

Feature Selection

Why we need Feature Selection

- Curse of dimensionality
 - If we have certain no. and getting optimal soln we increase feature After optimal # model of performance ↓
↓
why?
↓
Sparsity ⇒ on increasing dimension distance b/w two data points increases sparsity
- Computation Complexity
- Interpretability
 - difficult to infer use of ML model.

Techniques (Types) of FC

- Filter based Tech.
 - Variance threshold
 - Correlation based
 - ANOVA
 - chi sq
 - mutual info.
- Wrapper Methods
 - Exhaustive FS
 - Forward selection
 - Backward elimination
 - Recursive feature elimination.
- Embedded methods
 - Lasso, Ridge, Elastic → for linear data
 - DT, RF, GB → for non linear data.
- Hybrid technique
 - It comes in Hybrid

Filter based Feature Selection
Methods that use statistical to score each feature independently, then select a subset of feature based on these scores.

Called as "Filter" methods because they essentially filter out the features that do not meet some criterion.

Steps

- Deleting duplicate features
 - exactly same column with value

1) Variance

1) Variance Threshold

→ Constant: Variance = 0
Ex: B: 1, 1, 1, ...

→ Quasi's Constant
Variance → 0

Ex: B: 1, 1, ..., 999, 555, ..., 5
999 are 1 555 are 5

we will set up a variance threshold (ex: 0.05) if variance of each column < variance which will be dropped

→ before applying this we must firstly normalize or standardize these values. II

If data is standardized and normalized threshold should be in [0.1, 0.01]

When we should not apply variance threshold blindly?

1. ignores Target Variables:-
 • univariate method, evaluates features independently and doesn't consider the relationship b/w each feature that ~~high variance~~ and target variable.
 It may keep irrelevant feature that have high variance but not relationship with target or vice-versa.
2. ignores feature interactions:-
 A feature with low variance become very important (informative) when combined with another feature.
3. Sensitive to data scaling:-
4. Selecting threshold value is big challenge.

3. Correlation

find correlation betⁿ each pair of independent variable & we will discard features based on corr. coeff value is ~~below~~ \pm threshold (e.g. 0.8)

Disadvantages of correlation

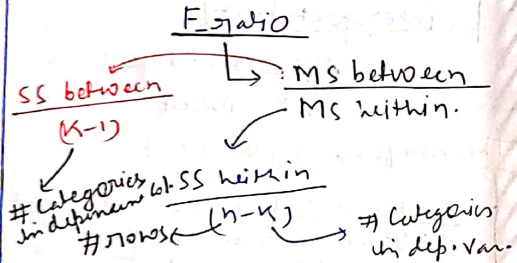
- (1) Linear Assumption
 It does not capture non linear behaviour while removing, that can be misleading.

- (2) Doesn't capture complex relationship:
 L may not capture relationship involving more than two variables.
- (3) threshold determination.
 L depending on feature & dataset.
- (4) Sensitive outliers:-

ANOVA

It applies only when all independent features \rightarrow numerical
 dependent features \rightarrow Catego^{>2}

to apply ANOVA
 L we have to calculate



Col: 1	y	S	L	w	8
1	S	1	2	3	
2	L				
3	w				
4	S				
5	w				

↑
unstacking.

group mean (\bar{x}) = mean of this entire col.
 \bar{x}_S = mean of S class after unstacking
 \bar{x}_L
 \bar{x}_w

$$SS_{within} = (1 - \bar{x}_S)^2 + (4 - \bar{x}_S)^2$$

Subtracting mean of each class from each observation of $\#$ belong to that class

$$(2 - \bar{x}_L)^2$$

we do for each class.

$$(3 - \bar{x}_w)^2 + (1 - \bar{x}_w)^2$$

$$SS_{between} = n_1(\bar{x}_S - \bar{x})^2 + n_2(\bar{x}_L - \bar{x})^2 + n_3(\bar{x}_w - \bar{x})^2$$

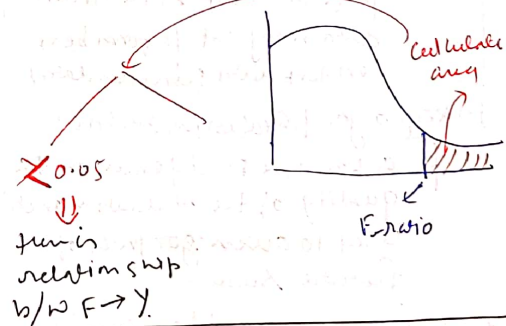
Category mean for that class

grand mean (mean of column)

observation under particular class..

$$\frac{(SS_w / (n-1))}{(SS_B / (K-1))} = F_{ratio} \text{ or } F_{stat}$$

It follows f-dist.



Chi sq
 it is used to find ~~best~~ selected columns based on categorical vs categorical
 i/p cols o/p cols

Let $f_1 \rightarrow Y$
 ex:- is there relationship b/w sex \rightarrow survived?

Step 1
(Contingency table)

	0	1
M	468	109
F	81	283

observed value

Step 2
(Expected value)

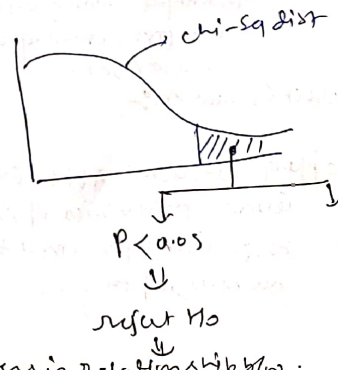
	0	1
M	195	355
F	120	221

$$\frac{(81 + 468) \times (81 + 283)}{\text{Total Sum.}}$$

Similarly we calculate for all cell.

Now we have to find diff b/w (observed value) vs (expected value)

$$\chi^2 = \sum_{i=1}^n \frac{(\text{observed val.} - \text{expected val.})^2}{(\text{expected val.})}$$



Disadvantage

Categorical Data only

- Can only be used with Cate. data
- not suitable for cont. data unless they have been discretized into categories

Independence of observations

- Chi-sq test assumes that the observations are independent of each other.
- This

Disadvantage of Filter based Selection

- we are only selecting features based on $X_i \rightarrow y$ relationship but we are not finding based on $X_1, X_2 \rightarrow y$ or other combinations of features relationship with y

↓ Solution

Wrapper method

It involve using a predictive model to score the combination of features.

Steps involved

- Subset generation → with several technique Adding one by one (or) removing, only one, or so on.
- Subset Evaluation:-
- Stopping criterion: After evaluating certain possibilities of subset, no further improvement then we will stop the process.

Wrapper method

- Exhaustive FS
- Forward Selection
- backward elimination
- Relative feature elimination

Exhaustive FS / best subset selection

- trying all the subsets along with Model development that subset whose accuracy is best, that subset will be selected in FS.

Disadvantage

- N cols → need $(2^n - 1)$ models to be train.

- Computational Complexity

- Risk of overfitting

- Feature combination that performs best on the training data may not perform best on test data (unless data)

- Req. a good evaluation metric.

- Exhaustive FS depends on the quality of the evaluation metric used to assess goodness of a feature subset.

* What is cross validation (CV)?

We split our data 5 times in train a test data set. we find accuracy for each split and at the end we will find avg score.

Sequential Backward Selection / Elimination

$F_1 \ F_2 \ F_3 \ F_4 \rightarrow$ Model

we develop 4 model by removing ~~each~~ each in each combination
model ↓
Select ~~the~~ having highest score.

3 model by removing each in each combination.

↓
select model having

$F_1 \ F_2 \ F_3 \ F_4 \rightarrow$ model $\rightarrow 0.89$

$f_1 \ F_2 \ F_3 \ F_4 \rightarrow$ model $\rightarrow 0.81$

$F_1 \ f_2 \ F_3 \ F_4 \rightarrow$ " " 0.71

$F_1 \ F_2 \ F_3 \ F_4 \rightarrow$ " " 0.91 (remove this feature model)

$F_1 \ F_2 \ F_3 \ f_4 \rightarrow$ Model $\rightarrow 0.65$

$f_1 \ F_2 \ F_4 \rightarrow$ model $\rightarrow 0.79$

$F_1 \ f_2 \ F_4 \rightarrow$ model $\rightarrow 0.81$

$F_1 \ F_2 \ F_4 \rightarrow$ model $\rightarrow 0.83$ (remove)

$F_1 \ F_2 \ F_4 \rightarrow$ model $\rightarrow 0.63$

$F_1 \ f_2 \ F_4 \rightarrow$ model $\rightarrow 0.53$

↓
Step 2

Now we will select best model

$f_1 \ f_2 \ f_3 \ f_4$	0.89
$f_1 \ f_2 \ f_4$	0.91
$f_1 \ f_2$	0.83
f_2	0.63

this is the best

Advantage

- faster
∵ for n columns we require $\frac{n(n+1)}{2}$ models to check
- We may miss best combination because after removing any feature, that feature can't be used in next combining iteration due to which it may loose best combination

Sequential forward selection

How we add in each iteration

Let's we have

F_1 F_2 F_3 F_4

Add iteration 1

f_1 f_2 f_3 f_4
 \downarrow \downarrow \downarrow \downarrow
 0.63 0.51 0.43 0.49

Add iteration 2

$F_1 F_2$ $F_1 F_3$ $F_1 F_4$
 \downarrow \downarrow \downarrow
 0.63 0.71 0.80

Add in iteration 3

$F_1 F_2 F_3$ $F_1 F_2 F_4$
 \downarrow \downarrow
 0.81 0.85

Add in iteration 4

$F_1 F_2 F_3 F_4$ \rightarrow 0.83

Now, select all iteration's best

$f_1 \rightarrow 0.63$
 $f_1 f_2 \rightarrow 0.71$
 $f_1 f_2 f_3 \rightarrow 0.85$ ✓ **best**
 $f_1 f_2 f_3 f_4 \rightarrow 0.83$ **final.**

Advantage

\rightarrow faster vs exhaustive
 $\rightarrow \frac{n(n+1)}{2}$ are models need to train.

Exhaustive $\rightarrow 2^n - 1$ models
 backward $\rightarrow \frac{n(n+1)}{2}$
 forward $\rightarrow \frac{n(n+1)}{2}$

Disadvantage

\rightarrow Here also we may skip best possible combination.
 Ex:- f_1 got selected in 1st iteration & it's getting clubbed f_2, f_3, f_4 then. (f_2, f_4) combination we are not checking.

Feature Selection using Embedded

Disadvantage of

Filter:- Missing mutual feature interaction which selecting features based on perform.

Wrapper:- slow

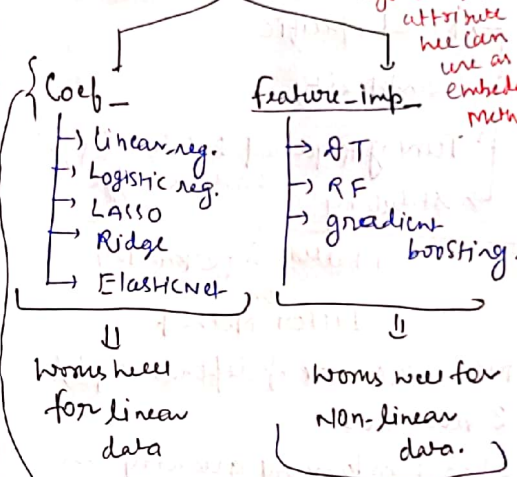
& missing some best combination.

\downarrow slow.

Embedded

works as while training only it gives feature importance to train that model. due to which it's faster as well as taking consideration of all possible combinations.

ML Algo



Ex:-

$$LPA = \beta_0 + (\beta_1) x y p a + (\beta_2) i q$$

there Co-eff only tells that by increasing their respective factor, how many units target variable will increase.

so, somewhat Co-eff of each feature tells their importance.

\downarrow

But it give static picture of feature importance of all assumption of suits very well.

- \rightarrow linearity
- \rightarrow independence
- \rightarrow Homoscedasticity
- \rightarrow Normality
- \rightarrow No Multicollinearity.

FS on Regularized Model

Regularized linear model are linear models that include a penalty term in the loss function during training.

The penalty term discourage the learning of a too complex model, which can help to prevent overfitting.

\downarrow

Ridge Lasso ElasticNet

\downarrow
 best to use for F.S.

\downarrow
 It keeps imp. feature's Co-eff as "Non-zero" when as many Co-eff. of other feature as "zero".

\rightarrow In all these algo there is an attribute called feature importance.

\rightarrow that will give importance of each feature under that Model.

Wrapper method's Recursive feature elimination

its kind of hybrid technique.

Wrapper embedded

$f_1 f_2 f_3 f_4 \rightarrow$ Model

remove it based on weak feature importance

$f_1 (f_2) f_3 \rightarrow$ Model

$(f_1) f_3 \rightarrow$ feature-imp

\therefore this way it is best feature.

We can do this same process to apply on rest feature to find 2nd best feature.

Mutual info Filter method

Advantages of Embedded technique

- Performance: they are generally more accurate than filter methods since they take the interactions b/w features into account.
- Efficiency: ↑ than wrapper \therefore they fit model only once
- Less prone to Overfitting: due to regularization, less prone to overfitting.

Disadvantage of Embedded

- Model Specific
- Complexity
- Tuning Required: Value of λ in Lasso & ridge
- Stability

Mutual information Filter Method.

MI is Measure of dependency b/w 2 variables

it is Fundamental quantity in information theory

$$MI = \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \frac{P(x, y)}{P(x) P(y)}$$

where,

$P(x, y) \rightarrow$ Joint Prob of x & y

$P(x) \rightarrow$ marginal Prob of x

$P(y) \rightarrow$ marginal Prob of y

$P(x, y)^{M=0} \rightarrow$ Probability of $y=0$ being $x=M$ or other diff cases based on value of x & y

X	Y
M	0
F	1
M	0
F	0
M	1

Cross table

	0	1
M	$\frac{2}{5}$	$\frac{1}{5}$
F	$\frac{1}{5}$	$\frac{1}{5}$
	$\frac{3}{5}$	$\frac{2}{5}$

$$P(x=M, y=0) = 2/5$$

Joint probability for this Condition.

to find MI sex we calculate each for each cells.

$$\begin{aligned} & \frac{2}{5} \log \left(\frac{2/5}{3/5 \times 2/5} \right) + \frac{1}{5} \log \left(\frac{1/5}{3/5 \times 2/5} \right) \\ & + \frac{1}{5} \log \left(\frac{1/5}{2/5 \times 2/5} \right) + \frac{1}{5} \log \left(\frac{1/5}{2/5 \times 2/5} \right) \end{aligned}$$

If $MI(\text{sex}) \uparrow$ then it higher important feature

Few points about MI

- It is Non-negative; always zero or +ve. with zero indicating that the variables are independent
- It is symmetric: $MI(X, Y) = MI(Y, X)$
- It can capture only kind of statistical dependency: Unlike correlation, which only capture linear relationship, it can capture any kind of relationship, include non linear ones.

How to deal with Numerical Variables?

→ unlike chi-sq. it is able to work on numerical data.

Ex:- age survived

35	0
60	1
52	1
16	1
27	0

generally it will convert into discrete or bins.

Disadvantage

1. Estimation Difficulty:- When the dimensionality of the data is high or the number of sample is low
2. Assumes Large Sample Size MI works best with large sample size
3. Computationally Intensive:-
4. Difficulty with Continuous Variables
5. No direct indication of the Nature of Relationships
6. Doesn't Account for Redundancy