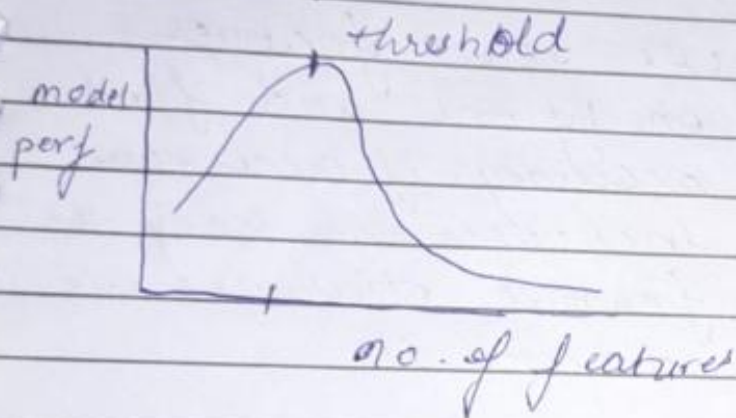# Feature Selection.

Choosing subset of features from all features.

## Why we need feature selection?

# Curse of Dimensionality :- If we have lot of features in our dataset, if we further increase no of features we will end up hampering our model growth Model will be learning from irrelevant data.

threshold



model perf

no. of feature

Bad features :-
1. irrelevant
2. redundant

Curse of Dimn
↓
dimensionality red$^n$

feature selection
a, b, c ... z / outp
5 best features
e, g, z, v, f → o/p

feature extra
a, b, c ... z
5 feature
az, ecq,
dm, a, e
↳ PCA

# Feature Selection Methods :-

① Filter Method

we have to filter out reqd features from rest of the features.

checks relevance of feature with output variable.

All features $\longrightarrow$ Select best subset $\longrightarrow$ Algo - ML

$\downarrow$

CHI square test
ANOVA test
correlation coeff

# �##② Wrapper Method
└ forward selection    └ backward elimination

ex Forward selection

A $\rightarrow$ ML $\rightarrow$ accuracy
model

AB $\rightarrow$ ML $\rightarrow$ accuracy ✓
↑ model

ABC $\rightarrow$ ml $\rightarrow$ accuracy ✗
✗ model

We select features one by one and find accuracy, if accuracy increases we keep the feature otherwise we discard it.

Backward Elimination
ABCDE          Test
$\downarrow$          { ANOVA
model          { CHI sq
               { Pearson coeff

Chi-squared $\rightarrow$ p-value
$P \leq 0.05 \rightarrow$ useful
$\downarrow$

③ Embedded Method
It learns the feature while building a model
Takes all PMC of features and check with model performs best.