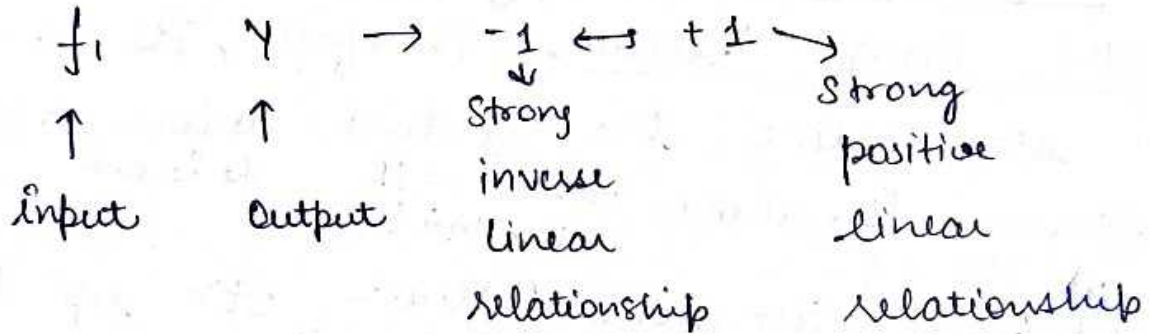


Correlation

Pearson Corr Coeff



f_1 f_2 f_3 - - - f_n y

$f_1 \rightarrow y$

$f_2 \rightarrow y$

$f_3 \rightarrow y$

\vdots

$f_n \rightarrow y$

Corr off $\rightarrow 0.3$

-0.3

\downarrow

feature X

$f_1 - f_2$

$f_1 - f_3$

$f_1 - f_4$

multicollinear \rightarrow

$f_1 - \textcircled{f_2} \rightarrow 0.9$ \rightarrow drop

$f_1 - \textcircled{f_3} \rightarrow 0.85$

Disadvantage

1. Linearity Assumption : Correlation measures the linear relationship between two variables. It does not capture non-linear relationships. If a relationship is non-linear, the correlation coefficient can be misleading.
2. Doesn't capture Complex Relationship :- Correlation only measures the relationship between two variables at a time. It may not capture complex relationships involving more than two variables. $x_1 \rightarrow x_2$ $x_1 \ x_2 \ x_3$ $\rightarrow Y$
3. Threshold Determination : Just like variance threshold, defines what level of correlation is considered "high" can be subjective and may vary depending on the specific problem or data set.
 $0.95 \rightarrow 0.9 \rightarrow 0.8$
4. Sensitive to Outliers :- Correlation is sensitive to outliers. A few extreme values can skew the correlation coefficient.

4. Anova

1. numerical \rightarrow Y
 f_1 categorical
 \downarrow
 more than 2 classes

2. numerical \rightarrow Y
 f_1 \downarrow
 numerical

Source of variation	Sum of squares (SS)	Degree of freedom (d.f)		
Between samples or categories	$n_1(\bar{x}_1 - \bar{\bar{x}})^2 + \dots + n_k(\bar{x}_k - \bar{\bar{x}})^2$ $(k-1)$		$\frac{SS_{\text{between}}}{(k-1)}$	$\frac{MS_{\text{betn}}}{MS_{\text{within}}}$
within samples or categories	$\sum (x_{i1} - \bar{x}_1)^2 + \dots + \sum (x_{ik} - \bar{x}_k)^2$ $i = 1, 2, 3, \dots$ $(n-k)$		$\frac{SS_{\text{within}}}{(n-k)}$	
Total	$\sum (x_{ij} - \bar{\bar{x}})^2$ $i, j = 1, 2, 3$ $(n-1)$			

$n = \text{rows}$

$k = \text{categories}$

f -statistic \rightarrow p -value \rightarrow $f \leftrightarrow Y$

f_i	Y
1	S
2	L
3	W
4	S
1	
1	

\bar{X}

grand mean

\Rightarrow

S	L	W
1	2	3
4	5	7
3	-6	-8
\bar{X}_S	\bar{X}_L	\bar{X}_W

SSw

$$(1 - \bar{X}_S)^2 + (4 - \bar{X}_S)^2 + (3 - \bar{X}_S)^2 + (2 - \bar{X}_L)^2 + (5 - \bar{X}_L)^2 + (6 - \bar{X}_L)^2 + (3 - \bar{X}_W)^2 + (7 - \bar{X}_W)^2 + (8 - \bar{X}_W)^2$$

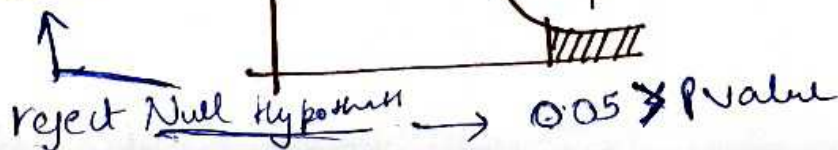
after calculate SSw = 50

$$SSb = 1(\bar{X}_S - \bar{X})^2 + 4(\bar{X}_S - \bar{X})^2 + 3(\bar{X}_S - \bar{X})^2 + 2(\bar{X}_L - \bar{X})^2 + 5(\bar{X}_L - \bar{X})^2 + 6(\bar{X}_L - \bar{X})^2$$

$$\frac{\frac{SSw}{k-1}}{\frac{SSbetw}{n-k}} = F\text{-ratio}$$

dist \rightarrow f-dist

$f_i \rightarrow Y \Rightarrow$ No relation



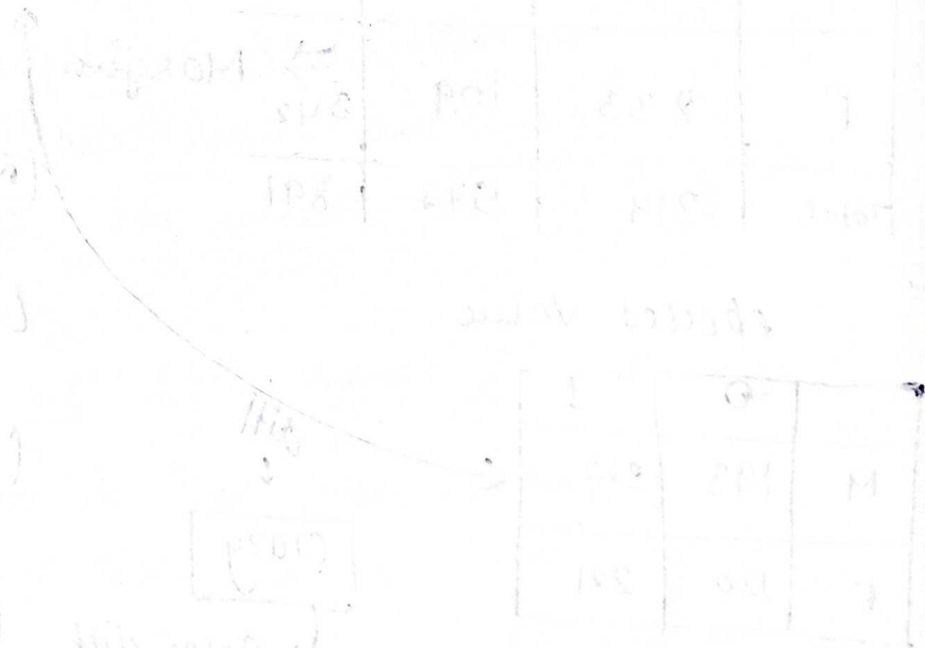
Disadvantage

1. Assumption of Normality: ANOVA assumes that the data for each group follow a normal distribution. This assumption may not hold true for all datasets, especially those with skewed distribution.
2. Assumption of Homogeneity of Variance: ANOVA assumes that the variances of the different groups are equal. This is the assumption of homogeneity of variance (also known as homoscedasticity). If this assumption is violated, it may lead to incorrect results.
- 3) Independence of Observations: ANOVA assumes that the observations are independent of each other. This might not be the case in datasets where observations are related (eg, time series, nested data).

4) Effect of Outliers :- ANOVA is sensitive to outliers. A single outlier can significantly affect the F -statistic leading to a potentially erroneous conclusion.

5) Doesn't Account for Interactions: Just like other univariate feature selection methods, ANOVA does not consider interactions betⁿ features.

f_1 longitude
 f_2 latitude



Chi-Square

f_1
 \downarrow
 category

Y
 \downarrow
 category

$f_1 \rightarrow 89$ relation \rightarrow form Contingency table \rightarrow frequency table
 \hookrightarrow Chi-square

observed value

Expected Value ideal

	0	1	
M	81	468	549
F	233	109	342
Total	314	577	891

	0	1
M	468	109
F	81	233

expected value

	0	1
M	193	355
F	120	221

diff

crazy

\hookrightarrow large diff

\hookrightarrow $[f_i \rightarrow Y] \rightarrow$ change in f_i and find variance Y

$$(M, 0) = \frac{314 \times 549}{891}$$

$$(F, 0) = \frac{314 \times 342}{891}$$

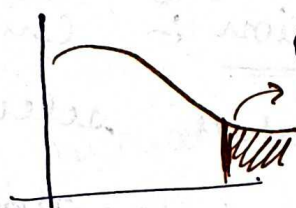
$$(F, 1) = \frac{342 \times 577}{891}$$

$$(M, 1) = \frac{577 \times 549}{891}$$

$$\text{Chi-squared stat} = \left(\frac{\text{Observed value} - \text{expected value}}{\text{expected value}} \right)^2$$

$$= \frac{(468-193)^2}{193} + \frac{(109-355)^2}{355} + \frac{(81-120)^2}{120} + \frac{(233-221)^2}{120}$$

let
Chi-squared = 15
stat



P-value \downarrow = Crazy \uparrow \uparrow
y diff

Disadvantage

1. Categorical Data Only :- The Chi-square test can only be used with categorical variables. It is not suitable for continuous variables unless they have been discretized into categories, which can lead to loss of information.
2. Independence of Observations: The Chi-square test assumes that the observations are independent of each other. This might not be the case in datasets where observations are related (eg. time series data, nested data)

3. Sufficient Sample size :- Chi-square test requires a sufficiently large sample size. The result may not be reliable. If the sample size is too small or if the frequency count in any category is too low (typically less than 5)

4. No Variable interaction :- Chi-square test, like other univariate feature selection methods, does not consider interaction between features. It might miss out identifying important feature that are significant in combination with other features.

f_1	f_2
-------	-------

$f_1 \rightarrow y$

Advantage and disadvantage

Advantage

1. Simplicity: filter methods are generally straight-forward and easy feature and select the top features based on this statistics.
2. Speed: These methods are usually computationally efficient. Because they evaluate each feature independently, they can be much faster than wrapper methods or embedded methods, which need to train a model to evaluate feature importance.
3. Scalability: Filter methods can handle a large number of features effectively because they don't involve any learning methods. This makes them suitable for high dimensional datasets.
4. Pre-processing step: They can serve as a pre-processing step for other feature selection methods. For instance, you could use a filter

method to remove irrelevant features before applying a more computationally expensive method, such as a wrapper method.

Disadvantage

1. Lack of Feature Interaction :- Filter method treat each feature individual and hence do not consider the interaction between feature. They might miss out on identifying important features that doesn't appear significant individual but are significant in combination with other feature.
2. Model Agnostic :- Filter method are agnostic to the machine learning model that will be used for the prediction. This means that the selected feature might not necessarily contribute to the accuracy of the specific model you want to use.
3. Statistical Measures Limitation :- The statistical measure used in these methods have their own limitation. For example,

Correlation is the measure of the linear relationship and might not capture non-linear relationship effectively. Similarly, variance based methods might keep feature with high variance but low predictive power.

4. Threshold determination :- For some methods, determining the threshold to select feature can be a bit subjective. For example, what constitutes "low" variance or "high" correlation might differ depending on the context or the specific dataset.