

Featureisation And Feature Engineering

most important

convert data to numerical vector

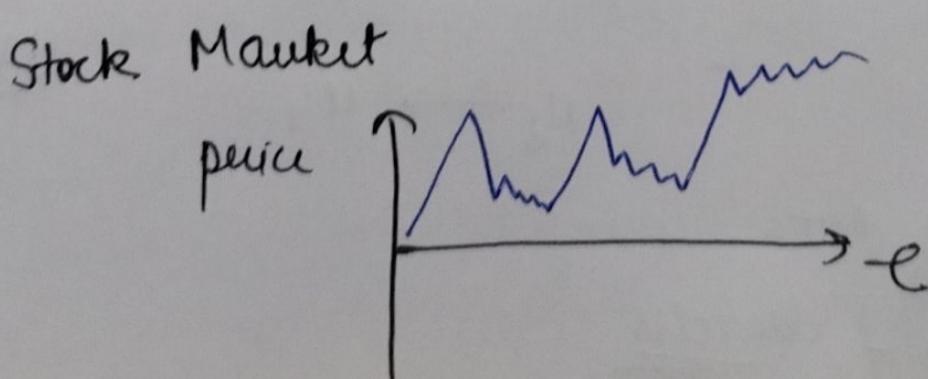
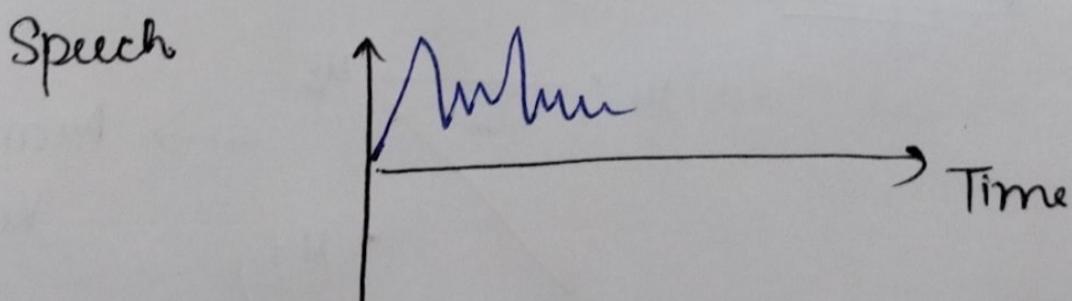
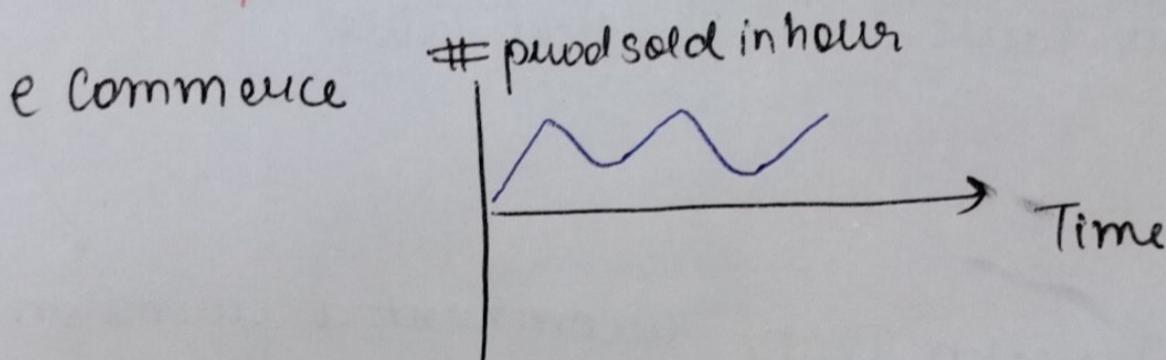
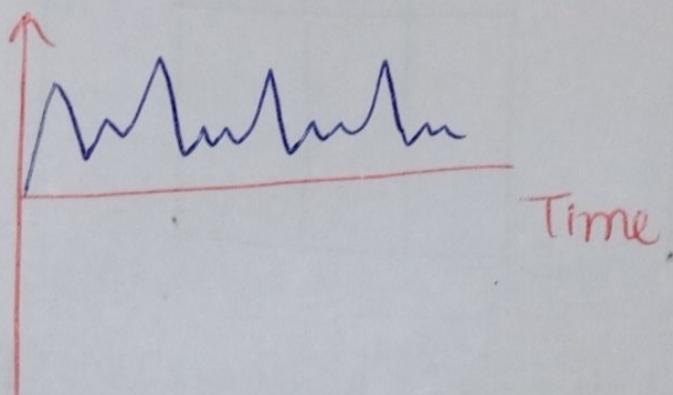
Text data: BOW, tfidf, avg w2v, tfidf w2v

↓
various featureization

categorical data: one hot encoding

mean response rate, domain specific

Time Series heart rate



Test: Sequence data

$w_1, w_2, w_3, w_4 \dots$

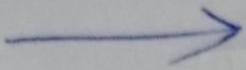


Image data Face detection, face recognition

X-rays, MRI scans, video ?

Image + Time
Series

Database Tables

C.Id	T ₁	locati

Numerical
features

Ans Table

C.Id	p.u.Id	Time

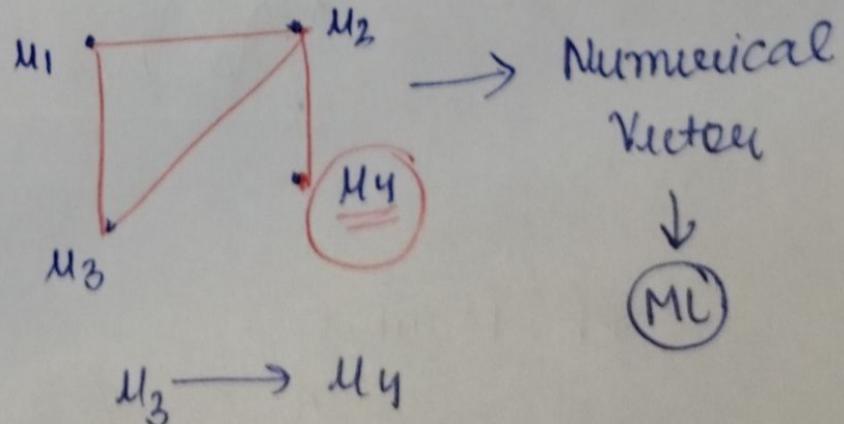
Purchase Table

	Recd	disc

Product Table

Graph Data recommend a friend on FB

[
→ vertex
→ edge]



$u_3 \rightarrow u_4$

↓
ML

Tons of types of data
→ feature extraction:- descriptors

→ cover commonly used featureization
(general purpose)

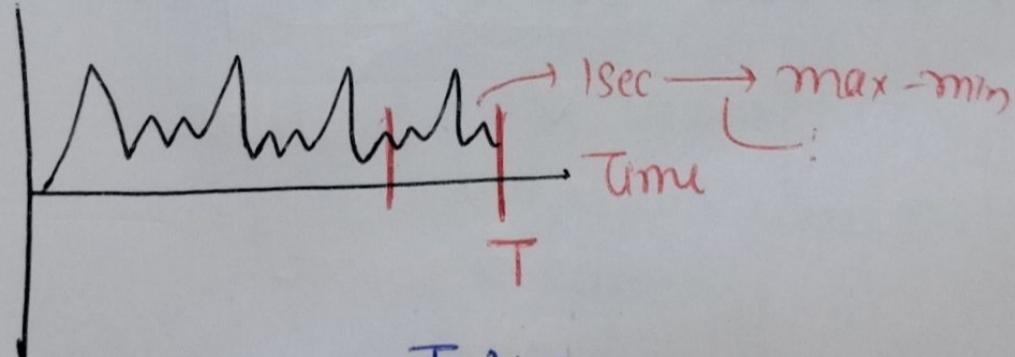
Moving Window

Simpliest featureization of Time Series

- ECG:
- ① Take a window of Time
 - ② Then compute mean, std dev - - -
median, quantile
 - Max, Min
 - Max - Min
 - Max / Min
 - local maxima & minima

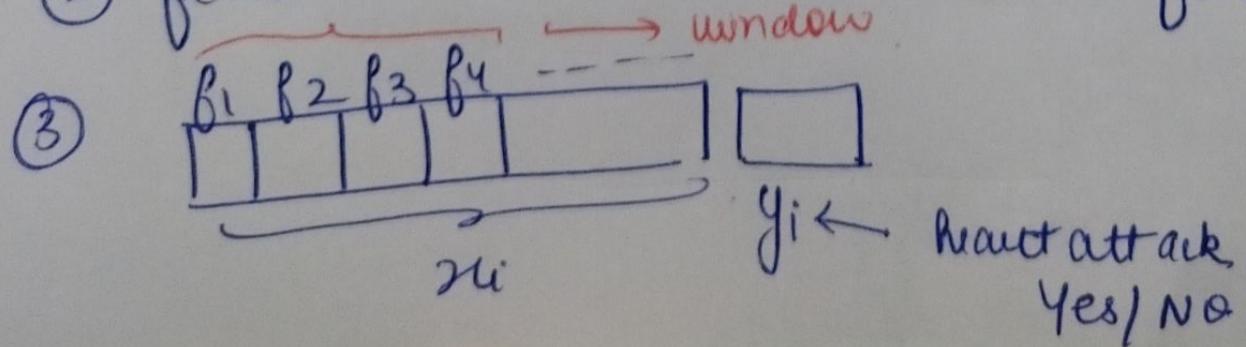
Zero Crossing (Times it crosses mean line)

ECG:



Task: predict cardiac heart attack in next 10 min

- ① Window-width
- ② features in window which are useful

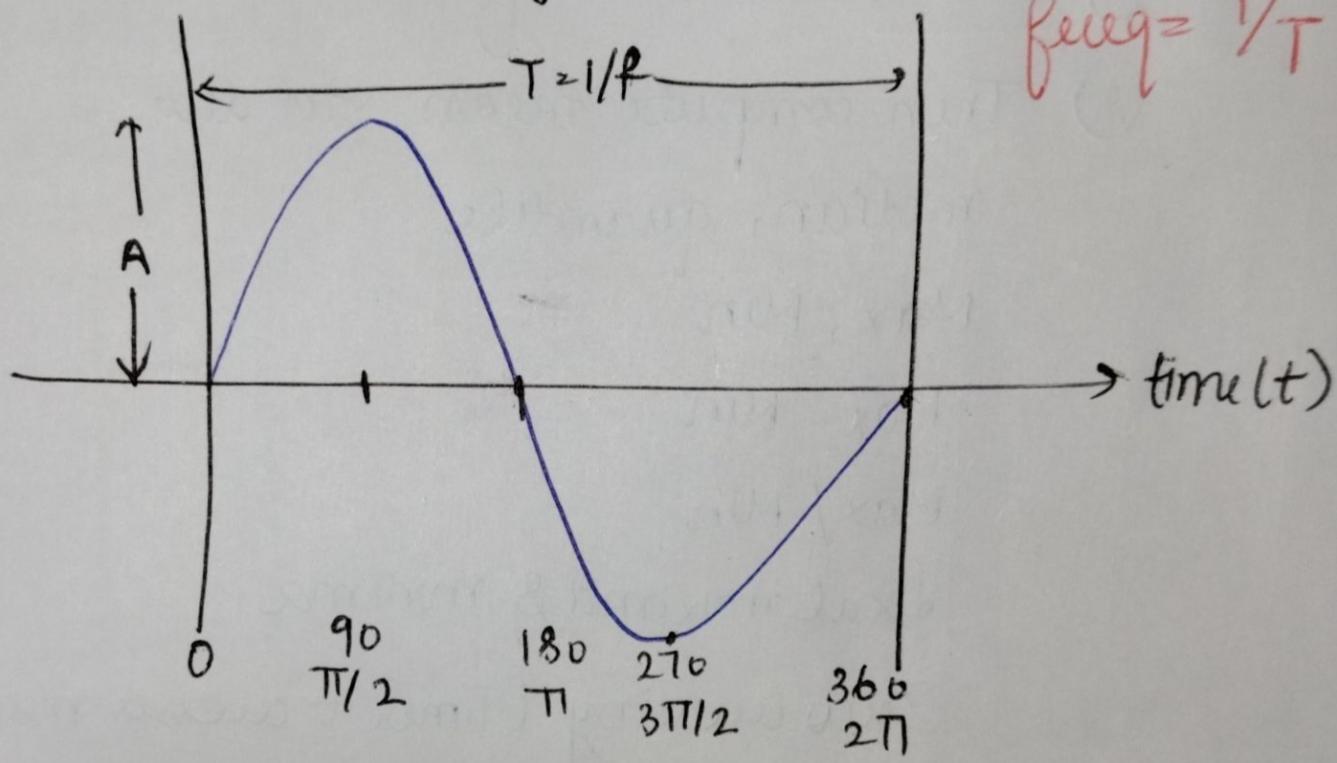


Fourier Decomposition/ Transform

method to represent Time Series

physics, applied math, communication, cs,
Signal processing

Oscillating Sine Wave



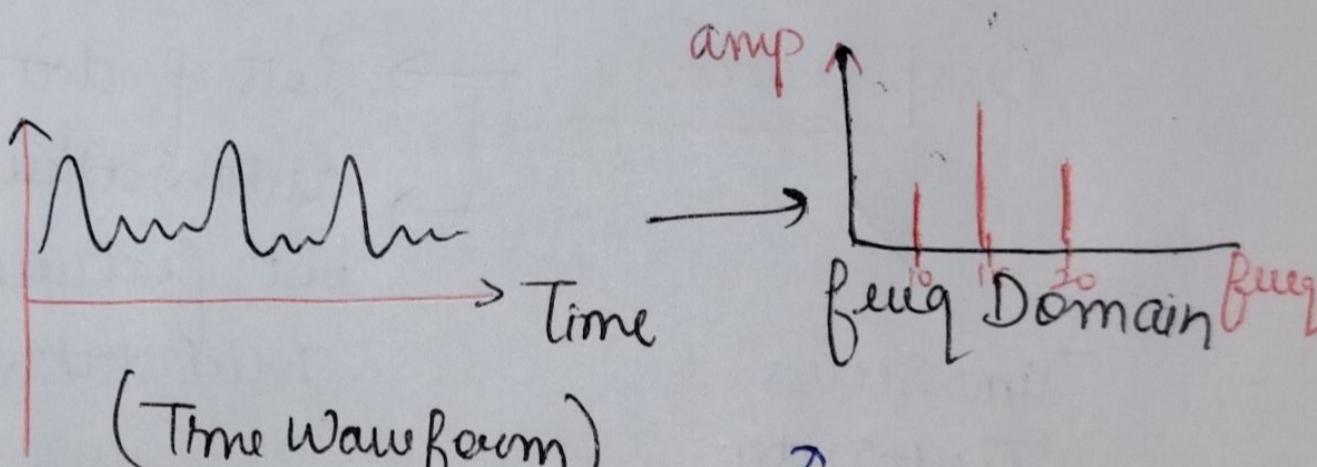
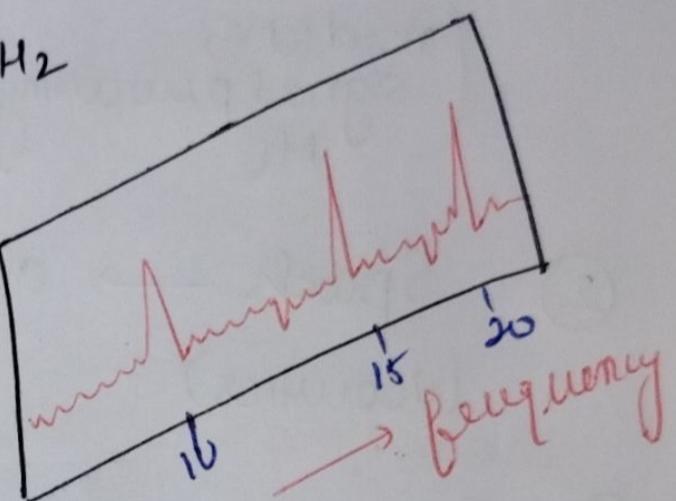
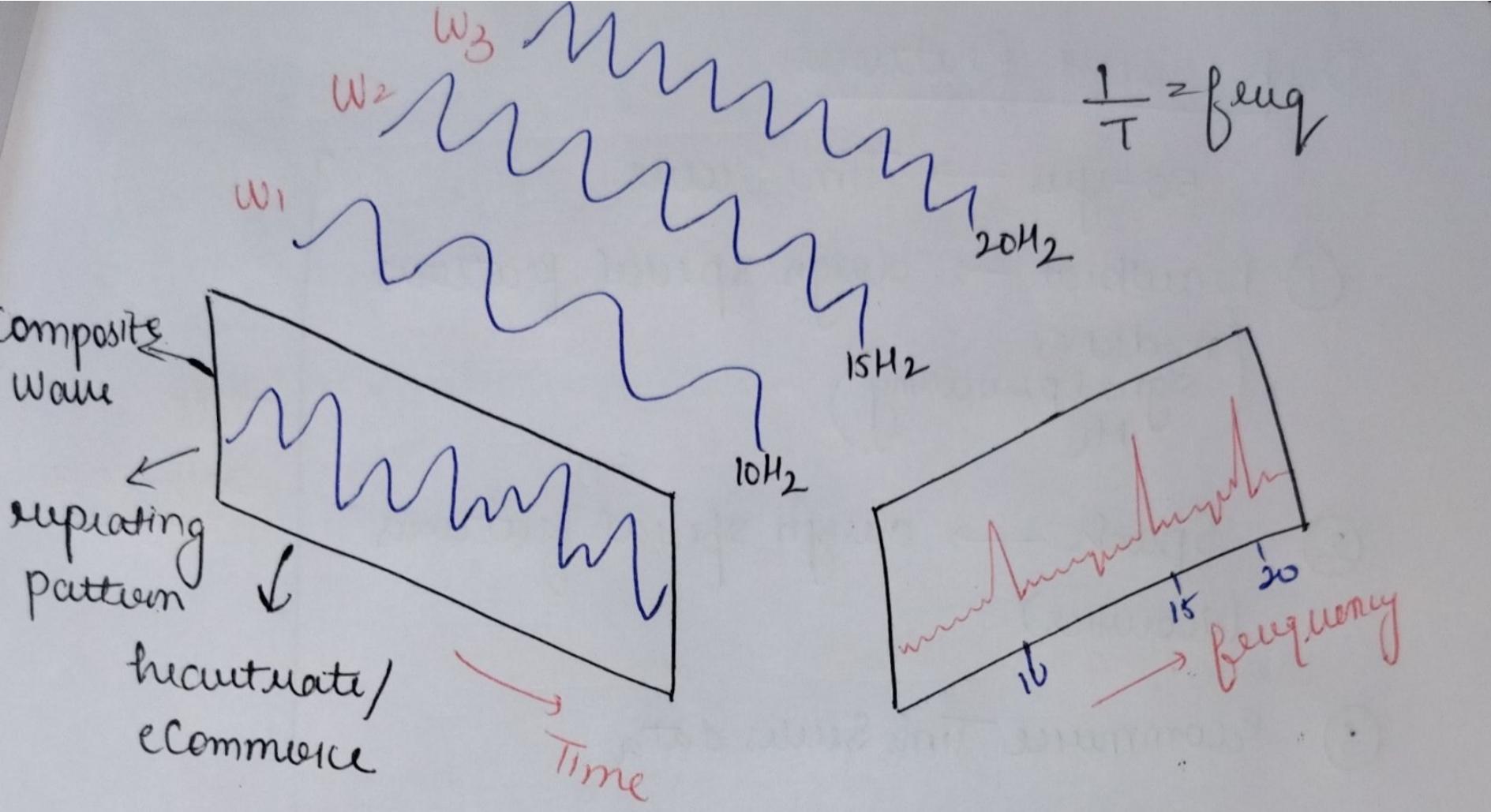
Human note: 60 - 100 beats per minute

↓ ↓
oscillates

1 to 1.6
beats per sec

↓

1 Hz to 1.6 Hz



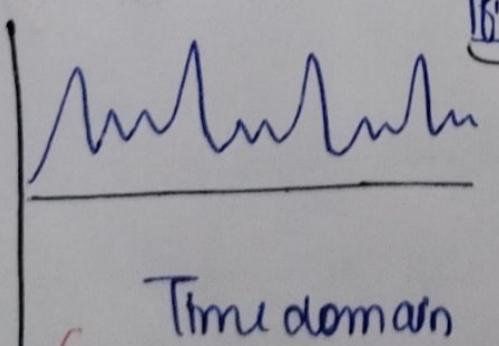
(Time Waveform)

Time domain

Fourier transform

useful whenever pattern
is repeating

$[P_1 | P_2 | f_3 | a_1 | a_2 | a_3]$ \rightarrow freq matrix



feature vector

FT $\xrightarrow{\text{24}}$

a_2
 a_3
 a_1

$\{f_1, a_1\}$
 $\{f_2, a_2\}$
 $\{f_3, a_3\}$

f_1 $10H_2$ f_2 $15H_2$ f_3
 $20H_2$

Deep learnt features

50+ yrs → Time Series

- ① Heartbeat → design special features
(medicine,
signal processing)
ML

Till
2012-2013

- ② Speech → design special features
(acoustics)

- ③ E-commerce Time Series data

Deep Learning

→ Lots of data
automatically learn the
best featureization for
your data

Time Series

Text Data

Image Data

5-6 years (2012-2017)

- { Google improved: Speech recognition
Siuu:
Song Recommendation } → deep learnt features

Image Histograms

Images: face, object, Scan, X rays, autonomous cars

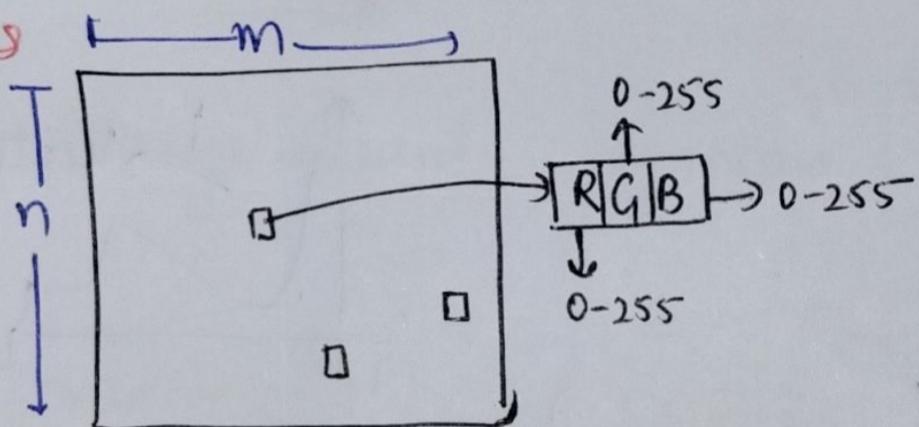


30+ years of research

2012: Deep learning

① Colour Histograms

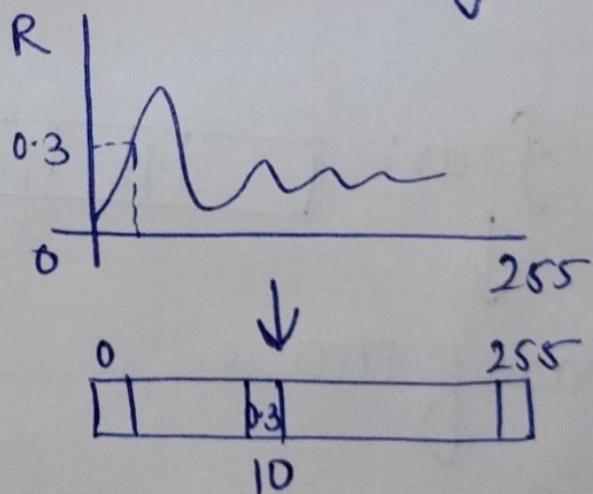
C



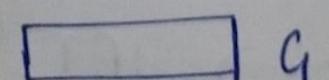
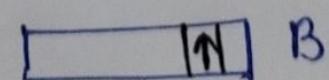
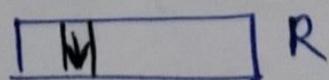
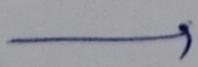
② Edge Histograms

① Red Values for each pixel

(n x m) data points \rightarrow histogram



② Green And Blue



Task: Sky or not

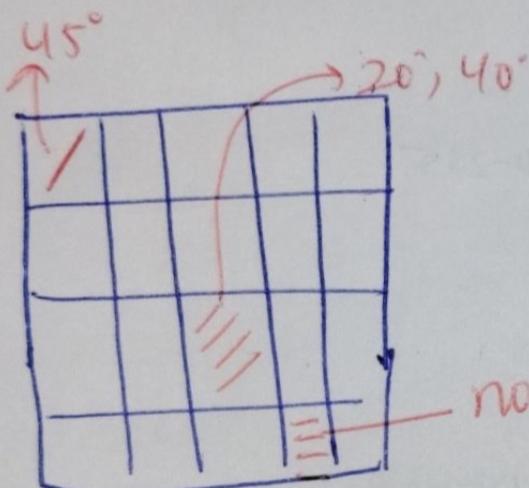
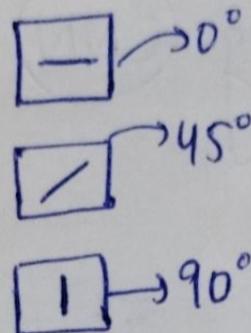
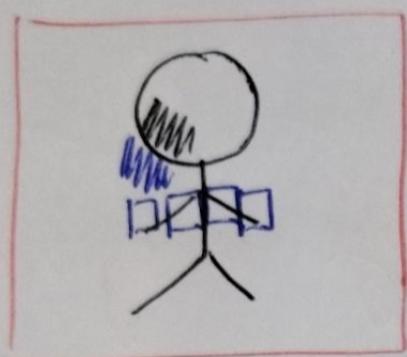
face detection: Skin colour \times



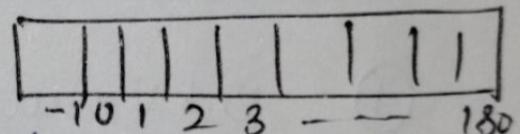
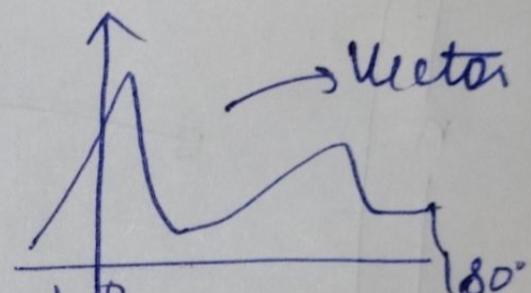
Colour hist $\rightarrow \times$ recognize shapes

Edge Histogram

(Image processing)



no edge $(-)$



each region: edge value / edge angle

↓
histogram

face:

Ram features

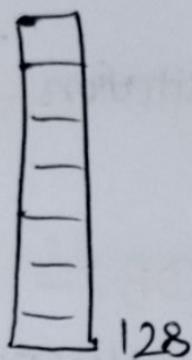
{ Colour hist
edge hist \rightarrow v.v basic & fundamental

→ SIFT (obj recog.)
→ CNN \rightarrow state of art

Scale Invariant Feature Transforms (SIFT)

- Image featureization
 - object detection (2000's)
 - creates vector for each point
128 dim (128 dim vector)

OpenCV



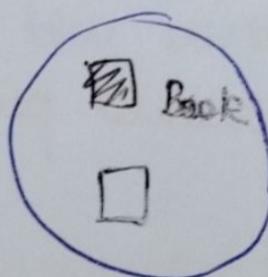
$$I = \{v_1, v_2, \dots, v_{10}\}$$

↑
128 dim

- Image Search (amazon.com) → mobile app

Take picture of anything
and search

Scale Invariance
Rotational
Invariance



← Search

Deep learning features

```

graph TD
    TS[Time Series] --> LSTM[LSTM]
    Images[Images] --> CNN[CNN]
    LSTM --> Featureize[featureize]
    CNN --> Featureize
    CNN --> AA[almost automatically]
    CNN --> BestFeatureization[best featureization (30+ yrs)]
    style Featureize fill:none,stroke:none
    style AA fill:none,stroke:none
    style BestFeatureization fill:none,stroke:none
    2012[2012] --> CNN
  
```

X-rays → lost of data/Images
↳ tumor see not → (CNI)

Relational Data And Featureization

eg

Cust Id	Cus Zipcode

Cust Table

Cust Id	PurId	Time
✓		
2	5	...

Cust viewing / visitation data

Cust Id	P.Id	Time

purchase data

PurId	Type	---

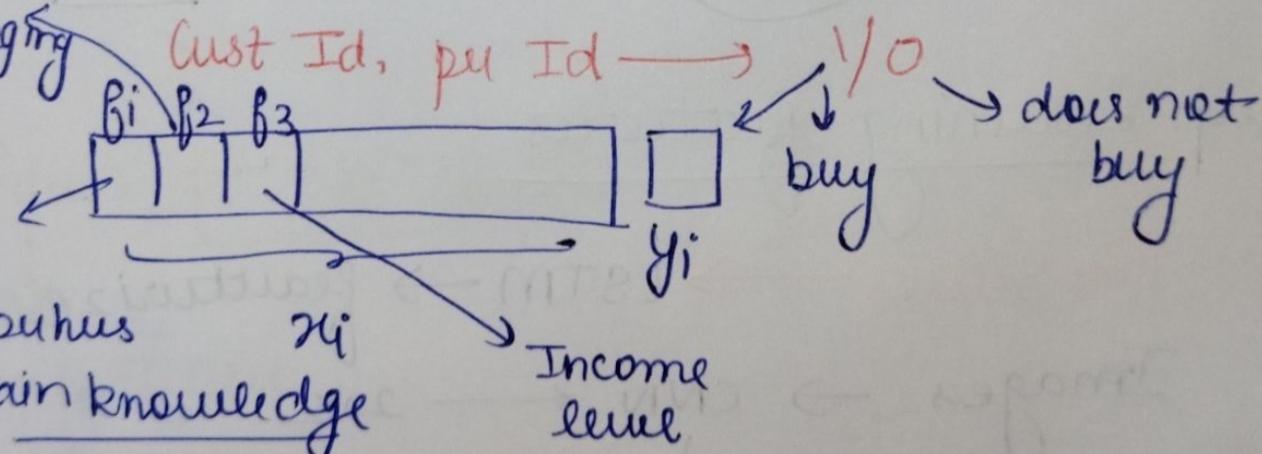
PurId. data

Relational DB:- Oracle,
MySQL, SQL Server

Task: predict if a customer could purchase product
in next seven days.

Cust Id visit

any pur. belonging
to same type



the cust
id viewed
per 10 pages in 2 hours

domain knowledge

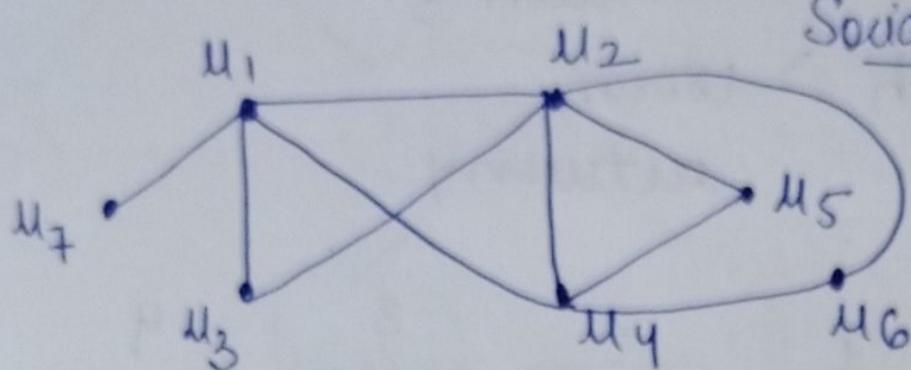
Relational data \rightarrow SQL

\downarrow
featureization

Graph data & featureization

e.g.

Data



Social graph

Facebook

U_i :- written (user)

edge (u_i, u_j) : friends

Task: Recommend new friends for a user (u₄)

ML

$$u_i, u_j \rightarrow 1/0$$

$$\checkmark \quad m_7 - m_4$$

} Recommend

$$M_7 - M_4 = u_1$$

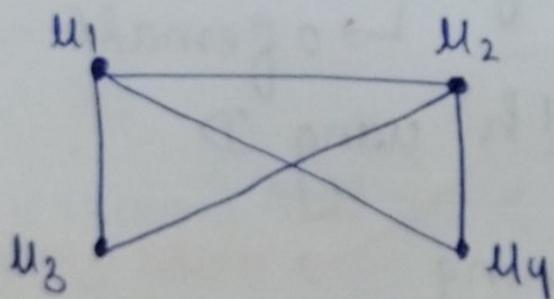
$$u_3 - u_4 := u_1, u_2$$

fi: # mutual funds

(f2) # paths b/w vertices u_3 & u_4

graph based features

(e.g.)



Feature Binning (buketing)

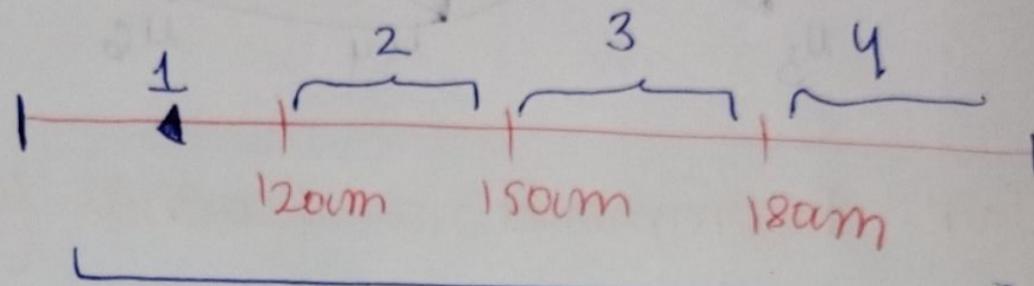
→ extension to Indicator variables

→ eg: height

} if height < 120cm
return 1

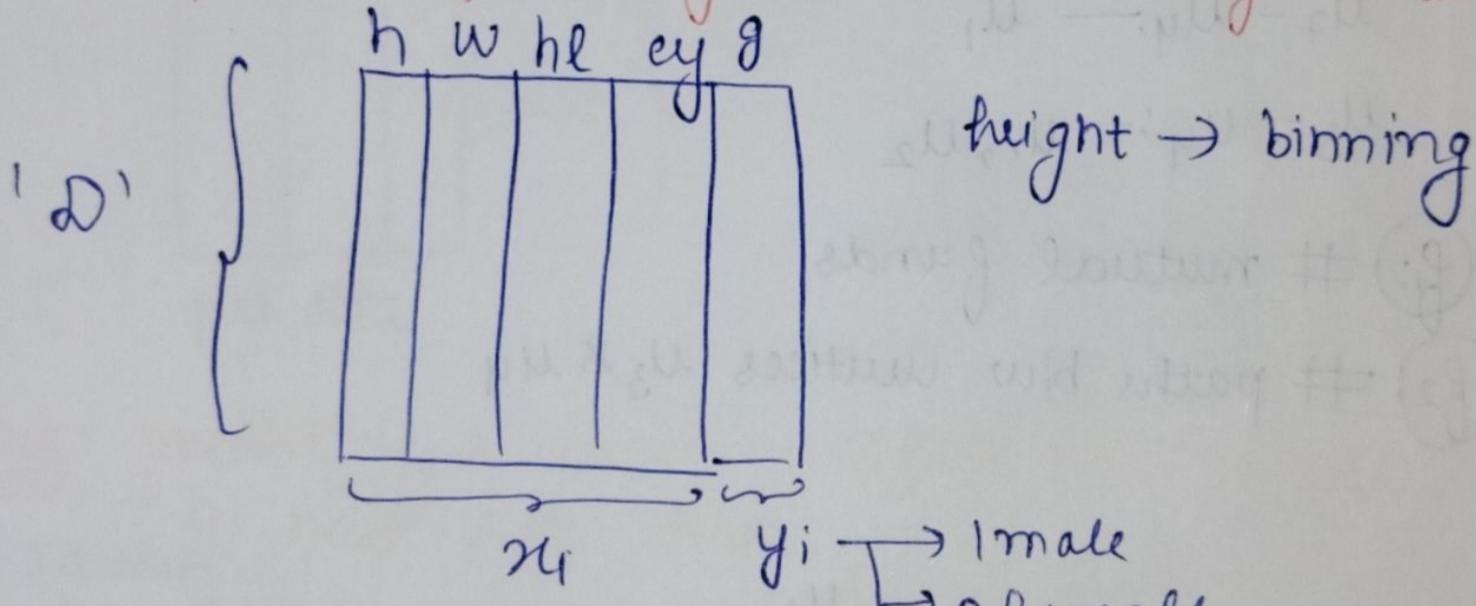
if height < 150cm And h ≥ 120cm
But

$\left\{ \begin{array}{l} \text{if } h < 180\text{cm and } h \geq 150\text{cm} \\ \text{entren 3} \\ \text{if } h > 180\text{cm} \\ \text{entren 4} \end{array} \right.$



4 bins (problem specific)

Task: predict gender given: $h, w, h, \text{eye colour}$

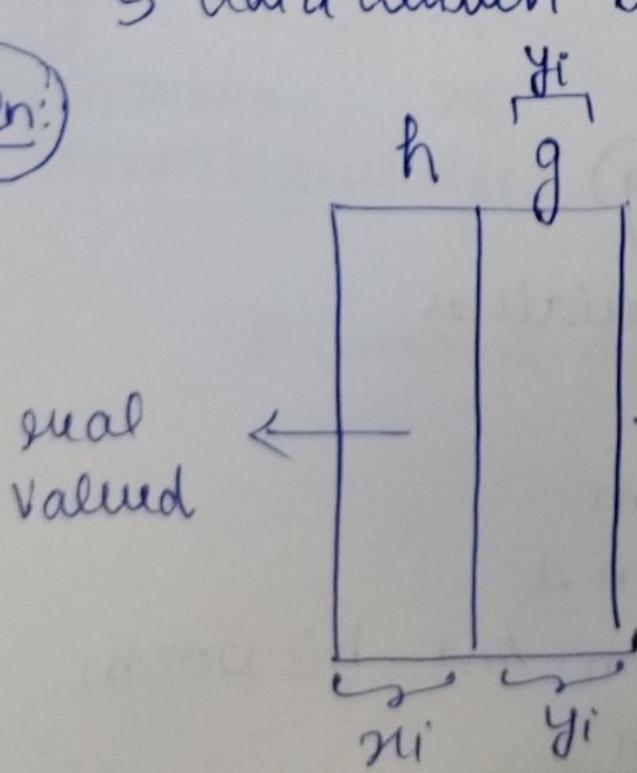


Challenge: binning of h using D

\hookrightarrow data driven binning \rightarrow make sense to use

y_i to perform binning

Soln:



Simple DT

$h < 130\text{cm}$ \rightarrow bind h using I_g

y
 $y_i = 0$

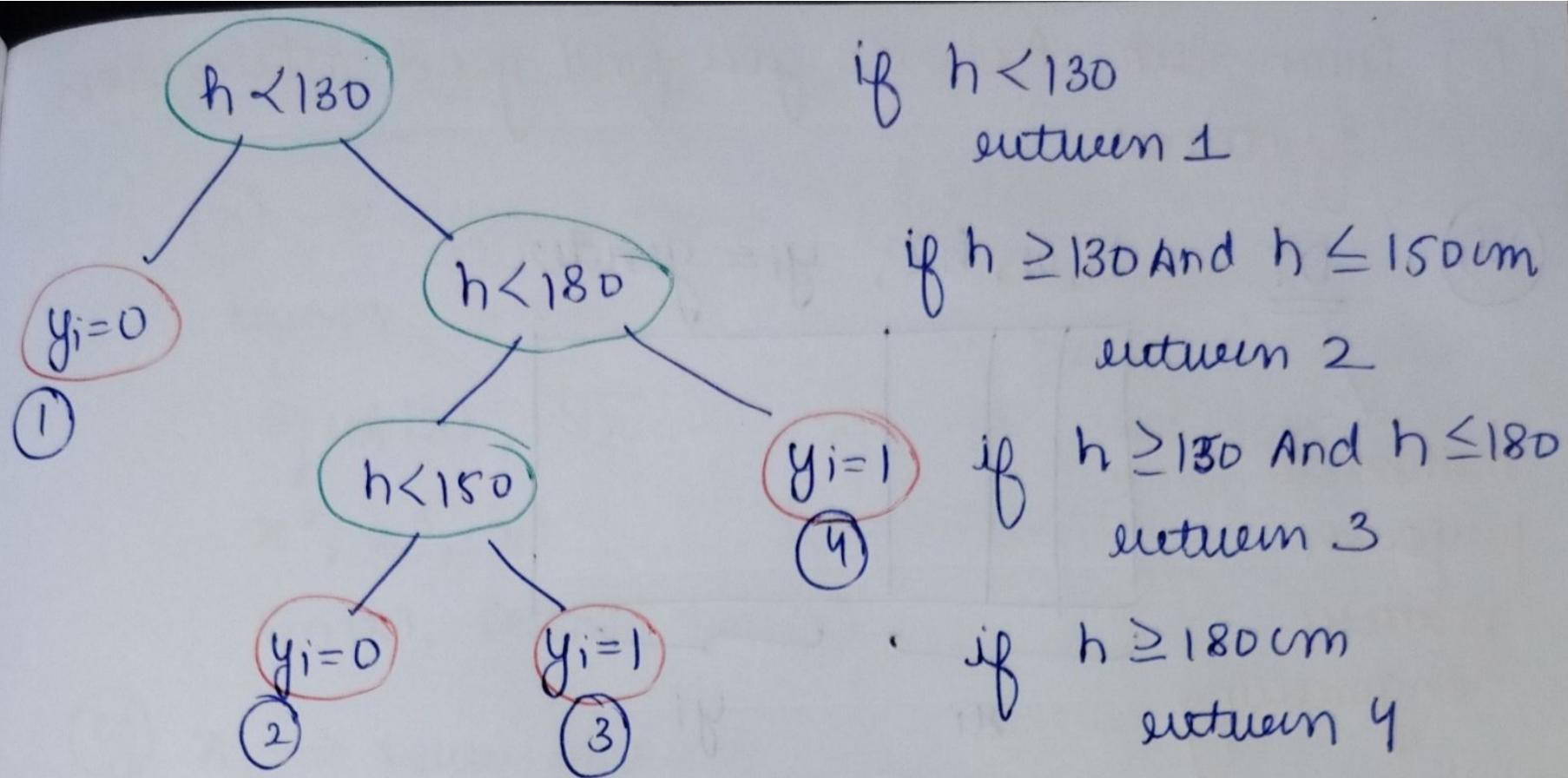
N

$h < 180$

y
 $y_i = 1$

y
 $y_i = 0$

$h < 150$



binned equal valued features using y_i 's & feature itself using DT

Interaction Variables

Task: predict gender

$\underbrace{h, w, hl, ec}_{x_i} \rightarrow \underbrace{\text{gender}}_{y_i}$

(eg) ① $(h < 150\text{cm}) \text{ AND } (w < 60\text{kg}) \rightarrow$ 2 way interaction feature

(eg) ② a) $h * w$ } math 2 way interaction feature
 b) $hl * w$

(eg) ③ $h < 150\text{cm} \text{ AND } w < 65\text{kg} \text{ AND } hl > 5\text{cm}$ } 3 way interaction feature

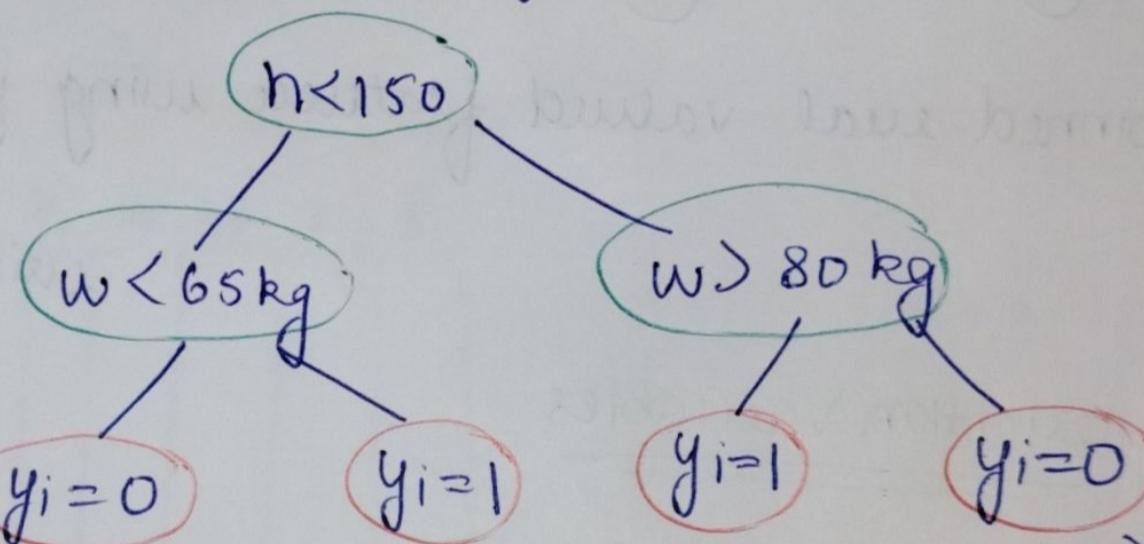
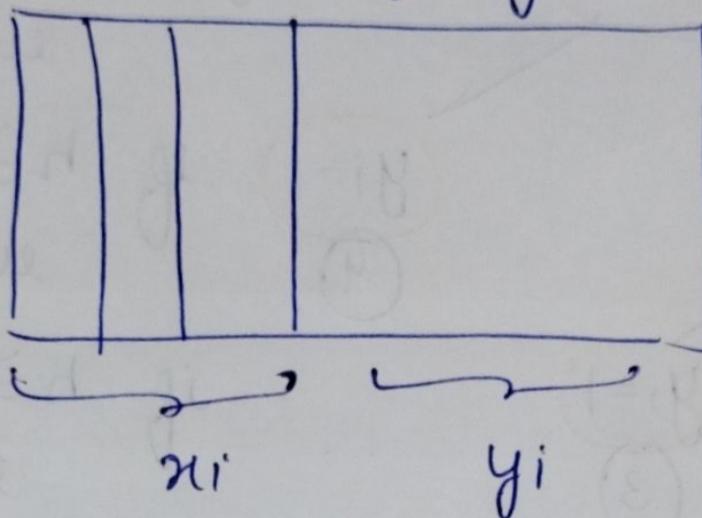
(Q) Given task, how do you find good interaction features

Soln

DT

$h, w, he, y_i = \text{gender}$

method to
perform
feature
engineering



$\underbrace{h < 150 \text{ AND } w < 65 \text{ kg}}$
new binary feature
(Interaction variable)

$\underbrace{h \geq 150 \text{ AND } w \leq 80 \text{ kg}}$
Corresponding to every leaf node create a
new feature.

Mathematical Transformations

- ② \rightarrow single feature (Q) what is the best transform
 $\log(n)$, e^n \rightarrow problem specific
 \sqrt{n} ; $\sqrt[3]{n}$ \rightarrow case study
 $n^2, n^3, n^4 \dots$ (polynomial)
 $\sin(n), \cos(n), \tan(n)$

(eg) $n \rightarrow$ power law dist

\downarrow
 $\boxed{\log(n)}$ (Gaussian dist)
 box contains transform

Model Specific Featureizations

(eg) $f_i \rightarrow$ power law dist $(\log(f_i))$

Logistic Reg. \rightarrow gaussian Naive Bayes



features are gaussian dist

(eg) $f_1, f_2, f_3, y \in \mathbb{R}$

$y \leftarrow f_1 - f_2 + 2f_3$ \leftarrow domain knowledge
 linear combination of f_i 's

DT may not work very well

linear models \rightarrow linear Regression

$$\begin{cases} w_1 = +1 \\ w_2 = -1 \end{cases} \quad w_3 = +2$$

eg ③ $\rightarrow y \rightarrow$ interaction of $f_1 \& f_2 \leftarrow$ Domain knowledge
 (gender) (n,w)

\downarrow DT / RF / GBDT tree based
 linear model (Log. Reg, SVM)

text

BOW \rightarrow linear models

\downarrow
 v. high dim-space \rightarrow hyperplanes $\xrightarrow{\text{SVM, Log. Reg}}$
 (+ve) (-ve)

Feature Orthogonality

* The more different / orthogonal the features are the better the model would be

$f_1, f_2, f_3 \xrightarrow{M} y$

~~$f_1 \xrightarrow{\text{corr}} y, f_2 \xrightarrow{\text{corr}} y, f_3 \xrightarrow{\text{corr}} y$~~

$f_1 \xrightarrow{h.\text{corr}} f_2$

f_2

$f_1 \xrightarrow{h.\text{corr}} f_3$

$f_2 \xrightarrow{h.\text{corr}} f_3$

f_1, f_2, f_3, f_4

$f_4 \xrightarrow{\text{corr}} y$; v. less corr with

$f_1 \xrightarrow{\text{corr}} y$

$f_2 \xrightarrow{\text{corr}} y$

$f_3 \xrightarrow{\text{corr}} y$

$f_1 \xrightarrow{\text{not corr}} f_2$

$f_2 \xrightarrow{\text{not corr}} f_3$

overall impat

$f_1, f_2, f_3 \xrightarrow{M} y$ (better)

how to create f_4

$f_4 \xrightarrow{\text{corr}} y$; less corr with f_1, f_2, f_3

idea

$f_1, f_2, f_3 \xrightarrow{M} y$

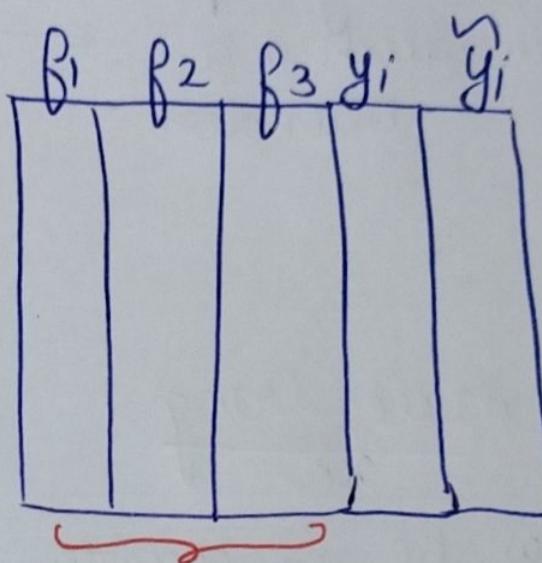
M:- \tilde{y}_i

(Q) how to assign new feature f_4

st

$f_4 \xrightarrow{\text{corr}} y_i$

less corr with f_1, f_2, f_3



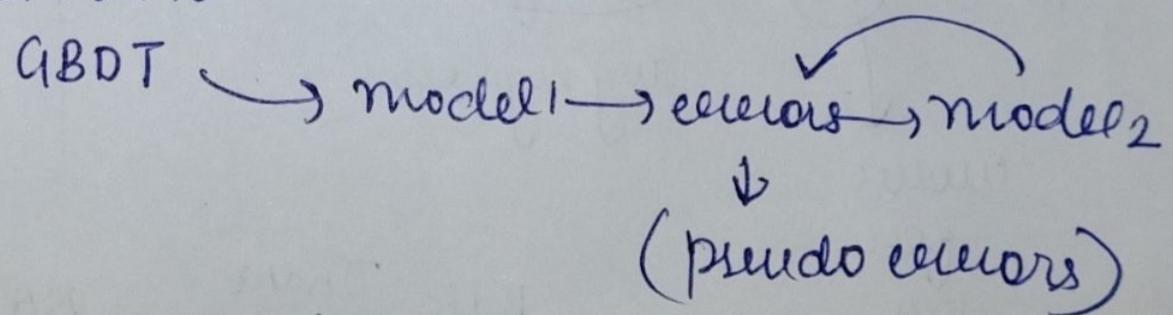
error:

$$y_i - \tilde{y}_i = e_i + x_i; \quad \left. \right\}$$

$f_4 \xrightarrow{\text{corr}} e_i$

Similar to

ABDT



$f_4 \rightarrow e_i$

avoid overfit \rightarrow Boosting

$f_4 \xrightarrow{\text{corr}} y_i - \tilde{y}_i$

such

\rightarrow less corr with f_1, f_2, f_3

Domain Specific featureization

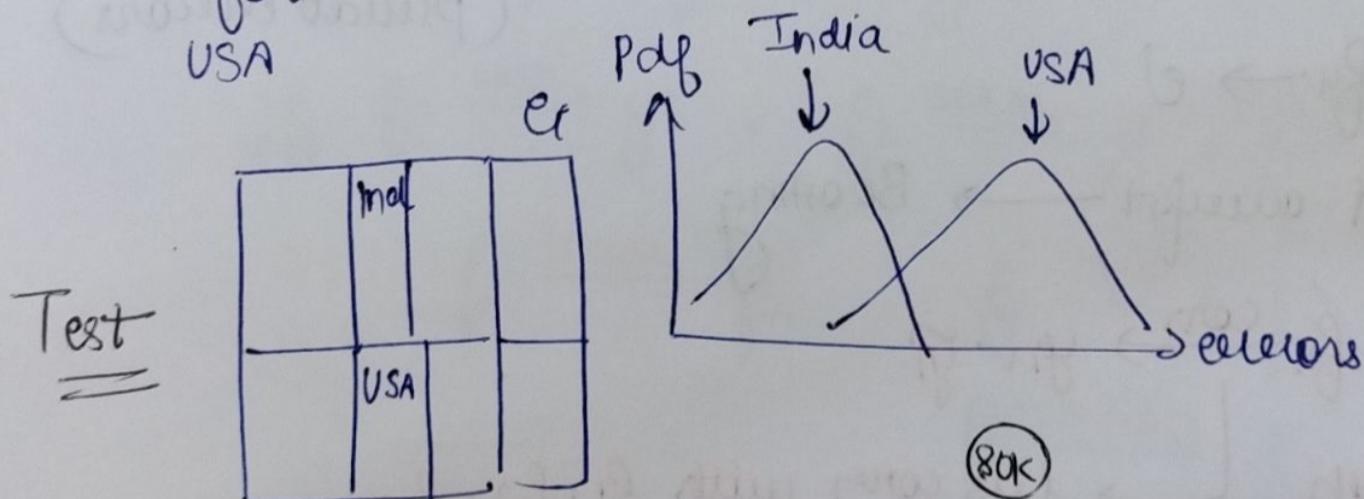
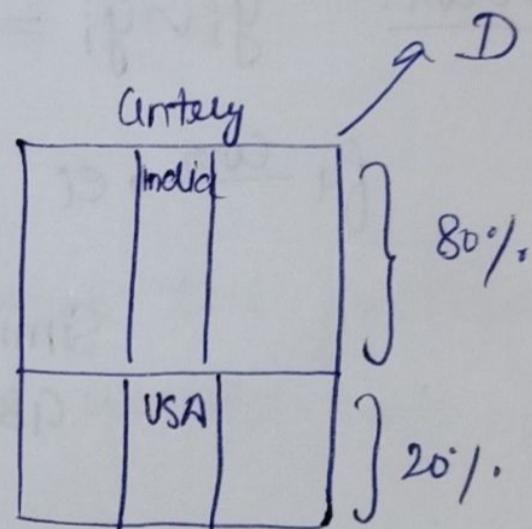
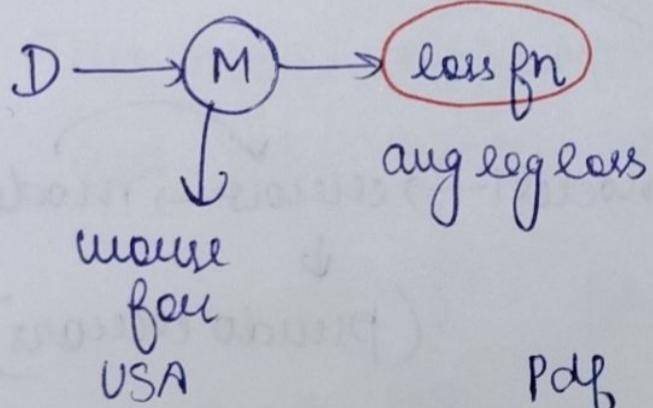
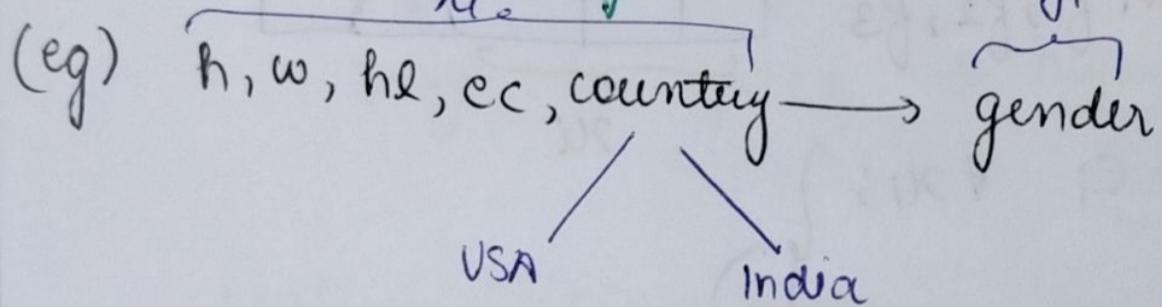
Heart attack → ECG data

doing no prior research is a blunder

impt to research existing featureizations by doc/specialists

create new features

Feature Slicing



- Slicing data on features
- ① cat₁, 8 cat₂ → different
 - ② sufficient # for each data

$D_{India} \rightarrow M_1$

$D_{USA} \rightarrow M_2$