# Types of Feature Selection

### Feature Selection

```
                    Feature Selection
        ┌──────────┬──────────┬──────────┐
        ↓          ↓          ↓          ↓
   Filter based  Wrapper   Embedded    Hybrid
   technique     method    technique   technique
```
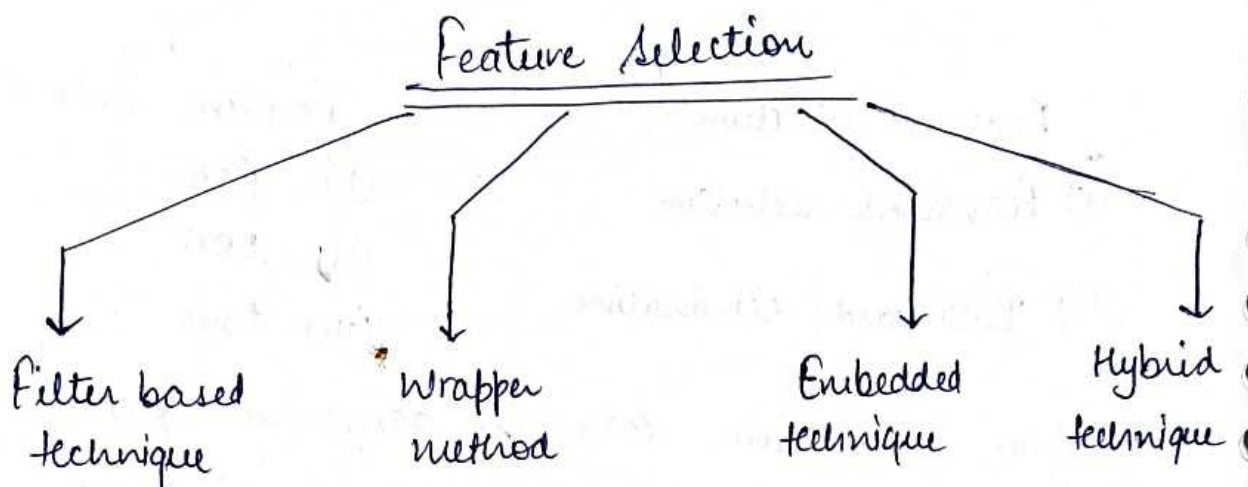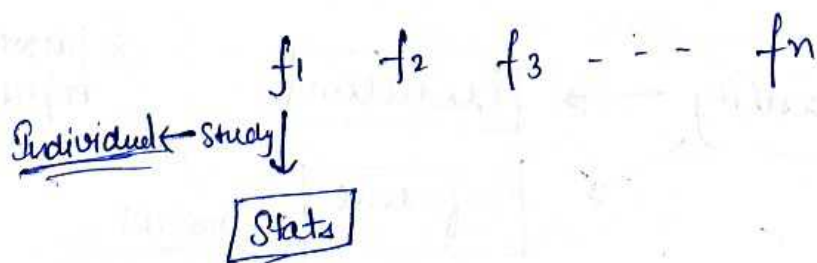
1.) **Filter based feature selection** :- Filter based feature selection techniques are methods that use <u>statistical measure</u> to <u>score each feature independently</u>, and then <u>select a subset of features based on these scores</u>. These methods are called "filter" methods because they essentially filter out the features that do not meet some criterion.

$$f_1 \quad f_2 \quad f_3 \quad - - - \quad f_n$$

Individual ← study ↓

[ Stats ]

## Techniques

- Variance Threshold
- Correlation

- Anova
- Chi square

- Mutual info

# 1. Check Duplicate Features

| $f_1$ | $f_2$ | $f_3$ | $f_4$ |
|-------|-------|-------|-------|
| 1 | 2 | 1 | 2 |
| 2 | 1 | 2 | 2 |
| 3 | 3 | 3 | 2 |

$f_1$ and $f_3$ → Same

dict → keep key →
    duplicate values

keys → original cols

values → dublicate

# 2. Variance Threshold

→ Constant

→ quasi constant

## Constant

— Drop

| A | B | C |
|---|---|---|
| 1 | 1 | X |
| 2 | 1 | N |
| 3 | 1 | X |

O variance → B column constant → always 1

B ————•———— Variance = 0

A —•———•———•— variance →

## quasi constant

Total 1000 rows

↳ 995 rows = 1
↳ 5 rows = 0

} variance ≈ 0 (quasi constant feature)

1) threshold = 0.1

2) fi → var  (check variance of every columns)

3) Drop those ∧variable which is less than

   feature variance

   threshold.

## Points to Consider

low variance ← $F_i$ → relate to y

$f_i$ → Y but not relate with Y

variance ↓ high

1. __Ignore Target Variable__ :- Variance threshold is
   individual a __univariate method__,
   meaning it __evaluates each feature independently__
   and doesn't consider the relationship between
   each feature and the target variable. This
   means it may keep __irrelevant features that
   have a high variance but no relationship
   with the target__, or discard potentially
   useful feature that have a low variance
   but a strong relationship with target.

2. __Ignores Features Interactions__ : Variance
   Threshold doesn't account for interaction between
   features. A features with a low variance
   may become very informative when combined
   with another feature.

   $f_1$ and $f_2$ relation↑↑
   variance↓ ← feature↑↑ → variance↑

3) <u>Sensitive to Data Scaling</u> : Variance Threshold
is sensitive to the <u>scale of the data</u>. If
features are <u>not on the same scale</u>, the
variance will <u>naturally be higher for features
with large values</u>. Therefore, it is important
to standardize the feature before applying
variance threshold. $10000 \leftarrow$ th $\qquad$ $f_2 \rightarrow 0.2$
<div style="text-align:center">minor chage $\rightarrow$ var$\uparrow$</div>

4) <u>Arbitrary threshold Value</u> :- It's up to the
user to define what constitutes a " low"
variance. The threshold is not always
easy to defined and the optimal value
can vary between datasets.

    0.1    or    0.01     or    0.2     is better