

## Module - 5

### Ch - Feature Engineering

→ Featureization & feature engineering:

Till now we have seen, how to convert text data to numerical data, categorical data to numerical data. But there are some other types of data which needs to be converted to numerical data.

- ① Time-series data → heart rate, stock price etc
- ② Image data
- ③ Video data.

We will learn how to deal with these type of data.  
Some other type of data is graph data.

There are tons of type of data.

Featureization is to convert that data to numerical data.

\* Moving windows for time series data:

This is simplest featureization of time series data.

In each window you can choose any polys

① Mean, std dev of window

② Median, Quantiles "

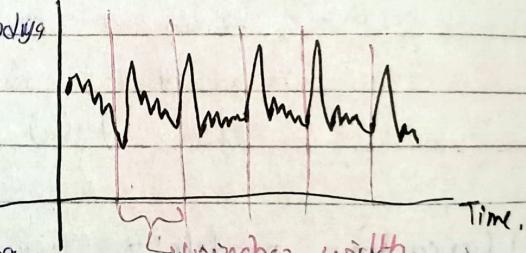
③ Max, min

④ No of local minima & local maxima

Step 1: Define window width

Step 2: Identify the useful features.

Step 3: put all the features as  $x_i$  & identify  $y_i$  according to task

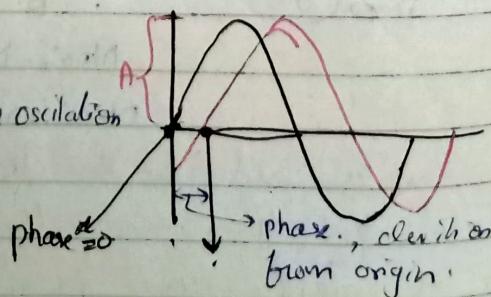


\* Fourier Decomposition / Transform:- It is another method to represent Time-series data.

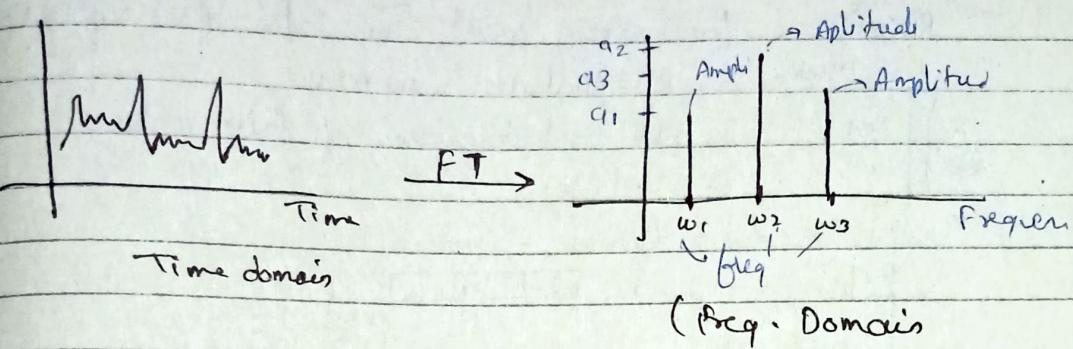
Amplitude:- height of wave.

period:- Time taken to finish one oscillation.

frequency-  $\frac{1}{\text{Time period}}$



If we have a repeating pattern, we can decompose it as sum of sin wave.



Some time it is useful to represent data in freq domain.

Now we get

freq	Amp
$f_1$	$a_1$
$f_2$	$a_2$
$f_3$	$a_3$

We can create a feature vector as  $[f_1, f_2, f_3, a_1, a_2, a_3]$

This type of representation is called Fourier representation of time-series data.

## \* Deep learnt features:-

Priorly, we had domain expert to identify the feature in data from time-series etc. They designed special features for each domain. But the problem is feature designed in one domain may not work in other data.

This was the common practice until 2012 - 2013. But now the thing changes.

Deep learning automatically (almost) learns the best featureization for data specially in time-series, text and image data. Only req. is it needs lots of data.

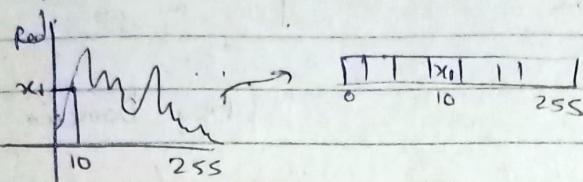
LSTM is one such model to automate featureization.

\* Image histogram:- we deal with lot of image in real world. Hence important to deal with image data. Today the best way to ~~learn~~ image featureization in Deep Learning (NN). But as of now we learn a old technique called image histogram.

There are 2 types of image histogram ① color histogram ② edge histogram

① Color histogram: we know every pixel have 3 values R G B, for every pixel we take Red values.

$n \times m$  rows  $\times n \times m$  columns  
Red value =  $n \times m$ ,  
plot a histogram of these values.



Similarly do for  $B \times 6$  colors. Hence for a image we have

3 vectors:

$x_i \rightarrow [R \quad G \quad B] \quad y_i$

② Edge histogram: each edge can have different direction, hence we can assign angle to each edge. Break up the big image into grid and ask a question,  $\oplus$  in a grid if there is a edge in cell with  $0^\circ, 45^\circ, 90^\circ, \dots, 360^\circ$ . If there are no edges assign the angle to that cell/region.  $\therefore$  for each region/cell we get a edge value / angle. Now plot a histogram of these values. And using the histogram convert it to vector.

## \* Scale Invariant feature Transform (SIFT)

This is a very popular image processing feature extraction technique. It is very useful in object detection. It identifies some key points of image as features and create a vector of 128 dimensions for each feature.

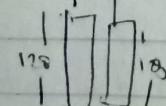
Scale invariant here means that

If the image is big or small

The key points will not change

SIFT is also rotation invariant.

$$[f_1 \ f_2 \ f_3 \ f_4]$$



$$\begin{bmatrix} f_1 & f_2 & f_3 \\ f_4 & f_5 & f_6 \\ f_7 & f_8 & f_9 \\ f_{10} & f_{11} & f_{12} \\ f_{13} & f_{14} & f_{15} \\ f_{16} & f_{17} & f_{18} \end{bmatrix}$$

\* Relational data & featurization : Some time data may be in discrete form table , we need to featurize them.

Ex:-	id zip code	ID, PIP, Time	CID, PIP, Time	PID, Product type
	Customer table	Customer viewing data	Purchase data	Product data

\* Task :- predict if a customer would purchase a product in the next 7 days.

customer, Product →  $\xrightarrow{\text{buy}}$  Not buy .

so how to we come up with features .

# CID viewed → Income level  
# CID visited any product belongs to same type as PID

In general we use SQL to feature relational data.

\* Graph - data & featurization :

In a graph based data, we need concepts from graph theory like no. of paths b/w two nodes, No. of neighbours etc.  
It depends on type of task.

\* Feature Engineering : Indicator Variable :-

Assume, we have height as feature - since height is a real-value feature , I can leave it as is . One other approach is we can change it to indicator variable :-

If  $h > 150 \rightarrow 1$

else  $0$

∴ height is now converted to a binary value feature i.e. indicator variable feature . Because the value of height indicates something that if  $h=1$  then height is more than 150 . It is a problem specific and there is no hard and fast rule .

Example ② categorical feature :- Country  
if Country = INDIA or Country = USA else,  
else down a ;

\* **Feature binning** :- Feature binning is a logical extension of indicator variable. we create bins here.

Ex:- If  $h < 120$

step 1

If  $h > 120 \& h < 150$

step 2

If  $h > 150 \& h < 180$

step 3

else step 4

we are creating 4 bins/category.

It is also called feature bucketing.

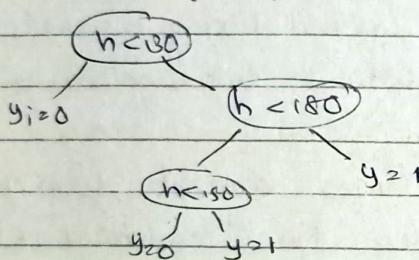
Question is How to get the right threshold?

Ans: getting right threshold is problem specific.

We can also use decision trees to make bins. of a particular feature. Suppose

There are  $N$  features and a class table, we want to ~~for~~ bin feature  $f_j$ . Then we will create a DT using only  $f_j$  and class table. (fully grown DT) using ~~some~~ ~~for~~  $f_j$  only.

Ex:-



"Now we can create bin.."

\* **Interaction variables** :- Assume our task is to predict gender using height, weight, hair length, eye color.

$h, w, hl, ec \rightarrow \text{gender}$ .

Ex. (1) If  $(h < 150)$  AND  $(w < 50kg)$  // Two way interaction variable

{ create a new feature  $f_5 = 1$  }

Q) we can use any mathematical operator instead of AND/Or we can also have 3-way, 4-way ... n-way interaction.

A) But the question is, given a task how to find good interaction feature.

Sol) Again we use Decision Tree; Ex:-  $h, w, hl \rightarrow \text{gender}$ . train a decision tree using these 3 features and for every leaf node of DT. we can create a interaction variable.

\* **Mathematical transforms**

Suppose we have a feature  $x$ , we can apply  $\log(x)$ ,  $e^x$ ,  $\sqrt{x}$

$f_3 = x^2 + y^2 - \sin(x) \cos(y) \tan(z)$

The question what's the best transform? Ans:- It is very very problem specific.

\* Model specific featurization: some type of featurization works on some model, some other type of " works better on one other model. Hence in that case we need model specific featurization. It depends on the model we are using.

\* Feature orthogonality:-

The more different features are, the better would our models be.

Let's we have 3 feature  $f_1, f_2, f_3 \rightarrow y$

If  $f_1$  ~~correlates~~  $f_2$ ,  $f_2$  ~~correlates~~  $f_3$ ,  $f_3$  ~~correlates~~  $f_1$ ; if features are highly correlated among themselves, the overall impact will be less.

Can we add a feature  $f_4$  which is correlated with  $y$  but not correlated with  $f_1, f_2, f_3$ ?

Ideas:- (1) build a model on  $f_1, f_2, f_3 \rightarrow y$

$M := \hat{y}_0$  (prediction by model)

(2) calculate errors:  $y_i - \hat{y}_0 = e_i^0$ , now while designing  $f_4$ , try to correlate  $f_4$  with  $e_i^0$  (error)

This method is very much similar to GBDT. (Boosting)

\* Domain specific featurization:- It is always upto research and study existing featurization that are designed by domain experts. Every domain is different and hence featurizing using domain knowledge is very important.

\* Feature slicing:- Let's assume we want to predict gender.

With  $h, w, \text{hair}, \text{eyecolor}, \text{Country} \rightarrow \text{gender}$

Like INDIA, USA

Suppose 80% of data is of Indians & 20% is of USA. If we try to build the model on this imbalanced data. Model tends to perform better on people from India because more data is available for Indians. In that case slice your data using country feature and train different models for each category.

D  $\rightarrow$  INDIA  $\rightarrow$  Model 1

D  $\rightarrow$  DUSA  $\rightarrow$  Model 2

The only criteria for feature slicing

- (i) Each category should behave differently,
- (ii) Each category should have enough data points.

#### CH- MISCELLANEOUS TOPICS

\* Calibration of Model:- Calibration of models means getting the probability of the output with itself. For example, suppose we have a 2-class classification model and it give output for  $x_2$ ,  $y_2=1$  or  $y_2=0$  but calibration means we need the probability of  $P(y_2=1/x_2)$  or  $P(y_2=0/x_2)$ .

So why do we need probability, probability is needed  
in log loss calculation.

\* Calibration plots:-  $f \rightarrow D_n$   $\xrightarrow{\text{model trained on } D_n}$

$$D_{CV} = \{(x_i, y_i)\}$$

$\hookrightarrow$  cross-validation data.

for each  $x_i$  in  $D_{CV}$ , send it to  $f$  and get  $\hat{y}_i$ :

$x_i, \hat{y}_i, y_i \rightarrow$  True class label.

$\hookrightarrow$  model predicted value

Step 1: Build the table.

Step 2: sort the table in increasing order of  $\hat{y}_i$

$x_i$	$\hat{y}_i$	$y_i$
$x_1$	$\hat{y}_1$	$y_1$
$x_2$	$\hat{y}_2$	$y_2$
$\vdots$	$\vdots$	$\vdots$
$x_n$	$\hat{y}_n$	$y_n$
		$y$

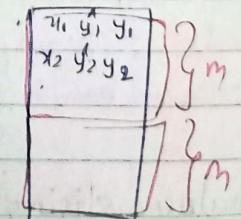
Step 3: break the table into pieces of size  $m$

Step 4: calculate average  $y_i \times \hat{y}_i$  in each piece.

$\hat{y}^{(i)} = \text{average of } i\text{th chunk}$

$y^{(i)} = \dots \dots \dots$

# of pieces =  $\frac{\text{size in } D_{CV}}{m}$

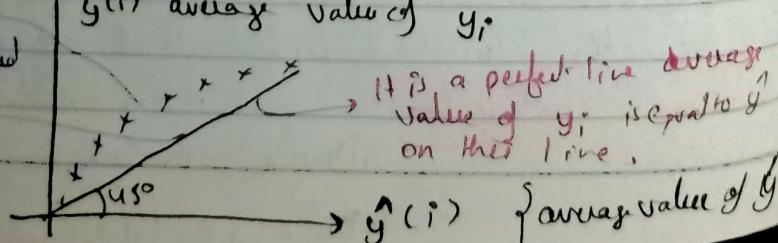


$$D_{ratio} = \{(\hat{y}^{(i)}, y^{(i)})\}$$

Using this calibration data we can plot calibration plots.

$y^{(i)}$  average value of  $y_i$

This is like this  
This is like unidirectional  
line. line a.



$\rightarrow$  It is a perfect line average  
Value of  $y_i$  is equal to  $\hat{y}^{(i)}$   
on this line.

$\rightarrow$   $\hat{y}^{(i)}$  average value of  $y$

Now the whole challenge of calibrations model comes down to we are given  $\hat{y}_i$ . predict  $y_i$ .

Platt Scaling / Sigmoidal calibration.

Task is : when we have a calibration dataset of  $(\hat{y}_i, y_i)$  and given a  $\hat{y}$ ; we need to predict a  $y$ .

$$\text{Dcalib} = \hat{y}_i, y_i \\ \xrightarrow{\text{Predicted by model}} p(y_i=1|x_i)$$

According to platt scaling,  $p(y_i=1|x_i) = \frac{1}{1 + \exp(-A(\hat{y}_i) + B)}$

Use the Dcalib data to find the perfect value of  $A$  &  $B$ .

\* But platt scaling only works when calibration plot looks like sigmoidal shape.

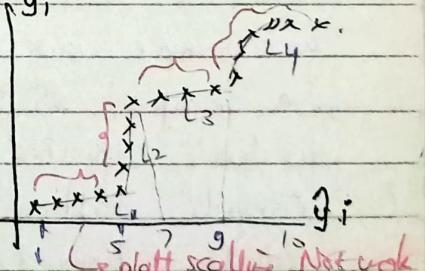
\* Isotonic Regression:- This technique works even if your calibration plot doesn't look like sigmoidal curve.  $y_i$

Concept of isotonic Regression is: Break the plot in small regions and fit a line in each region. These are called piece wise linear model.

If  $\hat{y}_i$  in  $(1, 5)$  we line  $L_1$

If  $\hat{y}_i$  in  $(5, 7)$  we  $L_2$

If  $\hat{y}_i$  in  $(7, 9)$  we  $L_3$



platt scaling Not work  
here bcz plot doesn't look like sigmoid

∴ how do we determine these boundaries and lines.

but we try to minimize a optimization problem

At the same time it also try to take care that we don't have too many lines.

\* Isotonic Regression needs large no. of data point in Dcalib

\* If Dcalib is very small then use simple platt scaling.

## \* Random Sampling Consensus (RANSAC)

It is statistical technique which is widely used in modern ML. If true to ans, can we build a robust model in the presence of outliers.

for simplicity let assume our model is linear regression.

In Lr. Reg we try to fit a line/plane in the data.

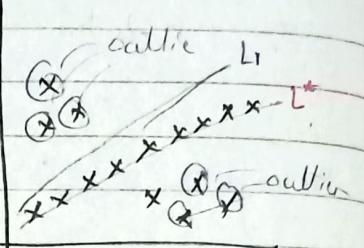
$$D_{\text{Tr}} = \{(x_i, y_i) | x_i \in R, y_i \in R\}$$

If we train a model on  $D_{\text{train}}$  we might end of getting  $L_1$  as our line because Lr. Reg try to minimize the squared loss for each point.

But problem is, some of the points are very far, then the sc-loss will be larger hence they tend to pull the line towards them. Hence  $L_1$  we got is a non-robust because it is impacted by outliers.

RANSAC is the solution.

- (1) Given  $D_{\text{train}}$ , get a  $D_o = \text{Random sample}(D_{\text{train}})$   
Build a Lr. model on  $D_o$ , let us get a line  $L_1$ . When we sample out data, probability of an outlier present in sampled data is very less. Hence we would get a less impacted line due to less No. of outliers.
- (2) Compute outliers based on  $L_1$  (by calculating dist of each points from  $L_1$ , points which have high dist are called outliers). Let us get a set  $O_1 = \{\text{outliers}\}$ .
- (3) Construct a next data set  $D_{\text{Train}}^1 = D_{\text{train}} - O_1$
- (4) Construct a dataset  $D_1 = \text{random sample}(D_{\text{train}}^1)$ , make a model  $L_2$  & find outliers  $O_2 \subseteq \{ \}$   
Repeat the step.  
Let us get a model  $L_1$  in 1st step &  $L_2$  in 2nd step. If they are very much similar then we will stop. Now we make a final model on  $D_{\text{Train}}^{1+}$ .



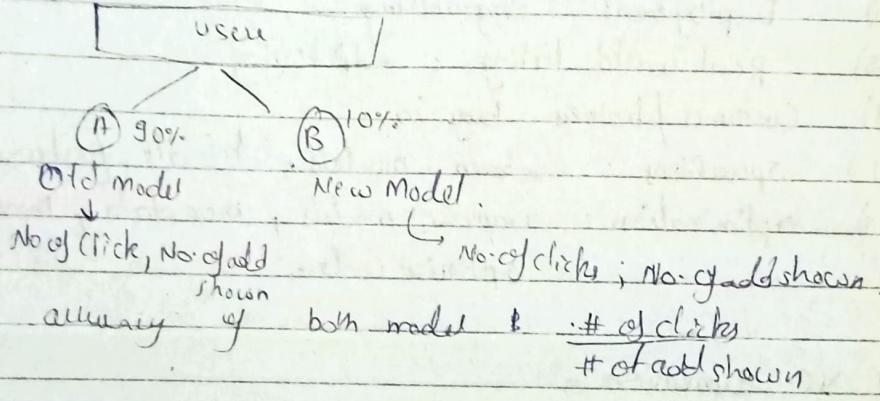
Retaining model periodically:-

A/B Testing :- It is also known as Bucket Testing or split run or controlled experiment.

Let assume we are building a model for Google to predict which ad to show for search query.

So how does A/B testing work here?

Suppose we have a old model and we have millions of user. Divide them in two parts.



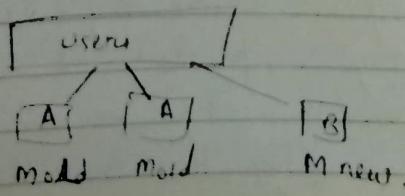
So we get accuracy of both model &  $\frac{\# \text{ of clicks}}{\# \text{ of add shown}}$

If accuracy A < accuracy B

Now we repeat this with ~~so~~ changing the percentage to 30, 50-50 etc. and if we are very sure then only we can deploy model + view.

But one thing to note is, Accuracy may not give correct interpretation. Hence instead of using Accuracy, we can use confidence interval. And if confidence interval of A & B are not overlapping then we can be sure about the result.

There is also a concept called N-N testing:-



## \* Data Science / ML project life-cycles:-

- 1) Understand business requirements :- define the problem, identify customers
- 2) Data acquisition :- ETL, SQL [DB, DW, log files, Hadoop]
- 3) Data preparation :- cleaning, pre processing.
- 4) Exploratory data analysis : plots, viz, hypothesis testing slice & dice data.
- 5) Modeling, Evaluation, & interpretation.
- 6) Communicate results : clear & simple, 1-page, 6-pages.
- 7) Deployment : engineering
- 8) Real world testing : A/B testing
- 9) Customer / business buy-in
- 10) Operations : retain - models, handle failures, process
- 11) Optimization : improve models, more data, more feature, optimize code.

## \* VC-dimension :-

Q) Is there any method to measure how powerful a class of models is

is there any numerical value that can be assigned to model which tells about power of class of model.

→ Linear model  
→ RBF-SVM  
→ Boundary etc.

VC-dimension is one such concept

Vapnik → Chernovnikov

VC-dimension is a theoretical concept, but it is not used much in practical scenario.

Let's take a linear models : line/plane/hyperplane are decision surface ex:- logistic Reg, linear SVM

Imagine we have only 2 points + & -

+ in whichever they are placed, we can always find a way to split the space to separate these points

Now, let's say we have 3 points, even in this case we can separate them using line/hyperplane.

Now what about 4 points? + - we can't  
- +

separate them.

There exist atleast one configuration s.t. ~~where~~ these 4 points cannot be shallowed/separated.

Hence VC dimension of linear model is 3

$$VC\text{-dim}(\text{Linear model}) = 3$$

$$VC\text{-dim}(\text{RBF-SVM}) = \infty$$

Theoretically RBF-SVM are super powerful but practically it is not the case.

## \* Productionizing models

We have the model  $f$ , in the runtime we keep getting  $y_2$  and we need to send  $y_2$  or  $\Pr(y_2=1)$  or  $P(y_2=0)$ .  
Productionizing models means running the model in real world.

There are 2 methods for this

(1) model persistence Using sklearn.

model persistence: store the model in HDFS and when you need it load it to ram and just give query points. We can persist a model in (i) string (ii) file:

sklearn stores model in a pickle file and we can save this .pkl file to disk.

(2) Custom implementation approach:

(a) Store all of the parameters of model to a file.  
(ex. logistic to store weights.)

So, why do we need it, because python is slow we can use the calculated parameters and create a calculator file using C/C++/Java which are faster languages.