

Towards Intelligence Machine:

Machine with human-like traits.

↳ casting intelligence into a machine
in human. { "sensors, control.
central nervous system" }
↳ cognitive aspects &
& decision-making.
(CNS).

Initially computer → store info.
cognitive process → develop new math tool and
h/w.

Biological neuronal process are complex and
difficult to understand.

Scientist can achieve only limited understanding
of brain. → Based on this new computer
theories under Neural Net has been evolving.
cognitive functions of brain → based on
relative grades of
information.

Our neurons accept info in the form of grades
rather than numbers.

The perceptions and actions → based on
the form of relative grades

Theory of fuzzy logic → provides notion of grades
emulation of cognitive functions.

Neural ^{Cognitive}
N/W ^{that} + Fuzzy } logic } → cognitive information.

Needs, Motivations and Rationale:

- 1. Collecting Information is very essential.
The collected info may be categorized on the basis of nature of experience
- 2. experimental data (example, sample, measurements, record, patterns of observations)
- 3. Structured human knowledge (experience, expertise, heuristics)
↳ expressed in (If-then rules.)
- 4. Eg: In man-machine control system, an experienced process operator employs → IF-THEN rules to control a process.
Operator estimates the important process variables and based on that it manipulates ctrl signal.
- 5. Industry process.

Soft computing | ML:

To solve decision making problems:

↳ variables related to the problem.

↳ Relation among the variables are expressed in terms of mathematical equations.

Decision → given by solution of These equations

Generalizing :

(2)

Recognizing similarity b/w different situations. This is very useful b/c we can use our knowledge in situations unseen earlier.

(4)

In ML tip is defined by learning task.

Part 1,

Different types of learning task appears in real-world application.

Classification learning → classification based on measurement / observation.

Regression learning → prediction.

clustering → Grouping the observations.

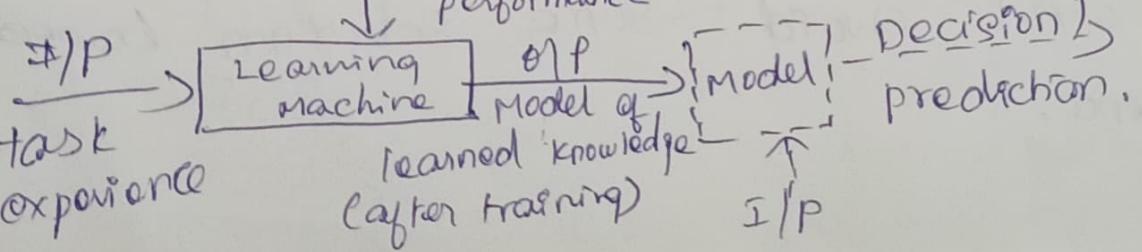
classification | regression → directed / supervised learning.

clustering → unsupervised learning.

Here experience → available in the form of data

I/P → Raw data can be in any form numeric, image, video etc.

O/P → Learned knowledge. → model. ↗ used for decision-making



Now-a-days ocean of data to be processed
Humans are unable to extract information from them.

So human-like ability need to be incorporated into the S/W. This is the essence of M/L.

With M/L a new mathematical theory has been emerged which is built on learning, memorizing, adapting and generalizing.

well - posed M/L problems:

M/L concerse with the question how to improve the performance of a task through experience?

To have a well-defined learning problem, we must identify following three features:

- * The learning task.
 - * The measure of performance
 - * The task experience.
- } → key concepts
learning from experience.

Important aspect of learning from experience is remembering, adapting and generalizing.

Remembering and Adapting:

Recognizing the situation and used the solution which worked last time for similar type of problems and avoid some solution based on past experience.

Cross-validation:

y2hbbhb. (3)

To overcome overfitting.

k-fold

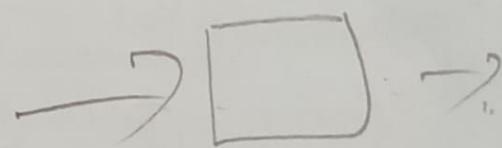
↳ Evaluate model for multiple times.

K → groups that given data sample can split.

Procedure:

1. shuffle the dataset randomly
2. split dataset into k-group.
3. For each group
 1. take a group as test dataset
 2. train with remaining group
 3. Fit the model & evaluate
 4. Retain the evaluation score and discard the model.
4. Summarize with evaluation score.

Ex:- 1, 2, 3, 4, 5, 6.



K = 3.

1, 3, 5, {2, 4, 6} →

	Train	Test	Score	Model based on these observations
F1	4, 3	F1	F1 Score	
F2	1, 6	F2	F2 Score	
F3	2, 5	F3	F3 Score	

Average Score

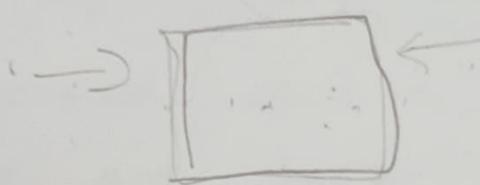
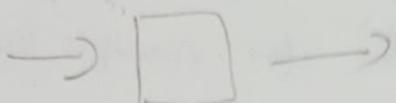
$X_{\text{train}} = C_1, C_2$

$Y_{\text{train}} = D_1$

`split(dataframe)` —> !

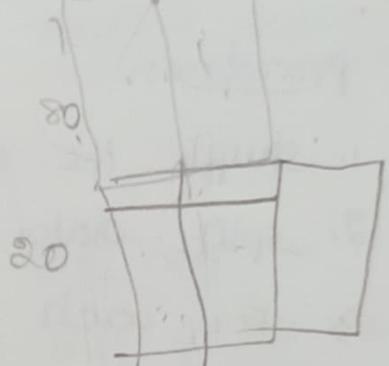
① shuffle.

⑤ 80%.



Condition-test

C_1	C_2	D_1
$\Rightarrow DF$		



$X_{\text{test}} = \underline{C_1, C_2}$

$Y_{\text{test}} = D_1$

$m_1 = \text{LR}_{\text{fit}}(X_{\text{test}})$

$m_2 = \text{LR}_{\text{fit}}(X_{\text{test}})$

$$\left\{ \begin{array}{l} \overline{m_2} = \text{LR} \\ \text{MAE} \\ \text{MSE} (\underline{Y_{\text{test}}, m_2}) \end{array} \right\} =$$

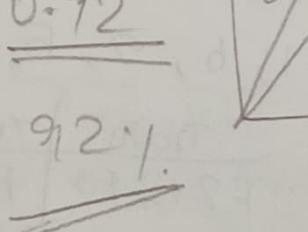
$R^2_s(C) = \underline{0.92}$

MR.

$$\begin{aligned} C &= a + b \\ &= \dots \\ &= a, b \end{aligned}$$

MAE

MSE



MR → optimal

0.92/↑

Application of ML in diverse fields:

(4)

N/L → ment for automatic learning from voluminous datasets.

mining knowledge from data.

Data mining → Google → gets info from billions of web pages. → list website → using PageRank algo.

Personalized product recommendations:

Eg:- Amazon → recommender system → ML.

based on data from social n/w.

→ to predict suitable content for the user.

Information collected based on rating

↳ explicit rating → given by users.

↳ Implicit rating → rating estimated from social media.

Computer n/w security:

system → to find unusual patterns of activity.

Privacy-preserving data mining assist in

protecting privacy-sensitive informations.

Medical diagnostics:

↳ Accurate diagnosis → challenging.

↳ Medical images. → major success.

ML + DL → more accurate diagnosis.

Finance Domain:
Helps to offer personalized services to the customers at low cost.
transactional fraudulent Predict frauds. → ML can scan large amount of data and identify if there is any behavior.
stock market forecasting:
prediction of stock index.
Here data is in the form of time series.
above two domain data are independent.
So here prediction is done using the past history.

Machine Vision:
Deep learning applied in major success.
Images were captured → converted to numerical values for the process.
Numerical values → pixel values
Crayon scale → matrix of pixel intensities
ML algo. ← vector representation
⇒ Biometric recognition, Medical image analysis,
handwritten digit recognition system.

Speech Recognition:
Signal processing technique → speech signal
↓
real values.
Virtual Personal Assistant → finding info when asked over voice.

Patient Index	Headache	Muscle-pain	Temperature	<u>Flu</u>
1	Yes	Yes	high	
2	No	No	low	
3			Normal	
4			Very-high	

Data are in Nominal representation \rightarrow acceptable if table size is small.

Large \rightarrow converted to some numerical value.

Direct / supervised learning \rightarrow outcome for each observation is known as a priori.

expressed by our \leftarrow attribute called decision attribute.

Information system of this kind are called decision systems.

In the above table Flu \rightarrow decision attribute. which categorize patient into \leftarrow two classes Flu \rightarrow Yes or Flu-No.

Other attributes are called condition attributes.

Features x_j	x_1, x_2, \dots, x_n	Decision y
Instances $s^{(i)}$		
s^1		
s^2		
\vdots		
s^N		

values can be numeric, categorical or discrete.

Training experience is available in the

controlling a robot with voice commands

Text Mining:

Identification of pattern in text data.

Eg: Natural Language Processing.

spam Email filtering.

Robotics and Automation:

ML → source of intelligence for robot.

Robots in Industrial automation, medical robots, military robots, so on.

Data Representation:

Experience → raw data.

Raw data → requires pre-processing.

this leads to information system that represents knowledge in raw data.

Data will be stored in datawarehouse.

Data warehousing provides integrated, consistent and cleaned data to machine learning algorithm.

Some application data → flat files → data tables.

Row → measurement / observation.

column → gives values of an attribute.

{ Row → instances, samples, record, pattern, object, cause, events.

{ Column → attributes & features.

↳ depends on the application.

(b)

Time series data:

Data structure of forecasting problems.

Time series data \rightarrow sequential. \rightarrow a sequence of observation is measured over time and each observation is indexed by t .

The measurement are taken at fixed time interval.

$$y(t) = f(y(t-1), y(t-2), \dots)$$

↓
O/P at t,
↓
past values of the signals.

NARMA \rightarrow (Nonlinear, Auto-Regressive, Moving Average) model.

$f(\cdot)$ \rightarrow very complex.

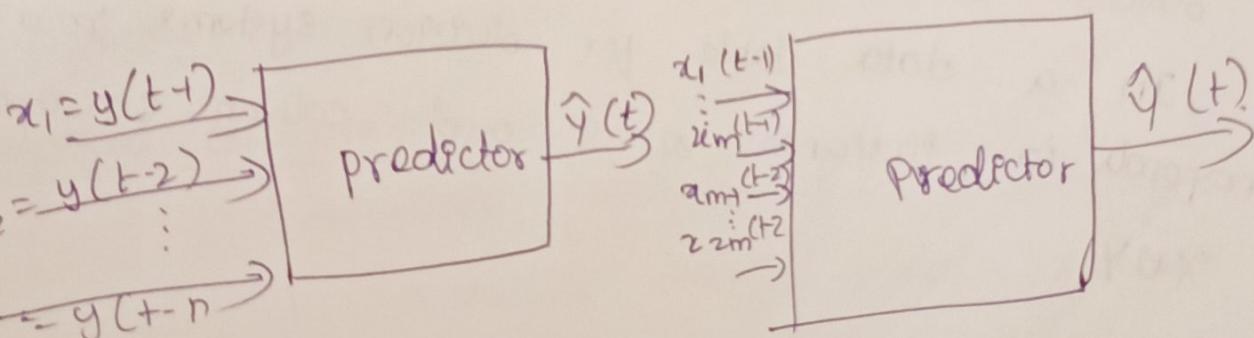
The number of lags L in the series is dependent on application.

$$\hat{y}(t) = f(y(t-1), \dots, y(t-L))$$

More past info \rightarrow prediction more accurate.

If less noise.

Time series \rightarrow Regression problem.



Univariant, Bivariant, Multivariate.
 Student, height
 ↓
 one variable 2 variable, 2 variables.

form of N examples.

$s^{(i)} \in S, i=1, 2, \dots, N$ where S is a set of possible instances.

Each instance $\in S$ with n features given by.

a vector

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = [x_1 \ x_2 \ \dots \ x_n]^T \in \mathbb{R}^n$$

the set X is a finite set of feature instances.

vector $x^{(i)}$ for all possible instances.

the pair (S, x) contributes the information system

where

$$\{s^{(i)}; i=1, 2, \dots, N\} \in S$$

$$\{x^{(i)}; i=1, 2, \dots, N\} \in X$$

$$\{x_j^{(i)}; j=1, 2, \dots, n\} = x^{(i)}. \quad v_{nj} = \{v_{1nj}, v_{2nj}, \dots, v_{dnj}\}$$

The tuple (S, X, Y) constitutes a decision system

where finite set of condition attributes $x \in X$, and

the decision attribute (output) $y \in Y$

In a data table for decision systems, row corresponds to instance of S and column corresponds to $X \cup Y$.

Structured / Unstructured Data:

Data modeling

establish, logical structuring.

↳ determine how data is stored, organized, and manipulated in DB.

Structured → Pre-defined data model.

↳ Relational model. → Tables.

↳ Easy to aggregate from various

easy location.
to analyse.

Unstructured → no data model / not organized in a predefined manner.

↳ text & multimedia data.

↳ images, videos, audio etc.

Having internal structure but it doesn't fit neatly in relational database.

Semi-structured. → Doesn't fit to data model but have some organizational properties that makes it easier to analyze.

Little processing → convert to the format acceptable by machines.

Eg: JSON [JavaScript Object Notation].

XML, CSV, HTML.

Power of Unstructured data:

80% → unstructured.

so non-relational DB → can be used to store unstructured data.

↳ Apache, MongoDB, Hadoop etc.

DL → tools for unstructured data } → analyze and search for some structural form in the data.

Text → will not have random letters as neighbours.
↳ follow grammar.

Images → nearby pixels are close to each other in values.

Speech → commonly used sounds → phonics.

Natural Language processing
Computer vision
Speech synthesis } → tool to capture some structures in unstructured data.

This is what done in ↳
DL.

Domain knowledge for productive use of ML:

ML knowledge + domain knowledge

Main for feature selection + feature extraction.

The key is to use features that
are computationally feasible
leads to good machine learning success and
reduce problem data into manageable data without
discarding valuable information.

compactness of feature vector is important } Reduce computation complexity, generalization.

Feature with less info } leads to poor generalization.

More data → More chances for errors.

↓
same range → no outliers

→ Beyond range → called outliers.

Some data set → large number of missing values.

Data loss ← Dropping ↙ Imputation → based on existing pattern.

know the importance of each attribute.
↓ In this case:

Domain knowledge places major role.

Forms of Learning:

↳ Supervised Learning

↳ Unsupervised Learning

↳ Reinforcement Learning

↳ Learning based on natural processes.: Evolution, Swarming and Immune System.

Supervised Learning:

uses labelled dataset.



Model will be developed and

tested with labelled dataset.

$$D = \{ s^{(i)}, y^{(i)} \}_{i=1, \dots, N}$$

$$g^{(i)} = x^{(i)} : \{x_1^{(i)}, x_2^{(i)} \dots x_n^{(i)}\}.$$

Two types of task in supervised Direct learning

→ Classification [pattern recognition]

→ Regression: [numeric prediction]

Unsupervised / Undirected Learning: [self-instructions].

No labels.

From the given set of p/p it will find a pattern and classify according to pt.

or similarity ↳ cluster analysis → very useful for data exploration
↓
Dimensional reduction algo. ↳ Understand structure of data

Reinforcement Learning:

↳ Learning the optimal behaviour in an environment to obtain maximum reward.

↳ Learn by trial and error method.

the two aspect are

trial-and-error

cumulative reward

Agent has to take the action and learn from the result of action taken.

Learning based on natural processes:

⑨

Evolutionary computation:

Driver towards to develop search and optimization methods that helps to solve complicated problems.

Genetic algorithm, evolution strategies, evolutionary programming and genetic programming will come under this.

Swarm Intelligence:

↳ Features of system of unintelligent agent with inadequate individual ability

1)

Intelligent behaviour ← collectively
→ Base for ant colony optimization. Finds shortest root for food.

2) Processing the knowledge → PSO. {particle swarm optimization}

Artificial Immune System:

↳ resist from diseases.

AIS → applied to solve pattern-recognition problems and cluster data.

Natural immune system → Match patterns of the cells and protect our body from foreign cells.

Supervised Learning:

Learning from observation.

$S \rightarrow$ set of observations/patterns.

$$s^{(i)} \in S; i=1, 2, \dots, n$$

A pattern is specified by n attributes/features. $x_j;$
 $j=1, 2, \dots, n$

where each feature has real values for a pattern

Domain of $\forall x_j \in \mathbb{R}$

Data pattern $s^{(i)}$ has feature set $\{x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}\}$
 where $x_j^{(i)} \in \mathbb{R}$.

Each pattern can be visualized with n
 numerical features as point in state space \mathbb{R}^n .

$$x = [x_1, \dots, x_n]^T \in \mathbb{R}^n$$

$X \rightarrow$ set of feature vector. $x^{(i)}$ for all patterns.
 ↳ representation space.

In supervised learning decision attribut

$y \rightarrow$ a prior

In classification problem $y \rightarrow$ given by the set

$$V_y = \{1, 2, \dots, M\}$$

O/P $y^{(i)}$ for $s^{(i)}$ takes values from the set V_y .

Thus in classification problem each
is associated with $y^{(i)} \in V_y$.
For regression.

$V_y \in R$
each $x^{(i)}$ has $y^{(i)} \in \overline{R}$.
 \Rightarrow State space.

Different patterns $x^{(i)}$ in $X \rightarrow$ have different
frequencies.

different probability distributions.

Each training examples $\xrightarrow{\text{trainer}}$ learning m/c
forward ~~op~~ $\xrightarrow{\text{m/c}}$
②

Training experience is in the form of data D.

$$D: \{x^{(i)}, y^{(i)}\}_{i=1,2,\dots,n}.$$

In general learning \rightarrow reliable \rightarrow training follow
distribution similar to unseen data.

Assume D is independently drawn and identically
distributed and a m/c is defined by a function

$$f: X \rightarrow Y \rightarrow \text{op. space} \quad \left. \begin{array}{l} \downarrow \\ \text{op. space} \end{array} \right\} \text{trained machine.}$$

Able to bind $\hat{y} = f(x)$ for a given pattern x.

To assess success of learning, \rightarrow evaluation
criterions are used.

Evaluation criterions

Based on decisions like errors, profits or loss, penalties or rewards etc.

In supervised learning \rightarrow minimization criterion potential loss.

From set of possible ML functions

optimal one \leftarrow choose \leftarrow
 \hookrightarrow which minimizes the loss.

ML function can be defined by

$f(x, w)$
 \hookrightarrow adjustable parameters.

A loss function

$L(y, f(x, w)) \rightarrow$ measure of error b/w actual output y and estimated o/p.

$$\hat{y} = f(x, w)$$

If $\hat{y} \rightarrow$ estimate o/p

$y \rightarrow$ actual o/p.

Loss associated with decision $f(x, w)$

$$E[L(y, f(x, w))] = \int_{x,y} L(y, f(x, w)) \underbrace{p(x, y)}_{\downarrow} dx dy = R(w)$$

Joint probability distribution of data.

(p.e) loss to calculated on entire space
of x and y .
 $(i|p)$ $(o|p)$

A risk \rightarrow predictable loss in decision theory.
 $R(w) \rightarrow$ Risk function.

Problem is to find decision function $f(x, w)$ against $p(x, y)$ that minimizes the risk function $R(w)$.

thus it is easy to design 'optimal' classifiers) numeric predictors if joint probability density function $\frac{p(x, y)}{\downarrow}$ is known.
very hard to find.

Empirical - Risk - Minimization:

Since $p(x, y) \rightarrow$ not known true risk given by $R(w) \rightarrow$ cannot be calculated.

With available info \rightarrow trying to calculate the risk.

Training dataset $\rightarrow x, y \rightarrow$ unknown probability distribution.
So, risk must be approximated with empirical risk, $R_{emp}(w)$.

$$R_{emp}(w) \triangleq \frac{1}{N} \sum_{i=1}^N L(y^{(i)}, f(x^{(i)}, w)).$$

(P-e) average over $p(x, y) \rightarrow$ replaced by
 average over training samples. (12)

In classification

$$L(y, f(x, w)) = \begin{cases} 0 & \text{if } y = f(x, w) = \hat{y} \\ 1 & \text{otherwise.} \end{cases}$$

classification error.

In regression.

$$L(y, f(x, w)) = \frac{1}{2} (y - f(x, w))^2$$

(P-e) difference b/w real and estimated value.
 $\cdot y - f(x, w)$.

But finding empirical risk is an ill-posed problem.

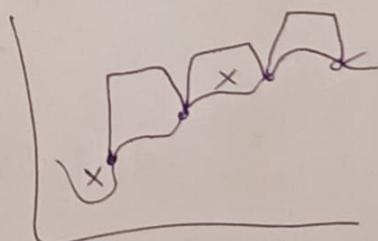
(P-e) problem may have a number of possible solutions.

The will work based on the assumption that unseen data are having similar distribution like training data. But this can not be guaranteed.

Real-world data may be incomplete, noisy and inconsistent.

Consider a regression problem.

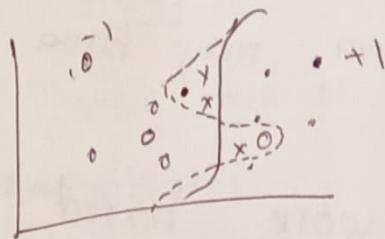
$$(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in X \times Y$$



Functions which interpolate those ~~values~~ data points will be zero for R_{emp} .

$X \rightarrow$ outside training data set.
 both are performing bad.
 thus, interpolate functions that results in
 zero empirical risk can mislead.
 There are many other approximating functions
 that will minimize empirical risk (training error)
 but not necessarily the true (expected) risk.
 consider the classification problem \Rightarrow two-class
 pattern recognition.

$$f(\cdot) : X \rightarrow \{\pm 1\} \text{ that is, } Y = \{\pm 1\}$$



~~class~~ line \rightarrow separating correctly,

consider test points \rightarrow not known
 but not used for training.

~~function~~ \rightarrow
 classifier \rightarrow zero training error not able
 to get all test point right.

\therefore minimization of empirical risk (i.e. training
 error) is thus not a solution for a learning
 task.

Minimization of true error \rightarrow Mahesh hidden
learning task not possible $p(x, w) \rightarrow$ unknown. (13)

\hookrightarrow Solution is productive inference.

Inductive Learning: \rightarrow Learning from examples.
 \hookrightarrow Generalized conclusion

Correlation

$$(x^{(i)}, f(x^{(i)})); i = 1, 2, \dots, N$$

returns function $h(x)$ that approximates $f(x)$

In statistical literature

approximation function

$h(x)$ is called hypothesis function.

$f(x) \rightarrow$ true function that maps input

space x to output space y .

But $f(x)$ is not known in real-world decision-making problem.

It is not easier to tell whether $h(\cdot)$ is a

good approximation of $f(\cdot)$.

\hookrightarrow Generalize well \rightarrow predict novel patterns correctly.

\therefore Generalization performance \rightarrow fundamental problem in inductive learning.

Judge using test dataset.

The off-training set error \rightarrow used to measure generalization performance.

Assumption \rightarrow Model hypothesis related to unseen patterns is one induced by observed training set.

Obtain general hypothesis \rightarrow empirical regularization
training sets. \leftarrow over
induce approximation \leftarrow
over mapping.

For a particular data set if a specific hypothesis
outperforms another, it is not necessary it will
work fine for all data sets.

BIAS and VARIANCE:

Random sample $\rightarrow D$ of n and measure training

error from the function.

$$R_{\text{emp}}(w) \triangleq \frac{1}{n} \sum_{j=1}^N (y_j - f(x_j, w))^2.$$

Let's denote error based on D and hypothesis h as $\text{error}_D[h]$.

For various samples (x_i, y_i) , we will get different $\text{error}_D[h]$.

$\therefore \text{error}_D[h] \rightarrow$ random variable.

For k experiment, the average over k :

$$\text{error}_D[h] = E_D \{ \text{error}_{Dj}[h] \}.$$

$E_D \{ \cdot \}$ \rightarrow expectation (or ensemble average).

Let's the regression function estimated is denoted $h(x; D_j)$ for training set D_j .

Due to random selection some approximation will be excellent ; some will be poor.

The natural measure of the effectiveness of the estimator can be expressed as its mean-square deviation

$$\text{error}_{D_j}[h] = [h(x:D_j) - f(x)]^2$$

The value of this quantity depends on the dataset D_j .

Average over complete ensemble of datasets as

$$\text{error}_D[h] = E_D \{ [h(x:D_j) - f(x)]^2 \}.$$

A non-zero error can arise for essentially two reasons.

1. It may be hypothesis function $h(\cdot)$ is, on average, different from regression function $f(x)$. This is called bias.

2. It may be hypothesis function is very sensitive to particular D_j , so that for a given x approximation error is larger for some dataset and smaller for other dataset. This is called variance.

variance measure level to which the hypothesis function $h(x; \theta_j)$ is sensitive to specific selection of data set.

Higher the complexity of hypothesis function, [i.e. more flexible with many parameters], lower is the approximation error.

Variance \rightarrow reflects capability of the trained model to generalize to other data samples.

Lower the variance \rightarrow estimate of $f(x)$ does not change much.

But higher the complexity of hypothesis, higher the variance.

Simpler model \rightarrow higher bias, lower variance.

Complex model \rightarrow lower bias, higher variance.

\therefore hypothesis that results in low bias and low variance is needed to minimize over all mean square error.

This is called bias-variance dilemma.

or bias-variance trade off.

In general finding optimal trade off is hard, but acceptable solution can be found.

$$\begin{aligned} \text{error}_{D_j}[h] &= [h(x; D_j) - f(x)]^2 \\ &= [h(x; D_j) - E_D\{h(x; D_j)\}] + \\ &\quad E_D\{h(x; D_j)\} - f(x)^2 \\ &= [h(x; D_j) - E_D\{h(x; D_j)\}]^2 + [E_D\{h(x; D_j)\} - f(x)]^2 \end{aligned}$$

$$(a+b)^2 = a^2 + b^2 + 2ab.$$

We take expectation of both sides over ensemble of dataset D.

$$\text{error}_D[h] = E_D\{[h(x; D_j) - f(x)]^2\}.$$

$$\begin{aligned} &= E_D\{[h(x; D_j) - E_D\{h(x; D_j)\}]^2\} + \\ &\quad \left(E_D\{[E_D\{h(x; D_j)\} - f(x)]^2\} - \right. \\ &\quad \left. 2 \downarrow \quad \downarrow \right). \end{aligned}$$

Here :

$$\begin{aligned} \text{variance} &\leftarrow E_D\{[h(x; D_j) - E_D\{h(x; D_j)\}]^2\} \\ \text{Bias} &\leftarrow E_D\{[E_D\{h(x; D_j)\} - f(x)]^2\}. \end{aligned}$$

Bias \rightarrow measures the level to which average hypothesis function differs from desire function $f(x)$.

Occam's Razor principle and overfitting

(1b)

avoldance:

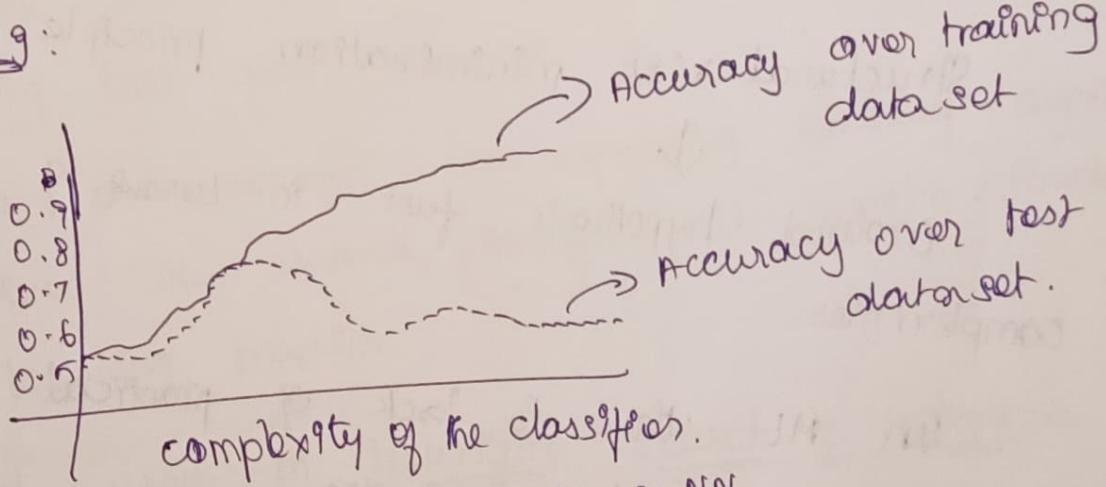
If there are two algorithms and both of them perform equally well on training set, then simpler algorithm can be expected to do better on a test set.

Simpler \rightarrow lesser parameters,
" training time
fewer attributes

" When a design with 'good enough' solution arrives, we not necessarily go for optimal one.

Occam's Razor principle suggest hypothesis function that avoid overfitting of train data.

overfitting:



[more no. of weights in NN,
more no. of nodes in tree etc.]

Accuracy get affected when test dataset lie in different state space compare to train dataset.

Heuristic search in Inductive Learning:

Int'l 34 MC Learning aim \rightarrow built a statistical model

22, 28, 29, 35, 46, ~~47, 48,~~
54.

good generalization

Success \rightarrow depends on \rightarrow hypothesis space complexity
① sample complexity.

② bias-variance trade-off.

Best generalization is achieved through the best compromise b/w requirement of small bias and small variance.

In order to find optimal bias & variance have to control complexity of hypothesis function.

Structural-risk minimization principle



ordered hypothesis funt in terms of their complexities.

In ML there is lack of practical tool for real-life practical complexities.

Bias-variance trade off is also theoretical result because lack of tools to find optimal bias and variance.

∴ Heuristic strategies ~~for like~~ 17
overfitting avoidance - regularization, early stopping,
pruning are some of heuristic tools explored
today.

Search through Hypothesis Space:

Heuristic \rightarrow cannot be precisely predefined.
uses trial-and-error method for good solution.
the search for hypothesis is done with the
following steps:

i) the search is first focused on class of hypothesis,
chosen for learning task. Prior knowledge and
experience are needed. Different hypothesis for different
learning task and dataset.

ii) For each member, corresponding learning algorithm
organize the search through all possible structures
of learning machine.

Techniques used in heuristic search to optimize
hypothesis complexity:

Regularization: \rightarrow promotes smoother function by
creating new criterion function not

only based on training error but also
algorithmic intricacy.

It avoid extremely complex hypotheses.
Now criteria can be written as

$$\tilde{E} = E + d\Omega$$

= error on training set + $\lambda \times$ hypothesis complexity.
 λ weight of penalty.
Regularization term.

when $\lambda = 0 \rightarrow$ no regularization. } causes overfitting
But have some variance }
↑ in bias.

when $\lambda \uparrow$, variance \downarrow .

Till some point, the decrease in λ
is useful b/e variance \downarrow without subsequent
↑ in bias.

But at certain point bias starts ↑.
and thus leads to underfitting.

$\lambda \rightarrow$ optimised using cross-validation.

Early stopping:

Training of learning machine \rightarrow iterative
decrease in error function.

During specific training session error \downarrow as iteration \uparrow
stopping iteration before attaining } \rightarrow affect the hypothesis
minimal error

ending:

Alternate solution to early stopping. (18)

Removing certain parts of model \rightarrow to prevent overfitting.

Goal is to improve model's performance, generalization and efficiency by reducing the complexity.

Evaluation of Learning System:

Some important aspect of evaluation are

Accuracy: Generalization capability is the measure of accuracy.

Robustness:

Machine should perform well in all circumstances including data with noises, irrelevant data, incomplete and with \hookrightarrow These all will present in real-time data.

Robustness based

Computational complexity and speed:

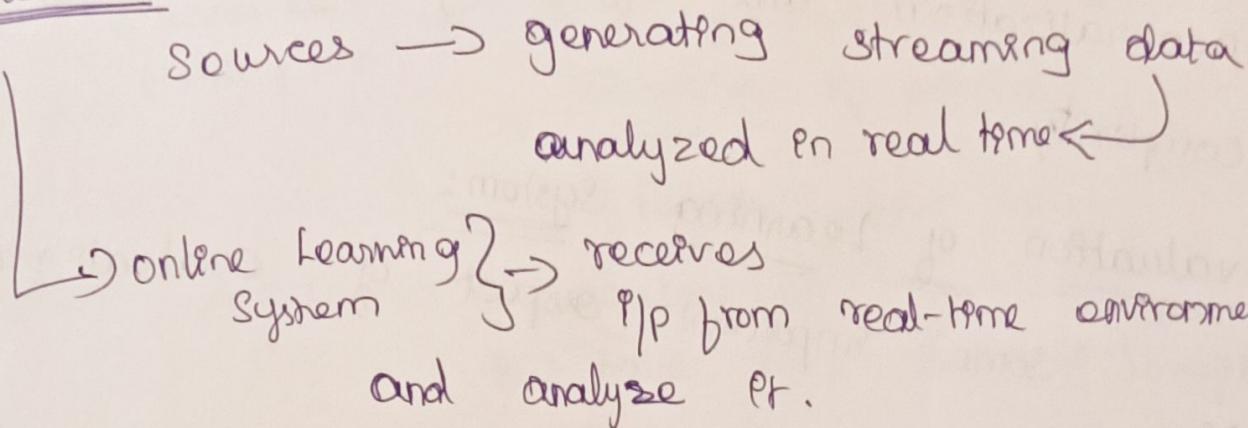
\hookrightarrow for efficiency of learning algo.

(i-e) how fast a machine can work and how much memory is needed by the machine.

Interpretability:

- ↳ It is subjective and difficult to eval.
- ↳ It is easy in decision tree,
but interpretability \downarrow when complexity \uparrow .

Online Learning:



Scalability:

High level of scalability is needed.

The assessment of scalability can be done with dataset size complexity.

Combining Multiple Models: Ensemble learning.

No Free Lunch \rightarrow no single learning algo. always gives most accurate learner.

Each algo \rightarrow model \rightarrow set of assumptions.
 this inductive bias $\begin{cases} \xrightarrow{\text{Error}} \text{if assumptions do not hold for data} \\ \xrightarrow{\text{less error}} \text{if assumptions work.} \end{cases}$

General model $\begin{cases} \xrightarrow{\text{choose the best performing}} \\ \xrightarrow{\text{combine them to produce ensemble of learned models.}} \end{cases}$

But suitably combining multiple models, accuracy can be improved.

Even though it is counter to Occam's razor [simplicity], due to ↑ predictive performance pt to be applied to both classification and regression tasks.

The goal of ensembling is to build a predictive model by combining the strength of base models.

The main aim is to find set of diverse learners that complement each other. And also the gain in overall success will not be achieved if [differences in their ↑ decision]

the base learners are not reasonably accurate in the domain of problem.

Maximizing overall accuracy of individual learners } more important
Requirement of diversity.

∴ Ensemble learning can be broken down into two tasks.

⇒ Generating a committee of base learners that complements each other.

⇒ combining opf of weak learners for maximum accuracy of the composite predictor, which leads to strong kernel.

The Base Learners:

Different algorithms to generate set of base learners.
 BTW makes different assumption about the data.
 → different predictors [Diversity].

Eg: For a given data set different algo can be applied

like NN
SVM
Decisiontree } → base learners. (or) same algo with different hyperparameters.

Generate base learners (or)
by training it with different subsets of training data.

NN-weight, hidden units ← Eg:

SVM - Kernel function
threshold in decision tree etc.

Combining base learners:

(20)

combining decisions } different
of various base learners } of P single prediction.

⇒ Bagging technique } → individual learners are created independently } → bag of base learners.

⇒ Boosting technique } → new model is impacted by the performance of earlier.

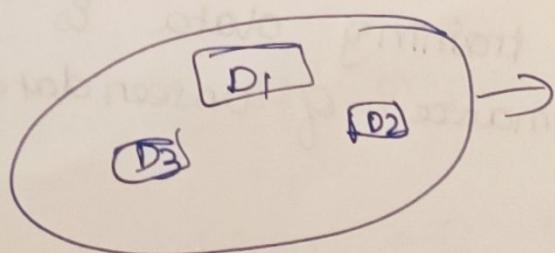
Model with accuracy > average } + new component learner

↳ join decision rules → high accuracy

performance boosted.

⇒ Random forest:

↳ Ensemble consist of Decision tree.



Group of DT forms the Random Forest.

Individual tree

↳ created with random choices of attribute at each node

To split.

Random forest built with } → bagging + Random attribute selection.
↳ combines model of same type [D.T].

In stacked generalization or stacking
different to base learners
complement each other.
↳ different learning algorithms.

But stacking is not used widely
difficult to analyze theoretically.

Decision tree → used more often for
classification problems
success rate is high
bagging / boosting

Estimating generalization errors:

Success of learning } → hypothesis space complexity
+
complexity of samples

Is the performance on training data is
the proper indicator of performance of unseen data.
↳ Ans is -ve.

To predict the performance
of learning model } → test data are used
not used for building
the model.

Some learning models

models } \rightarrow 2 stages

(21)

↳ Building the elementary base structure.

→ optimize the parameter in the structure.

two separate sets of data may be required in the structure.

In such cases three datasets are considered

training data → create the structure of learning model

validation data \rightarrow optimize the parameters.

test data. \rightarrow compute the error rate of the optimized model.

Each of the dataset should be selected independently.

D large set of } \rightarrow training + test data

If properly trained \rightarrow best model.

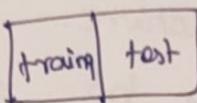
Matrices → test data target } compared with
feature
- predicted value ↗

Ib dataset → training + validation + test
↓
used only for optimization.

Holdout method and Random subsampling:

training + testing.

↓
1/3rd of data.



time series data →

earlier data
↓
training

} latter part
of testing.

consider classification problem:

Each class in the full dataset should be represented in both training and testing.

Eg: Image → cat, dog.

In training the target should include both cat, dog.

similarly in testing also.

To ensure this random samples are selected for training and testing dataset.

This is known as stratification.

safeguard ↓
irregular representation.

More general way → Random subsampling.

↓
Holdout is iterated k times

Average accuracy is considered.

This will be useful if data set $\xrightarrow{\text{split}} 50:50$
But it is possible to take more than 22%
50% of data \rightarrow training.

so a simple variant of holdout cross-validation becomes more powerful method.

K-fold cross-validation:

Given data $\rightarrow k$ folds

D_k where $k=1 \dots K$.

Training and testing done k times.

Iteration #	Train	Test
1	$D_2 \cup D_3 \dots D_K$	D_1
2	$D_1 \cup D_3 \dots D_K$	D_2
:		
K	$D \cup D_2 \dots D_{K-1}$	D_K

If stratification is used, it is known as

stratified K-fold cross-validation.

For each $K \rightarrow$ average accuracy. The

one with least error may be deployed.

$K=10 \rightarrow$ standard number.

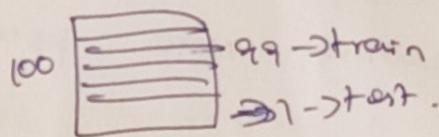
\hookrightarrow But not generalized

\hookrightarrow But practically used.

Leave-one-out Cross-validation.

Exceptional case of K-fold cross-validation.

Each iteration \rightarrow Single sample is left out \rightarrow test set.



Overall accuracy \Rightarrow average of N iterations

Computation $\uparrow \rightarrow$ b/c whole learning process iterated N times.

Small dataset ✓

Large dataset ✗

No stratification can be applied.

Bootstrapping:

\hookrightarrow Sampling with replacement.

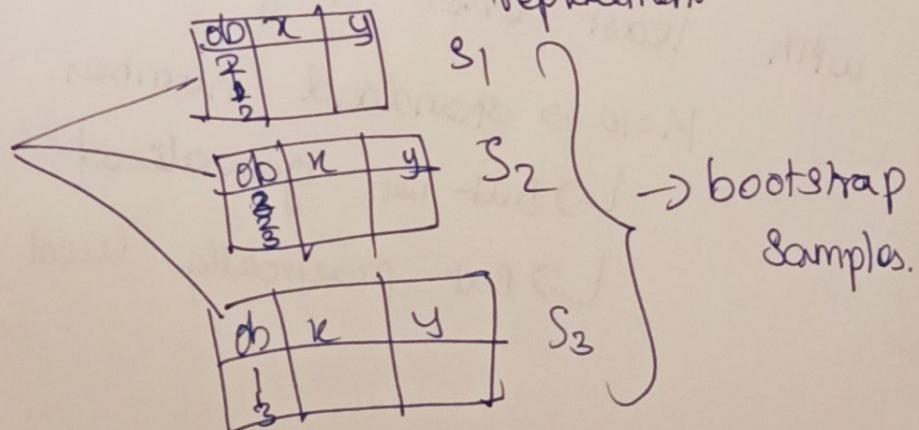
Earlier techniques \rightarrow samples \rightarrow never replaced.
 or
repeated.

Most learning \rightarrow repeated samples \rightarrow impact on outcome.

Aim of bootstrapping \rightarrow form sample dataset by replacement

ob	x	y
1	x ₁	y ₁
2	x ₂	y ₂

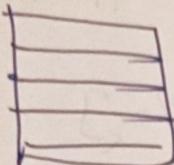
S



(any) bootstrapping technique.

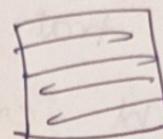
popular - 0.623 bootstrap.

(23)



N instance \rightarrow sample

DS



N instances.

Eg: $N = 10$.

10 \rightarrow observation.

In that sample(10) is created with a replacement.

Ex: 1, 3, 4, 3, 2, 5, 6, 5, 10, 4. \rightarrow Train set

Here 3 \rightarrow 2 times

4 \rightarrow 2 times.

5 \rightarrow 2 times.

[7, 8, 9.] \rightarrow Test set.

If this is repeated several times on an average 63.2% of data \rightarrow training
36.8% of data \rightarrow testing.

Probability that particular instance will be picked $\rightarrow \frac{1}{N}$

not picked $\rightarrow 1 - \frac{1}{N}$

During whole cycle \rightarrow picked $\rightarrow \frac{N}{N}^N$
not picked $\left[1 - \frac{1}{N}\right]^N$

$N \rightarrow$ Large \rightarrow probability approaches $e^{-1} = 0.368$ [base of natural algo].

thus for large dataset

↳ test set \rightarrow 36.8% of instances
train set \rightarrow 63.2% " "

In 10-fold cross validation \rightarrow 90% \leftarrow training.

To compensate, bootstrap combines training error with test error and gives final error estimation as

$$\text{Error estimate} = 0.632 \times \text{Error given by test instances} + 0.386 \times \text{Error given by train instances.}$$

then whole procedure is repeated and results are averaged.

It works best for small dataset.

Metrics for assessing Regression (Numeric Prediction)
Accuracy.

(24)

$$f = X \rightarrow Y; f(x) = y.$$

Task is to find $h(x)$ that explains underlying data, i.e. $h(x) \approx y$ for all samples (x, y) .

In statistics \rightarrow function description } \Rightarrow regression.
of data

In multivariate function $f(x)$ is determination of an approximating function $h(x, w)$
 \downarrow weight vector.

Accuracy \rightarrow measured on $\underbrace{\text{test data set}}_{\substack{\downarrow \\ \text{not used in model training.}}}$

Estimating error for predictor using holdout and random subsampling, cross validation and bootstrap methods.

Several alternative metrics can be used to assess the accuracy of numeric prediction.

Mean Square Error: (MSE)

Most commonly used metric.

For MSE \rightarrow no statistical info are used.
Mean of training set is considered.

$$MSE = \frac{1}{N} \sum_{i=1}^N \frac{(y^{(i)} - h(w, x^{(i)}))^2}{\square} \rightarrow$$

average squared deviation
of predicted value from
true value.

Root Mean square error: (RMSE)

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(w, x^{(i)}))^2}$$

↓
More clearly related to individual error.

Sum of Error Squared:

$$\text{Sum of Error Square} = \sum_{i=1}^N (y^{(i)} - h(w, x^{(i)}))^2.$$

(Expected) Mean square error:

When I/P and O/P are random variables
this metric can be used.

$$(\text{Expected}) \text{Mean square error} - E \left[\sum_{i=1}^N (y^{(i)} - h(w, x^{(i)}))^2 \right].$$

E → statistical expectation operation.

Mean Absolute Error:

$$MAE = \frac{1}{N} \sum_{i=1}^N |y^{(i)} - h(w, x^{(i)})|$$

↓
Average deviation of predicted value from true value

Suppose model 1 $\rightarrow h_1(w, x)$ works well for 95% and
totally differs on other 5%.

{ model 2 $\rightarrow h_2(w, x)$ small error over full range of data.

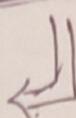
model 2 will be selected. \rightarrow Amount of variability is considered.

MSE \rightarrow exaggerate \rightarrow outliers.

MAE \rightarrow All errors are treated equally according to their magnitude.

Metric for assessing classification (Pattern Recognition)

Accuracy:



(25)

Recognizing patterns
in data.

ENL ^{Pattern recg means} → Assigning label to given I/P pattern.

Different metrics are used to assess the accuracy of a classifier.

Misclassification Error:

Traditional classification algorithm aim to minimize the no. of errors made during classification.

Misclassification error = $\frac{\text{Number of data points for which } (y^{(p)} - \hat{y}^{(1)}) \neq 0}{N}$

→ work well for situations where class tuples are more or less evenly distributed.

When classes are imbalanced → poor performance b/c all misclassifications should not be treated equally.

Log Loss/Cross Entropy: It is a cost function that loss function → evaluates how well algo models data set.

Prediction $\begin{cases} \text{not good} \rightarrow \text{high value} \\ \text{good} \rightarrow \text{low value.} \end{cases}$

Classification accuracy → count of $\begin{cases} \text{no. of predictions} \\ \text{equal to actual value.} \end{cases}$

Log loss \rightarrow takes uncertainty of prediction base on how much predictor varies from actual label.

Log loss \downarrow robustness \uparrow .

Likelihood function:

Consider \rightarrow observations with corresponding class labels. parameters $\xrightarrow{\text{estimated}}$ maximum likelihood approach.

It is computed as follows.

For each data point \rightarrow compute likelihood funt.

Then likelihood funt \rightarrow compute combination of likelihood \oplus for all the data points.

Estimate set of parameters $\xrightarrow{\text{which}}$ maximize the likelihood function.

Let $X = \{x^{(1)}, \dots, x^{(N)}\} = \{x^q, q=1 \dots N\}$.

The probability of predicting y is given by.

$P(y^{(q)} | x^{(q)}, w)$ $\xrightarrow{\text{parameters governing the distribution}}$

For whole $x^{(q)}$; $q=1 \dots N$, $y^{(q)}$ for each $x^{(q)}$ is

$$\prod_{q=1}^N P(y^{(q)} | x^{(q)}, w).$$

Likelihood function can be defined by Here data fixed

$$L(w|x) = \prod_{q=1}^N P(y^{(q)} | x^{(q)}, w). \text{ Parameters vary.}$$

Goal \rightarrow find w that maximizes $L(w)$. (2b)

$$(1) w^* = \arg \max_w \{L(w|x)\}$$

Often $\log(L(w|x))$ will be tried to increase because it is analytically easier.

$\therefore w^*$ that maximize log-likelihood will also maximizes likelihood.

$$\begin{aligned} \log(L(w|x)) &= \log \left[\prod_{i=1}^N P(y^{(i)} | x^{(i)}, w) \right] \\ &= \sum_{i=1}^N \log P(y^{(i)} | x^{(i)}, w) \end{aligned}$$

For binary classification problem.

$$y^{(i)} \in \{0, 1\}.$$

$$\text{If } P(y^{(i)} = 1 | x^{(i)}) = \hat{y}^{(i)}$$

$$\text{then } P(y^{(i)} = 0 | x^{(i)}) = 1 - \hat{y}^{(i)}$$

Combining these two function.

$$P(y^{(i)} | x^{(i)}, w) = (P(y^{(i)} = 1 | x^{(i)}))^{y^{(i)}} (P(y^{(i)} = 0 | x^{(i)}))^{1-y^{(i)}}$$

Here only one of the two classes will be active depending on the value of $y^{(i)} \in \{0, 1\}$.

$$\Rightarrow P(y^{(i)} | x^{(i)}, w) = (\hat{y}^{(i)})^{y^{(i)}} (1 - \hat{y}^{(i)})^{1-y^{(i)}}$$

This gives.

$$\begin{aligned}\log(L(w|x)) &= \log \left[\prod_{i=1}^N (g^{(i)})^{y^{(i)}} (1-g^{(i)})^{1-y^{(i)}} \right] \\ &= \sum_{i=1}^N \log \left[(g^{(i)})^{y^{(i)}} (1-g^{(i)})^{1-y^{(i)}} \right] \\ &= \sum_{i=1}^N \log(g^{(i)})^{y^{(i)}} + \log(1-g^{(i)})^{1-y^{(i)}} \\ &= \sum_{i=1}^N [y^{(i)} \log g^{(i)} + (1-y^{(i)}) \log(1-g^{(i)})]\end{aligned}$$

↑
This function is maximized to get w.

Average of loss over entire dataset

$$\boxed{\text{log loss} = -\frac{1}{N} \sum_{i=1}^N [y^{(i)} \log g^{(i)} + (1-y^{(i)}) \log(1-g^{(i)})]}$$

[log loss is straightforward of log-likelihood function with maximization transferred to minimization.

$$= \underset{i=1}{\cancel{\sum}} - [y^{(i)} \log g^{(i)} + (1-y^{(i)}) \log(1-g^{(i)})]$$

log loss \downarrow when predicted probability is high.

Cross entropy: \rightarrow From information theory.

log loss & cross entropy \rightarrow both can be used interchangeably.

multi class problem.

(27)

Class = M.

$$\text{target output} = \{y_1, y_2, \dots, y_M\} = \{y_{av}\}_{v=1}^M$$

q^{th} Sample \rightarrow target value $y_{av}^{(q)}$
predicted value $p(y_{av}^{(q)} | x^{(q)}) = \hat{y}_{av}^{(q)}$

Common way of representing target output

↳ one-hot vector.

For M classes the target o/p are represented as $M \times 1$ vectors.

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

Sample point that belongs to class 1, has 1 in

the 1st row and others are 0.

$$(i.e) y_1^{(q)} = 1, y_2^{(q)} = 0, \dots, y_m^{(q)} = 0.$$

However, predicted output may give probability value for all M outputs, with

$$\sum_{v=1}^M g_{av}^{(q)} = 1$$

For k^{th} target class, a good classifier will have a high value (close to 1) for k^{th} o/p and almost negligible value for all other o/p.

The probability of predicting class $y^{(i)}$ is given by

$$P(y^{(i)} | x^{(i)}, w); y^{(i)} \in \{y_1^{(i)}, y_2^{(i)}, \dots, y_M^{(i)}\}.$$

On receiving the whole series, the probability of predicting $y^{(i)}$ for each $x^{(i)}$ is

$$\prod_{i=1}^N P(y^{(i)} | x^{(i)}, w).$$

∴ Likelihood function $L(w|x) = \prod_{i=1}^N P(y^{(i)} | x^{(i)}, w)$

For the given $x^{(i)}$ only one cell in the target o/p will be active.

$$\therefore \text{Predicted Probability} = P(y^{(i)} | x^{(i)}, w) = \prod_{q=1}^M (P(y_{qv}^{(i)} | x^{(i)}, w))^{y_{qv}^{(i)}}$$
$$= \prod_{q=1}^M (\hat{y}_{qv}^{(i)})^{y_{qv}^{(i)}}.$$

For $\hat{y}_1^{(i)}$ \Rightarrow class label ($y_1^{(i)}=1, y_2^{(i)}=0, \dots, y_M^{(i)}=0$) and so on.

$$L(w|x) = \prod_{i=1}^N \prod_{q=1}^M (\hat{y}_{qv}^{(i)})^{y_{qv}^{(i)}}.$$

Log likelihood is given by

$$\log(L(w|x)) = \sum_{i=1}^N \sum_{q=1}^M y_{qv}^{(i)} \log \hat{y}_{qv}^{(i)}.$$

Maximization is transferred to minimization and average over entire dataset is taken.

$$\text{log loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{q=1}^M y_{qv}^{(i)} \log(\hat{y}_{qv}^{(i)}).$$