

Module 5

classmate

Date _____

Page _____

ch1 - Featurization & Feature Engineering

Video: Introduction

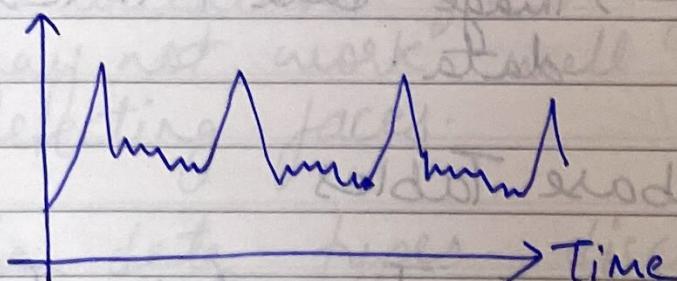
Featurization means converting any sort of data to numerical vectors.

Feature engineering means modifying the numerical vector so it works well with a model.

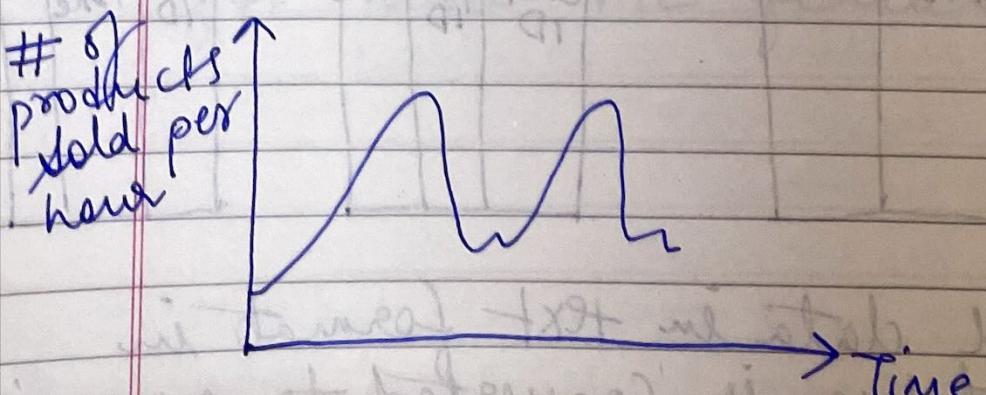
e.g. - BOW, TFIDF, Avg. W2V, etc are various featurizations of text data.

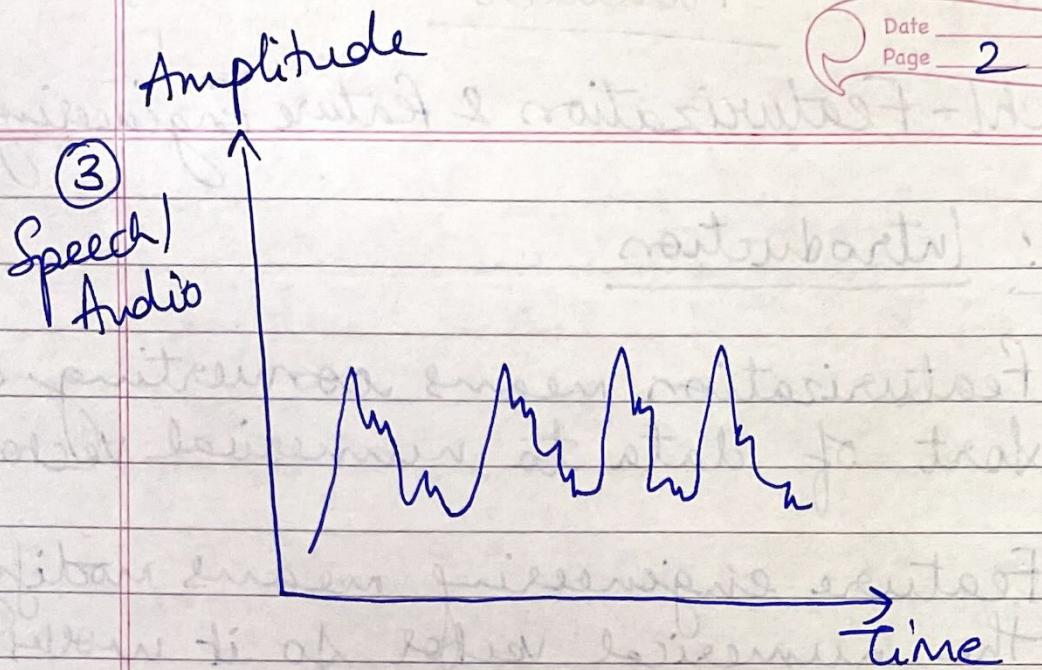
I Examples of time series data :-

① Heart Rate



② e-commerce





Other examples would be ~~the~~
Stock market, etc.

II Image data → Face detection,
Face recognition, MRI Scans, etc

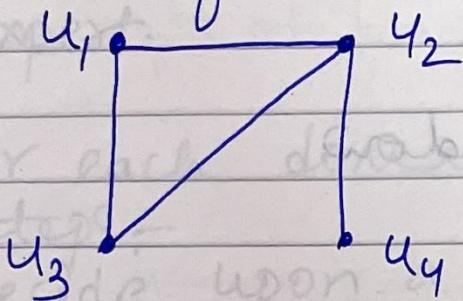
video → Image data + Time Series
data.

III Database Tables

T ₁ Customer Table	T ₂ purchase table	T ₃ product table
custID	custID prID Time	prID Price Di

All data in text format in database is converted to numeric features before supplying to ML model

IV Graph data \rightarrow recommend a friend on Facebook.



There are tons of types of data. Researchers have spent decades for getting proper featurizations.

30+ yrs. of research for image featurizations only.

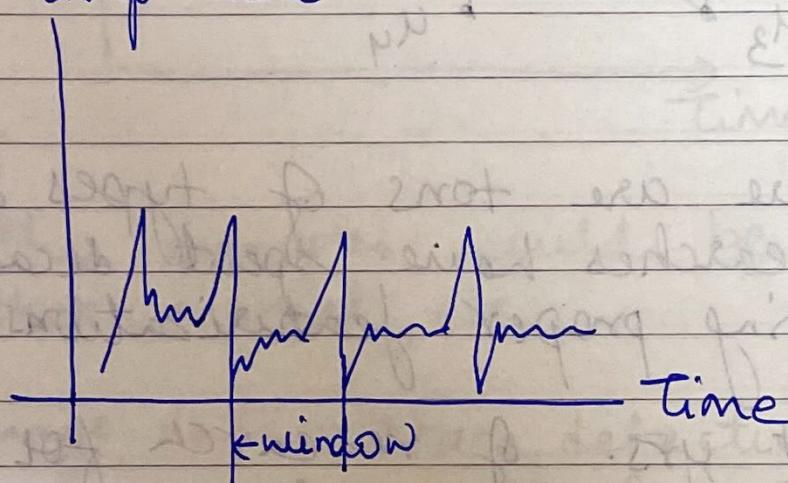
Featurizations for x-ray images may not work well with detecting faces.

The data types discussed here are only some of the many types, refer data and research for more ML featurizations.

Video 2: Featureization of Time Series

Data

- 1) Sliding Window Amplitude



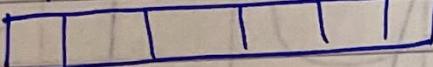
Right window width should be taken care of and it depends upon domain and application specific.

- We can do various types of calculations upon all the values belonging to the window like
- 1 mean, std-dev
 - 2 median, quantiles
 - 3 min, max
 - 4 max - min
 - 5 max ÷ min
 - 6 # local minima or maxima
 - 7 # mean crossing
 - 8 # zero crossing

which all out of these is suitable is decided by domain expert.

For each domain do :-

Steps :-

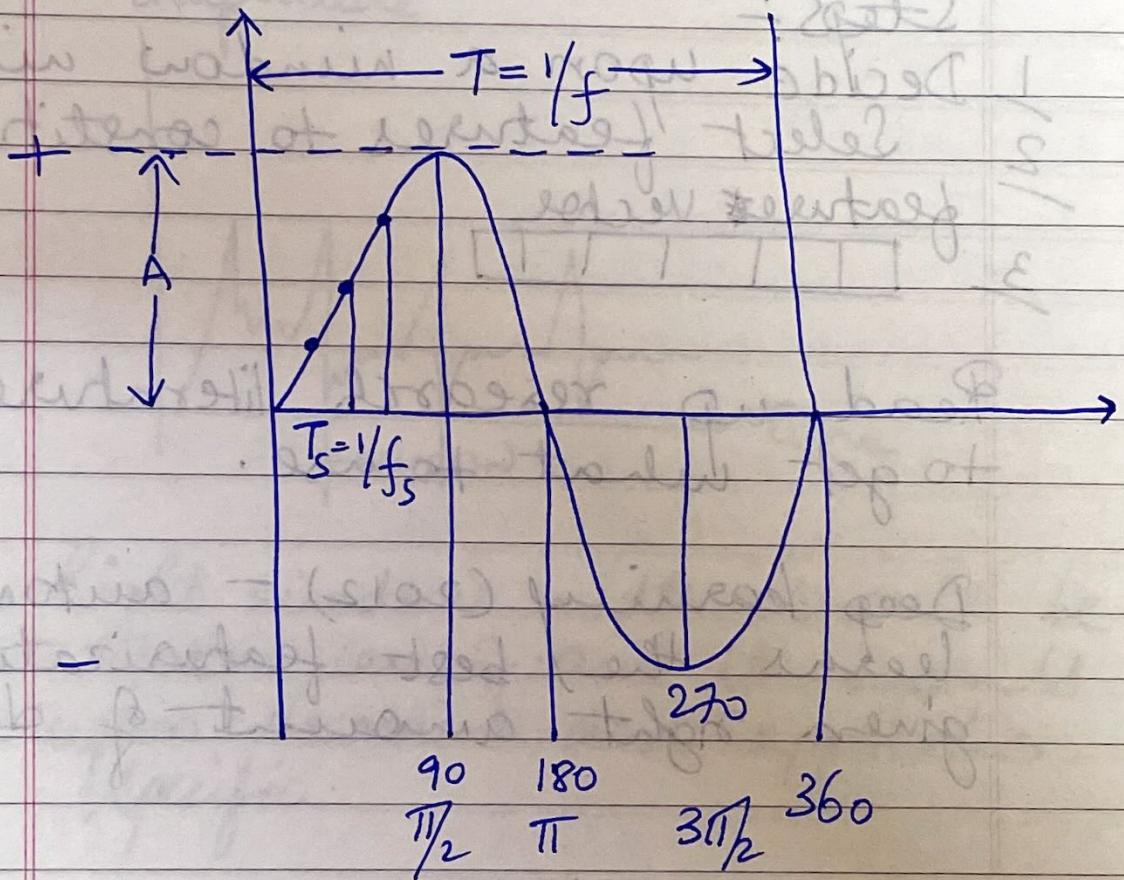
- 1 Decide upon a window width
- 2 Select features to constitute feature vector
- 3 

Read up research literatures to get what to use.

Deep learning (2012) - automatically learns the best featureizations given right amount of data.

Video3: Fourier Decomposition / Transform

It is a method to represent Time Series data.

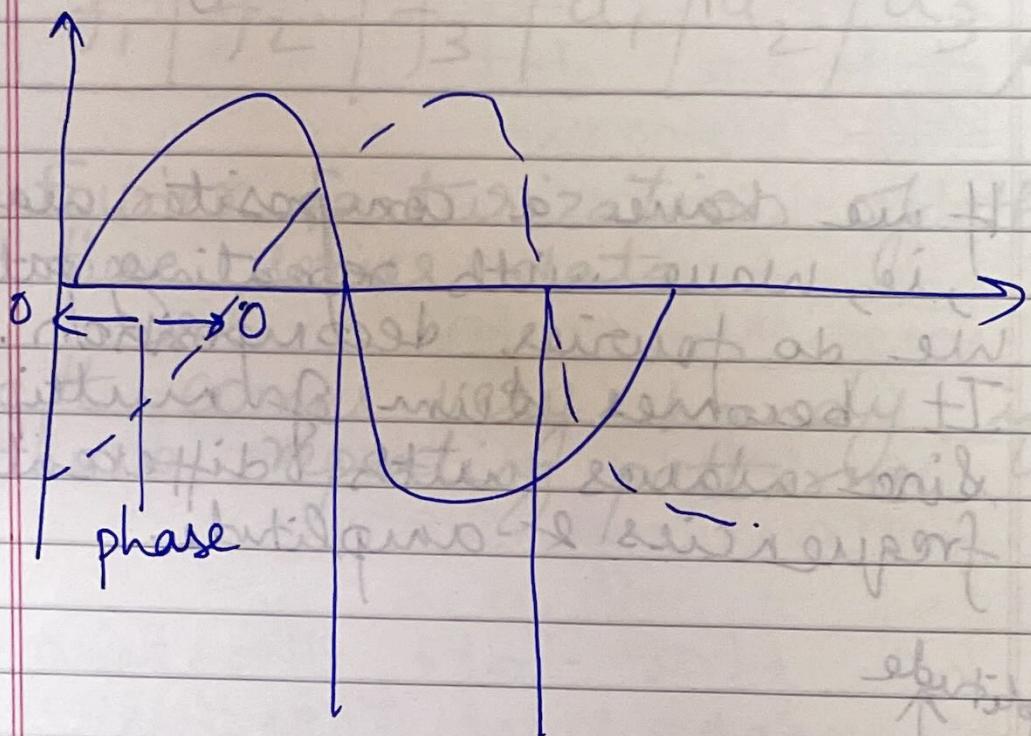


T = period
 f = frequency
 A = amplitude.

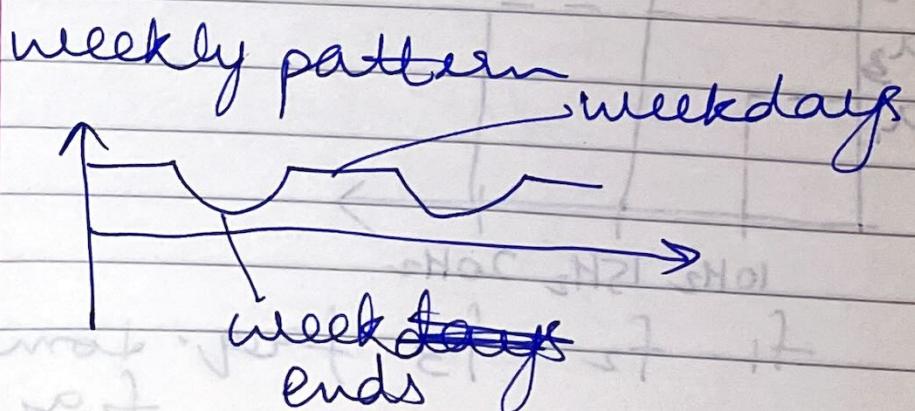
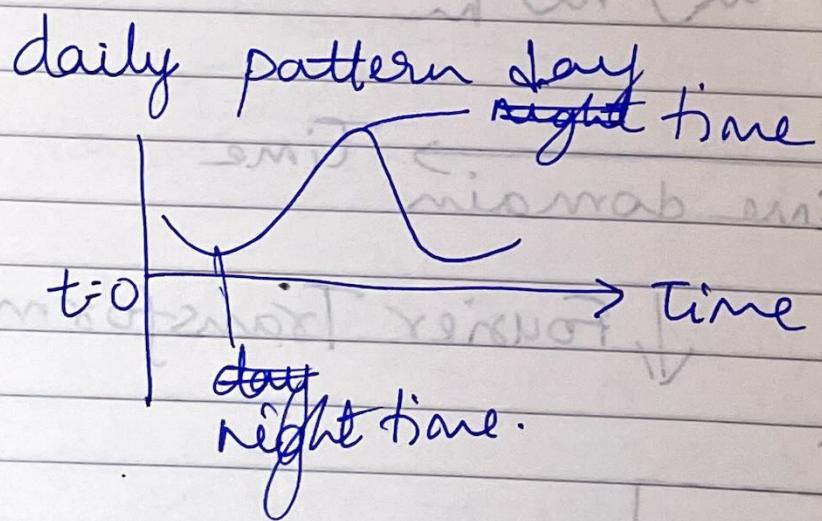
$$\text{freq} = \frac{1}{T}$$

1 oscillation/sec = 1 Hz freq.
 2 " " " = 2 Hz "

Phase of a wave -

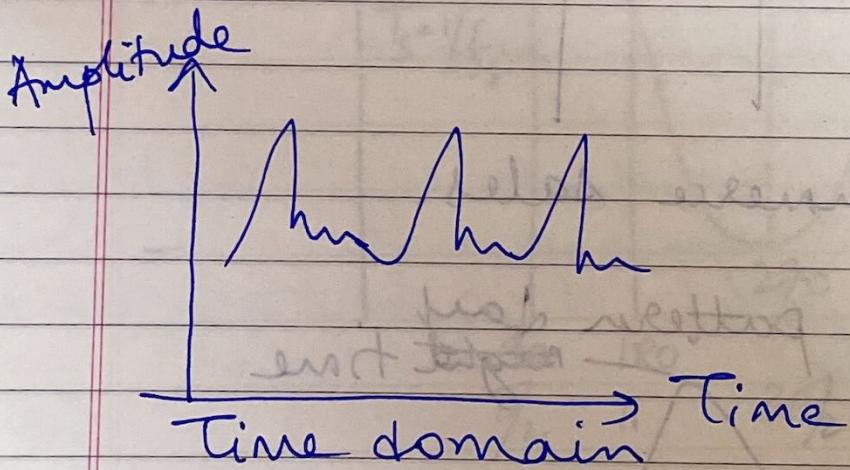


eg- e-commerce sales

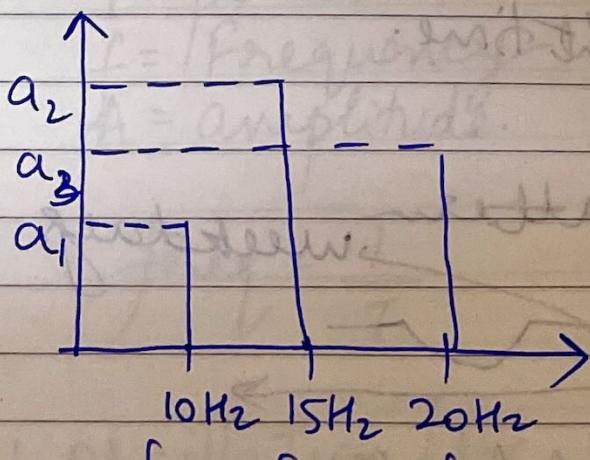


Similarly annual pattern
 \rightarrow sleep PPS during vacation

If we have a composite wave, ie, wave with repeating pattern we do Fourier decomposition. It becomes sum of multiple sine waves with different frequencies & amplitude.



\downarrow , Fourier Transform



$f_1 \quad f_2 \quad f_3$ freq. domain
 $a_1 \quad a_2 \quad a_3$

$$f_2, a_2$$
$$f_3, a_3$$

Feature Vector is like :-

f_1	f_2	f_3	a_1	a_2	a_3
-------	-------	-------	-------	-------	-------

Fourier featurization of the time series data

It comes very handy if we have repeating patterns.

Video 5: Image Histogram

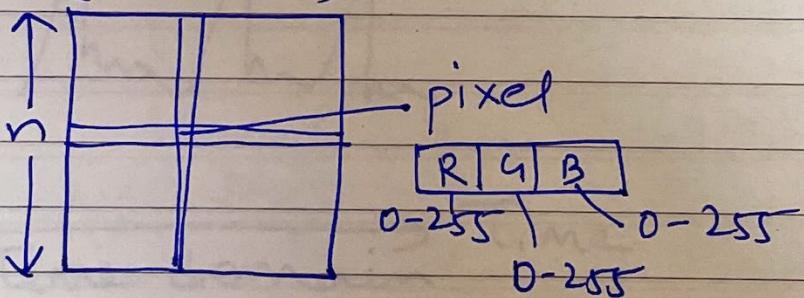
2 types of ~~featureizations~~ ^{Image} histograms of image data:-

- 1 ~~Image~~ ^{Colour} histogram
- 2 Edge histogram

1 ~~Image~~ ^{Colour} histogram (Colour)

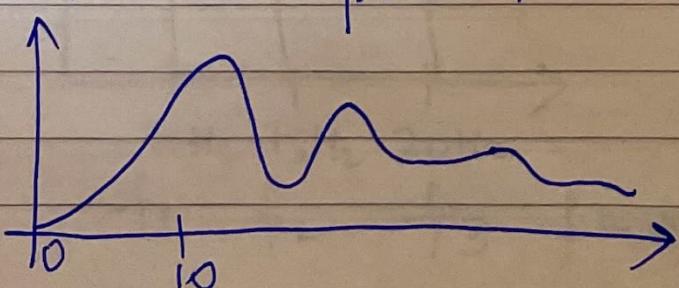
The colour histogram in the following steps:-

We first take the RGB value of each pixel like this:-



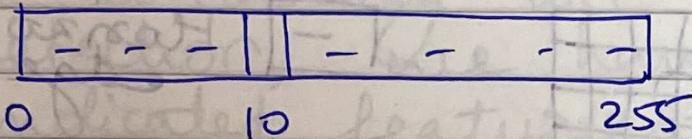
The pixel has RGB values in the range 0-255.

So we collect red values from each of the $n \times m$ datapoints pixel and plot the histogram



This gives the frequency of occurrence of a particular intensity of red colour in the entire picture.

The histogram can then be converted to a vector.



Similarly we have vectors for Green & Blue colours as well.

We can concatenate them to get the resultant feature vector.

Such featurization can be used for tasks like: Check whether there is sky or not in the given image.

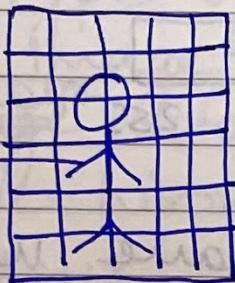
If there is sky, there will be high % of high intensity blue colour.

But colour histograms cannot recognize shapes.

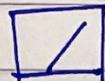
For recognizing shapes, we have Edge Histograms.

2 Edge Histograms

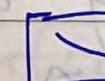
Edge
in
Image



Horizontal
edge of 0°



45°



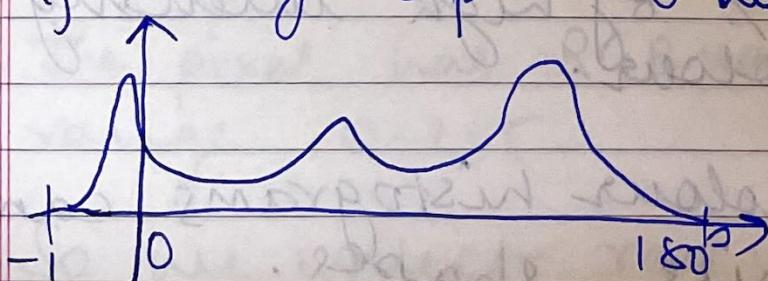
135°



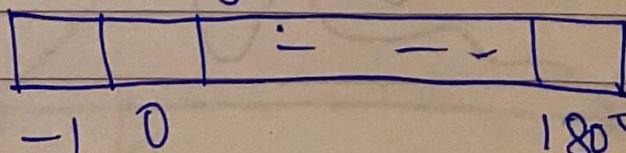
90°

Image is broken down into grid as in above and edge angle for each pixel is assigned.

For each region we calculate edge value or edge angle and plot histogram. If no edge is present we give -1°



Histogram ~~has~~ is converted to vector of size 181



These 2 histogram featurizations are very rudimentary. For something like face detection, we have more complicated featurization called "Haar" features.

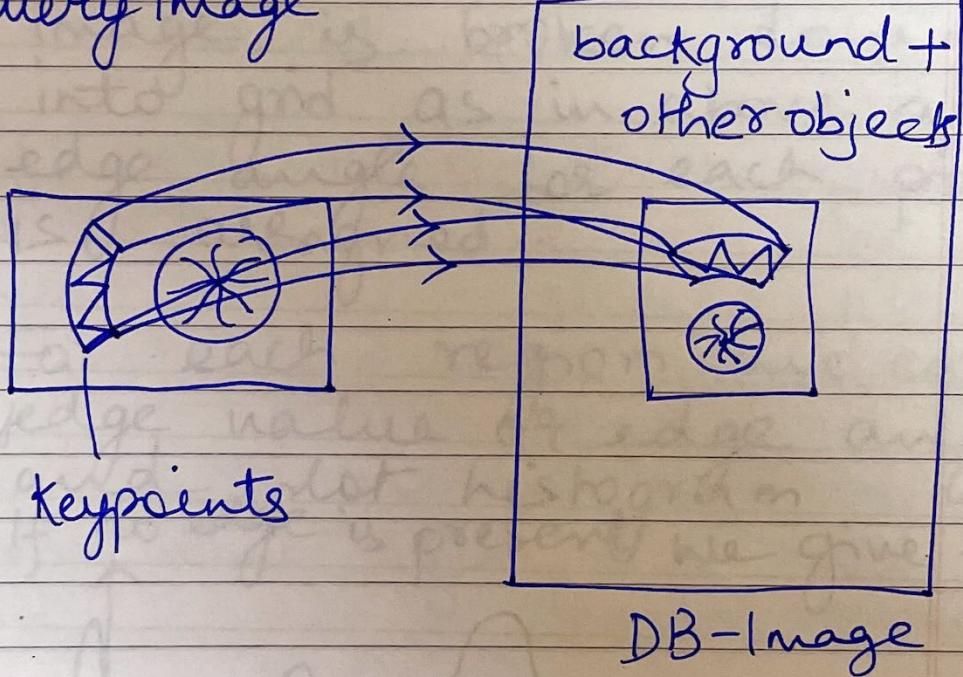
Another library called OpenCV can be used to get SIFT features off a given image.

Video 6: SIFT (Scale Invariant Feature Transform)

Extremely useful for detecting objects in an image.

Variants of SIFT are used in image search eg - Amazon.

Query Image



SIFT detects keypoints mostly corners. For every keypoint, it creates a 128-dimension vector. These keypoints help in detecting an object in the image.

Scale Invariance in SIFT means, even if query image is larger or vice versa than DB image, it will still detect.

Similarly, we have Rotational Invariant?

Another library called OpenCV can be used to get SIFT features of a given image.

Video 8: Relational Data & Feature Extraction

e.g. → Customer Table Customer Viewing/Visiting Purchase Table

Cust ID	Zip Code

Cust ID	Prod ID	Time

Cust ID	Prod ID	Time

Product data

Prod ID	Prod Type

This database can be featurized by extracting using pandas, Python, SQL & domain knowledge.

eg → Task is to predict if a customer would purchase a product in the next 7 days.

customer + product → $| \begin{matrix} \text{ID} \\ \text{ID} \end{matrix} |$
buy \downarrow
buy \downarrow

We have to use domain knowledge to convert the given (cust ID + prod ID) to a feature vector for the given task.

Some relevant features:-

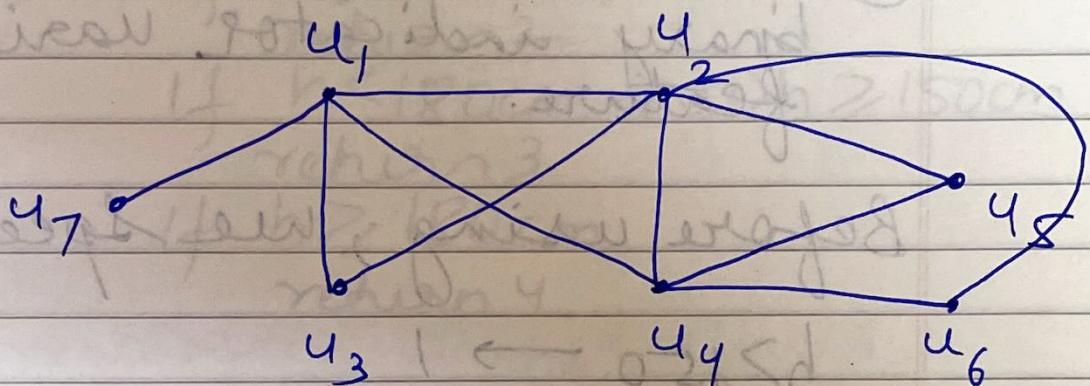
1 f_1 : # of times the cust ID visited the page of given prod ID in last 24 hrs.

2 f_2 : # of times the cust ID visited any product of same type as prod ID in last 24 hrs.

f_3 : income level

f_4 : Zip codes \rightarrow signifies urban/rural areas.

Video 9: Graph data & featurization



Facebook/Social networking websites use social graphs.

u_i : vertex (users)

(u_i, u_j) : edge (friendship)

Task: Recommend features would be

1 f_1 : # of mutual friends b/w u_4 & u_i $\forall i \neq 4$

2 f_2 : # of paths b/w u_4 & u_i $\forall i \neq 4$.

Video 10: Indicator Variable

e.g. Height can be used

- ① Directly as a real valued feature
- ② Can be converted to binary indicator variable feature.

Before using, we specify:-

$$h > 150 \rightarrow 1$$

$$h \leq 150 \rightarrow 0$$

why we may used to do so & decide threshold are problem specific

It may refer male as 1 & female or anything else.

e.g. 2 - Categorical feature like Country :-

if (country = India) OR
(country = USA) 1

return 1

else return 0

VII: Feature Binning (Bucketing)

This is a logical extension to indicator variable.

eg - If $h < 120\text{cm}$
return 1

If $h \leq 150\text{cm}$ AND $h > 120\text{cm}$
return 2

If $h < 180\text{cm}$ AND $h \geq 150\text{cm}$
return 3

If $h \geq 180\text{cm}$
return 4.

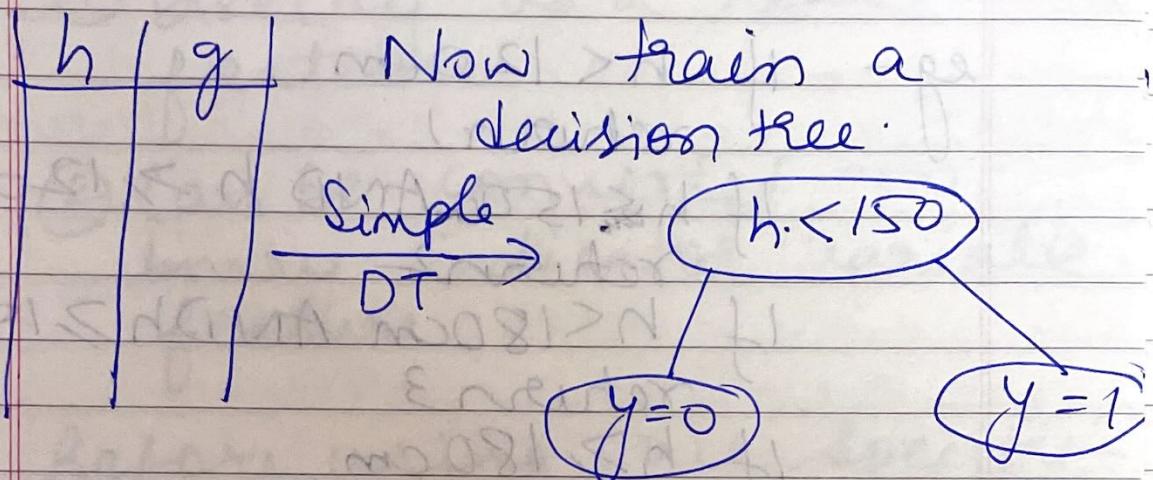
Interview Question: Predict gender given weight(w), height(h), hair length(hL) and eye colour (ec)

D _{Train}	h	w	hL	ec	g

Challenge - Perform data driven binning on feature 'height' using

D_{Train}

Soln: Take h & g & make a new table.



Many Kaggle competitions, this trick comes handy.

V12: Interaction Variables

e.g. $h, w, hL, ec \rightarrow \underbrace{\text{gender}}_{y_i}$

Say from domain knowledge we know if $(h < 150) \& (w < 60)$ female is high probability.

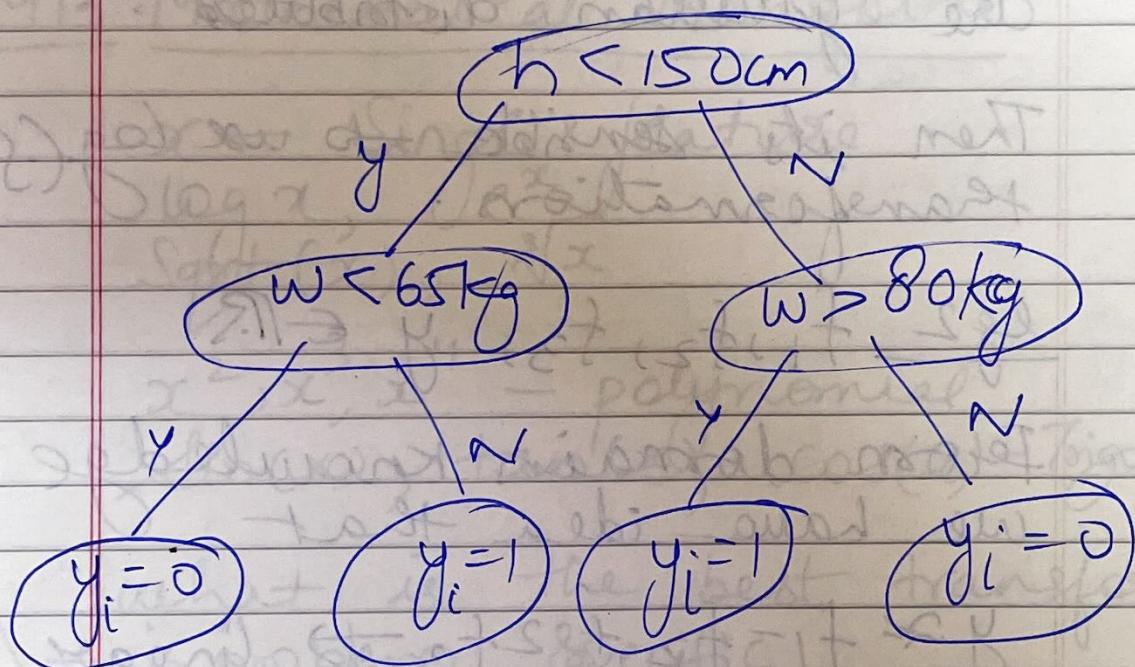
Then we can create a new 5th feature by 2 way logical interaction.

$$f_5 = (h < 150) \text{ AND } (w < 60)$$

Similarly we can also have

$f_6 = h \neq w$. - 2 way interaction feature.

Again DT can be used to create nice interaction feature.



$$(h < 150\text{cm}) \text{ AND } (w < 65\text{kg})$$

Similarly for each leaf node we can have an interaction variable.

And these interaction features are guaranteed to be helpful, because DT was constructed to predict gender using information gain.

Model Specific featurization

V1.3: Mathematical Transforms

eg) $f_1 \rightarrow$ has powerlaw distribution.

We want to apply Logistic regression same as Gaussian Naive Bayes assuming features are gaussian distributed!

Then it's sensible to use $\log(f_i)$ transformation.

$$\log f_1, f_2, f_3, y \in \mathbb{R}$$

From domain knowledge we have idea that

$$y \approx f_1 - f_2 + 2f_3 \rightarrow \text{linear combination of } f_i's$$

Then linear models are powerful if used as compared to others.

e.g. Linear Regression or Linear SVM.

$$\text{where } w_1 = 1, w_2 = -1$$

$$w_3 = 2$$

fits decently.

eg 3 From domain knowledge
 say we know
 $y \xrightarrow[\text{on } f_1, f_2]{\text{depends}} \text{Interactions of } f_1, f_2$

Then better use DT / RF / GBDT

M13 Mathematical Transforms

x — Single feature

$\log x$, e^x
 \sqrt{x} , $\sqrt[3]{x}$

x^2, x^3, x^4 — polynomial

$\sin(x), \tan(x), \cos(x)$ — Trigonometric

— what is the best transform
 — is problem specific

eg-elfs

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

V15: Feature Orthogonality

The more different/orthogonal (or \perp) the features are, the better the model would be.

So whenever adding a new feature, keep in mind that the new feature should be correlated to the o/p - ' y ' & very less correlated to other features.

e.g. $f_1, f_2, f_3 \xrightarrow{\text{Model}} \hat{y}_i$ & we

also have y_i = actual labels.

How to design a new feature f_4 such that

$f_4 \xrightarrow{\text{correlated}} y_i$ &

$f_4 \xrightarrow{\text{less correlated}} f_1, f_2, f_3$

$$\text{Solve } f_4 = y_i - \hat{y}_i$$

Try to design f_4 correlated to e.g. so f_4 related to \hat{y}_i but

less correlated to other features.
It is similar to boosting.

But do remember to avoid overfitting.

V16: Domain Specific featurizations

e.g - Predicting heart attack using ECG data.

Always important to research and study existing featurizations by Doctors/Spl.

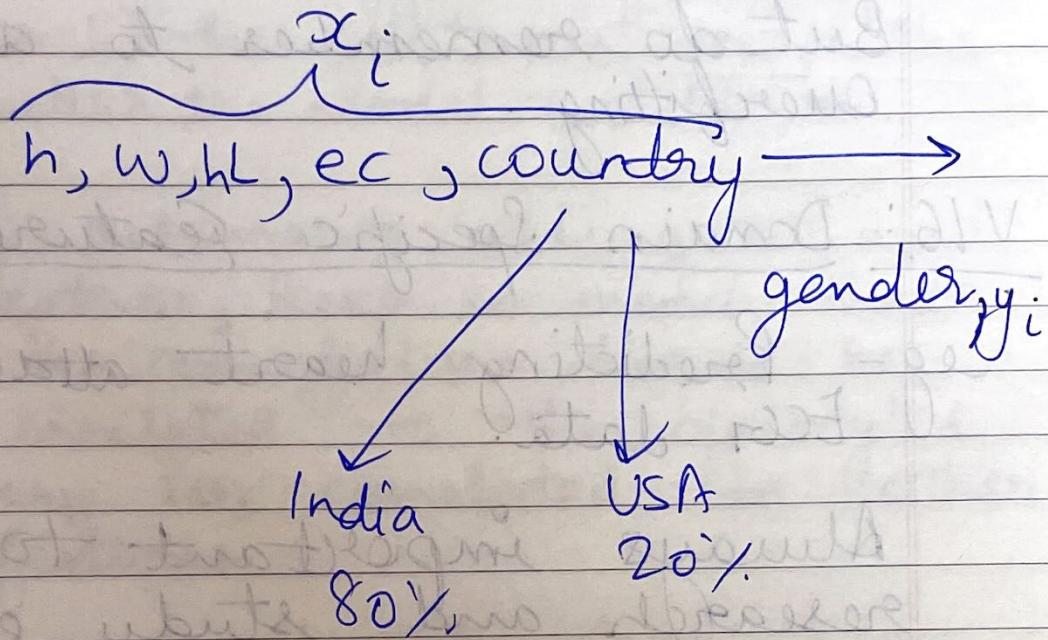
Can create new feature out of these existing using ML hacks from experience.

blog.kaggle.com/winner-solution

Has details of the feature engineering they have done.

V17: Feature Slicing

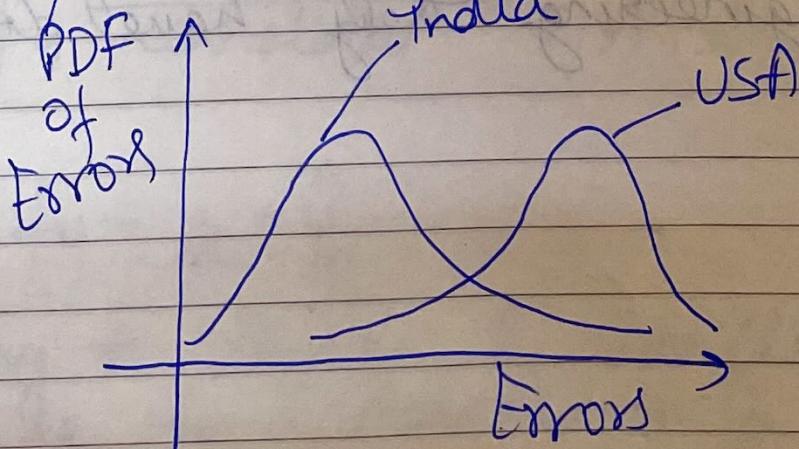
e.g.



Due to majority of datapoints belonging to India, model performs

- 1) better on Indian test data
- 2) worse on USA test data

So plot PDF of Errors



Strategy → Split D to

D_{India}

and train M₁

D_{USA}

and
train M₂

This idea of slicing the dataset and training separate models helps if:

- 1) Both category have different behaviour
- 2) Sufficient data points for each category.

Real World Example: Mobile & Desktop users in Amazon have different types of traffic.

So feature slicing is done.