

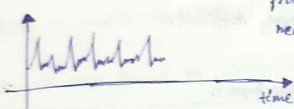
Featurization & Feature Engineering

- * Featurization → converting any sort of data to numerical vectors.
- * Feature engineering → modify the numerical vector so it works well with a model (hacks)

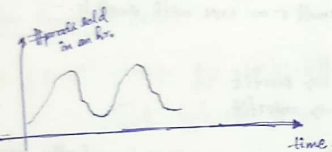
Ex → BoW, tf-idf, avg W2V etc are various featurizations of text data.

(i) Examples of time series data

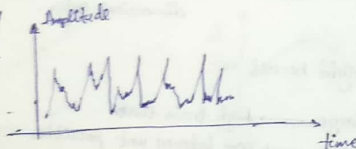
① heart-rate



② temperature



③ speech / audio



Other examples would be stock market, etc.

(ii) Image data → face-detection, face recognition, MRI-scans etc.

video → Image data + time series data

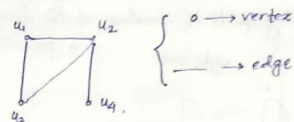
(iii) Database tables

T ₁		T ₂			T ₃		
custID	doc	custID	prod	time	P-ID	P-Val	Discount

converted to numerical features before supplying to ML.

In the end, ML algos can only process numbers. So, everything needs to be converted to that before processing.

(iv) Graph data → recommend a friend on FB.



There are tons of types of data.

Researchers have spent decades for getting proper featurization.

* 30+ yrs of research for image featurization only.

* featurization for x-ray images maynt work well with detecting faces.

* cant read up all the featurizations to all kinds of data
So, lookup when required.

Featurization of time series data

① Sliding window

Window-width → domain & application specific.

We can do various types of calculations upon all the values belonging to the window like

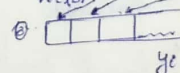
- ① mean, std-dev
- ② median, quantiles
- ③ min, max
- ④ max-min
- ⑤ max/min
- ⑥ # local minima or maxima
- ⑦ # mean crossing
- ⑧ # zero crossing.

* Which all out of these is suitable is decided by domain expert.

Steps

① decide upon a window width.

② select features to constitute feature vector



Read up research literatures to get what to use.

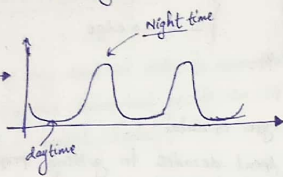
Drop-learning (2012) - automatically learns the best featurization given right amt of data.

② Fourier Decomposition / Transform →

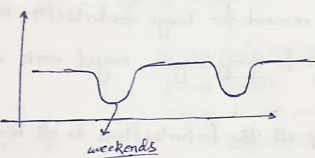
Whenever we have some repeating pattern, looking at the frequency is a good idea.

ex → e-commerce

daily pattern →



weekly pattern →

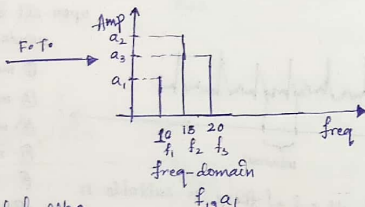
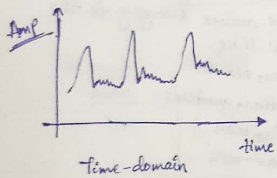


Similarly, annual pattern → dip during vacations.

Composite wave (with repeating pattern)

Fourier decompose →

sum of multiple sine waves with different frequencies & amplitudes.



Feature vector → somewhat like

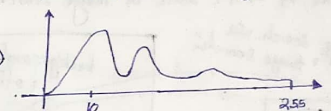
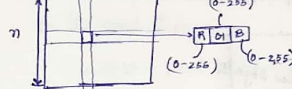
$[f_1 \ f_2 \ f_3 \ a_1 \ a_2 \ a_3]$

fourier featurizatⁿ of time series data.

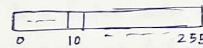
comes very handy if we have repeating patterns.

Featurization of Image data

② Image histograms (colour histogram) → collect red values for each of the $(m \times n)$ pixel & plot the histogram



The histogram can then be converted to a vector.



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

Similarly we'll have vectors for green & blue colors.

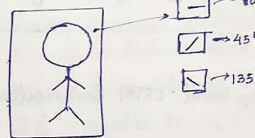
We can concatenate them to get the resultant feature vector.

* Such featurizatⁿ can be used for tasks like:— check whether there is sky or not in given image.

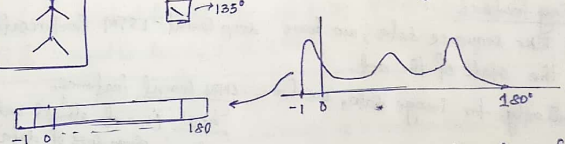
→ If there is sky, there'll be high percentage of high intensity blue color.

* But colour histograms can't recognise shapes.

edge histogram



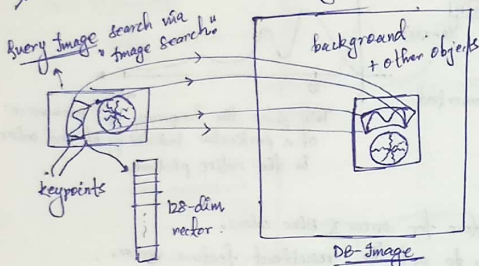
So, image is broken down into pixels as in above and edge angle for each pixel assigned.
No edge = -1.



For these two histogram featurizations are very rudimentary. For something like face detection, we have more complicated featurization — "haar" features.

(b) Scale Invariant Feature Transform (SIFT)

- * extremely useful for detecting objects in an image.
- * variant of SIFT used in image search - Amazon



→ It detects keypoints, — mostly edges corners.

for every keypoint, creates a 128-dim vector.

These keypoints help in detecting an object in the image.

SIFT → scale invariant → Even if query image is larger than DB-Image or vice-versa, it'll still detect.

Similarly, Rotation invariant ✓

SPECV → can be used to get SIFT features of a given image.

Deep Learning Features

Just like sequence data, we have deep learnt LSTM featurization as the state of the art.

Similarly, for image data, we have LNN learnt features.

best (learnt almost automatically given lots of data).

Relational data & featurization

custID	custZipcode

cust table

custID	PID	Time

customer View / visitatn

custID	PID	Time

Purchase table

PID	PType

Prod data

Can featurize DB-data by extracting using pandas, python, SQL & domain knowledge.

ex → task is to predict if a customer would purchase a product in the next 7-days.

custID + prodID → 1/0 — doesn't buy, buys.

→ We have to use domain knowledge to convert the given (custID + prodID) to a feature vector for the given task.

Some relevant features —

- ① f_1 → #times the custID visited the page of prodID in last 24hrs.
- ② f_2 → #times the custID visited any prod of same type as prodID in last 24hrs.
- ③ f_3 → income level.
- ④ f_4 → zipcodes — signifies urban/rural areas.

Graph-data & featurization.



social graph.

u_i :- vertex (users).

(u_i, u_j) :- edge (friendship).

task :- recommend new friends to a user u_i , say u_1 .

Some useful features would be

- ① f_1 — #mutual friends b/w u_1 & u_i $\forall i \neq 1$.
- ② f_2 — #paths b/w u_1 & u_i $\forall i \neq 1$.

Feature engineering: Indicator variable.

ex → height can be used — directly as a real valued feature, or converted to binary indicator variable before using.

ex → categorical feature → country
if (country = India) or (country = USA) return 1
else return 0.

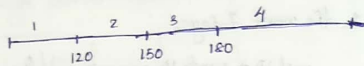
ex → $h > 150 \rightarrow 1$
 $h \leq 150 \rightarrow 0$

why we may need to do so & decide threshold are problem specific
may refer male (1) & female or anything.

Feature Binning (Bucketing)

→ logical extension to ~~feature~~ indicator variable.

Ex → if $h < 120$ cm
return 1
else if $h < 150$
return 2
else if $h < 180$
return 3
else return 4



Interview Ques → predict gender given — weight (w), height (h), hair length (hl), eye color (ec).

D_{train}

h	w	hl	ec	g
				✓
				✓
				✓

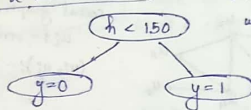
challenge → perform data driven binning on feature 'height' using D_{train}

So! → take h & gender (g) & make a new table.

h	g

Now train a decision tree.

simple
DT



DT finds the threshold with max inf. gains.

* Many kaggle competitions this trick comes handy.

Interaction Variables

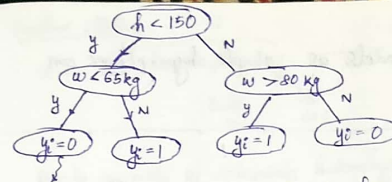
Ex → $h, w, hl, ec \rightarrow$ gender
 $x_i \rightarrow y_i$

say fm domain knowledge we know if $(h < 150) \& (w < 60)$
→ female high probability.

Then, we can create a new 5th feature by 2-way ~~and~~ logical interaction.

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

similarly, we can also have $f_6 = h * w$ — math 2-way interaction feature.
Again, DT can be used to create nice interaction features.



$(h < 150 \text{ cm}) \text{ AND } (w < 65 \text{ kg})$ similarly for each leaf node we can have an interaction variable.

* And those interaction features are guaranteed to be helpful because DT was constructed to predict gender using information train.

Mathematical Transforms

x — single feature.

$\log x, e^x$

$\sqrt{x}, \sqrt[3]{x}$

x^2, x^3, x^4 — (polynomial).

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

what is the best transform?
→ problem specific.

Ex → if ~~feature~~ x_i has power law distributions, then taking $\log(x)$ transform as it'll have a gaussian distributions.

Model Specific Featureization

Ex → f_1 → has powerlaw distribn

We want to apply Logistic Regression — Same as Gaussian Naïve Bayes (assumes features are gaussian distributed).
Then its sensible to use $\log(f_1)$ transformation.

Ex → $f_1, f_2, f_3 \rightarrow y \in \mathbb{R}$.

fm domain knowledge we have idea that

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

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

linear regression with $w_1 = 1, w_2 = -1, w_3 = 2$ — fits decently.

Ex → fm domain knowledge, say we know y depends on interaction of f_1 & f_2 .
Then better use DT/RF/GBDT

Ex: BOW \rightarrow high dimensions.

Then better to use linear models as simple hyperplanes can separate data

Feature Orthogonality

* The more different/orthogonal 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.

Ex $\rightarrow f_1, f_2, f_3 \xrightarrow{\text{Model}} \hat{y}$ & we also have y - actual labels.

How to design a new feature f_4 so that $f_4 \xrightarrow{\text{correlat}} y$ & $f_4 \xrightarrow{\text{less correl}} f_1, f_2, f_3$

Solⁿ $\rightarrow e_i = y_i - \hat{y}_i$ Try to design f_4 un-related to e_i
So, f_4 related to y_i but less correlated to other features.

Similar to boosting

* Remember to avoid overfitting.

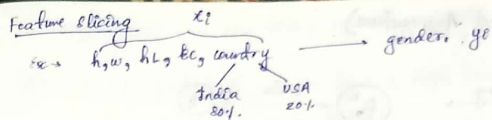
Domain Specific Featureizatⁿ

Ex \rightarrow Predicting heart attack using ECG data

* Not reading these at all is a big blunder. $\left\{ \begin{array}{l} \text{Always important to research and study} \\ \text{existing featureizations by doctors/specialists.} \end{array} \right.$

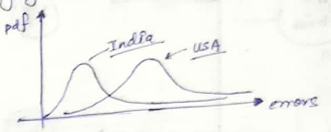
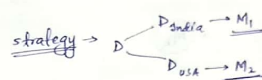
Can create new feature out of these existing using ML hacks from experience.

Blog: kaggle.com/winner-solution - details of the feature engineering they have done.



Due to majority of datapoints belonging to India, model performs

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



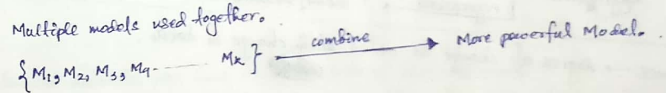
This idea of slicing the dataset and training separate models helps if

- ① both category have different behaviour.
- ② sufficient data points for each category.

Real world example \rightarrow Mobile & Desktop users in Amazon have different types of traffics.
So, feature slicing done.

Ensemble Models

Multiple models used together.



4 types of ensembles

- ① Bagging (Bootstrapped Aggregation).
- ② Boosting
- ③ Stacking
- ④ Cascading

Mostly used in Kaggle as well as Real world scenarios.

Key Aspect \rightarrow More different the models are, the better & powerful it becomes.