

Feature Selection and Engineering

[SHALA2020.github.io](https://shala2020.github.io)

Learning objectives

- List advantages and disadvantages of feature selection and engineering
- List basic feature selection method types
- Write the formulae for a few feature selection metrics, e.g. t-test, AIC
- Write the steps for a few feature selection procedures, e.g. FS, BE
- Write the objective of LASSO and elastic-net
- Write the steps for PCA
- List some common features used for images, speech, and text

Why select features?

- Reduce overfitting
- Reduce confusion
- Reduce collinearity
- Reduce training time
- Simplify interpretation

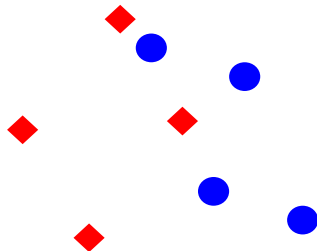
Disadvantages

- Inadvertently throw away useful information, because all selection methods have their own assumptions and biases

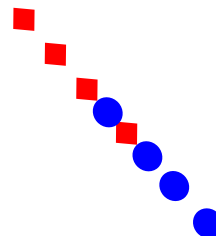
Compare two features



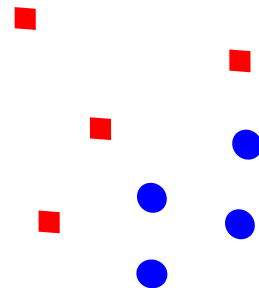
One of the
dimensions may
be useless



One of the
dimensions may
be useless

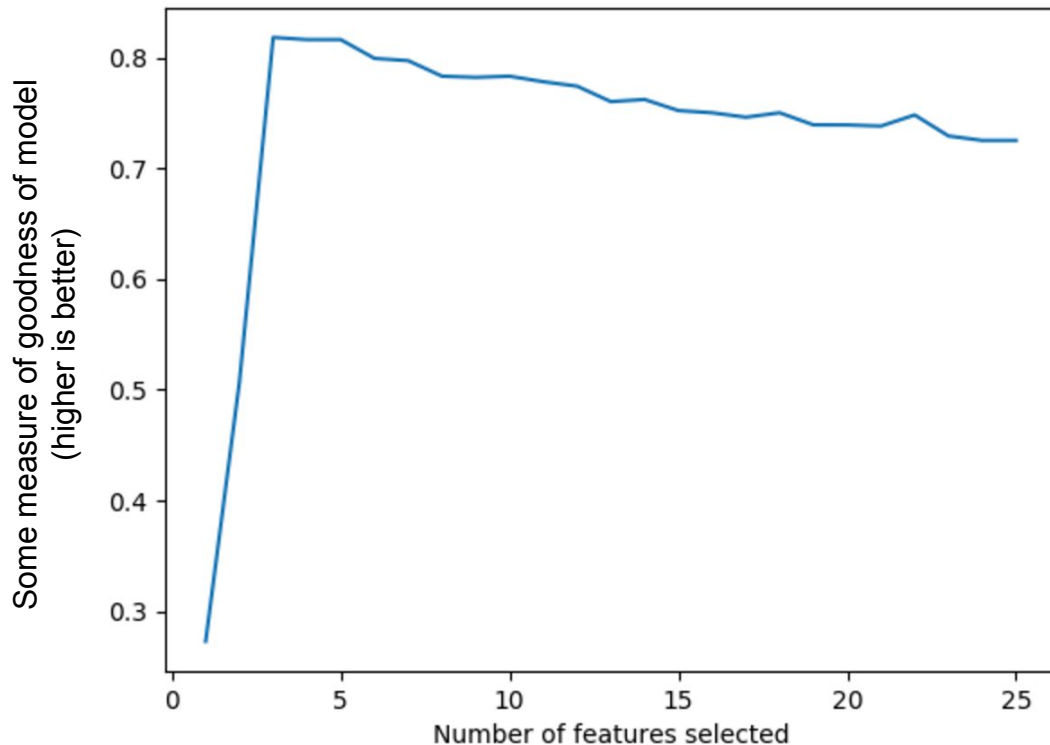


Rotating the
axes may be
useful



Rotating the
axes may be
useful

Usually, there is an optimal number (and subset) of features



Filtering

Basic idea: Remove “useless” or “redundant” features

- What makes a feature useful?
- What makes it non-redundant?

Algorithm:

- Loop through features x_i
 - Compute measure of utility and non-redundancy $m(x_i)$
- Sort features based on measure m
- Pick the top- k features

Filtering measures for regression

- Pearson correlation

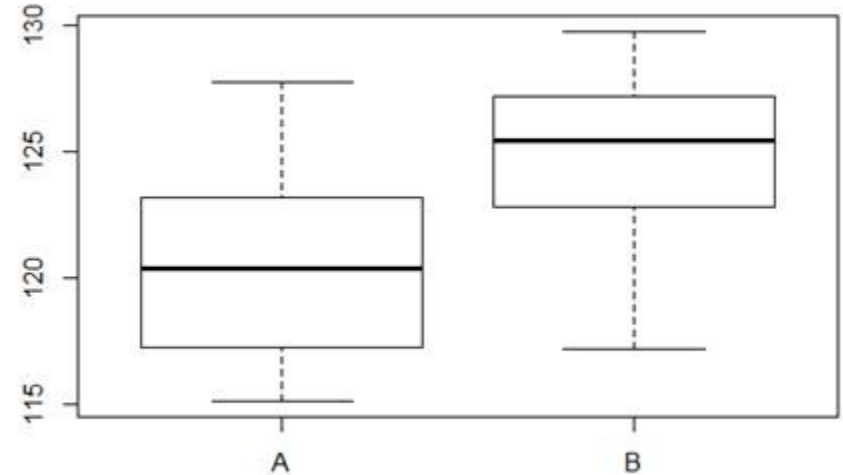
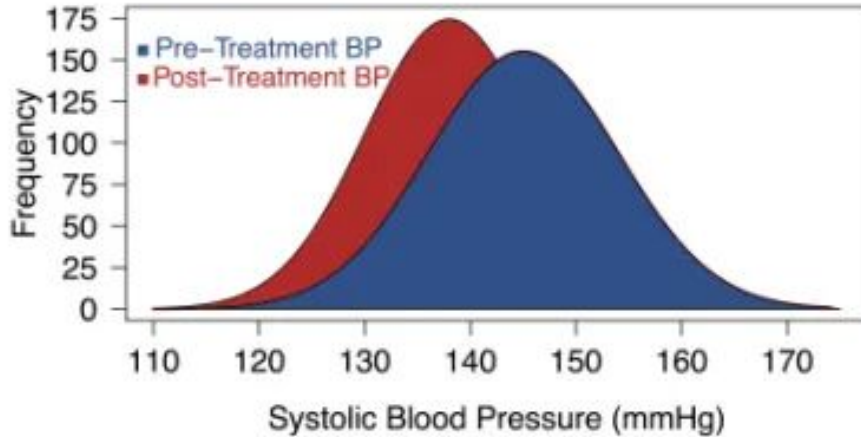
$$\frac{Cov(x, y)}{\sqrt{Var(x)}\sqrt{Var(y)}}$$

- Spearman correlation

$$\frac{Cov(rank(x), rank(y))}{\sqrt{Var(rank(x))}\sqrt{Var(rank(y))}}$$

Filtering measures for classification

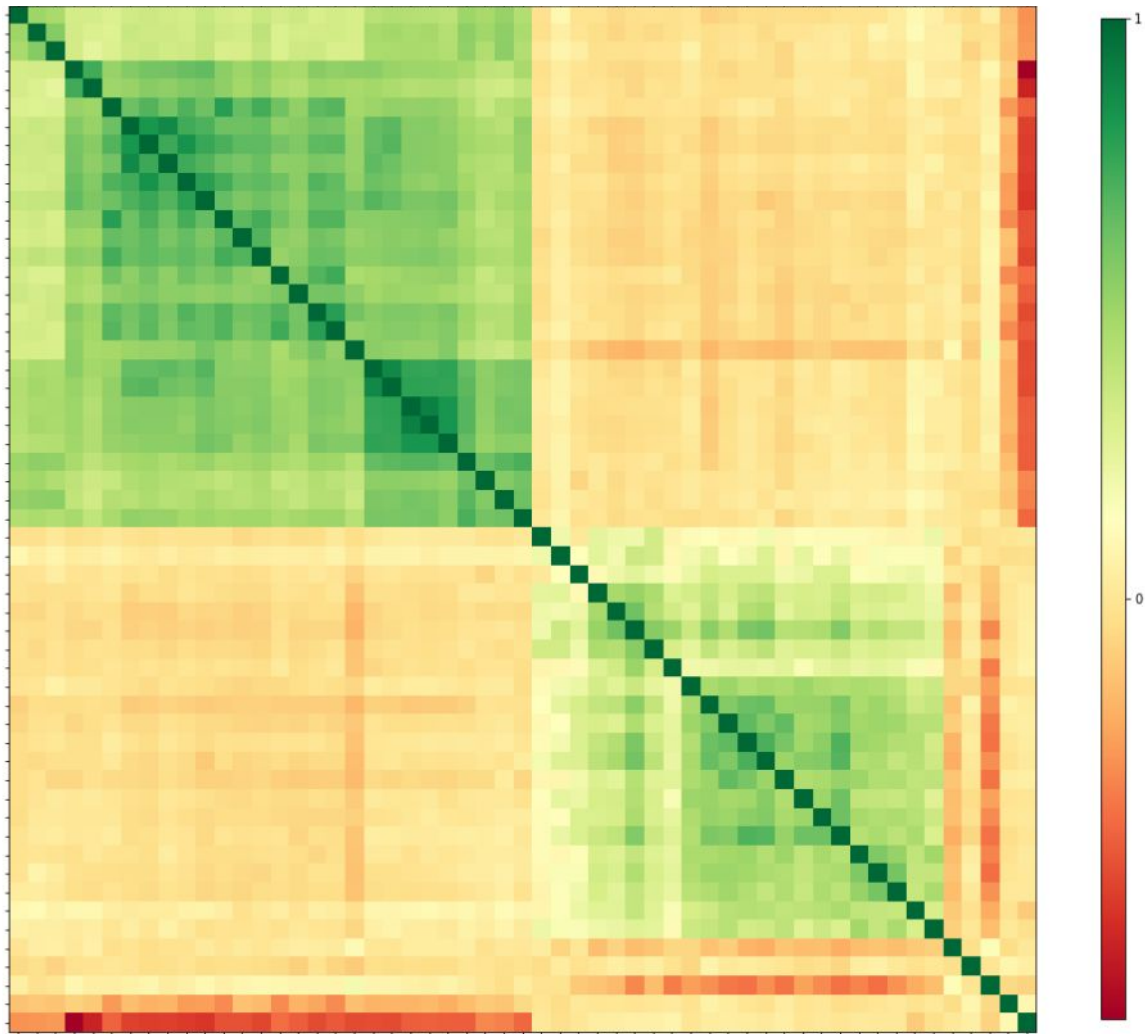
- Hypothesis testing such as t-test, Wilcoxon rank test etc.



Correlation-based clustering removes redundant features

One can choose a representative variable of each cluster

The selected variable can be based on interpretability



Wrapper methods

Basic idea:

- Add or remove features
- And measure model accuracy

Model selection criteria:

- K-fold CV
- AIC
- BIC

K-fold cross-validation

	Held-out		Training		
Result 1	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5
Result 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5
Result 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5
Result 4	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5
Result 5	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5
Overall result					

AIC and BIC

AIC - Akaike Information Criterion

- $2k - 2 \ln(L)$; where k is number of parameters, L is likelihood

BIC - Bayesian Information Criterion

- $k \ln(n) - 2 \ln(L)$; where k is number of parameters, n is number of samples, L is likelihood

Forward selection

- Start with no variables in set S and all variables in set A
- For $i = 1$ to d
 - for $j = 1$ to $d-i+1$
 - Compute a measure m of adding variable $x_{A(j)}$ from A to the model
 - Select the variable with the best measure
 - If the change in measure m meets some criteria
 - Remove it from A and put it in S
 - Else exit

Backward elimination (also RFE)

- Start with all variables in set S and no variables in set A
- For $i = 1$ to d
 - for $j = 1$ to d
 - Compute a measure m of removing variable $x_{S(j)}$ from S to the model
 - Select the variable with the best change in measure
 - If the change in measure m meets some criteria
 - Remove it from S and put it in A
 - Else exit

Forward selection and backward elimination compared

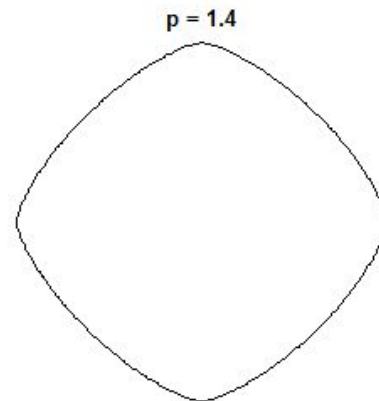
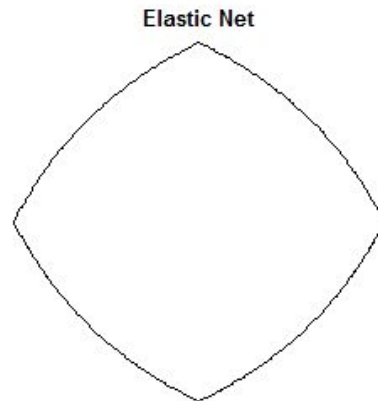
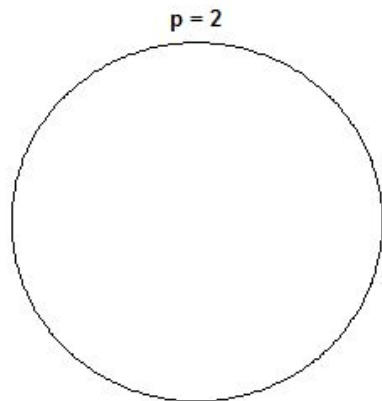
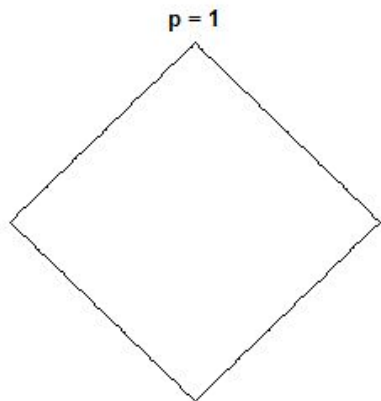
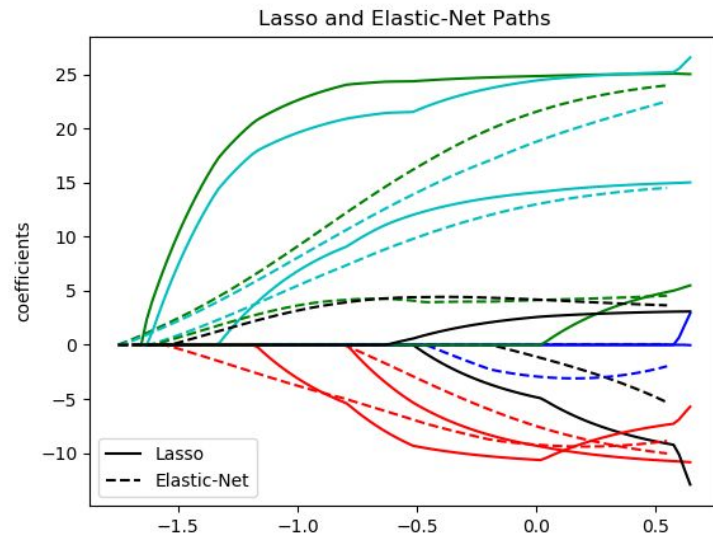
- Both are greedy methods
- In FS two variables individually may be uninformative and thus not considered, but together may be informative
- In BE we may start with a very complicated model initially itself

Regularization for feature selection

- LASSO (L1 penalty over weights)

$$L_p \text{ norm is } (\sum_i |w_i|^p)^{1/p}$$

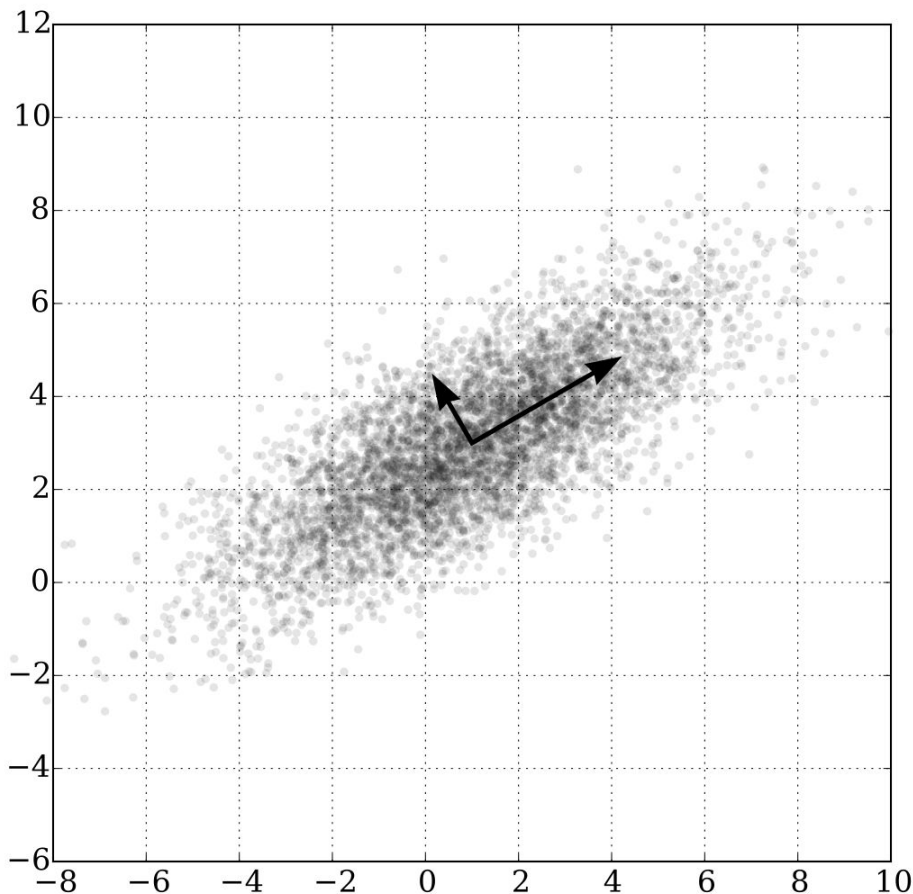
- Elastic-net (L1+L2 penalty over weights)



Unsupervised feature reduction using PCA

Principal component analysis

- Start with D -dimensional data
- Compute covariance matrix Σ
- Eigen decompose $\Sigma = U\Lambda U^T$
 - U contains eigenvectors (principal directions) and Λ is a diagonal matrix of eigenvalues
- Select $d < D$ directions (eigenvalues) with the highest eigenvalues
- Project data to those d dim.s



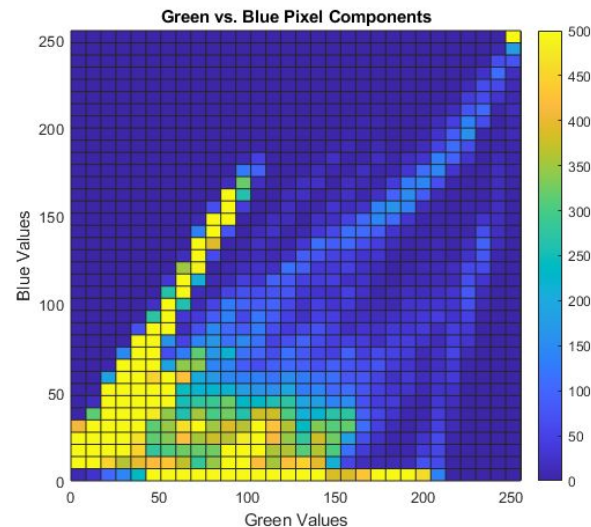
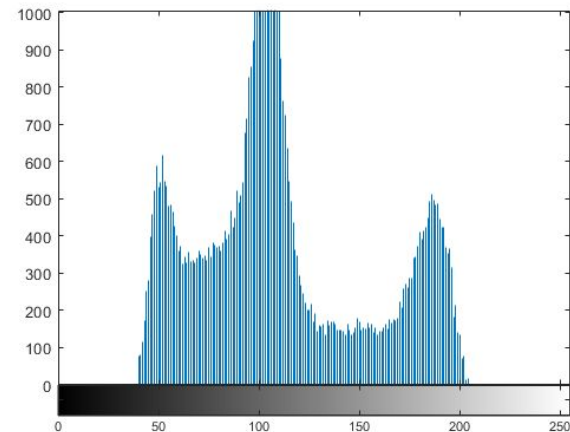
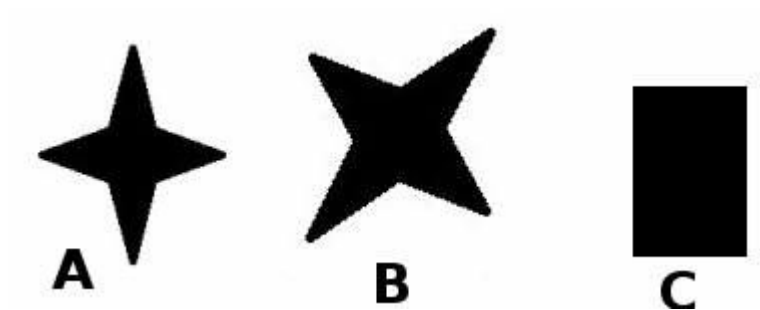
How to engineer features?

Features based on domain knowledge and statistics for:

- Images
- Audio
- Text

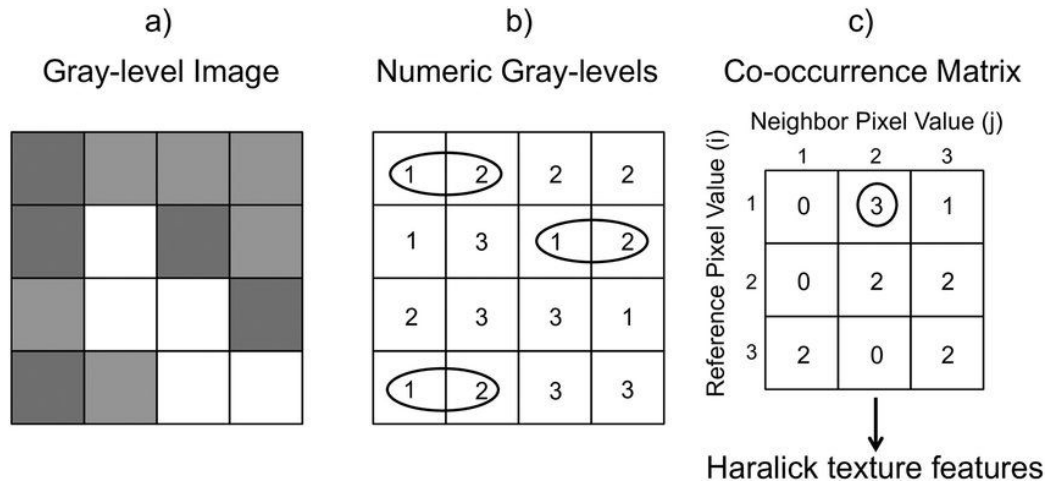
Features for images

- Intensity histogram
 - Its mean, median, mode, skew, kurtosis
- Color histogram
 - Bivariate histogram
- Hu invariant moments



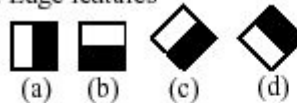
Features for images

GLCM

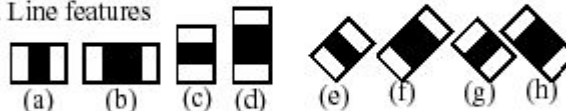


Haar features

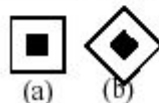
1. Edge features



2. Line features

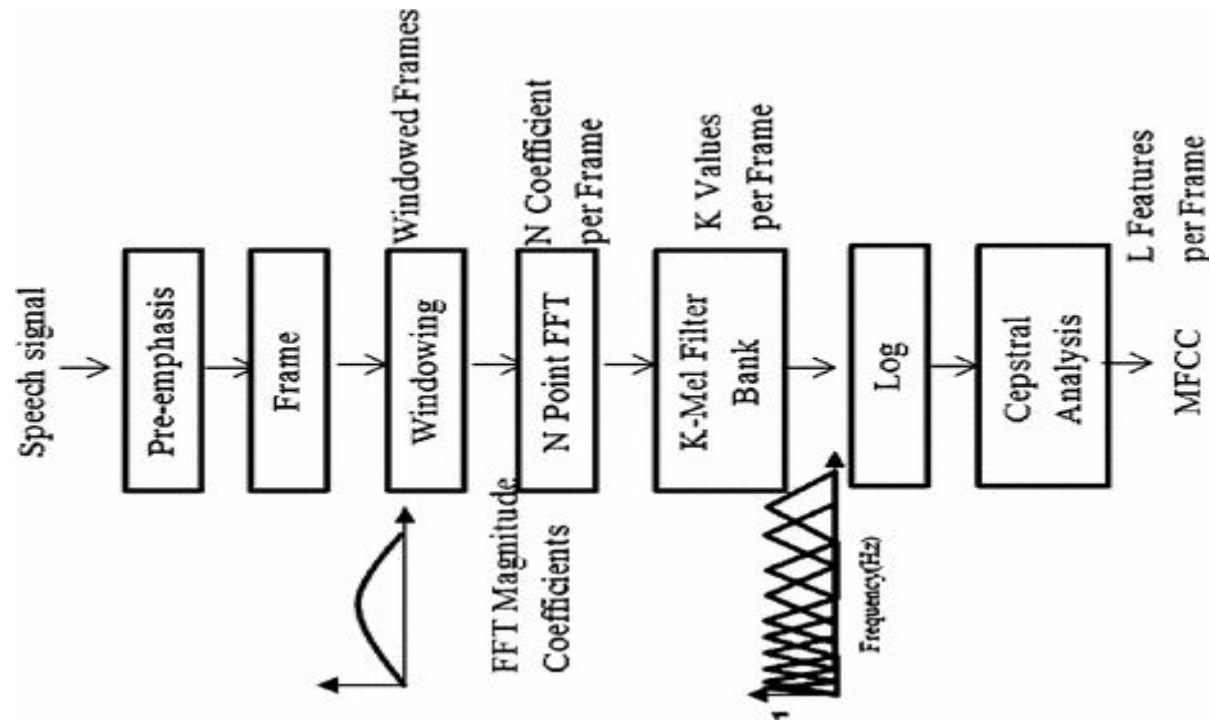


3. Center-surround features



Features for audio

MFCC features



Features for text

- Bag-of-words and Frequency of n -grams

- “India posted a score of 256/8 in their allotted 50 overs in the third and deciding ODI of the series. Virat Kohli was the top-scorer for men in blue with a classy 71, while Adil Rashid and David Willey picked up three wickets each”

- Counts
 - $\begin{matrix} \text{India} \\ \text{Posted} \\ \text{Score} \end{matrix} \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$
- The counts can be normalized
- The words can be standardized
 - Score
 - Scorer
- What about uninformative words?

- Term frequency – inverse document frequency
- TF $f_{t,d}$ is the count of term t in document d
 - Usually normalized in some sense
 - $\text{tf}(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$

- TF-IDF

- IDF penalizes terms that occur often in all documents, e.g. “the”
 - $\text{idf}(t, D) = \log \frac{|D|}{1 + |\{d \in D : t \in d\}|}$
- TF-IDF is $\text{tf}(t, d) \times \text{idf}(t, D)$
- Form a vector of TF-IDF for various terms
 - Which terms?