

Lecture - 01

1/9/24

Machine Learning:

Feature Selection: आनकाश्वरी कोल \rightarrow selective cols.

always correlation के से col
इसमें कहा optimal होता है।

Generalization: Model predictor की अनुमति
वाला जल्दी prediction करता है।

\nexists not generalized
 \exists overfitted for Gulshan.

Neural network \rightarrow existed before 2000.

but data was scarce,
जल्दी एवं popular नहीं थी।

Computational power (उपर्युक्त)
Neural networks are data hungry.
We have now data.

Data cloud
Data Generator,

Principal Component Analysis:

vector direction change করে মান

ক্ষেত্র কর্তৃত পরি স্ব. ঘনত্ব dimension

redundant এবং যাই, then feature/Vector
space কর্ম কোর্স।

(Feature selection)

1 EN মাঝে অন্ধিকার করে হতে পারে

FNA

baseless for

model এর ব্যবহারে

100% এরও বেশি \rightarrow K-Nearest Neighbors

এটা কোথা থেকে

1 EN এর সাথে কোথা

ব্যবহার করে আপনার কাজে

এটা এখন এখন এখন

ব্যবহার কর

নির্দেশ কর

ML: Experience \rightarrow performance \uparrow

training data \uparrow performance \uparrow

Supervised Learning:

Model build करते तो Ans भी (target Variable/label) दिया जाता है।

— supervised learning

Regression

"मात्रा वाली prediction."

काम इसी

→
target Variable is a
continuous Variable"

Classification

"Student औ level
term prediction.

Here target Variable
is a discrete
variable."

binary classification

covid - yes/no

Multi-
class
classi-

fication
prob.
cell - healthy/
unhealthy / damaged.

Unsupervised learning

Clustering

Target Variable अनुदान नाही ।

Data अजून Essence विस्तारात नाही ।

Reinforcement learning

learning of human beings.

Win/Loss/Rewards.

Artificial Intelligence

Perception

knowledge

Environment

Perceptual input

Agent

Intention

Effect

Motor output

Sensory input

Memory

Lecture - 02

2/9/24

Label = Target Variable.

Unsupervised

Data \rightarrow Label নাই, We want to extract essence.

Data Similarity.

Reinforcement learning:

objective function

States

Win/Loss Model.

for কিন্তু State G traverser করে profit/loss

analysis, কোথা optimal এর দিয়া সোজাল্য
চেষ্টা কৃতা।

Datasets:

Supervised

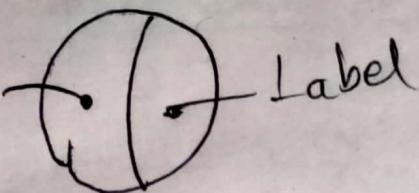
i) Training Set: প্রক্রিয়াজনযোগ্য
info + label মাঝে।

Features + Labels মাঝে।

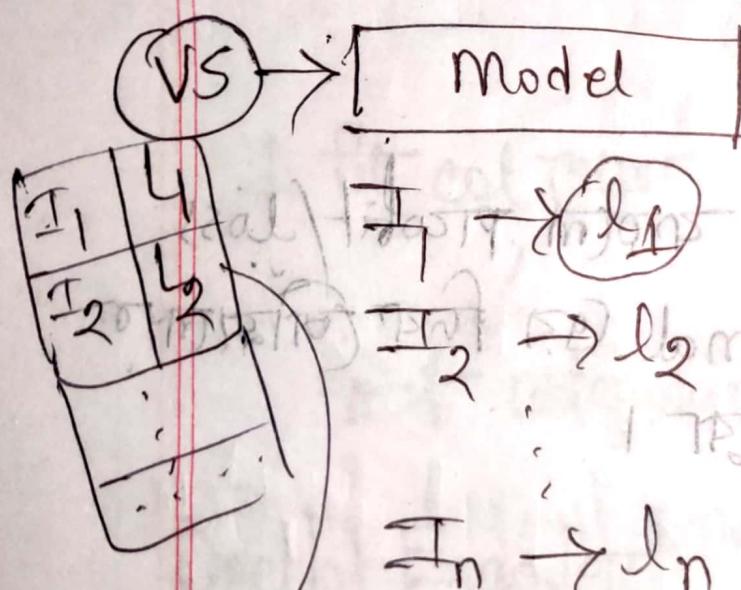
ii) Validation Set: Not used in training.
Performance analysis এ।

~~cycle~~

Features



$f(\text{feature info})$
= label



Ground
Truth
Labels

Their performance
analysis

1. PFA + PMD

2. Prioritization

1. Performance

Model আলোকণ্ঠ উপায় :

1) Training Data ↑
2) Hypothesis, h . } house info from price
} predict করতে পারে।
Horizon broadens করতে হবে।

3) Hyperparameters tuning.

Validation test আলোন করলে Hypothesis
অসম্ভব complex করতে হতে পারে।

Model, Validation Data এর essence নিম্ন

ফলস্বরূপ এটিৰ performance 85%.
জ্ঞান significant ফল।

Test Dataset: Model building এ কোনো আবেদ্ধ
কৈবল্য use কৰা মাৰ্যাদা।

5% data নিম্ন test ফলাফল। Model দুধ
ফলস্বরূপ data চৰি ॥ ৩৫ labels
জৰুৰ নাই, তাত্ত্বিক এই test set
ফোল ফিলে হবে - Ideal

papers publish \rightarrow always Independent
নতুন data না আসা difficult test data ফলে কোম্পানির নথি
loophole খোঁজে।

But industry'তে নতুন data আসতেই হারা।
So, উৎপন্ন ফলে দেখা যাবে নাহি।

'purely unseen' - myth!

Rule of Thumb: 65 — 20 — 15

70 — 20 — 10

TS — VS — Test Set

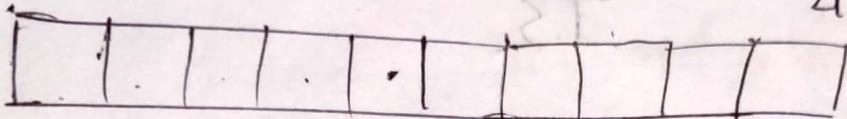
কেবল data জিতে যে test + training কি?

Validation Test \rightarrow H.P. tuning

Independent Test \rightarrow performance reporting.

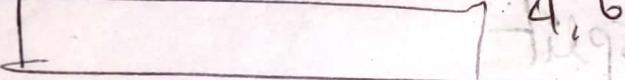
Training & testing: K-fold CV

1



4:6 (stratified)

2



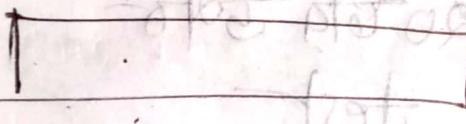
4:6

100

60% \rightarrow neg

40% \rightarrow pos

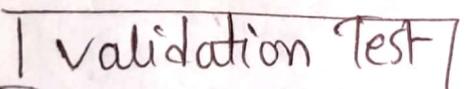
9



4:6

Balanced

10



4:6

Multi-class

No "n" & balance

10 fold

Each has 10 data points

neg - pos
50 - 50

49 - 51 ✓

0 - 1 - 2

33% - 34% - 33%

90% data for training

→ 10% Validation

→ 10% Validation

Folds [1 2 3 4 5 6 7 8 9] [10]
VS
T.S.

Folds [1 2 3 4 5 6 7 8 9 10]
VS
T.S.

pure 80 - 20 split

↓
training (80% data)
↓
Used 20%
↓
test

K-fold

$K = 10$

90% training (

overall 100% test data

for Hyperparameters $\langle h_1 \rangle \dots \langle h_5 \rangle$

for $K = 1 \dots 10$

{

90% পরিষেবা
ফর্মুলার করা

for some h_i

So, we are so happy with
 h_i .

নতুন ফিল্ড প্রেডিক্ট করাতে তান h_i ব্যবহার
করা যাবে ?

customer কে কোথায় কাজ করা যাবে ?.

Final Model

h_1 optimized

h_2 optimized

h_1, h_2 ৯০%
training data
fit $f_{h(1)}$ করা যাবে।
 $f_{h(1)}$ এবং $f_{h(2)}$ এর মধ্যে
১০০% ফিল্ডে বেশি
যাওয়া হবে।
বেশি করা যাবে।
বেশি করা যাবে।

train $f_{h(1)}$ এবং $f_{h(2)}$ এর
model

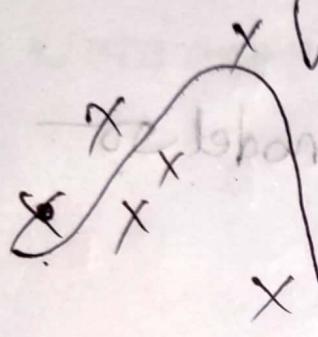
expected to be
generalized,

$f_{h(1)}$

Lecture 3

- Hypothesis \rightarrow এই func ফলে input \rightarrow output
mapping রেখোক করা হয় ক্ষয়গতি ।

২ Linear - Very Simple Model



polynomial

যেটো simple, তাদে bias সোফি।

জো তাৰ নিষেব limit এৰ মাঝিটো বেশি কিম্বা
ভিন্ন ক্ষয়গতে পাইৱল।

প্রতিটা func এৰ inherent bias থাই।

ওৱে মাঝীৰ ও ধূম-চৰকাৰ কৃতে পাইলাম।

- slide ৬ এ ৫ টা hypothesis এৰই datapoints same.

যেটো sophisticated / polynomial degree \uparrow
overfit ক্ষয়গতে তাম উতে।

- Degree 12 টে unnecessary peak দৈখা রাখে।

okham's Razor : parsimony

যথোন্ত জ্ঞান আৰু simple model এ, অধিক সুপ্ৰসাৰ
model

Bias एवं opposite = Variance.

Generalized Model & Variance

— अलग डेटेसेट पर यहां, अलग फूलों के लिए।

II sample एवं डेटा depend करे model तो
output highly Varied करता है →

Variance high.

final
representative
line नहीं

Degree - 12 (9)

बहुवर्षीय कम्पनी

मार्केट अलग करता है।

Variance high.

linear

Variance low.

II Sampling तो diff तो यह model output
highly Varied तो → Highly Variance Model.

- Model heavily biased \rightarrow
 100 test points \rightarrow 60% performance
 200 " " \rightarrow 62% "
 200 " " \rightarrow 63% "
 Variance \rightarrow data यांत्रिक लालनाडी
 Model एवं अल्पवर्णना / Model
 एवं विशेषज्ञ कोष्ठक प्राप्ति
 एवं data फैले जाते |
- We have to trade off betⁿ bias & variance
- Decision tree \rightarrow Russell Norvig एवं book एवं
 example
- Ordering of features to build a decision tree
 \rightarrow correlation analysis.
- Information gain
- Balanced dataset - 60% yes, 60% no.
 \rightarrow Type कृपया यहां subsample करनावा उद्धीश्य
 distribution ratio same. $\boxed{50-50}$
- \rightarrow pattern fit कृपया check करना। $\boxed{2:4} \quad \boxed{4:12}$
 एवं लालनाडी। $\boxed{2 \ 5 \ 9 \ 10}$

Where there is uncertainty, there is message.

$-\log p$ info probability
or $\log \frac{1}{p}$ info value

Attribute property \rightarrow Entropy

for a set Random Variable info content (weighted)

$$H(V) = \sum_{k=1}^K P(v_k) \log \left(\frac{1}{P(v_k)} \right)$$

$$\checkmark \text{ for } K \text{ different values} = \sum_{k=1}^K -P(v_k) \log P(v_k)$$

Entropy ক্ষেত্রে কমান্তর পাই।

Lecture -04

8/9/24

Uncertainty \uparrow Statement Value \uparrow
Probability \downarrow

Information is inversely proportional to
probability

Entropy \downarrow Uncertainty \downarrow

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

true const probability

Remainder (A)

→ A Variable select करते हैं -

मध्य सबसे ज्ञान प्राप्ति अपनी

P और नया probability

$$\begin{cases} 0 \times \log 0 = 0 \\ 1 \times \log 1 = 0 \end{cases} \quad \left. \begin{array}{l} \text{इस तरह एक guaranteed} \\ \text{वर्ते } \end{array} \right\}$$

Tree - 8
 Gain(A) = $B \left(\frac{P}{P+n} \right)$ - Remainder(A)
 2
 1
 Gain prob of A
 entropy
 Gain prob of B
 entropy

As A is gain(A) ↑ P or C or I
 ↓ Information ↓ Entropy

0	0	0	0	0	0
0	0	0	0	0	0

$$0.5 \log \frac{1}{2} + 0.5 \log \frac{1}{2}$$

$$= -\frac{1}{2}$$

Feature Ranking by Information Gain.

Recursive elimination of features

Data Collection

Pre-processing

Feature Engineering

Feature Selection

Domain knowledge
needed to know
what are important.

Performance Metrics

self-study

AUROC

AUPR

Confusion Matrix: gives stat

		Actual class	predicted class	stat	
Row	Col			TP	FN
Pos	Neg	Pos	Pos	TP	FN
Neg	Pos	Neg	Neg	FP	TN

$$P = \text{pos}, N = \text{neg}$$

dataset में कितने बिल्कुल सही

$P = 20$, dataset में 20 बिल्कुल pos classed points हैं।

Neg class 200 नीत इन्हें।

Prediction positive
वास्तव के बहुमत
तथा वास्तव के फल

TN

) को predict करा गया

या

predict करा गया

अदृष्ट समिन्द्रियों

परिप्रेक्ष्य

प्राप्ति - लोग

गणना

गणना

FN



negative predict

false

फ़ल

prediction

प्राप्ति = N - 209 = 9

अब अन्तरा बताएं।

अब अन्तरा बताएं।

(T+class)

multiclass G, diagonal परिवर्तन नहीं होता।

Accuracy: युवा dataset G के भूलीं corrrect prediction कराएं लाभार्थी + 9T

$$Acc = \frac{TP + TN}{P + N}$$

$$\boxed{\begin{array}{l} P = 99 \\ N = 1 \end{array}}$$

Model (xP)

Data balanced का रेल Accuracy misleading.

out of 25
with Recall

$S_n \rightarrow$ positive class go Accuracy.

$$S_n = \frac{TP}{TP + FN}$$

(sensitivity)

S_p (specificity) = neg class go Accuracy

$$= \frac{TN}{(TN + FN)}$$

Data balanced or unbalanced

Accuracy important

$S_n = 20\%$, $S_p = 65\%$. ✓ good model

$S_n = 20\%$, $S_p = 40\%$. ✗ bad model.

Precision: - অস্তি হা অফল তা precise হয়।

① অস্তি positive টা pos ofml।

$$\text{Prec} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

Limitation: neg নিষে ফোল কথা নাই।

*** MCC

Unbalanced dataset এর জন্য অনেক

অনুমতি measure হ।

If dataset has no P
point, precision 0 হয়ে মান।

but S_p উচ্চ ধৰণ আলো রাখ।

(if Model is good ofc)

F1 score:

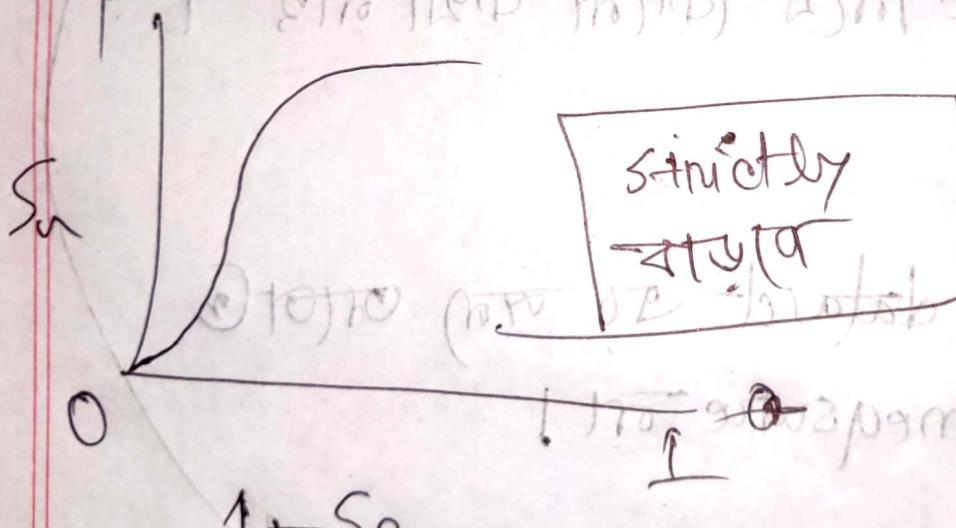
harmonic series formula

$$\frac{1}{\sum_{i=1}^K \frac{1}{x_i}}$$

into the form = maximum

AUROC: / AUC

Area Under Receiver Operating Curve.



1 - Sp

- data balanced না, বেটাম্য AUROC

- একে হস্তান্তি prediction করলা

(কোকো বিল্ডার্স)

Linear Regression

why linear? $w_0 + w_1 x$

$$y = c + mx$$

$$h(x) = w_0 + w_1 x$$

↓ ↓
we want to
get w_0, w_1 .

known

$$\frac{\sum_{K} |y - h(x)|}{K}$$

Mean absolute
error অনুপর্যুক্ত
বিরুদ্ধ

$$\frac{\sum_{K} (y - h(x))^2}{K}$$