} (y_d - \sigma(w \cdot x_d))^2$$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_d \frac{\partial E}{\partial g_d} \cdot \frac{\partial g_d}{\partial w_i} \rightarrow \sigma'(w \cdot x_d) x_{id}$$

$$\sum_d (y_d - \hat{g}_d) \cdot \left(\frac{\partial}{\partial w_i} (y_d - \sigma(w \cdot x_d)) \right)$$

$$= \sum_d (y_d - \hat{g}_d) / g_d \cdot x_{id} \cdot \hat{g}_d (1 - \hat{g}_d)$$

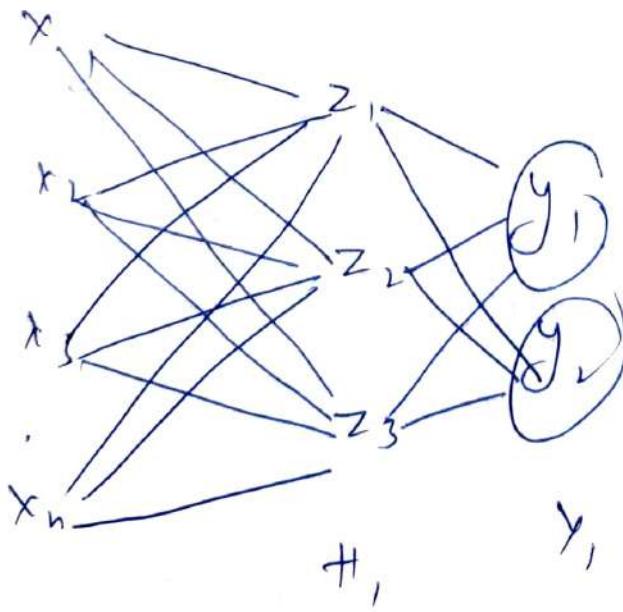
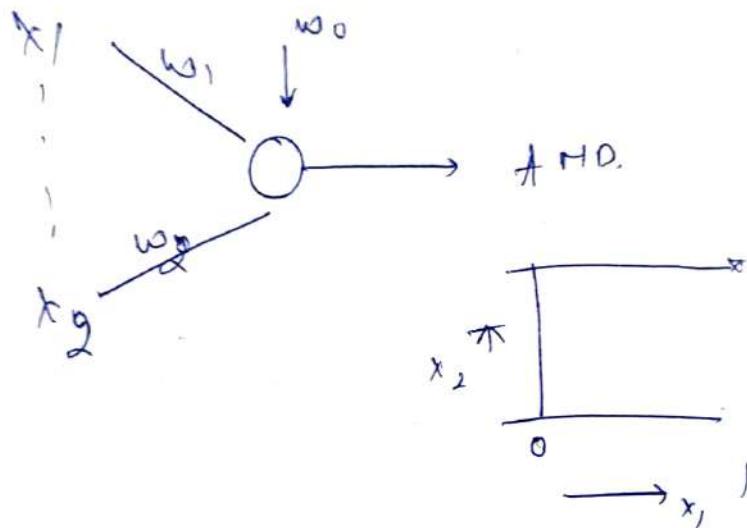
$$\Delta w_i = \eta \sum_d (y_d - \hat{g}_d) g_d (1 - g_d) \cdot x_{id}$$

Chapter 18

Multilayer Neural Network

Limits of Perception

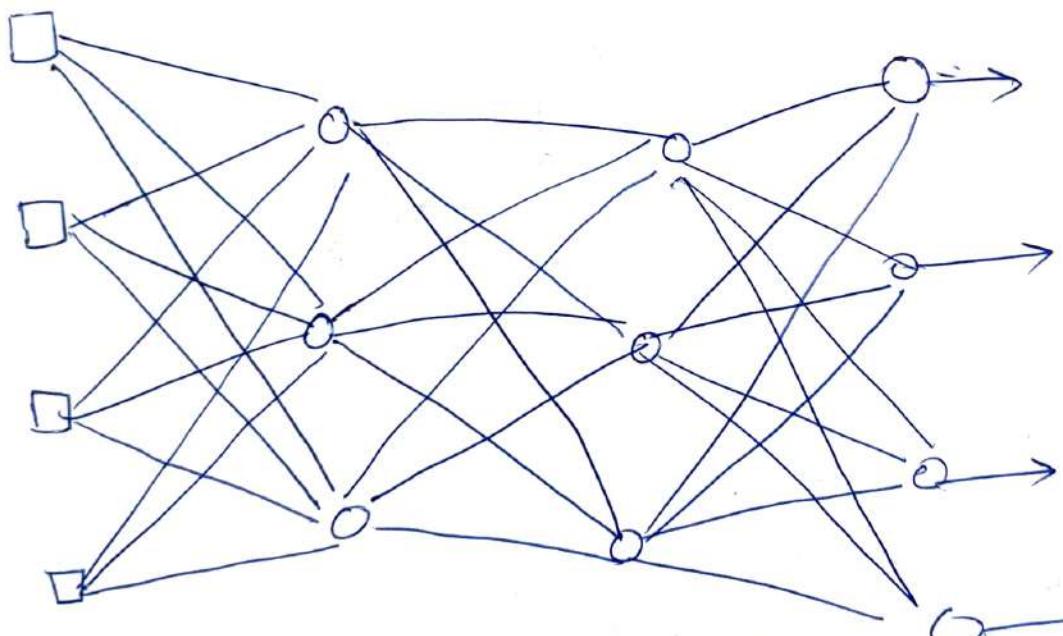
- Perceptions have a monotonicity property.
If a line has the weight, activation can only
rise as corresponding input value rises. (irrespective
of other input values)
- Can't represent where input interactions can
cancel one another's effect (e.g. XOR)
- Can represent only linearly separable functions



Power / Expressiveness of multilayer networks

- Can represent interactions among inputs
- Two layer networks can represent any Boolean functions and continuous functions (within a tolerance) as long as the no. of hidden units is sufficient and appropriate activation function used
- Learning algorithms exist, but weaker guarantees than perceptron learning algorithm

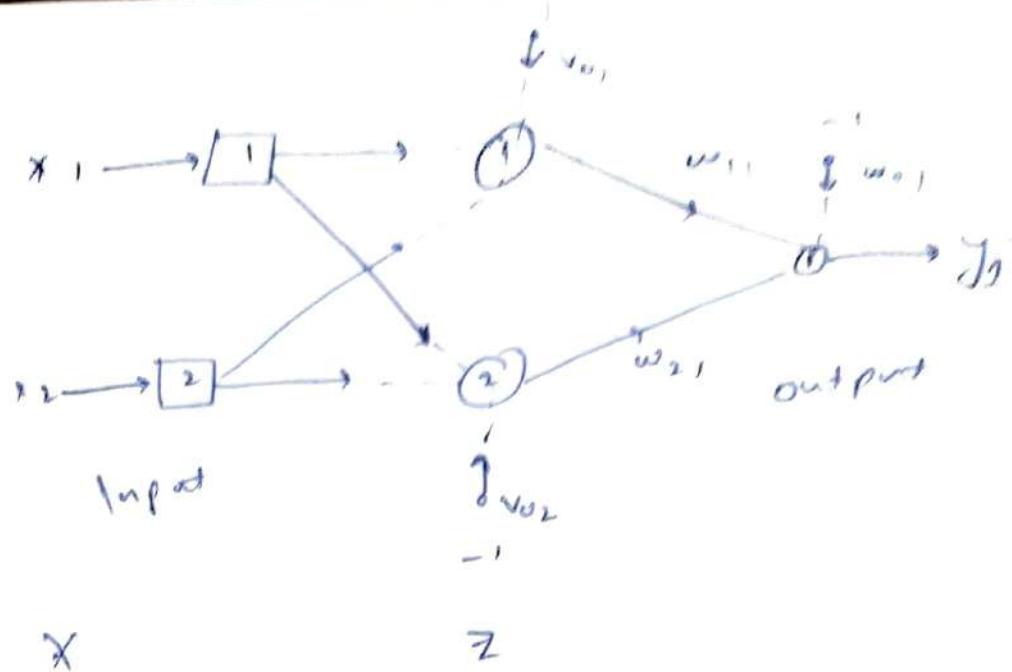
Two-layer back-propagation



The back-propagation training algorithm

Step 1: Initialisation

set all the weights and thresholds
levels of the network to random no. uniformly
distributed inside a small range



Backprop

initialization

forward computing

- apply on input vector x to input unit
- compute activation/output vector z on hidden layer
- compute the output vector y on output layer

$$z_j = \varphi(\sum_i w_{ij} x_i)$$

$$y_k = \varphi(\sum_j w_{jk} z_j)$$

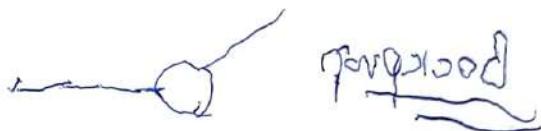
y is the result of the computation.

Learning for BP nets

- update of weights in w (by output and hidden layers)
 - delta rule!

- Not applicable to updating V (b/w input and hidden)
 - don't know the target values for hidden units Z^1, Z^2, \dots, Z^P
 - solution: propagates errors at output units to hidden units to drive the update of weights in V (again by delta rule) (error propagation learning)
 - Error backpropagation can be continued

do



From Output x we can see
and receive weights from
Input neurons

Derivation

- for one output neuron, the error function is

$$E = \frac{1}{2} \sum (y - \hat{y})^2$$

- for each unit j , the output O_j is defined

$$O_j = \varphi(\text{net}_j) = \varphi\left(\sum_{k=1}^n w_{kj} O_k\right)$$

The input net to a neuron is the weighted sum of output O_k of previous n neurons

- finding the derivative of the error

$$\frac{\partial e}{\partial w_{ij}} = \frac{\partial e}{\partial o_j} \frac{\partial o_j}{\partial w_{netj}} \frac{\partial w_{netj}}{\partial w_{ij}}$$

$$\frac{\partial w_{netj}}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \left(\sum_{k=1}^n w_{kj} o_k \right) = o_i$$

$$\frac{\partial o_j}{\partial w_{netj}} = \frac{\partial}{\partial w_{netj}} \varphi(w_{netj}) = \varphi'(w_{netj})(1 - \varphi(w_{netj}))$$

~~$$\frac{\partial E}{\partial o_j} = \text{error}$$~~

$$\frac{\partial E(o_j)}{\partial o_j} = \frac{\partial E(w_{net_1}, w_{net_2}, \dots, w_{net_l})}{\partial o_j}$$

$$\frac{\partial e(o_j)}{\partial o_j} = \frac{\partial E(w_{net_1}, w_{net_2}, \dots, w_{net_l})}{\partial o_j} = \sum_l \left(\frac{\partial e}{\partial w_{net_l}} \frac{\partial w_{net_l}}{\partial o_j} \right)$$
$$= \sum_l \left(\frac{\partial e}{\partial o_l} \frac{\partial o_l}{\partial w_{net_l}} \cdot w_{jl} \right)$$

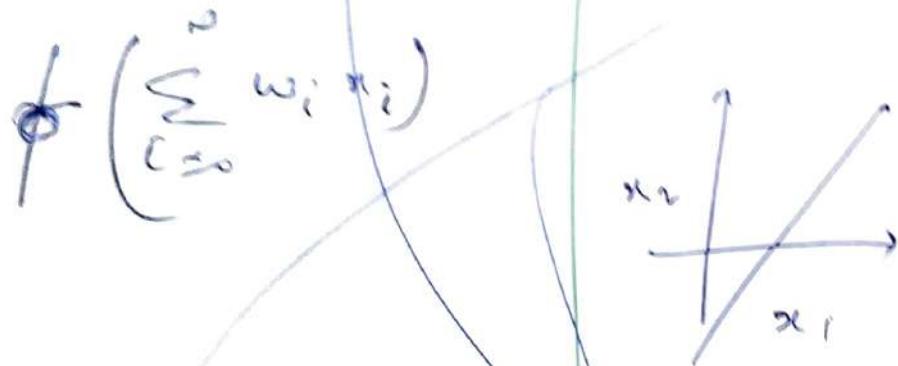
$$\delta_j = \delta_l = \frac{\partial e}{\partial o_l} \frac{\partial o_l}{\partial w_{net_l}} \cdot w_{jl}$$

$$= (\delta_j \cdot y_j) \circ_j (1 - y_j)$$

$$\sum_i (\delta_j \cdot w_{ji}) \circ_j (1 - y_j)$$

Chapter 28 Neural Networks and Backpropagation Algorithm

single layer Perceptrons learn decision boundaries



Chapter 30: Deep neural network

Chapter 32 Introduction to Computational Learning theory

Goal of learning theory

- To understand what kind of tasks are learned
- what kind of data is required for learnability
- what are the requirements of learning algorithms?

- To develop and analyze models that provable meet some requirements
- Develop algorithms that provable meet some criteria
- Prove guarantees for successful algorithms

→ Learnable

what to learn?

what type of data

Resource requirement - space, time

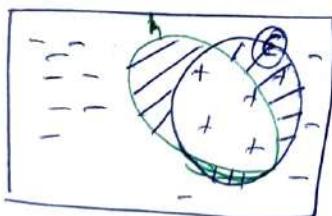
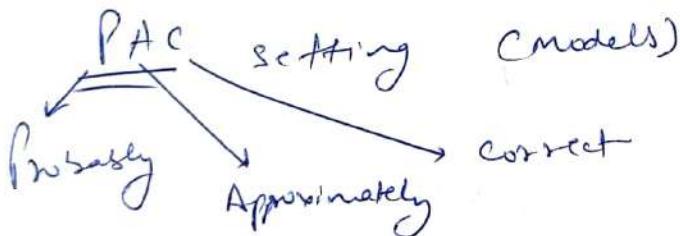


→ Develop algorithm

→ Prove guarantees

→ Design algorithm

→ confidence → Generalization ability



$$C \oplus h$$

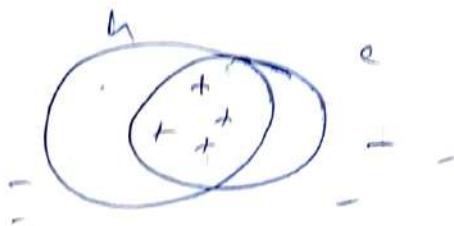
$$\text{Probability } Pr(P_C(c \oplus h) \leq \varepsilon) \geq (1-\delta)$$

$$0 \leq \varepsilon \leq \frac{1}{2}$$

$$0 \leq \delta \leq \frac{1}{2}$$

Intuitively about concept learning task.

- Given:
 - instance space $X = \mathbb{R}^d$ or $X = \{0, 1\}^d$
 - Distribution D over X
 - Target function c
 - Hypothesis Space H
 - Training Examples $S = \{(x_i, c(x_i))\}_{i=1}^m$ drawn from D .
- Determine
 - A hypothesis $h \in H$ s.t. $h(x) = c(x)$ for all x in S ?
 - A hypothesis $h \in H$ s.t. $h(x) = c(x)$ for all x in X ?
- An algorithm does optimization over S to find hypothesis h .
- Goal: find h which has small error over D .

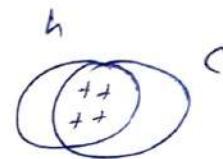


Instance of X

Consistent hypothesis

computational learning theory

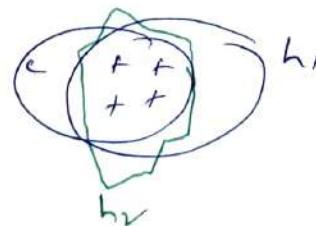
- can we be ~~certain~~ certain about how the learning algorithm generalizes?
- we would have to see all the examples
- Inductive inference -
generalizing beyond the training data is impossible unless we add more assumption
(e.g. priors over h)
we need a bias!



Instance space

function Approximation

- How many labeled examples in order to determine which of the $2^{|H|}$ hypotheses is correct one?
- All $2^{|H|}$ instances in X must be labeled!
- Inductive inference: generalizing beyond the training data is impossible unless we add more assumption
(e.g. bias)



$$H = \{h: X \rightarrow Y\}$$

$$|H| = 2^{|X|} = 2^{|H|}$$

M boolean features

2^M possible instance

$h_i \quad 2^{2^M}$ possible partitions of instance space

Error of a hypothesis

The true error of hypothesis h , with respect to the target concept c and observation distribution D is the probability that h will misclassify an instance drawn according to D .

$$\text{error}_D(h) = \Pr_{x \sim D}[c(x) \neq h(x)]$$

In a perfect world, we'd like the error to be 0

Bias Bias: Fix hypothesis space H .

c may not be in $H \Rightarrow$ find h close to c

A hypothesis h is approximately correct if

$$\text{error}_D(h) \leq \epsilon$$

h measure the error of h on S

$\text{error}_S(h) \approx$ True error of h

$\text{error}_S(h) = \Pr_{n \sim S}[c(n) \neq h(n)]$

$h \in H$

$c \notin H$

h close to c

$P(c \in H)$

$\text{error}_D(h) \leq \epsilon$
 ↓
 approximately correct

PAC model

- Goal: h has small error over D
- True error: $\text{error}(h) = \Pr_D[h(x) \neq c^*(x)]$
- How often $h(x) \neq c^*(x)$ over future instances drawn at random from D .
- But, can only measure:
 Training error: $\text{error}_D(h) = \frac{1}{m} \sum_i I(h(x_i) \neq c^*(x_i))$
 How often $h(x) \neq c^*(x)$ over training instances
- Sample complexity: bound $\text{error}_D(h)$ in terms of $\text{error}_D(h)$
- PAC learning concern efficient learning
 - we would like to prove that
 - with high probability on (efficient) learning algorithm will find a hypothesis that is approximately identical to the hidden target concept
- we specify two parameters, ϵ and δ and require that with probability at least $1 - \delta$, system learns a concept with error at most ϵ

Sample complexity for supervised learning

Theorem

$$m \geq \frac{1}{\epsilon} \left[\ln(H) + \ln\left(\frac{1}{\delta}\right) \right]$$

Labeled examples are sufficient so that with prob.

1 - δ, all $h \in H$ with $\text{error}_D(h) \geq \epsilon$ have $\text{error}_{\mathcal{D}}(h) > 0$

- Inversely linear in ϵ
- logarithmic in $|H|$
- ϵ error parameter: D might place low weight on certain parts of the space
- δ confidence parameter: there is small chance the examples we get are not representative of the distribution.

Theorem: $m \geq \frac{1}{\epsilon} \left[\ln(H) + \ln\left(\frac{1}{\delta}\right) \right]$ Labeled eg
are sufficient so that with prob. 1 - δ, all $h \in H$ with
 $\text{error}_D(h) \geq \epsilon$ have $\text{error}_{\mathcal{D}}(h) > 0$

Proof! Assume K bad hypotheses $H_{\text{bad}} = \{h_1, h_2, \dots, h_K\}$

with $\text{err}_D(h_i) \geq \epsilon$

fix h_i . Prob. h_i consistent with first training

example is $\leq 1 - \epsilon$. Prob. h_i consistent with first m training examples is $\leq (1 - \epsilon)^m$.

- Prob. that at least one h_i consistent with D in training examples

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

calculated value of m so that $|H|(1-e)^m \leq \delta$

- Use the fact that $1-x \leq e^{-x}$, sufficient to set

Sample complexity $\frac{|H|}{\epsilon} e^{-\epsilon m} \leq \delta$ An anti-hypothesis space realization
PLC: How many examples suffice to guarantee small error w.h.p.

Theorem

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

Labeled examples are sufficient so that with prob. $1-\delta$, all $h \in H$ with $\text{err}_D(h) \geq \epsilon$ have $\text{err}_S(h) \geq 0$

statistical learning way!

with probability at least $1-\delta$, all $h \in H$ s.t. $\text{err}_D(h) = 0$ we have

$$\text{err}_D(h) \leq \frac{1}{m} \left[\ln(|H|) + \ln\left(\frac{1}{\delta}\right) \right]$$

$$P(\text{consist}(H_{\text{bad}}, D)) \leq |H|e^{-\epsilon m} \leq \delta$$

$$e^{-\epsilon m} \leq \frac{\delta}{|H|}$$

$$-\epsilon m \leq \ln\left(\frac{\delta}{|H|}\right)$$

$$m \geq \left[-\frac{\ln(\delta)}{|H|} \right] / \epsilon \quad \text{flip inequality}$$

$$m \geq \left\lceil \frac{\ln \frac{1/\delta}{\epsilon}}{\delta} \right\rceil / \epsilon$$

$$m \geq \left\lceil \frac{2 \ln \frac{1}{\delta} + \ln \frac{1}{\epsilon}}{\delta} \right\rceil / \epsilon$$

Chapter 33 Sample Complexity: finite
Hypothesis space

- For a single hypothesis to have misleading training error

$$\Pr[\text{error}_D(f) \leq \epsilon + \text{error}_D(f)] \leq e^{-2n\epsilon^2}$$

- we want to ensure that the best hypothesis has error bounded in this way
 - so consider that any one of them could have large error

$$\Pr[\exists f \in H \text{ such that } \text{error}_D(f) \leq \epsilon + \text{error}_D(f)] \leq |H|e^{-2n\epsilon^2}$$

- from this we can derive the bound for the no. of samples needed

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

Sample complexity - finite hypothesis spaces

consistent case

Theorem

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

Labeled examples are sufficient so that with prob $1 - \delta$, all $h \in H$ with $\text{err}_D(h) \leq \epsilon$ have $\text{err}_S(h) \leq \epsilon$

Inconsistent case

what if there is no perfect h ?

Theorem: After m examples, with probability $\geq 1 - \delta$, all $h \in H$ have $|\text{err}_D(h) - \text{err}_S(h)| < \epsilon$ for

$$m \geq \frac{2}{2\epsilon^2} \left[\ln(|H|) + \ln\left(\frac{2}{\delta}\right) \right]$$

Sample complexity : example

c: conjugation of n Boolean functions. $\ln(C)$ learnable?

$$|H| = 3^n$$

$$m \geq \frac{1}{\epsilon} \left(n \ln 3 + \ln\left(\frac{1}{\delta}\right) \right)$$

- Concrete examples:

- $\delta = \epsilon = 0.05$, $n = 10$ gives 280 examples

- $\delta = 0.01$, $\epsilon = 0.05$, $n = 10$ gives 312 examples

- $\delta = \epsilon = 0.01$, $n = 10$, gives 1560 examples

- $\delta = \epsilon = 0.01$, $n = 50$ gives 5981 examples

- Result holds for any consistent learner, such as find

H = Conjunction of Boolean literals.

$x_1, \bar{x}_1, \dots, x_n, \bar{x}_n$

$x_1 \wedge \bar{x}_2 \wedge \dots \wedge x_n$

$x_1 \wedge \bar{x}_2 \wedge \bar{x}_3$

$\bar{x}_1 \wedge x_2$

Concept Learning

$h = (? , \text{and}, \text{high}, ?, ?, ?, ?)$

Indicates that Aido enjoys his favorite sport on cold day with high humidity

most general hypothesis: $(?, ?, ?, ?, ?, ?, ?)$

most specific hypothesis: $(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset)$

Find-S Algorithm

① Initialize h to the most specific hypothesis in H

② For each training instance X
for each attribute constraint a_i in h
if the constraint a_i in h is satisfied by x
Then do nothing
Else replace a_i in h by next most
general
Output hypothesis h

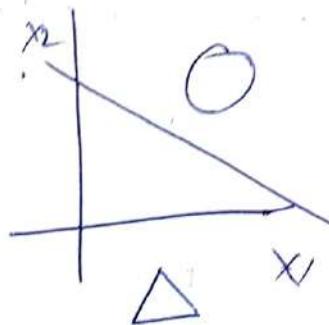
Hypothesis Sample complexity of learning
Autotopy Number of hypothesis

$$m \geq \frac{1}{c} (Cn)^{d+1}, \ln(C) \leq \left\lceil C^{(d+1)} \left(\frac{1}{c}\right)\right\rceil$$

Chpt 3 VC dimension

Sample complexity: In finite hypothesis spaces

- Need some measure of the expressiveness of infinite hypothesis space
- The Vapnik-Chervonenkis (VC) dimension provides such a measure, denoted $VC(H)$
- Analogous to $\text{Inf}(H)$, there are bounds for sample complexity using $VC(H)$

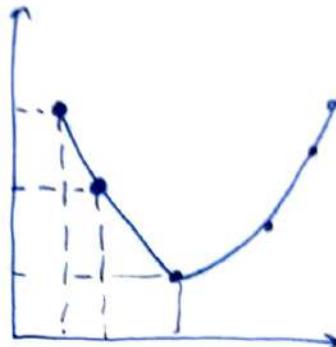
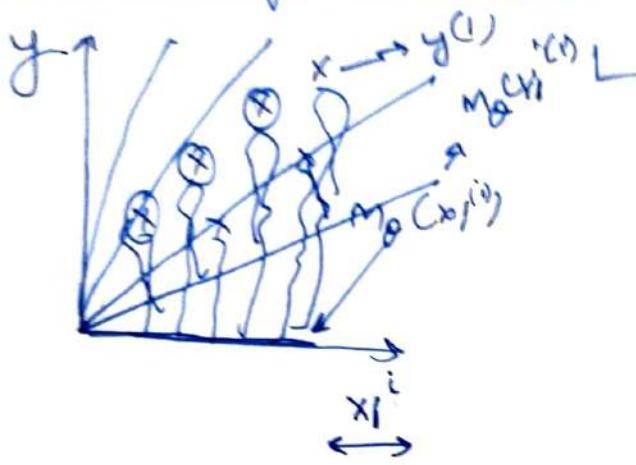


$$\text{Loss L} = \underset{\theta_0, \theta_1}{\operatorname{argmin}} \sum_{i=1}^n (M_\theta(x_i^{(1)}, y^{(1)})^2)$$

Least square cost function

Defining a Good model (selection of parameters)
How to select optimal θ_0, θ_1 ?

Curvature of the Loss function



In this curve θ ,
there is single minimum point



$$y = \theta_0 + \theta_1 x$$

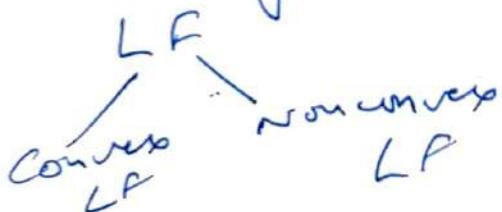
If intercept is not there then line pass throughs origin

$\sum (M_\theta(x_i^{(1)}, y^{(1)})^2)$ we are changing the value of θ

too 2D point of 3D graph is called contour

Generalized Loss function (L_f)

Non convex surface comprises of local minima
and global minima



two ways to obtain
optimal of θ_0 & θ_1

- (a) numerical method
- (b) gradient descent

Numerical methods

$$L = \frac{1}{2m} \sum_{i=1}^m (y_i^{(i)} - \theta_0 - \theta_1 x_i^{(i)})^2$$

$$= \frac{1}{2} \sum_{i=1}^m (\theta_0 + \theta_1 x_i^{(i)} - y_i^{(i)})^2 \quad \textcircled{1}$$

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i^{(i)})^2 + (y_i^{(i)})^2 - 2(\theta_0 + \theta_1 x_i^{(i)}) y_i^{(i)}$$

$$L = \frac{1}{2m} \sum_{i=1}^m \left[\theta_0^2 + \theta_1^2 (x_i^{(i)})^2 + 2\theta_0 \theta_1 x_i^{(i)} + (y_i^{(i)})^2 - 2\theta_0 y_i^{(i)} - 2\theta_1 x_i^{(i)} y_i^{(i)} \right]$$

$$\frac{\partial L}{\partial \theta_0}, \frac{\partial L}{\partial \theta_1} = ?$$

$$\frac{\partial L}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m \left[2\theta_0 + 2\theta_1 x_i^{(i)} - 2y_i^{(i)} \right]$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m \left[2\theta_1 (x_i^{(i)})^2 + 2\theta_0 x_i^{(i)} \right] - 2x_i^{(i)} y_i^{(i)}$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m \left[2\theta_1 (x_i^{(i)})^2 + 2\theta_0 x_i^{(i)} \right] - 2x_i^{(i)} y_i^{(i)}$$

to obtain values

Equating the derivatives to zero and corresponds to minimum of function L

$$\frac{\partial L}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m \left[2\theta_0 + 2\theta_1 x_i^{(i)} - 2y_i^{(i)} \right] = 0$$

$$\theta_0 + \frac{2}{2m} \sum_{i=1}^m \theta_1 x_i^{(i)} - \frac{1}{m} \sum_{i=1}^m y_i^{(i)} = 0$$

$$\theta_0 = \frac{1}{m} \sum_{i=1}^m y_i^{(i)} - \frac{\theta_1}{m} \sum_{i=1}^m x_i^{(i)}$$

$$\hat{\theta}_0 = \bar{y} - \theta_1 \bar{x},$$

To ensure minima we have to check the 2nd derivative

$$\frac{\partial^2 L}{\partial \theta_0^2} = \frac{(m+1)}{2} > 0 \quad (\because \text{minima is guaranteed})$$

optimal value of $\hat{\theta}_1 = ?$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m [2\theta_1(x_i^{(i)})^2 + 2\theta_0 x_i^{(i)} - 2x_i^{(i)}y_i^{(i)}] = 0$$

$$\theta_1 \frac{1}{m} \sum_{i=1}^m (x_i^{(i)})^2 + \frac{1}{m} \sum_{i=1}^m (\bar{y} - \theta_0 \bar{x}_i) x_i^{(i)} - \frac{1}{m} \sum_{i=1}^m x_i^{(i)} y_i^{(i)} = 0$$

$$\theta_1 \frac{1}{m} \sum_{i=1}^m (x_i^{(i)})^2 + (\bar{y} - \theta_0 \bar{x}_i) \frac{1}{m} \sum_{i=1}^m x_i^{(i)} - \frac{1}{m} \sum_{i=1}^m x_i^{(i)} y_i^{(i)} = 0$$

$$\theta_1 \left[\frac{1}{m} \sum_{i=1}^m (x_i^{(i)})^2 - \bar{x}_i \bar{x}_i \right] = \frac{1}{m} \sum_{i=1}^m x_i^{(i)} y_i^{(i)} - \bar{y} \bar{x}_i$$

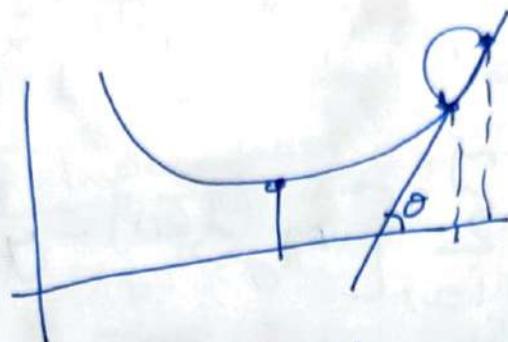
$$\hat{\theta}_1 = \frac{\frac{1}{m} \sum_{i=1}^m x_i^{(i)} y_i^{(i)} - \bar{y} \bar{x}_i}{\frac{1}{m} \sum_{i=1}^m (x_i^{(i)})^2 - \bar{x}_i \bar{x}_i}$$

$$\hat{\theta}_0 = \bar{y} - \theta_1 \bar{x}$$

$\hat{\theta}_0, \hat{\theta}_1 \rightarrow \text{optimal values}$

Since Second method gradient descent

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \frac{\partial L}{\partial \theta}$$



$$L = \frac{1}{2m} \sum_{i=1}^m (M_\theta(x^{(i)}) - y^{(i)})^2$$
$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i^{(i)} - y^{(i)})^2$$

2nd method.

$$\frac{\partial L}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m [2\theta_0 + 2\theta_1 x_i^{(i)} - 2y^{(i)}]$$
$$= \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i^{(i)} - y^{(i)})$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (M_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_0^{\text{new}} = \theta_0 - \alpha \frac{\partial L}{\partial \theta_0} = \theta_0 - \alpha \sum_{i=1}^m (M_\theta(x^{(i)}) - y^{(i)})$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m [2\theta_1 (x_i^{(i)})^2 + 2\theta_0 x_i^{(i)} - 2x_i^{(i)} y^{(i)}]$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{m} \left\{ \sum_{i=1}^m [\theta_0 + \theta_1 x_i^{(i)}] - y^{(i)} \right\} u_i^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m \{ [\theta_0 + \theta_1 x_i^{(i)}] - y^{(i)} \} u_i^{(i)}$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m \{ m_\theta(x_i^{(i)}) - y^{(i)} \} u_i^{(i)}$$

$$\theta_1 := \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m \{ m_\theta(x_i^{(i)}) - y^{(i)} \} x_i^{(i)}$$

Gradient Descent Algorithm (single feature)
repeat till minimum is reached

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m \{ m_\theta(x_i^{(i)}) - y^{(i)} \}$$

$$\theta_1 := \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m \{ m_\theta(x_i^{(i)}) - y^{(i)} \} u_i^{(i)}$$

Vectorization Techniques (Motivation)

Gradient descent Algorithm

repeat till minimum is reached

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m \{ m_\theta(x_i^{(i)}) - y^{(i)} \}$$

$$\theta_1 := \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m \{ m_\theta(x_i^{(i)}) - y^{(i)} \} u_i^{(i)}$$

↙ input?

What is vectorization Technique?

$$Z = \sum_{i=1}^n A_i B$$

$$Z = \vec{A}^T B$$

$$A = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad B = \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}$$

MATLAB
Python
library

Vectorized Implementation of cost function

$$M_\theta(x^{(i)}) = \vec{\theta}_0 + \vec{\theta}, x^{(i)}$$

$$\text{Let } \vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \in \mathbb{R}^{2 \times 1}, x^{(i)} = x_i$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix} = \begin{bmatrix} x_1^{(1)} \\ x_1^{(2)} \\ \vdots \\ x_1^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

$$\tilde{x} = [1 \ x] = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & x_1^{(3)} \\ 1 & x_1^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times 2} \rightarrow \begin{array}{l} x^{(2)} = [1 \ x_1^{(2)}] \\ x^{(2)} = x_1^{(2)} \end{array}$$

$$Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

$$L = \frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$$L = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2$$

$$L = \frac{1}{2m} \sum_{i=1}^m \left[[x_i^\theta - y_i]^2 \right]$$

single vectorized

Double vectorized (no summation in eq. A)

so we can write, $M_\theta(\theta) = \underline{\underline{x}}^\theta \frac{\underline{\underline{\theta}}}{2m}$,

$$\text{Loss} = \frac{1}{n} (\underline{\underline{x}}^\theta - \underline{\underline{y}})^\top (\underline{\underline{x}}^\theta - \underline{\underline{y}})$$

↳ similar to eqn. A

$$\begin{aligned} &= \frac{1}{m} \left[(\underline{\underline{x}}^\theta - \underline{\underline{y}})^\top (\underline{\underline{x}}^\theta - \underline{\underline{y}}) \right] \\ &= \frac{1}{m} \left[(\underline{\underline{\theta}}^\top \underline{\underline{x}} - \underline{\underline{y}}^\top) (\underline{\underline{x}}^\theta - \underline{\underline{y}}) \right] \\ &= \frac{1}{m} \left[\underline{\underline{\theta}}^\top \underline{\underline{x}}^\theta - \underline{\underline{\theta}}^\top \underline{\underline{x}}^\theta - \underline{\underline{y}}^\top \underline{\underline{x}}^\theta + \underline{\underline{y}}^\top \underline{\underline{y}} \right] \end{aligned}$$

$$L = \frac{1}{m} \left[\underline{\underline{\theta}}^\top \underline{\underline{x}}^\theta - 2 \underline{\underline{\theta}}^\top \underline{\underline{x}}^\theta + \underline{\underline{y}}^\top \underline{\underline{y}} \right]$$

$$\frac{\partial L}{\partial \theta} = \frac{2}{m} \underline{\underline{x}}^\theta \underline{\underline{\theta}} - \frac{2}{m} \underline{\underline{x}}^\theta \underline{\underline{y}} = 0$$

$$\underline{\underline{x}}^\theta \underline{\underline{\theta}} = \underline{\underline{x}}^\theta \underline{\underline{y}}$$

$$\hat{\theta} = (\underline{\underline{x}}^\theta \underline{\underline{x}})^\top \underline{\underline{x}}^\theta \underline{\underline{y}}$$

Disadvantage of normal method
If large no. of features are there, then $(\underline{\underline{x}}^\theta \underline{\underline{x}})^\top$
will be a big square matrix,

$$Y = \begin{pmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ -(x^{(3)})^T \\ \vdots \\ -(x^{(m)})^T \end{pmatrix} \Rightarrow Y = \begin{pmatrix} Y_1^{(1)} & Y_2^{(1)} & Y_3^{(1)} & \dots & Y_n^{(1)} \\ Y_1^{(2)} & Y_2^{(2)} & Y_3^{(2)} & \dots & Y_n^{(2)} \\ Y_1^{(3)} & Y_2^{(3)} & Y_3^{(3)} & \dots & Y_n^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ Y_1^{(m)} & Y_2^{(m)} & Y_3^{(m)} & \dots & Y_n^{(m)} \end{pmatrix} \in \mathbb{R}^{m \times 1}$$

Design matrix

$$X^{\circ} = [1 | X] = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & \dots & x_n^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & x_3^{(3)} & \dots & x_n^{(3)} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & x_3^{(m)} & \dots & x_n^{(m)} \end{pmatrix} \in \mathbb{R}^{m \times (n+1)}$$

$$\text{Loss} = \frac{1}{m} (X^{\circ} \theta - Y)^T (X^{\circ} \theta - Y)$$

$\underbrace{m \times (n+1)}_{(n+1) \times 1} \quad \underbrace{m \times (n+1)}_{(n+1) \times 1} \quad \underbrace{m \times 1}_{(n+1) \times 1}$

multiple feature

$\leftrightarrow \quad \leftrightarrow$

$1 \times m \quad m \times 1$

Gradient Descent Algorithm (Multiple features)
 repeat till maximum is reached achieved

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m [M_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m [M_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)}$

$\{j = 0 + n\}$

↓
features

Vectorized formulation

Using Normal update

$$\theta = \underbrace{(X^T X)}_{(n+1) \times (n+1)}^{-1} X^T y$$

$(n+1) \times (n+1)$

→ slow if n is large $O(n^3)$

Using gradient Descent

$$\theta := \theta - \alpha \frac{1}{m} \sum_{i=1}^m (X^T \theta - y) \quad \begin{matrix} m \\ (n+1) \times n \\ n \times (n+1) \end{matrix}$$

Gradient

Descent

Normal equation

- α is required

- no iteration is required

- slow if features are well

- Required

- need iteration

- work well if features are large

Any drawback with this cost function?

$$\text{Loss}(L) = \underset{\theta_0, \theta_1}{\operatorname{argmin}} \frac{1}{2m} \sum_{i=1}^m (y_i - (\theta_0 + \theta_1 x_i))^2$$

Least fit square overfit to outliers.

Solutions Least Absolute Deviation
→ we use absolute error $|y_i - (\theta_0 + \theta_1 x_i)|$ instead of squares

- still treat negative & positive error equally, but do not exaggerate

- emphasize the importance of large errors

Least Absolute Deviation cost

$$\text{Loss}(L) = \underset{\theta_0, \theta_1}{\operatorname{argmin}} \frac{1}{2m} \sum_{i=1}^m |y_i - (\theta_0 + \theta_1 x_i)|$$

$$A - B \Rightarrow |A| + |B|$$

Comparison b/w least square & Absolute Deviation

Weighted Regression (Dealing with Duplicates)

$$L = \frac{1}{2^m} \sum_{i=1}^m (M_0(x^{(i)}) - y^{(i)})^2$$

$$\# [(M_0(x^{(1)}) - y^{(1)})^2 + \underset{\beta_1, \text{ term}}{+} + (M_0(x^{(1)}) - y^{(1)})^2]$$

$$+ [(M_0(x^{(2)}) - y^{(2)})^2 + \underset{\beta_2, \text{ term}}{+} + (M_0(x^{(2)}) - y^{(2)})^2]$$

$$+ [(M_0(x^{(3)}) - y^{(3)})^2 + \underset{\beta_3, \text{ term}}{+} + (M_0(x^{(3)}) - y^{(3)})^2]$$

β_{nterm}

$$L = \frac{1}{\beta_1 + \beta_2 + \dots + \beta_m} \sum_{i=1}^m \beta_i [(M_0(x^{(i)}) - y^{(i)})^2]$$

$$= \cancel{\frac{1}{(\beta_1 + \beta_2 + \dots + \beta_m)}} \sum_{i=1}^m \beta_i [M_0(x^{(i)}) - y^{(i)}]^2$$

$$= \frac{1}{\beta_1 + \beta_2 + \dots + \beta_m} \sum_{i=1}^m \beta_i [x^{(i)} - y^{(i)}]^2$$

$$20 [M_0(x^{(5)}) - y^{(5)}] + 10 [M_0(x^{(4)}) - y^{(4)}]$$

Taking inverse will slow down the speed. See in next section

It's also ill-conditioned

$$\frac{\partial L}{\partial \theta} = \frac{2}{m} x^T (x^T \theta - y)$$

update: $\theta := \theta - \alpha \frac{2}{m} x^T (x^T \theta - y)$

faster speed compared to eqn. (5) and Σ is not there,
complete matrix multiplication.

target (y)

Linear Regression with Multiple Features

Feature = y
sample = $x_1 \ x_2 \ x_3 \ x_4$

$x_1^{(1)} = 2017 \quad | \quad x_2^{(5)} = 0$ there is no feature

$x_2^{(2)} = 1 \quad | \quad x_4^{(3)} = 5.3$

$$x^{(24)} = \begin{bmatrix} 2017 \\ 1 \\ 1.75 \\ 6.1 \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

Parameters

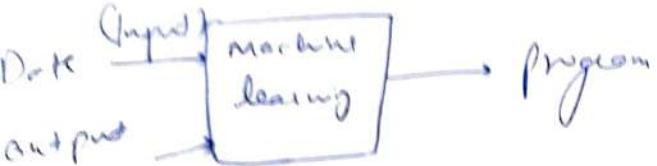
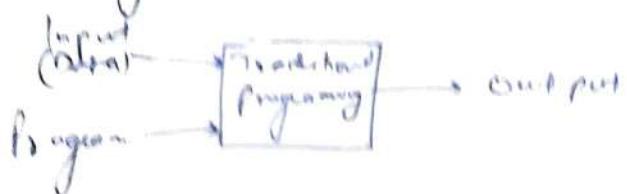
$$y_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \theta_3 x_3^{(i)} + \dots + \theta_n x_n^{(i)}$$

$$\text{Let } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{(n+1 \times 1)}$$

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

Chp 1: Introduction to Machine Learning

It is the field of study that gives computer the capability to learn without being explicitly programmed.



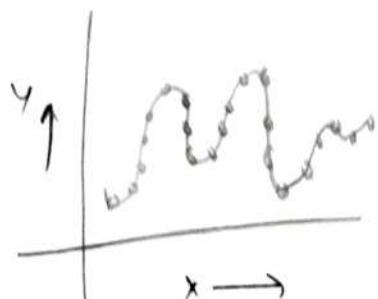
$$\text{Error margin} = \text{Actual output} - \text{Predicted output}$$

Chp 2: Classic & Adaptive machines

① Classic / Non adaptive system

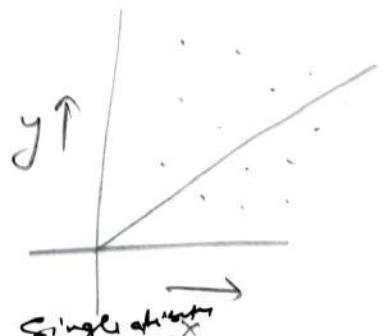
Chp 3: Basic of training and testing

Chp 4: Overfitting and Underfitting



Four attributes
Sphere
Play
Color
Radius = 5cm

- Line try to cover every points.
best fit



Single attribute
Ball
Sphere
- Features are less
- no difference b/w ball and orange

Chp 5 - Feature Selection techniques

1. filter methods optimal feature
- IG (Information Gain method)
- Chi-Square test
- correlation coeff

2. Wrapper methods

- Recursive feature elimination
- Genetic Algorithms

3. Embedded methods

- Decision tree One vs All

Chp 6: Principle Component Analysis

$$x \mid y \quad C = \begin{pmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{pmatrix} \quad \begin{array}{l} \text{One vs one} \\ \frac{n(n-1)}{2} \end{array}$$

$$\text{cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$x - \bar{x} \quad (x - \bar{x})(y - \bar{y})$$

$$y - \bar{y} \quad (y - \bar{y})(y - \bar{y})$$

$$y - \bar{y} \quad x - \bar{x} \quad y - \bar{y} \quad (x - \bar{x})(y - \bar{y})$$

Chp 8: $C - \lambda I = 0$

$$\begin{pmatrix} 0.6165 & 0.6157 \\ 0.6154 & 0.7165 \end{pmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$$\begin{bmatrix} 0.6165 - \lambda & 0.7161 \\ 0.6154 & 0.7165 - \lambda \end{bmatrix}$$

$$\therefore \lambda^2 - 1.332\lambda + 0.067 = 0$$

$$\lambda_1 = 0.0490$$

$$\lambda_2 = 1.2840$$

$$CV = \lambda_2$$

$$\begin{bmatrix} 0.6165 & 0.6154 \\ 0.6154 & 0.7165 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 0.0490 \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}$$

$$0.6165 x_1 + 0.6154 y_1 = 0.0490 x_1$$

$$0.6165 x_1 + 0.7165 y_1 = 0.0490 y_1$$

$$\therefore 0.5174 x_1 = -0.6114 y_1$$

$$\therefore 0.6154 x_1 = -0.6074 y_1$$

$$x_1 = -1.0845 y_1$$

$$\begin{bmatrix} -1.0845 \\ 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} -1.7614 \\ 1.7614 \end{bmatrix} + 1$$

$$= \frac{1.47517}{1.47517} = \frac{-1.0845}{1.47517}$$

$$= \begin{bmatrix} -0.7357 \\ 0.6778 \end{bmatrix} \quad \frac{1}{1.47517}$$

$$x_2 = 0.92195 y_2$$

$$\underbrace{0.92195}_1 = 0.849941$$

$$\therefore \begin{bmatrix} 0.6778 \\ 0.7357 \end{bmatrix} = \sqrt{1.47517} = 1.2601$$

Ch 9

Regression Analysis

- Dependent & Independent Variable
- Outlier
- Multicollinearity
underfitting and overfitting



Ch 10 : Linear and Logistic Regression

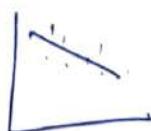
Dependent variable
continuous in nature

Linear Relationship

$I \rightarrow I$, & $I \rightarrow O$

$I - O \& I \rightarrow I$
multiple linear

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + \epsilon$$



- Dependent variable in binary
- I (One) & O (Zero)
- Goal is to find the best fitting model for I & O variable relationship
- Independent variable can be continuous & binary
- Also called Logit & R
- Used in medicine
- deals with probability \rightarrow
- means the relation depends on independent variables

Chp 11: Confusion matrix

		Predicted		P.
		No	Yes	
Actual	No	50	19	69
	Yes	5	71	76
		S	S	N.

$$\text{Recall} = \frac{\text{TP}}{\text{actual Yes}}$$

$$= \frac{100}{105} = 0.95$$

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{Total}} = \frac{(100 + 50)}{165} = 0.71$$

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

or

$$\frac{\text{FP} + \text{FN}}{\text{Total}} = 0.09$$

$$\text{Precision} = \frac{\text{TP}}{\text{Predicted Yes}} = \frac{100}{110} = 0.64$$

Ch13 Curse of dimensionality

Ch14 Many missing values
 One hot coding During coding
 after threshold after n/r.

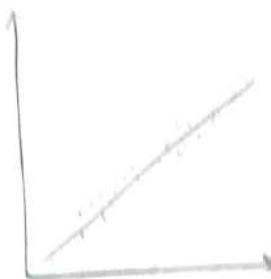
Ch15 Regress

Ch16 Logistic regression

$$y = \frac{1}{1+e^{-x}} \quad \text{sigmoid function}$$

30 - 50 data point

Chp 17 Polynomial Regression



$$y = \alpha_0 + \alpha_1 x_1$$

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4$$

$$\boxed{y = \alpha_0 + \sum_{i=1}^m \alpha_i x_i}$$

0 degree poly : $y = \text{constant}$

1 degree poly : $y = mx + c$

2 degree poly : $y = ax^2 + bx + c$

$$y = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \alpha_3 x^3 + \dots + \alpha_n x^n$$

$$\boxed{y = \alpha_0 + \sum_{i=1}^m \alpha_i x_i + F_p}$$

Chp 18 Ridge Regression

Regularisation \rightarrow

- Ridge

- Lasso

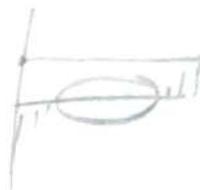
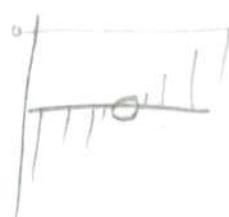
$$\text{Ridge } L = \text{Loss} + \alpha \|w\|^2$$

$$\text{Ridge } L = \text{Loss} + \alpha \|w\|^2$$

Lasso Regressor

$$\text{Lasso R} = \text{Loss} + \alpha \|w\| \text{ penalty}$$

$$\|w\| = |w_1| + |w_2| + |w_3| + \dots + |w_n|$$



magnitude of coefficient exactly zero

— Feature selection

Chp 20: Elastic Net regression

hybrid version of Ridge and Lasso

$$\text{Ridge R} = \text{Loss} + \alpha \|w\|^2$$



$$\text{Lasso R} = \text{Loss} + \alpha \|w\|$$

$$\text{Elastic Net R} = \text{Loss} + \underline{\alpha_1 \|w\|^2} + \underline{\alpha_2 \|w\|}$$

Chp 21 Bayes theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Prior

Posterior

likelihood

marginal

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

=====

Chp 22 Naive Bayes Classification Algorithm

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

37

Chp 23 Back propagation of error + nn, h

Stanford lecture 4

$$f(x, w) = w_x \quad \text{score function}$$

$$L_i = \sum_j f(y_i) \max(0, s_j - s_{y_i} + 1) \quad \text{sum loss}$$

$$L = \frac{1}{n} \sum_{i=1}^n L_i + \sum_k w_k^2 \quad \text{data loss + regularization}$$

want $\boxed{\nabla_w L}$

$$\frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \text{sigmoid function}$$

$$\frac{e^{-x}}{(1 + e^{-x})^2}$$

Always check the gradients
with respect to w's should
not be the same shape

Lecture 2:

Introduction to TensorFlow
Dr. Mehdi Hasheminezhad (Google)

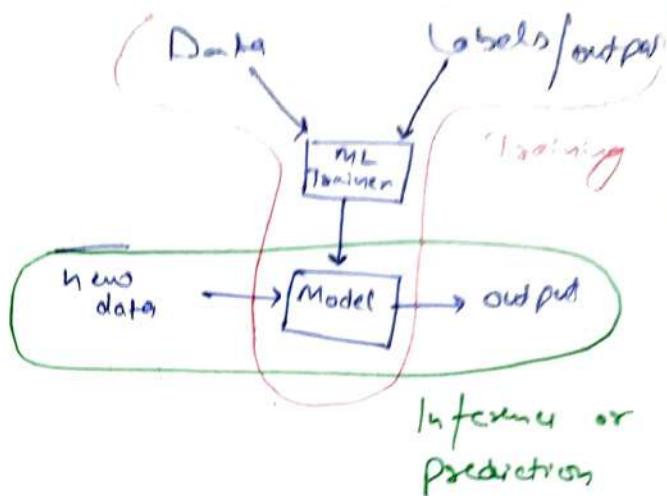
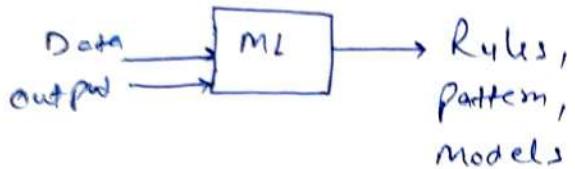
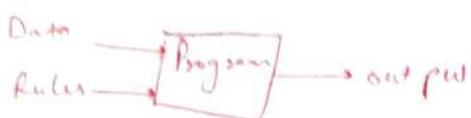
Add two no.

$$a+b$$

$$\begin{matrix} 8 & 7 & 2 & 8 \\ \boxed{8} & \boxed{7} & \boxed{2} & \boxed{8} \end{matrix}$$

func $f(a, b)$

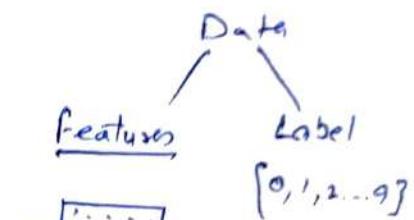
return $(a+b)$



1. Data :

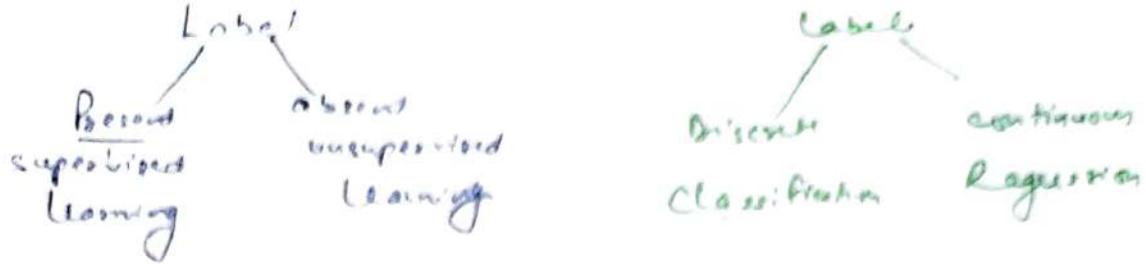
eg Handwritten images $x^{(j)}$

$$D = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$$



features \rightarrow
 - House price prediction
 # bedrooms $x_1^{(i)}$
 + Area (sq.foot) $x_2^{(i)}$
 + Distance from school $x_3^{(i)}$
 Label \rightarrow Price of House $x^{(i)}$: $x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}$

$$D_{\text{all}} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$$



feature : # bedrooms, area

(i) Numerical feature

(ii) Categorical

city : ["Mumbai", "Delhi", "Chennai"]

	F_Mumbai	, F_Delhi	, F_Chennai	
Mumbai	1	0	0	One-hot-encoding
Delhi	0	1	0	
Chennai	0	0	1	

Lecture 3 :-

Dr. Ashish Tendulkar (Google)

① Data preprocessing

$$x_j : \text{mean} : \mu_j \quad \frac{x_j - \mu_j}{\sigma_j} \quad Z\text{-score} : -3 \text{ to } +3$$

std. dev. : σ_j

$$\min_j, \max_j \quad \frac{x_j - \min_j}{\max_j - \min_j} \quad 0 \text{ to } 1$$

log transformation

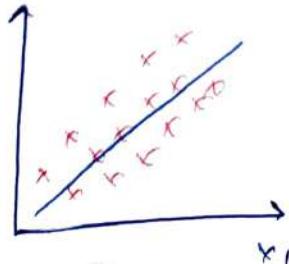
Model

(i) Linear regression

$$y = b + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$(b, w_1, w_2, \dots, w_n) \leftarrow \text{Parameters}$

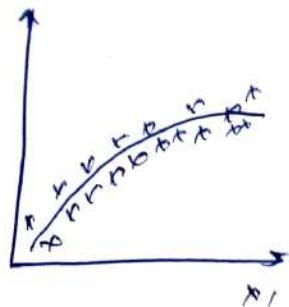
$$y = b + \sum w_i x_i$$



m - features : x_1, x_2, \dots, x_m

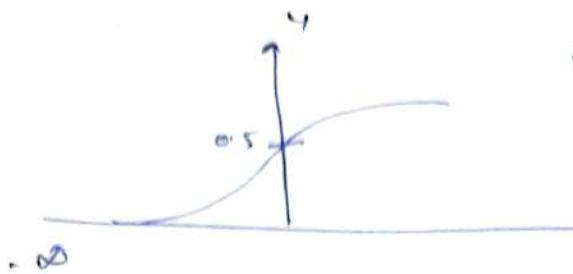
$$y = mx + c$$

$$y = b + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$



② Logistic regression:

$$Z = b + w_1x_1 + w_2x_2 + \dots + w_m x_m$$



Sigmoid function

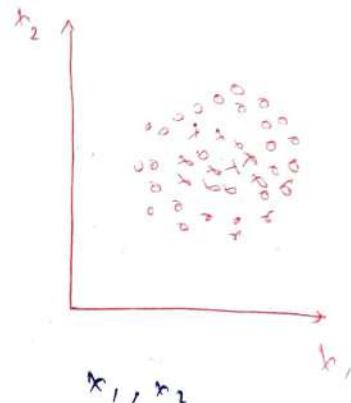
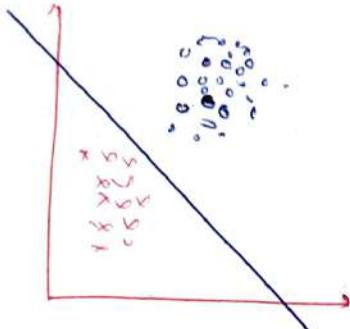
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$Z = 0, \sigma(Z) = 0.5$$

$$Z = \infty \rightarrow \sigma(Z) \rightarrow 1$$

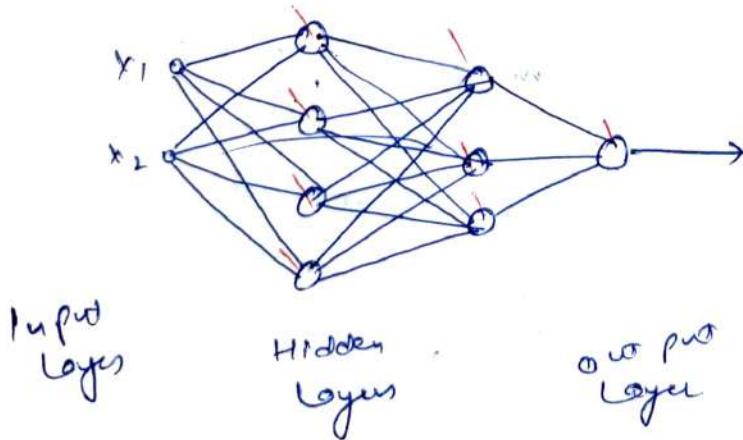
$$Z = -\infty \leftarrow \sigma(Z) \rightarrow 0$$

$$P_{\theta}(y^{(i)}=1|x^{(i)}) = \frac{1}{1 + \exp[-(b + w_1x_1 + \dots + w_mx_m)]}$$

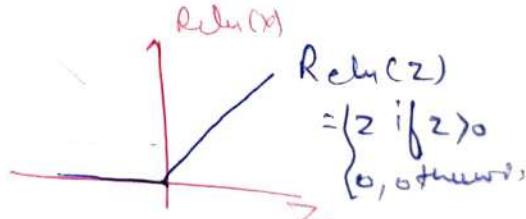
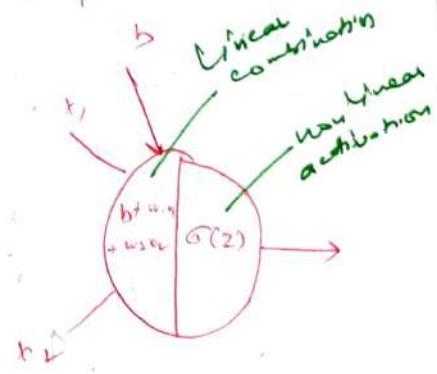


$$x_1^L, x_2^L, b, z_L$$

Feedforward neural network



features =



Lecture 4 Dr. Ashish Tadumadka Loss function.

Linear regression

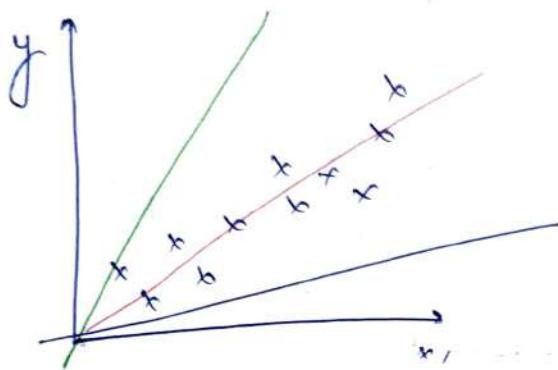
$$b + w_1x_1 + w_2x_2 + \dots + w_mx_m$$

$$h_{w,b}(x) = b + \sum_{i=1}^m w_i x_i$$

Logistic regression

$$P(y=1/x) = \frac{1}{1+e^{-(z)}}$$

$$z = b + w_1x_1 + w_2x_2 + \dots + w_mx_m$$



$$y = b + w_1x_1$$

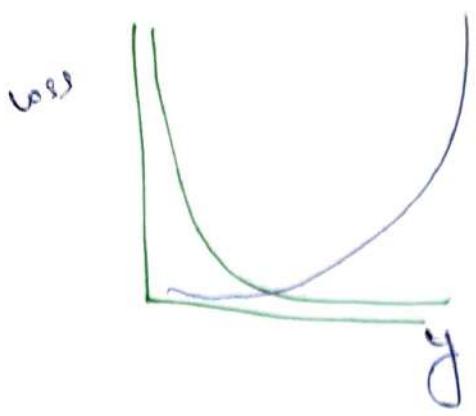
$$y = w_1x_1$$

Loss function

$$J(w, b) = \frac{1}{2} \sum_{i=1}^n (h_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2} \sum_{i=1}^n ((b + w_1x_1^{(i)} - y^{(i)})^2)$$

actual (y)	Predicted (\hat{y})	
1	0	error
0	1	error
1	1	ok (no-error)



$$f(y) = -y \log(p)$$

$$y=0 \quad -(1-y) \log(1-p)$$

cross entropy loss

$$= -y \log(p) - (1-y) \log(1-p)$$

for $y=1$

$$= -\log(p)$$

for $y=0$

$$-\log(1-p)$$

Lecture 5 :- Gradient Descent

loss function

$$J(w, b) = \frac{1}{2} \sum_{i=1}^n (h_{w,b}(x^{(i)}) - y^{(i)})^2$$

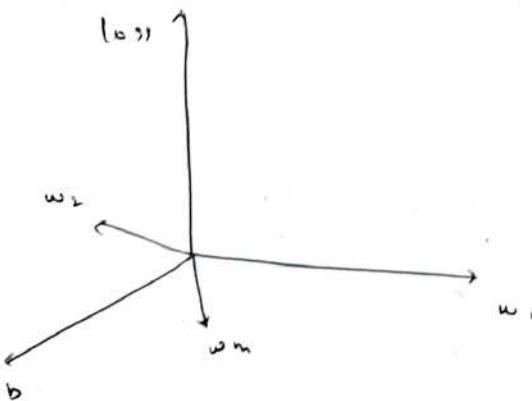
predicted value actual value

$$y = b + w_1 x_1 + w_2 x_2 + \dots + w_m x_m$$

$(b, w_1, w_2, \dots, w_m)$

m- features
(m+1) parameters

(b, \vec{w})



Data

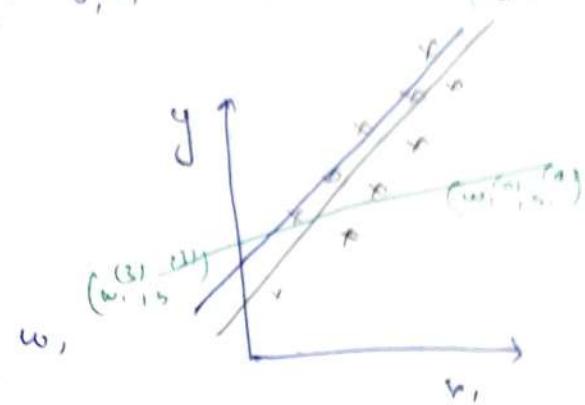
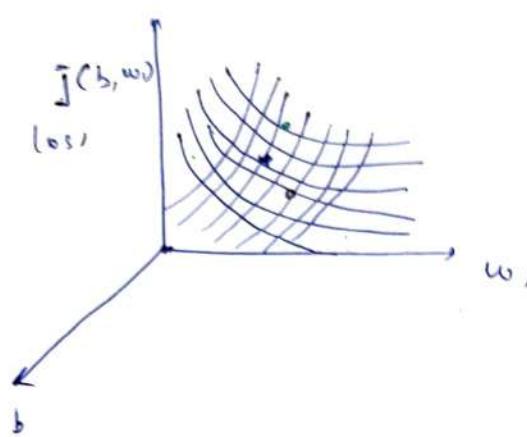
$$\text{Data} \rightarrow \{x^{(i)}, y^{(i)}\}_{i=1}^n$$

$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$$

$$x^{(i)} : x_i, y^{(i)} \in \mathbb{R}$$

$$\text{Model: } y = b + w_1 x_1$$

$$\text{Loss: } J(b, w_1) = \frac{1}{2} \sum_{i=1}^n (h_{b, w_1}(x^{(i)}) - y^{(i)})^2$$



Learning problem

find w, b such that loss is minimized

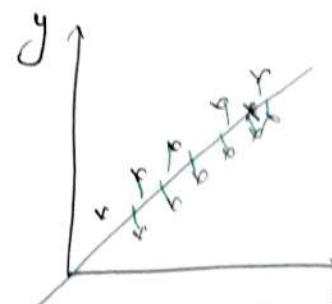
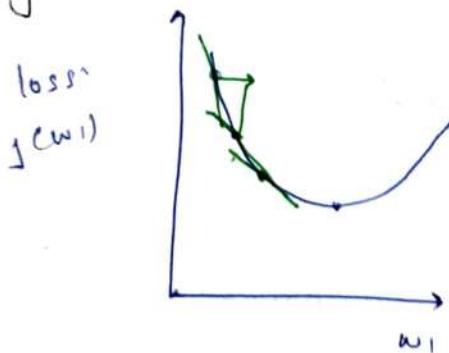
$$J(w, b)$$

$$y = b + w_1 x_1 ; b = 0$$

find out w_1 st. $J(w_1, b)$ is minimized

$$y = w_1 x_1$$

"Gradient Descent"



- ① Randomly initialize w_1
- ② Loop $J(w_1, b)$ calculate predictions
- ③ Gradient calculation

Learning ended

$$\textcircled{4} \quad w_{\text{new}} = w_{\text{old}} - \alpha \cdot \text{gradient}$$

| finds shorter and
shorter and
shorter)

Gradient Descent b, w_1, w_2, \dots, w_n

$$\textcircled{5} \quad \text{Randomly initialize } y = b + \sum_{i=1}^m w_i x_i$$

$$b, w_1, w_2, \dots, w_m \quad J(b, w) = \frac{1}{2} \sum_{i=1}^n (h_{w,b}(x^{(i)}) - y^{(i)})^2$$

⑥ Repeat until convergence

⑦ Predict $\hat{y}^{(i)}$ for each data point in training

⑧ calculate loss $J(b, w)$

⑨ calculate gradient of $J(b, w)$

$$b^{(\text{new})} = b^{(\text{old})} - \alpha \cdot \text{gradient } b$$

$$w_1^{(\text{new})} = w_1^{(\text{old})} - \alpha \cdot \text{gradient } w_1$$

$$w_n^{(\text{new})} = w_n^{(\text{old})} - \alpha \cdot \text{gradient } w_n$$

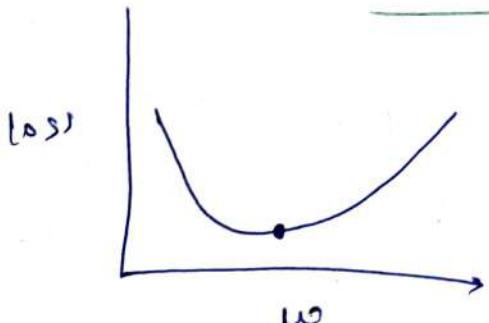
⑩ Update

$$b, w_1, w_2, \dots, w_n$$

simultaneously

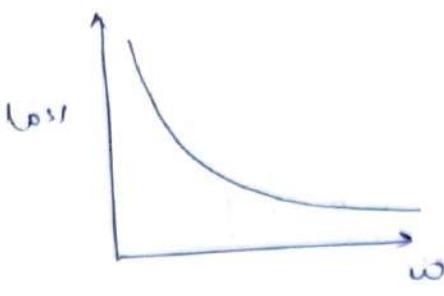
go to step (ii)

Lecture 6: Gradient descent Variation



① No change in parameter value
gradient = 0
 $w_1^{(\text{new})} = w_1^{(\text{old})} - \alpha \cdot \text{gradient}$

② No. of steps



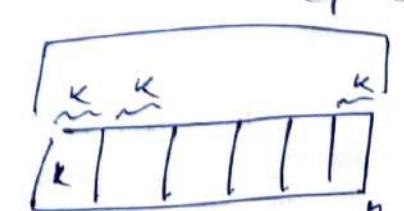
Gradient Descent of Cost func.

$$J(w, b) = \sum_{i=0}^n (h_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$\begin{aligned} \frac{\partial}{\partial w_j} J(w, b) &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{i=0}^n (h_{w,b}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2} \cancel{\sum_{i=1}^n} h_{w,b}(x^{(i)} - y^{(i)}) \frac{\partial}{\partial w_j} (h_{w,b}(x^{(i)} - y^{(i)})) \\ &= \sum_{i=1}^n (h_{w,b}(x^{(i)} - y^{(i)})) \cdot x_j \\ &= \sum_{i=1}^n (h_{w,b}(x^{(i)} - y^{(i)})) \cdot x_i \\ \frac{\partial}{\partial w_j} J(w, b) &= \sum_{i=1}^n (h_{w,b}(x^{(i)} - y^{(i)})) \cdot x_j \\ \frac{\partial}{\partial b} J(w, b) &= \sum_{i=1}^n (h_{w,b}(x^{(i)} - y^{(i)})) \cdot 1 \\ \frac{\partial}{\partial w_j} J(w, b) &= \sum_{i=1}^n (h_{w,b}(x^{(i)} - y^{(i)})) \cdot x_j \end{aligned}$$

Batch gradient descent "n" data points

$$\sum_{i=1}^k (h_{w,b}(x^{(i)} - y^{(i)}) \cdot x_j^{(i)})$$



k: batch size
mini-batch grad.

1 epoch = $\frac{n}{k}$ batch iterations

$K=1$ · stochastic G · D

SGD

G · D

1

n

($b^{(i)}, y^{(i)}$)

① Draw a bunch of K examples ($b^{(i)}, y^{(i)}$)

② Randomly initialize parameters

Repeat until converge

- obtain prediction from model

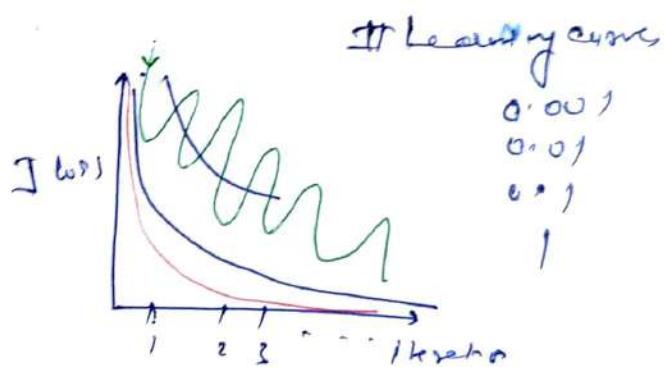
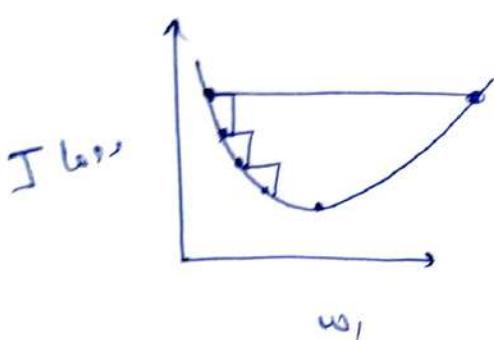
- calculate loss on the batch

- calculate gradient of loss w.r.t parameters

- $w^{(\text{new})} = w^{(\text{old})} - \alpha \cdot \text{gradient}$

- update simultaneously

How do I know if my model is learning?



Learning rate:

① Small

$$w^{(\text{new})} = w^{(\text{old})} - \alpha \cdot \text{gradient}$$

② Larger

③ Too large

Lecture 7: Model selection and evaluation

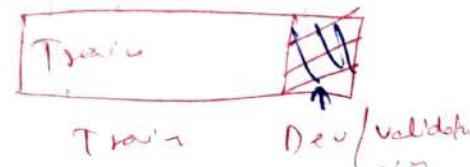
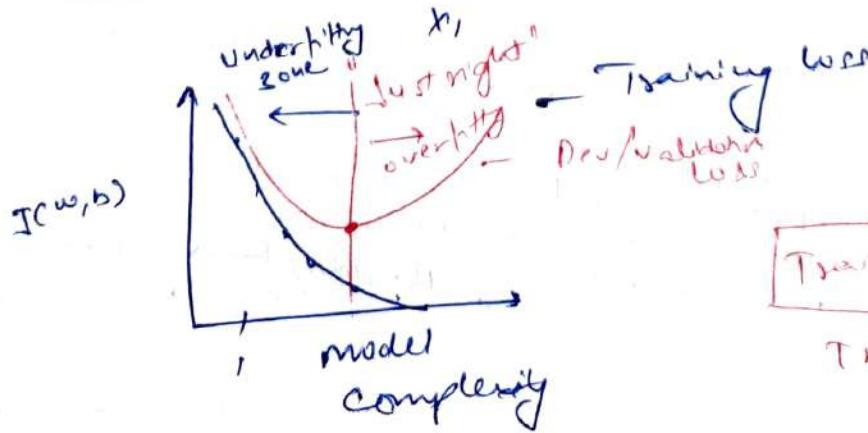
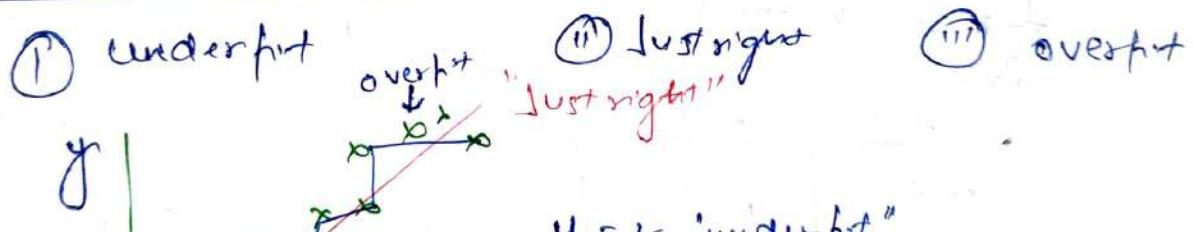
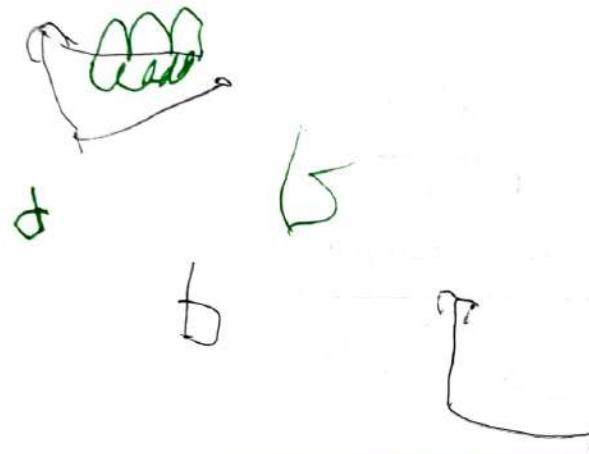
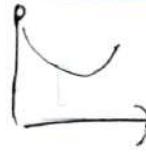
Regression

ms e
MAE



$T_{\text{train}}^{(20\%)}$ $T_{\text{test}}^{(20\%)}$

- Classification
- # Precision, recall -
 $f=1$ score score
 - # Accuracy /
 - # ROC PR-curves
-



Fitting underfitting

(1) x_1, x_2

$$x_1^2, x_2^2, x_1 x_2$$

(2) Reduce λ to fix underfitting

Fitting overfitting

(1) Get more data

(2) Reduce model complexity

"regularization"

$$\mathcal{J}(w, b) + \frac{\lambda \text{model complexity}}{L_2 - \text{regularization}} \sum_{i=1}^m \|w_i\|_2^2$$

$$\begin{aligned} L_1 - \text{regularization} & \sum_{i=1}^m \|w_i\|_1 \\ & = \|\sum_{i=1}^m w_i\|_1 \end{aligned}$$

	More Data	More complexity
Underfit	X	Decrease
Overfit	✓	most complex Inferior less complexity

Confusion matrix

Predicted		Actual	
		+1	-1
+1	TP	FN	
	FP	TN	

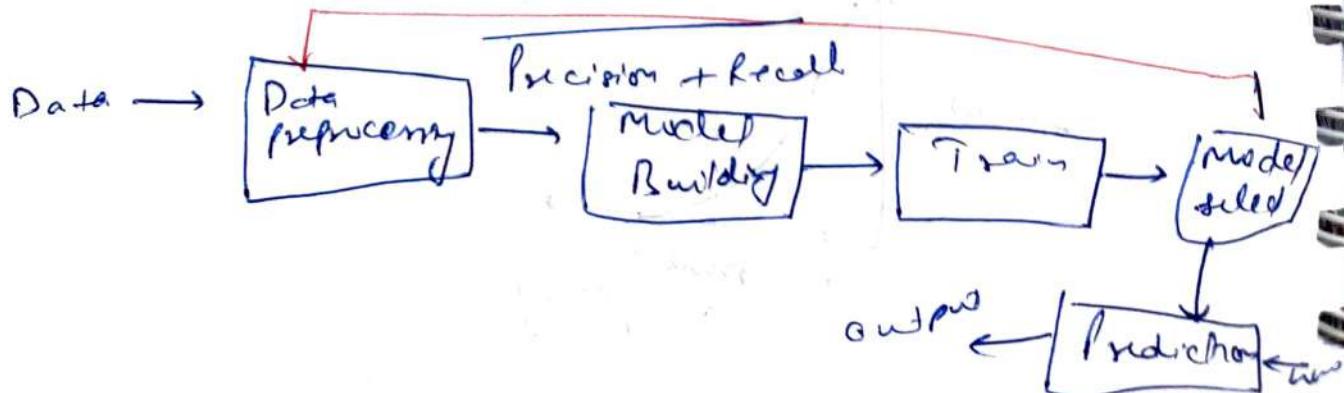
$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

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

$$F_1 - \text{score} = 2$$

Precision + Recall

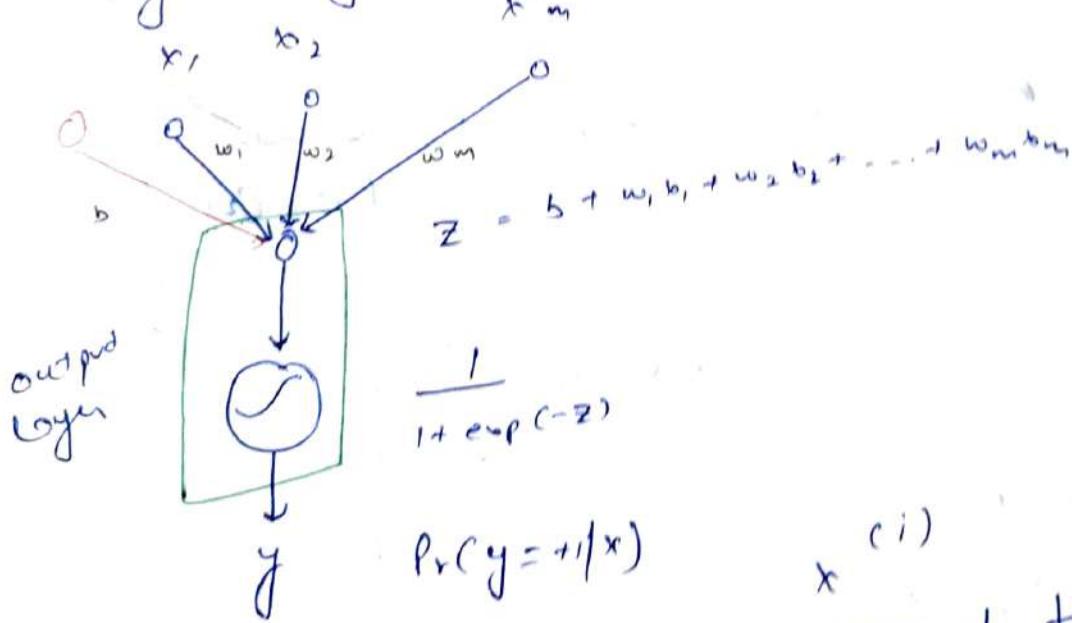


- ① Training data: $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$
- ② Model $y = h(x) + b$ with $\hat{y} = g(h(x))$
- ③ Cost $J = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$ (Predict - actual)² / m
- ④ Train $G.O.$ → gradient descent
- ⑤ Evaluate mse / MAE / linear regression

Lecture 8: Machine learning visualization
CA Neural network (playground) \rightarrow graph
Tinker with a Neural network

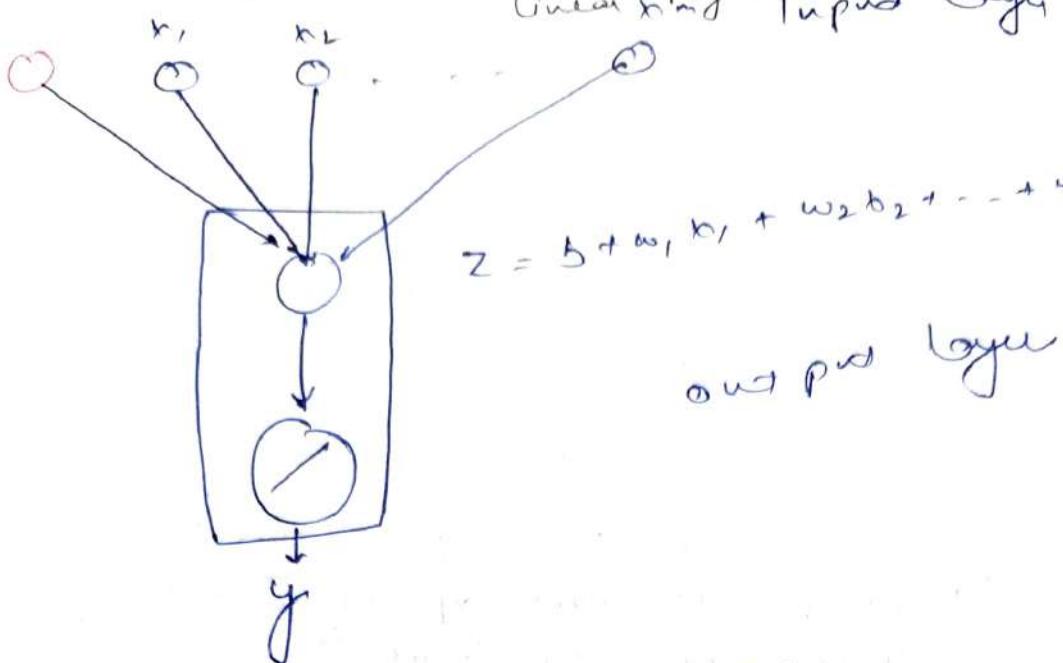
Lecture 9: Deep Learning Refresher

logistic regression

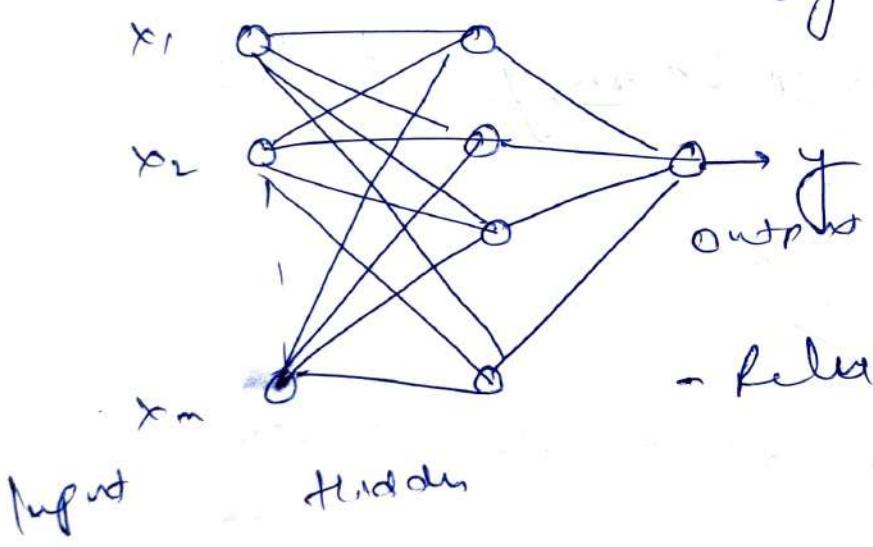


$$\Pr(y=+1|x)$$

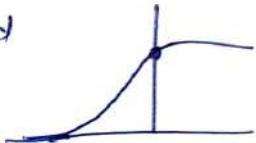
$x^{(i)}$
 m : # of features
 $m+1$ nodes (input)



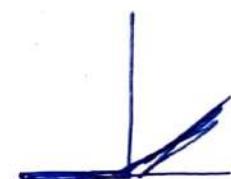
Building blocks of DL



- sigmoid

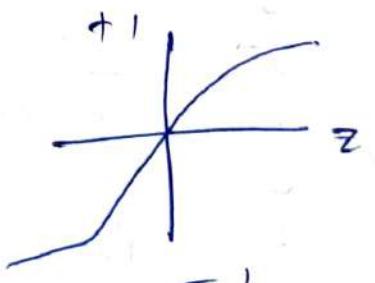


- ReLU



$$\text{ReLU}(z) = \begin{cases} z & z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- Tanh



$$z = b + w_1x_1 + w_2x_2 + \dots + w_nx_n$$

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

Row

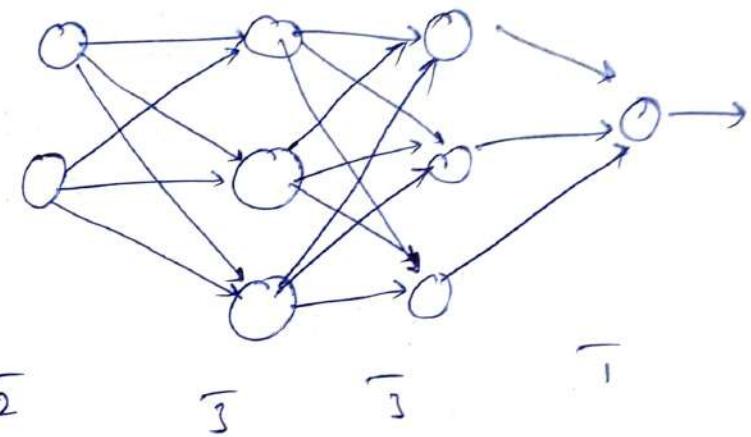
$$y = \text{sigmoid}(z_2 w_1^{(2)} + z_1 w_2^{(2)} + b)$$

$$z_1 = \text{relu}(b + w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2)$$

$$z_2 = \text{relu}(b + w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2)$$

feed forward NN

Model :-



Loss : MSE for regression.

Lecture 10 Introduction to Tensors

Scalar

$$x = \text{np.array}(10)$$

Vector 1D tensor

$$v = \text{np.array}(4, 1, 5, 9)$$

$v.$ shape

Matrix (2D tensor) $\begin{bmatrix} [1, 2, 3], \\ [4, 5, 6], \\ [7, 8, 9] \end{bmatrix}$

3D tensor $x = np.array([[[1, 2, 3],$
 $[4, 5, 6],$
 $[7, 8, 9]],$
 $[[1, 2, 3],$
 $[4, 5, 6],$
 $[7, 8, 9]]])$

$x.shape$
 $(3, 3, 3)$

$x.dtype$

Tensor in MNIST datasets

from future import absolute, division, print_function,
unicode_literals

! pip install tensorflow = ~~2.0.0-beta1~~ 2.0.0-beta1
Import tensorflow as tf

[] mnist = tf.keras.datasets.mnist

$(x_{\text{train}}, y_{\text{train}}), (x_{\text{test}}, y_{\text{test}}) = \text{mnist.load_data}$

- selecting a single datapoint
 $\underline{x} = \underline{x}_{\text{train}}[0]$

x_1

Select multiple

Select data point from 10 to 100 (100 is not included)

$$x_{\text{train_slice}} = x_{\text{train}}[10:100]$$

Print ($x_{\text{train_slice}}$)

Here we explicitly select two points

$$x_{\text{train_slice}} = x_{\text{train}}[10:10, 1, 1]$$

Here we explicitly select data according to our needs by specifying
the first and last elements

$$x_{\text{train_slice}} = x_{\text{train}}[10:100, 0.5:0, 0.1:0]$$

$$x_{\text{train_bs_slice}} = x_{\text{train}}[:, 14:, 14:]$$

Data batches

first batch - first 128 examples. Each batch
has 128 examples

$$\text{batch}_1 = x_{\text{train}}[:128]$$

Next batch - next 128 examples

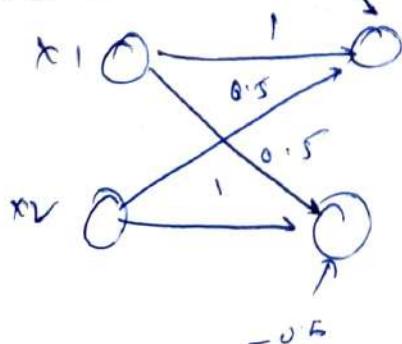
$$\text{batch}_1 = x_{\text{train}}[128:256]$$

Lecture 11 Mathematical foundation of deep learning

Key operation in neural networks

tf.keras.layers.Dense(128, activation='relu')

$$\text{output} = \text{relu}(\text{dot}(\text{w}, \text{input}) + b)$$



$$[1, 0.5], [0.5, 1]$$

2×2
2d tensor with shape(2, 2)

Input tensor

$$[x_1, x_2]$$

$$[[x_1^{(1)}, x_2^{(1)}], [x_1^{(2)}, x_2^{(2)}]]$$

2-d tensor (2, 2)

Bias vector $[2, 0.5]$

1-d tensor with shape

out put = np. maximum (0., z)

Element wise addition

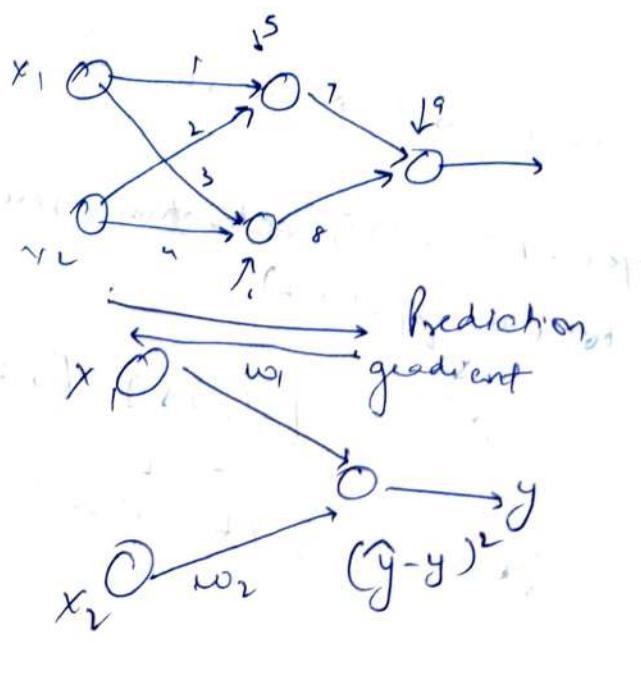
Broadcasting

$x_1 = \text{np. random. rand}(32, 10)$

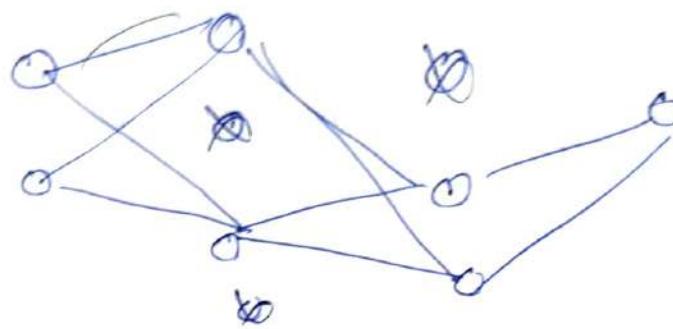
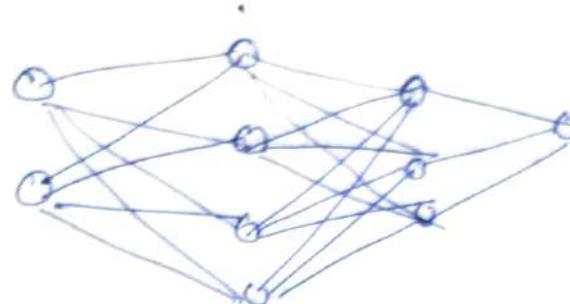
$x_2 = \text{np. random. rand}(10,)$

Reshaping

Transpose $x = \text{np. transpose}$



Dropout 0.2 or more
0.5 or 1.0



0.5 2.

Lecture 12 A Building data pipelines for tensorflow

Text

structured data

Image

Time series

spontaneous

shuffle Dataset.shuffle()

Preprocessing

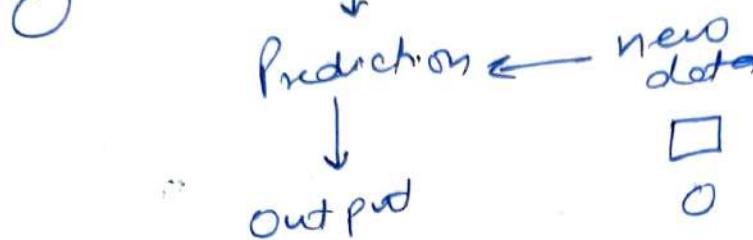
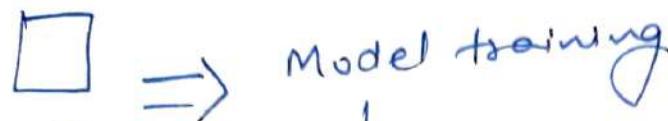
Tf. keras tf. estimator

Lecture 12 B

Load csv with tf. data

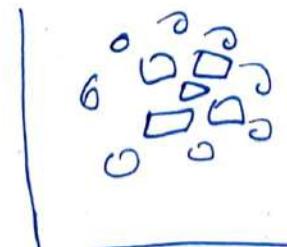
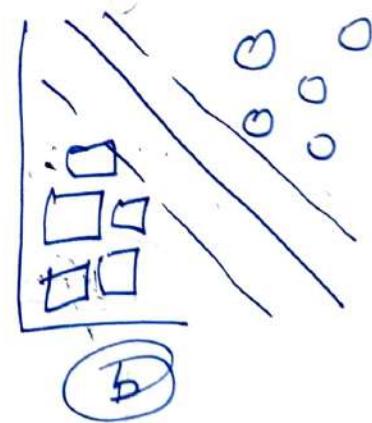
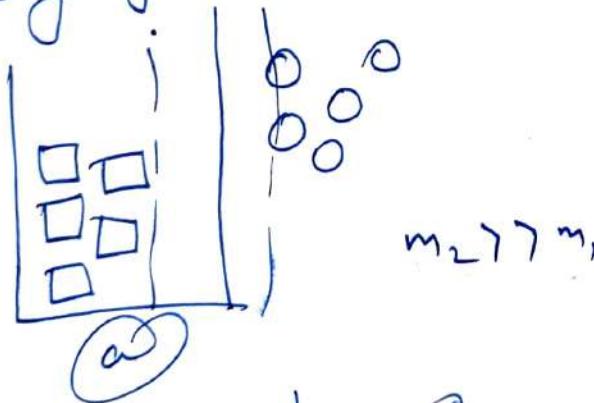
Support vector Machine (SVM)

- supervised learning and classification



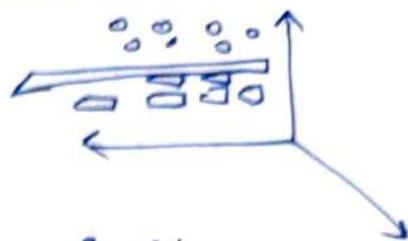
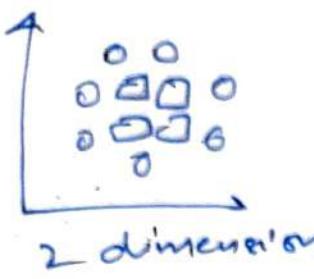
- find the nearer point in one class which is very near to other class, then draw line which touch the one point and \parallel to the hyperplane
- support vector :- the point from which two \parallel lines to hyperplane drawn, that point called support vector

why hyperplane decision



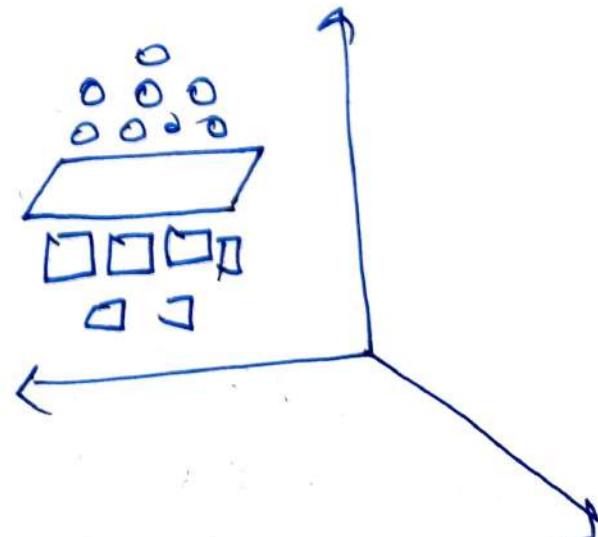
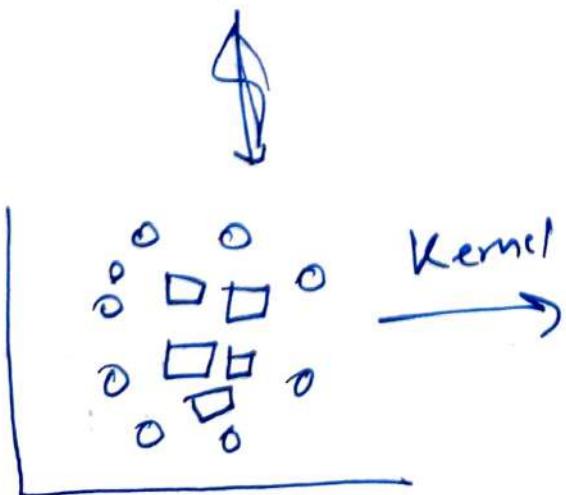
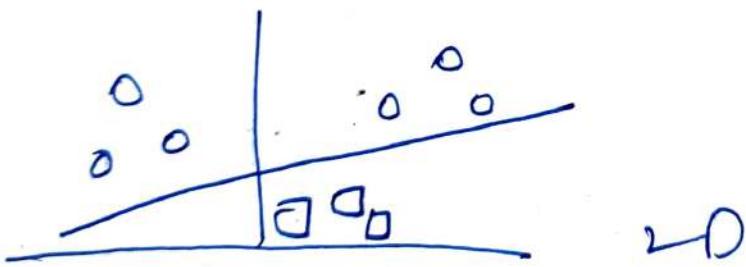
- condition is that we need maximum width of the margin so in (b) which is angled and (a) is perpendicular, angled margin distance width is more that why this called maximal margin hyperplane should get selected

- Miss classification minimum
 - Maximal margin hyperplane
- Non Linear SVM and Kernel function



Low dimension feature space
Kernel
High dimension feature space

~~0 0 0 □ □ □ 0 0 0~~ 1 Dimension
↓ Kernel.



Decision tree Algorithm

- select one attribute which is known as target attribute
- then detect the information gain of that target attribute

$$I \cdot G = -\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$E(A) = \sum_{i=1}^v \frac{P_i + N_i}{P+N} I(A|P_i, N_i)$$

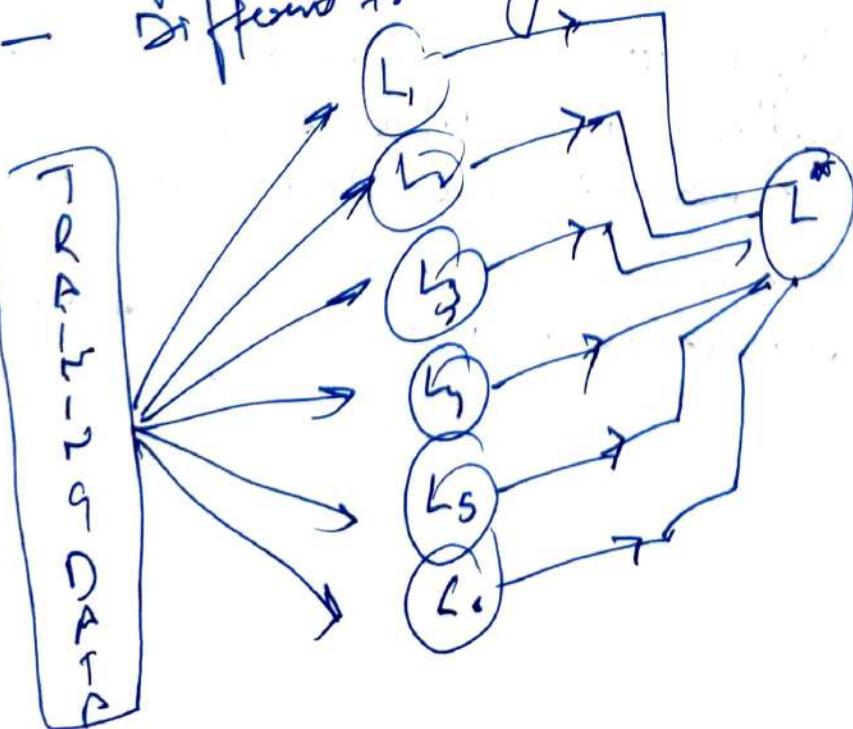
$$\text{Gain} = I \cdot G - E(A)$$

$$\log_2 x = \frac{\log_{10} x}{\log_{10} 2}$$

Root node which has maximum Gain information
and each obs with decrease the value

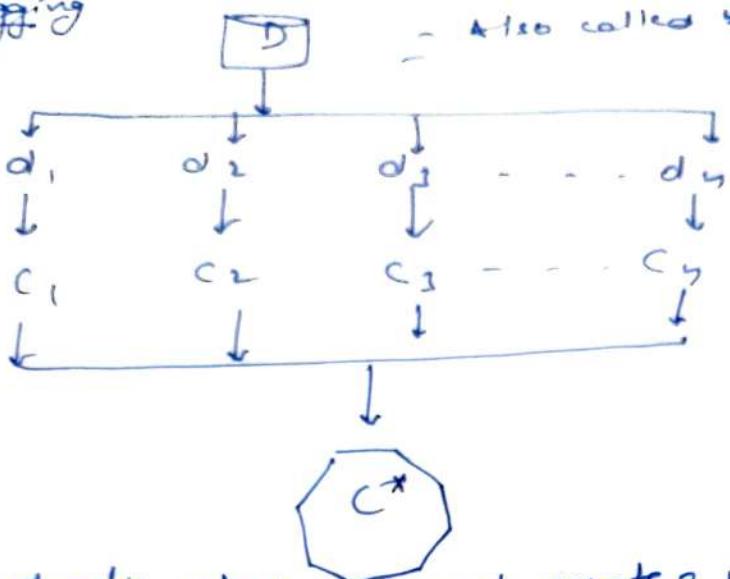
Ensemble learning

- Different Algorithms
- Different training data set



Ensemble Methods

Bagging



- Also called bootstrapping aggregation

- Randomly select data and create a new data set several times i.e. $d_1, d_2, d_3, \dots, d_n$ with replacement i.e. sampling with replacement

$d_1, d_2, d_3, \dots, d_n$ called bootstrap sample

- these classifiers are trained. $c_1, c_2, c_3, \dots, c_n$
- combine these classifiers and generate ensemble classifier

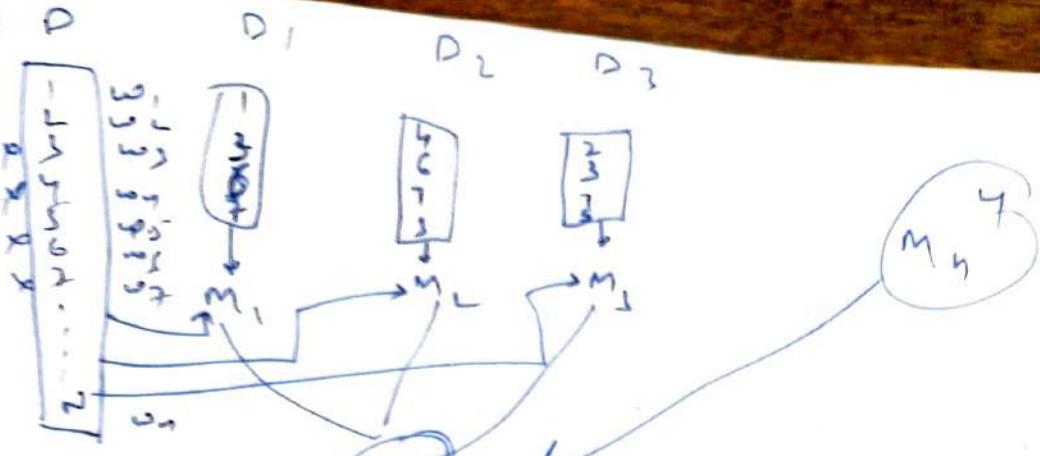
→ error rate is low compare to individual classifier in ensemble classifier and accuracy is more in ensemble classifier

- If we give x tuple to all classifier then individual classifier give different results then to decide which is the true answer we using voting.

technique is used in case of classification

- Voting problem

Boosting



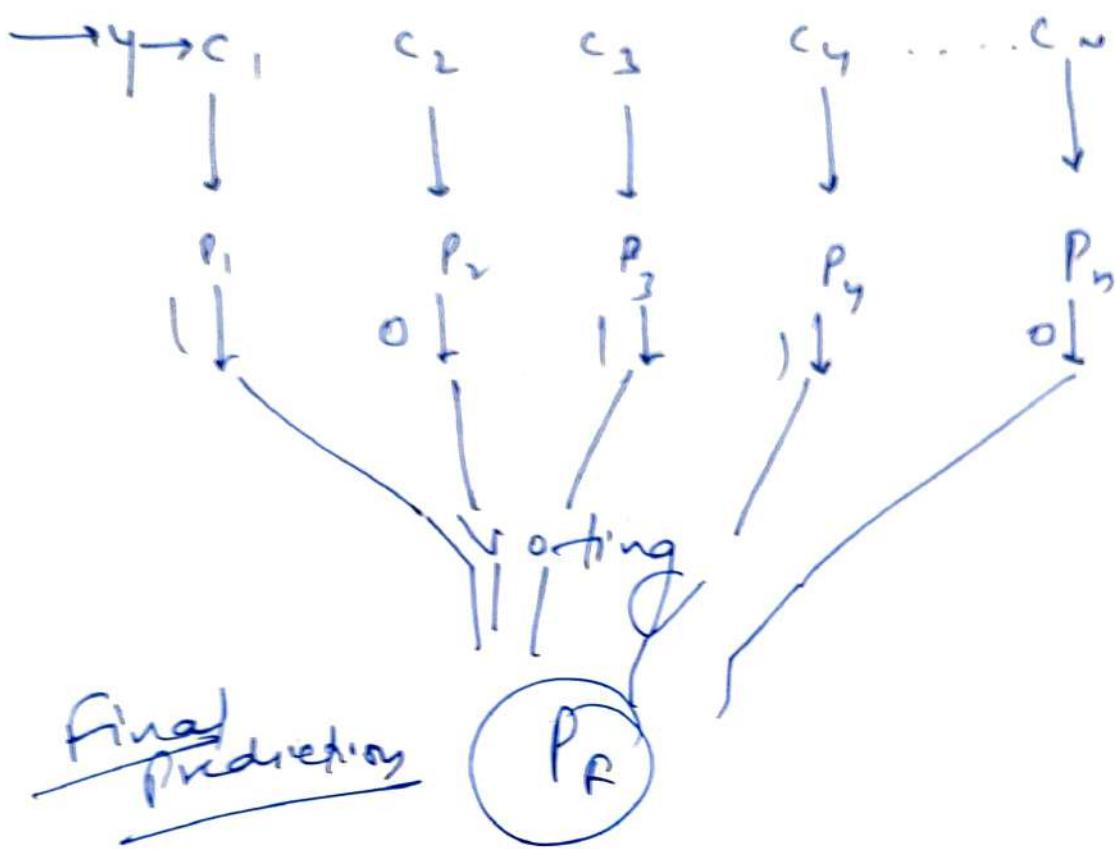
- n instances and individual instances associated with different weights and all weights are equal,
- There will be some misclassification
- Before to random sampling update the weight particularly those which are wrongly classified in previous model
- Weights of wrongly classified now bias higher degree than chance of selecting wrongly classified instance in the new model is more
- If all models are created then combine all models and create a strong model by combining all weak model
- use voting method for test data

Voting classifier Classification

① Hard voting also called majority voting
 On Training data we obtain n classifiers then every classifier feed test data
 so all classifier give their prediction after prediction voting is used for final prediction

- majority of voting decide which is the correct answer for classification

Training data



Soft voting depending on probability
range from 0 to 1

$$c_1 \quad 0.9, \quad 0.1$$

$$c_2 \quad 0.8, \quad 0.2$$

$$c_3 \quad 0.4, \quad 0.6$$

$$\text{Final Prediction} = \frac{0.9 + 0.8 + 0.4}{3} = 0.7$$

$$\frac{0.1 + 0.2 + 0.5}{3} = 0.3$$

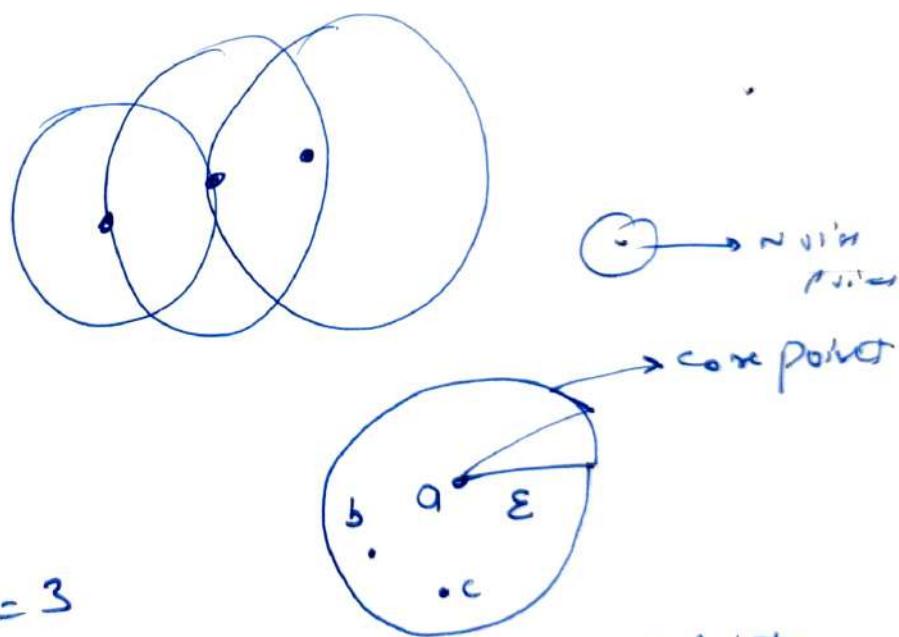
1, 2, 3 class

Random forest It is like ensemble classifier which use decision tree

we have original data and from that create a bootstrap data (BD)

- select subset of total variables DBSCAN "density Based spatial clustering of Application with noise"
- A is better than B and C
- Voting counting is performed, voter with majority of votes attribute win

DBSCAN



noise point :

- minimum points = 3
- In a circle
- Boundary Point :- neighbourhood of core point
- Directly density reachable

K mean Algorithm

Euclidean distance

$$\sqrt{(x_0 - x_c)^2 + (y_0 - y_c)^2}$$

- Agglomerative
 - Divisive
- Down to up
up to Down

```
fig, ax = plt.subplots(figsize=(8, 8))
```

```
sns.distplot(dataset.salary)
```

```
dataset['salary'].fillna(dataset['salary'].median(), inplace=True)
```

```
salary_dataset = salary - dataset.dropna(how='any')
```