

Feature Engineering:

ଏହି ଏବେ ରେଜିମିନ ଲାଗିପାଇଛି; ଏହି ଏକାକି ଧ୍ୟାନ, ରେଜିମିନ ଡିଜିଟାଲ ରେଖାକ ପରିମାଣ କରିବାକୁ ଦେଇବାକୁ କଥା ହେଉଥିଲା ଏବେ ରେଖା କାହାର କୁଳମାତ୍ର କାହାର ରେଖାକୁ,

ରେଫରେନ୍ସ: ଯୁଦ୍ଧ ଅନୁଷ୍ଠାନ ଏହି ଡିଜିଟାଲ ପ୍ରଦାନ କରିବାକୁ ରେଜିମିନ କିମ୍ବା କୁଳମାତ୍ରା, ବିଶ୍ଵାସ କରିବାକୁ ରେଜିମିନ କିମ୍ବା କୁଳମାତ୍ରା କରିବାକୁ ଏବାକୁ ଏହାର ବିଶ୍ଵାସ କରିବାକୁ ଏହାର ବିଶ୍ଵାସ କରିବାକୁ ଏହାର ବିଶ୍ଵାସ କରିବାକୁ ଏହାର ବିଶ୍ଵାସ କରିବାକୁ

Feature Engineering ଏବା ଫେଚ୍ ଏବା ଫେଚ୍ ଏବା ଫେଚ୍ ଏବା ଫେଚ୍
Feature Transformation ଏବା ଫେଚ୍ ଏବା ଫେଚ୍

→ ରେଜିମିନରେ ଆପ୍ରାପ୍ତ କ୍ଲାସ କ୍ଲାସ

ମୁହଁର୍ବଦ୍ଧ,

→ ଏହି ରେଡିଆ ରେଖା କିମ୍ବା କାହାର ଗାଲା,

→ କିମ୍ବା କିମ୍ବା ବାବୁ କାହାର predict

କାହାର ଗାଲା,

Distance Based Algorithm

Feature Transformation includes:

- Scaling and normalization
- Logarithmic Transformation
- Polynomial Features
- Dimensionality Reduction
- Encoding categorical values
- Handling missing values
- Tokenizing text data
- Stemming and lemmatization

Why is Feature Transformation important?

- Model performance
- Computational Efficiency
- Improved Interpretability
- Handling Non-linearity

Regression Algorithm

Continuous value predict করা

যাবে

(1) Linear Reg

(2) Rig "

(3) Lam " "

(4) Elastic Net "
(Rig + Lam)

এইচ্যু 3 Tree based regression
and Non-linear Algorithm

Classification algorithm list:

(Yes/No, Cat/Dog), (Par/Fail)

① Logistic Regression

② k - Nearest Neighbors

③ Naive Bayes

Tree based classifiers,

Kernel based Algorithm.

Euclidean distance:

→ नियमित दूरी का ज्ञान

→ दो बिंदुओं के बीच दूरी,

(L₂ distance) $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

→ Use करना ML में data

प्राप्ति इसका एक उपयोग है।

उपयोग, दो बिंदुओं के बीच दूरी का ज्ञान

→ एक दूरी का ज्ञान

$$F = P + Q$$

→ दो बिंदुओं के बीच दूरी का ज्ञान

लोगिक और गणितीय रूप से

दो बिंदुओं के बीच दूरी का ज्ञान

बहुत सारे विकल्प थे

Manhattan distance : (L1 Norm)

5 ft. (15 m) दूरी तक

6 मीटर लंबाई तक
Taxicab
distance 3 मीटर.

⇒ for diagonal 6 मीटर,

points: $x_1(2, 3)$

$y_1(5, 7)$

$$\therefore \text{M.D.} = |5-2| + |7-3| \\ = 3+4 = 7$$

Robotics : वाहन यात्रा
दूरी

→ Image processing : 5 ft pixels.

→ Path finding games : यात्रा

Tetris यात्रा गेम.

Minkowski distance

একটি মাপিয়ে ফর্মুলা দেখ
E. D. or M. D. = $\sqrt[p]{\sum_{i=1}^n (x_i - y_i)^p}$ দিয়ে

বড়,

Suppose if $p = 1 \rightarrow MD$ [MD formula 75 এমি
MD formula 90 এমি] $\rightarrow ED$ [MD formula 90 এমি]

if $p = 2 \rightarrow ED$ [MD formula 90 এমি]

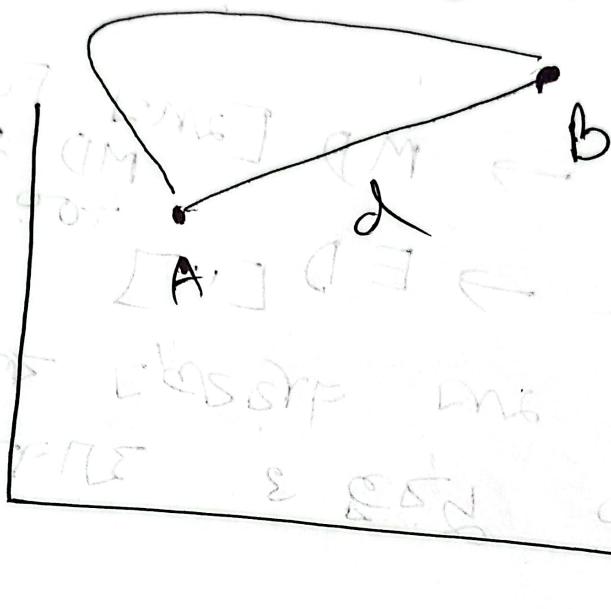
$p \geq 2$ এমি গুরুত্ব করে। $\sqrt{2}$
ধৰণ: $\sqrt[3]{3}$ হয়।

Formula : Distance = $\left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$

এমন, x_i, y_i $\sqrt{2}$ হয়।

প্রদর্শন = $\sqrt{2^p + 2^p} = \sqrt{2^{p+1}}$,
 $(1, 2, 3, \dots, n)$

- Pattern Recognition
- Clustering Algorithm
- Distance Tuning



Note: $p_{\geq 0}$ is not valid in Minkowski distance metric.

$$\text{Ex: } \lim_{p \rightarrow 0} (x_i - y_i)^p = 0$$

$$\text{Ex: } \lim_{p \rightarrow 0} \frac{1}{p} = \text{undefined}$$

Hamming distance

Length of string
 position 1
 position 2
 Hamming distance = 2.

Karim

Str 1 : "Karim"
 Str 2 : "Karm" (not same)

Str 2 : "Wards" (not same)

Str 2 : "wards" (not same)

∴ Hamming distance = 1

- Error detection & correction
- DNA / protein sequence
- Cryptography → RSA (3rd year)
- ATM (4th)

Feature Transformation

Techniques

- ① Normalization
- ② Standardization
- ③ Log Transformation
- ④ Robust Scaler
- ⑤ (Max) absolute scalar

1. Normalization (Min-max scaling)

Data \rightarrow 0 (223, 1, 28, 315)

Ques:

Formula:

$$x - x_{\min}$$

$$x_{\text{scaled}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Ex:

Suppose height 120 cm
= 180 cm

150 cm

$$\frac{150 - 120}{180 - 120} = \frac{30}{60} = 0.5$$

height = 0.5 \times 200 cm,
120 cm is 200 cm

223, 315

, floor board 1756 -

2. Standardization (Z-score calling)

Data \rightarrow जबकि scale \rightarrow 10.
में mean = 0, SD = 1

Formula: $X' = \frac{X - \mu}{\sigma}$

$$X \text{ scaled} = \frac{X - \mu}{\sigma} \quad \text{below } X$$

Ex: Given $\mu = 3.5, \sigma = 0.5$

उत्तराः $GPA = 3.0$, $SD = 0.5$

$$\frac{3.0 - 3.5}{0.5} = \frac{-0.5}{0.5} = -1.0$$

जहाँ -1.0 "standard deviation above average"
Good result.

3. Log Transformation:

Data ~~is~~ ~~not~~ ~~normally~~ ~~distributed~~ ~~(skewed)~~ onto the log scale.

(Skewed) ~~onto~~ ~~the~~ ~~log~~ ~~scale~~ ~~(normal)~~

and profit one variable

Formula: $y = \log(x)$

x transformed

Ex: annual orders 55 $\rightarrow \log(100)$

salary 1,000 $\rightarrow \log(1000)$
= 3

10,000 $\rightarrow \log(10000)$

20,000 $\rightarrow \log(20000)$

30,000 $\rightarrow \log(30000)$

Data \rightarrow difference
and 2150 61

2000 \rightarrow 2000
2000 \rightarrow 2000

4. Robust Scalar:

outlier ($\text{अनोर्मल वैल्युट्स}/(h_i^2)$) में से कम से कम 15%, median तथा interquartile range का अन्तर्वाला

Ex:

\bar{x} और s नियन्त्रित करने के लिए, income का अन्तर्वाला बढ़ाव दिया जाता है। इसका उद्देश्य यह है कि अन्तर्वाला की विफलता के कारण अंतर्वाला का मैडल अपेक्षित मैडल से भिन्न हो जाए। इसका उद्देश्य यह है कि अंतर्वाला की विफलता के कारण अंतर्वाला का मैडल अपेक्षित मैडल से भिन्न हो जाए।

उदाहरण के लिए, अंतर्वाला का मैडल अपेक्षित मैडल से भिन्न हो जाए।

उदाहरण के लिए, अंतर्वाला का मैडल अपेक्षित मैडल से भिन्न हो जाए।

5. Max Absolute Scalar:

for Data \rightarrow तारिखों के 2NN
for. तरीका \rightarrow $[-1, 1]$ scale

2. ~~non-monotonic~~ \rightarrow ~~discrete~~

Formula:

$$x_{\text{scaled}} = \frac{x}{|x_{\text{max}}|}$$

Ex:

Feature : $[-50, 0, 50, 100]$

Max = 100

Scal. : $250 \cdot \frac{1}{100} = 2.5$ 3NN
 $[-0.5, 0, 0.5, 1]$