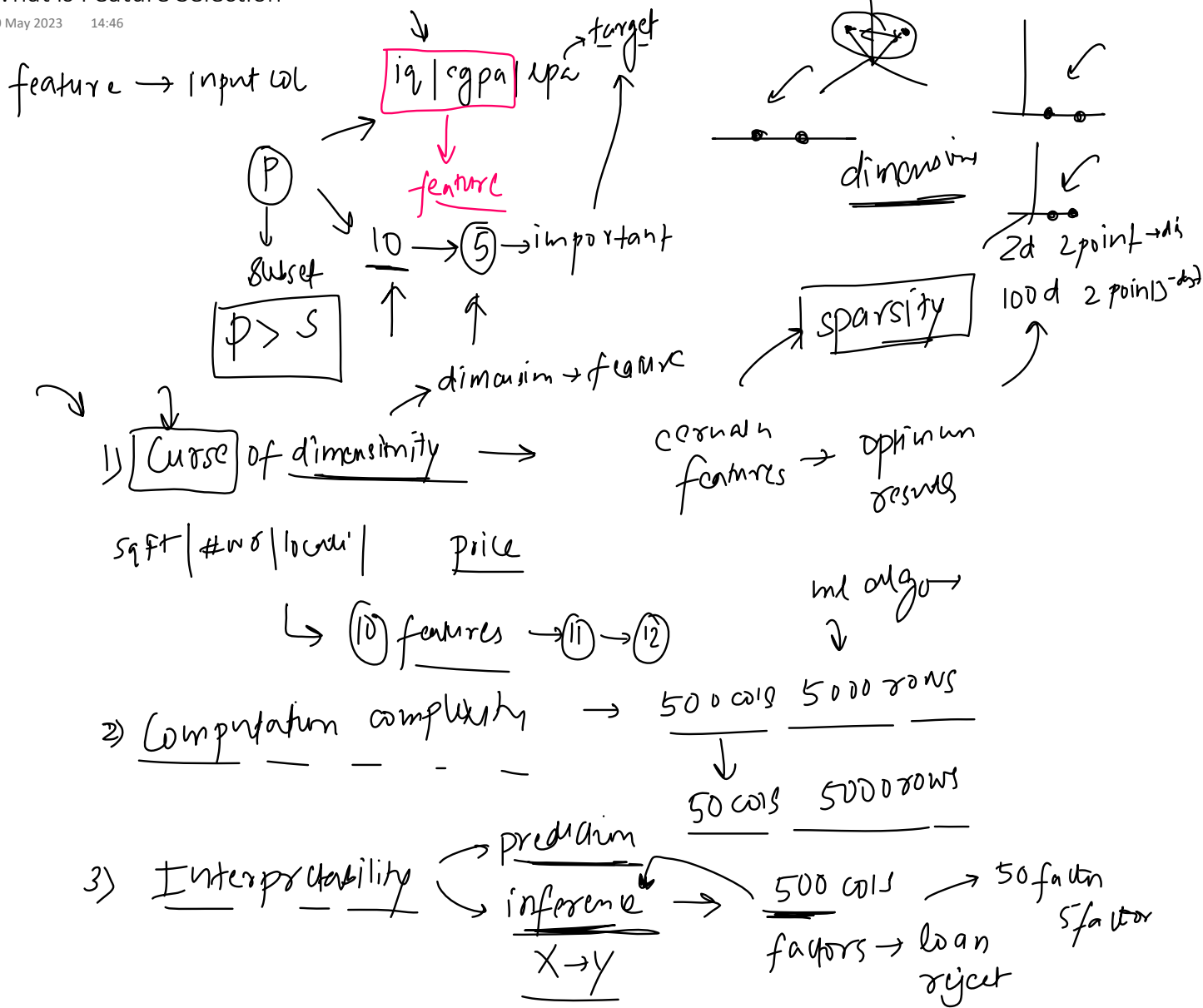


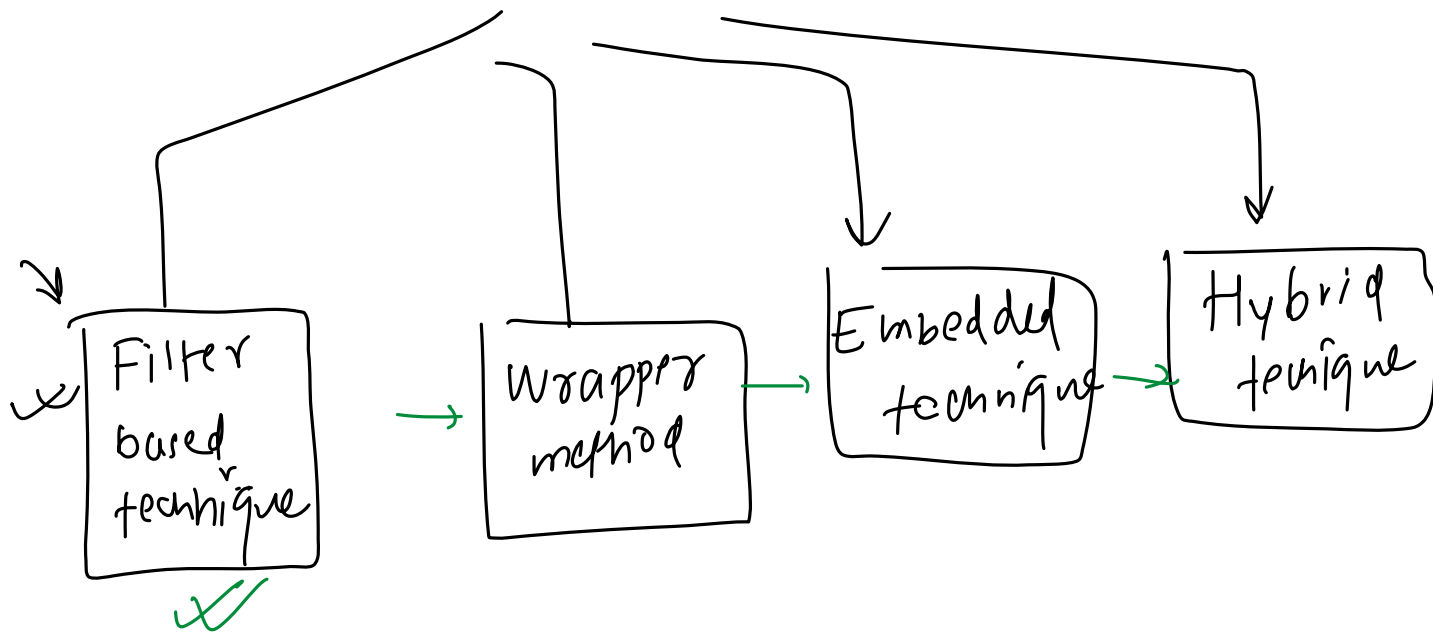
What is Feature Selection

10 May 2023 14:46



Types of Feature Selection

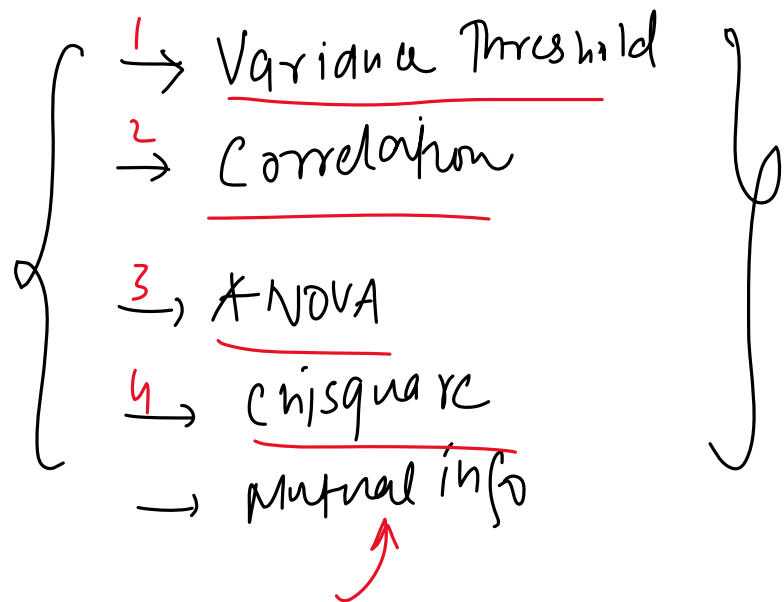
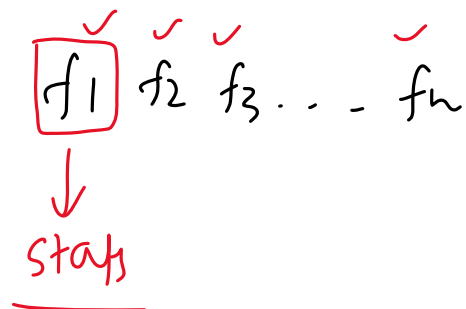
10 May 2023 14:47



Filter Based Feature Selection ←

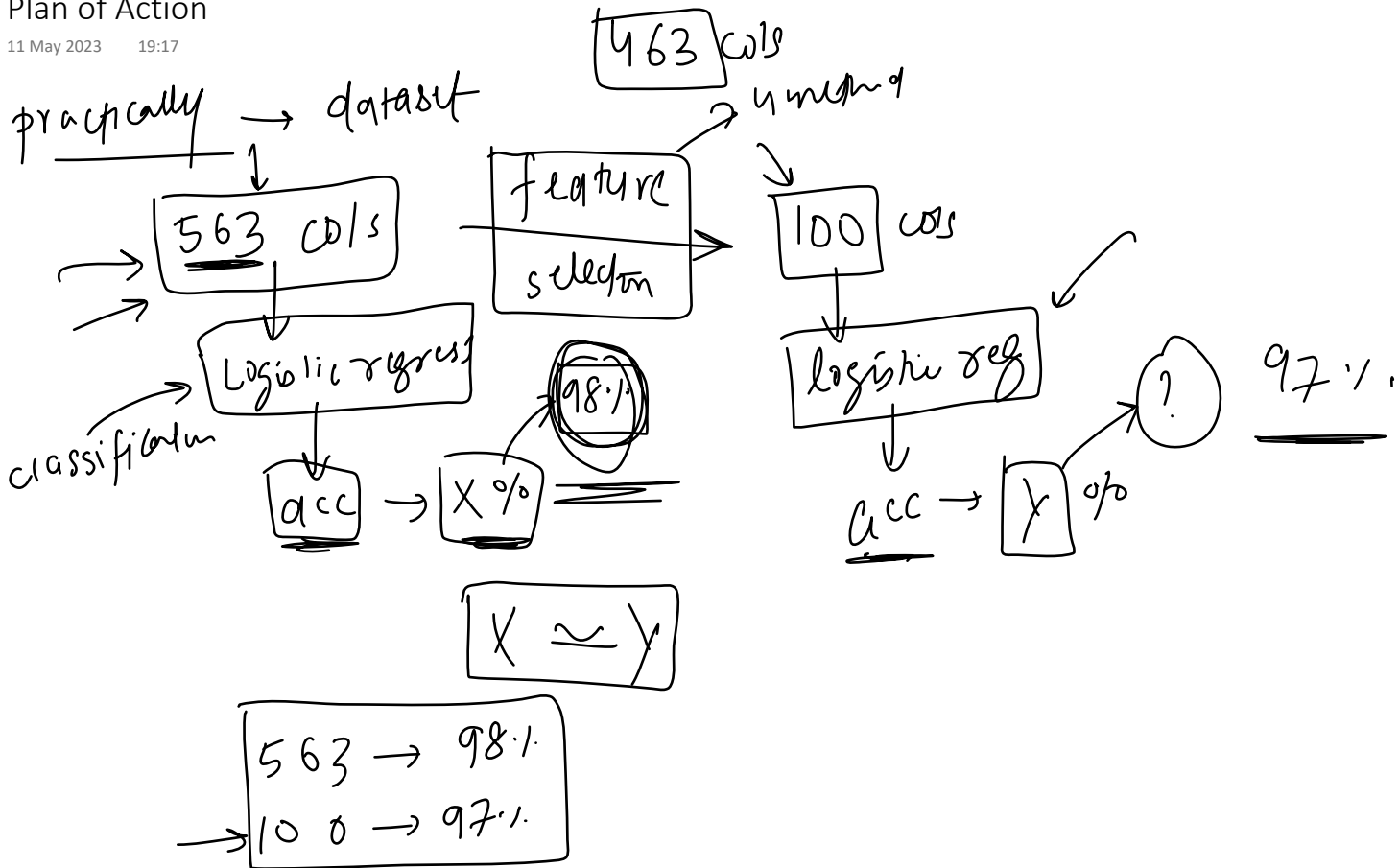
10 May 2023 14:47

{ Filter-based feature selection techniques are methods that use statistical measures to score each feature independently, and then select a subset of features based on these scores. These methods are called "filter" methods because they essentially filter out the features that do not meet some criterion.



Plan of Action

11 May 2023 19:17



1. Duplicate Features

10 May 2023 14:47

f_1	f_2	f_3	f_4	f_5	output
1	2	1	2	3	X
2	1	2	2	3	N
3	3	3	2	3	N

same

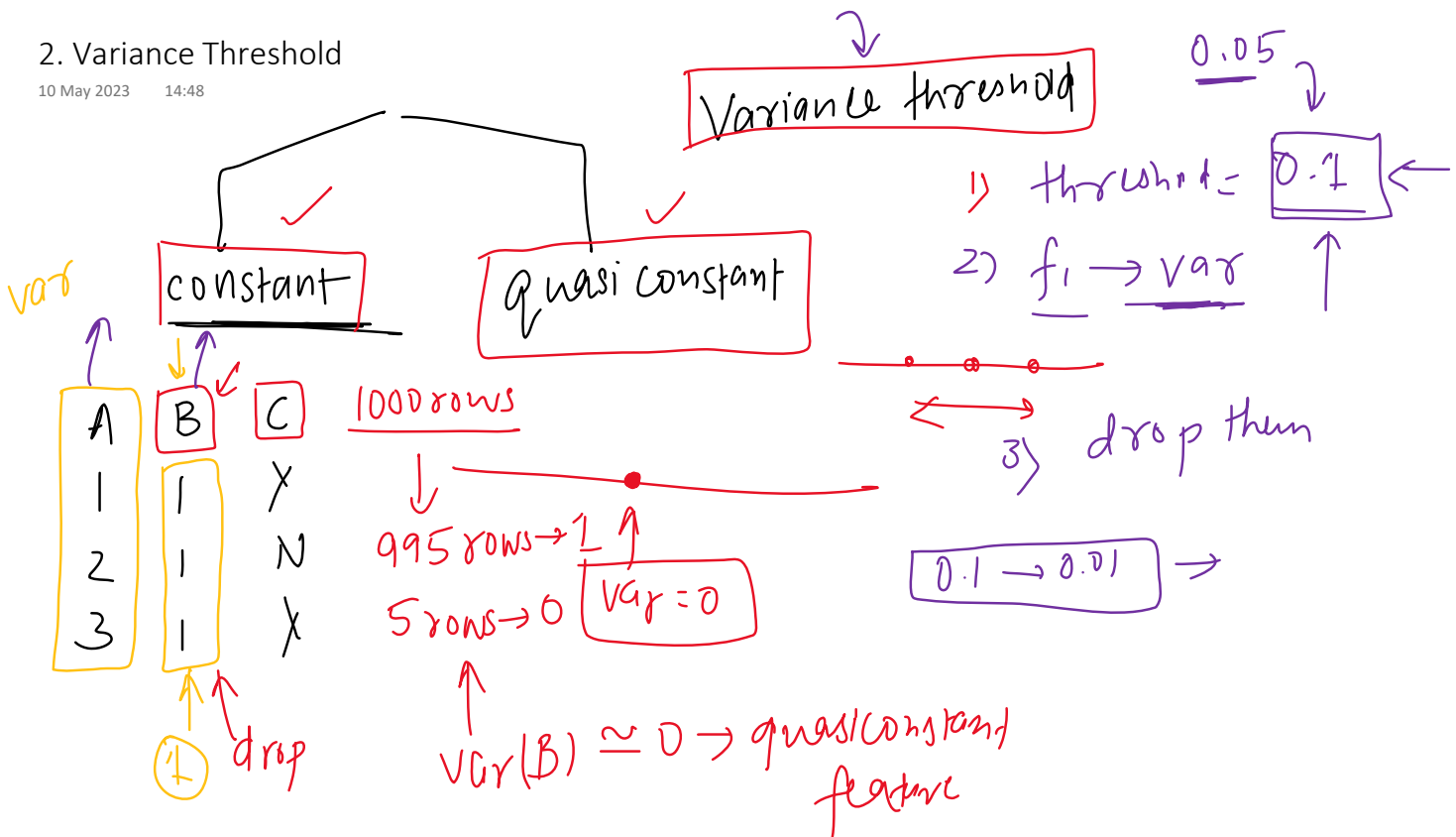
561 → 540

↓
21 duplicate

dict → keep keys → del → values
key → original cols
values → duplicate

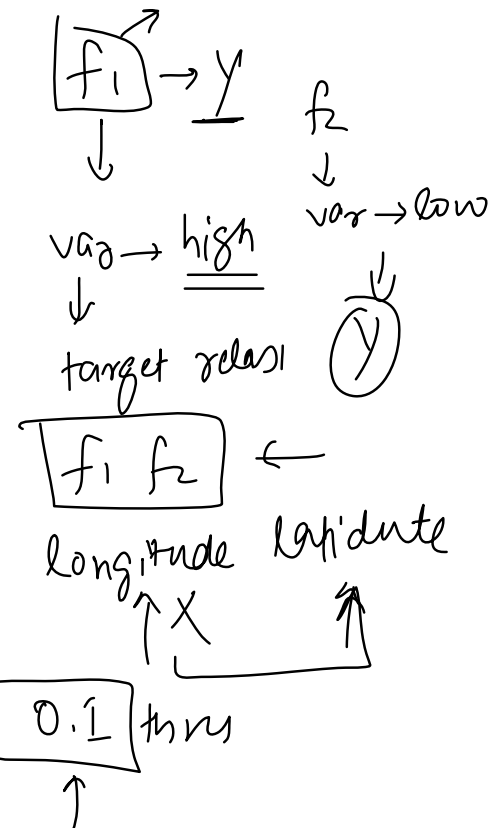
2. Variance Threshold

10 May 2023 14:48



Points to Consider

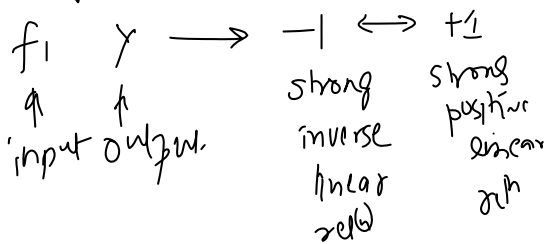
- 1. Ignores Target Variable:** Variance Threshold is a univariate method, meaning it evaluates each feature independently and doesn't consider the relationship between each feature and the target variable. This means it may keep irrelevant features that have a high variance but no relationship with the target, or discard potentially useful features that have a low variance but a strong relationship with the target.
- 2. Ignores Feature Interactions:** Variance Threshold doesn't account for interactions between features. A feature with a low variance may become very informative when combined with another feature.
- 3. Sensitive to Data Scaling:** Variance Threshold is sensitive to the scale of the data. If features are not on the same scale, the variance will naturally be higher for features with larger values. Therefore, it is important to standardize the features before applying Variance Threshold.
- 4. Arbitrary Threshold Value:** It's up to the user to define what constitutes a "low" variance. The threshold is not always easy to define and the optimal value can vary between datasets.



3. Correlation

10 May 2023 14:48

pearson corr coeff



$X \rightarrow Y$
0.9

-0.9 0
 $X \rightarrow Y$ $X \leftrightarrow Y$

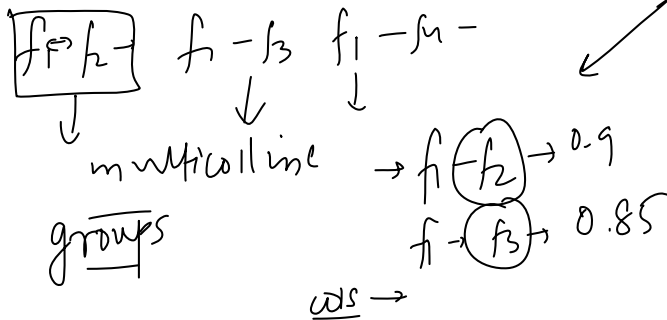
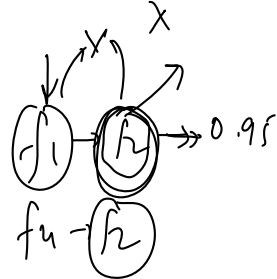
$f_1 f_2 f_3 \dots f_n y$

$f_1 \rightarrow y \rightarrow$
 $f_2 \rightarrow y \rightarrow$
 $f_3 \rightarrow y \rightarrow$
 \vdots
 $f_n \rightarrow y \rightarrow$

cutoff $\rightarrow \begin{cases} 0.3 \\ -0.3 \end{cases}$
feature X

brute force
corr

0.95
 $f_1 f_2$



561 \rightarrow 541 \rightarrow 341 \rightarrow 152 \rightarrow 100
duplizer var thr corr ANOVA

Disadvantages

- Linearity Assumption:** Correlation measures the linear relationship between two variables. It does not capture non-linear relationships well. If a relationship is nonlinear, the correlation coefficient can be misleading.
- Doesn't Capture Complex Relationships:** Correlation only measures the relationship between two variables at a time. It may not capture complex relationships involving more than two variables.
- Threshold Determination:** Just like variance threshold, defining what level of correlation is considered "high" can be subjective and may vary depending on the specific problem or dataset.
- Sensitive to Outliers:** Correlation is sensitive to outliers. A few extreme values can significantly skew the correlation coefficient.

$f_1 \leftrightarrow f_2$ $\begin{bmatrix} f_1 & f_2 & f_3 \end{bmatrix} \leftarrow y$

$0.95 \rightarrow 0.9 \rightarrow 0.8$

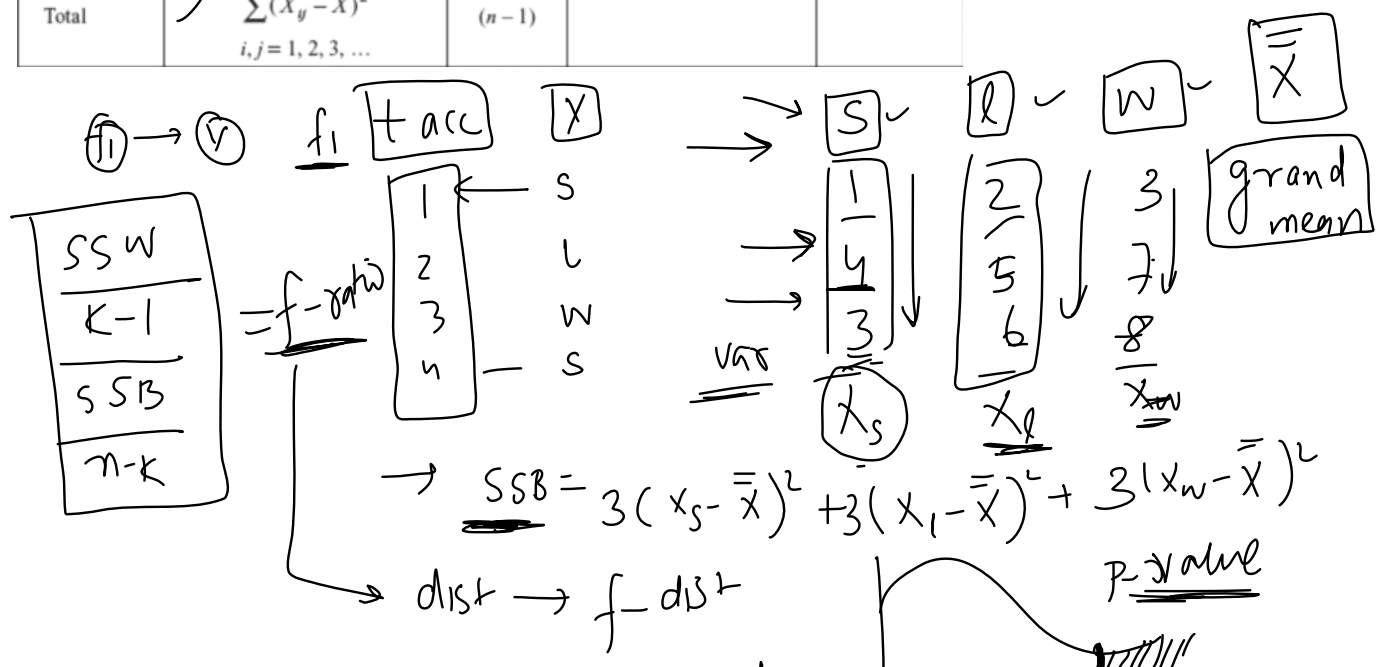
10 May 2023 14:49

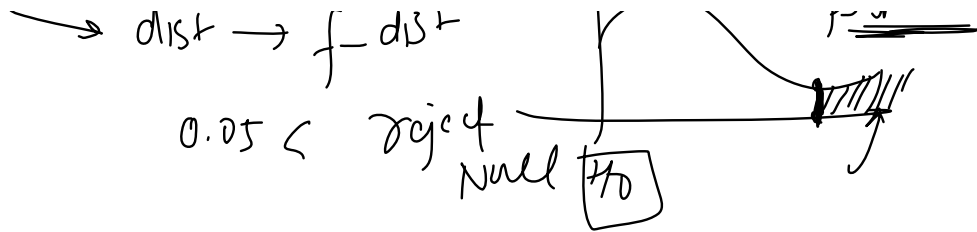
$f_1 \rightarrow y$ relⁿ $k-1$ $n-k$ $\eta \rightarrow \# \text{ rows}$ $k \rightarrow \# \text{ categories}$ $\textcircled{6}$
 \uparrow strong \rightarrow Hypothesis 1 way ANOVA

f -statistic \rightarrow p-value \rightarrow num $f \rightarrow y$
 \uparrow
 $f_1 \ f_2 \ \dots \ f_{152}$
ANOVA
 \rightarrow $(f_1) \rightarrow (y) \rightarrow$ no rel \rightarrow sklearn
 \rightarrow cat > 2

$$\underline{SSW} = 50$$

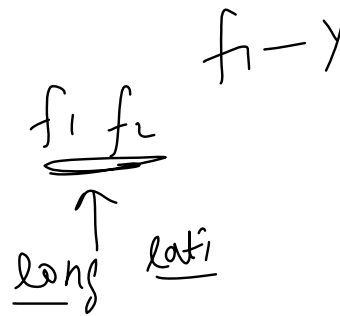
$$(1-x_s)^2 + (4-x_s)^2 + (3-x_s)^2 + (2-x_l)^2 + (5-x_l)^2 + (6-x_l)^2 + (3-x_w)^2 + (7-x_w)^2 + (8-x_w)^2$$





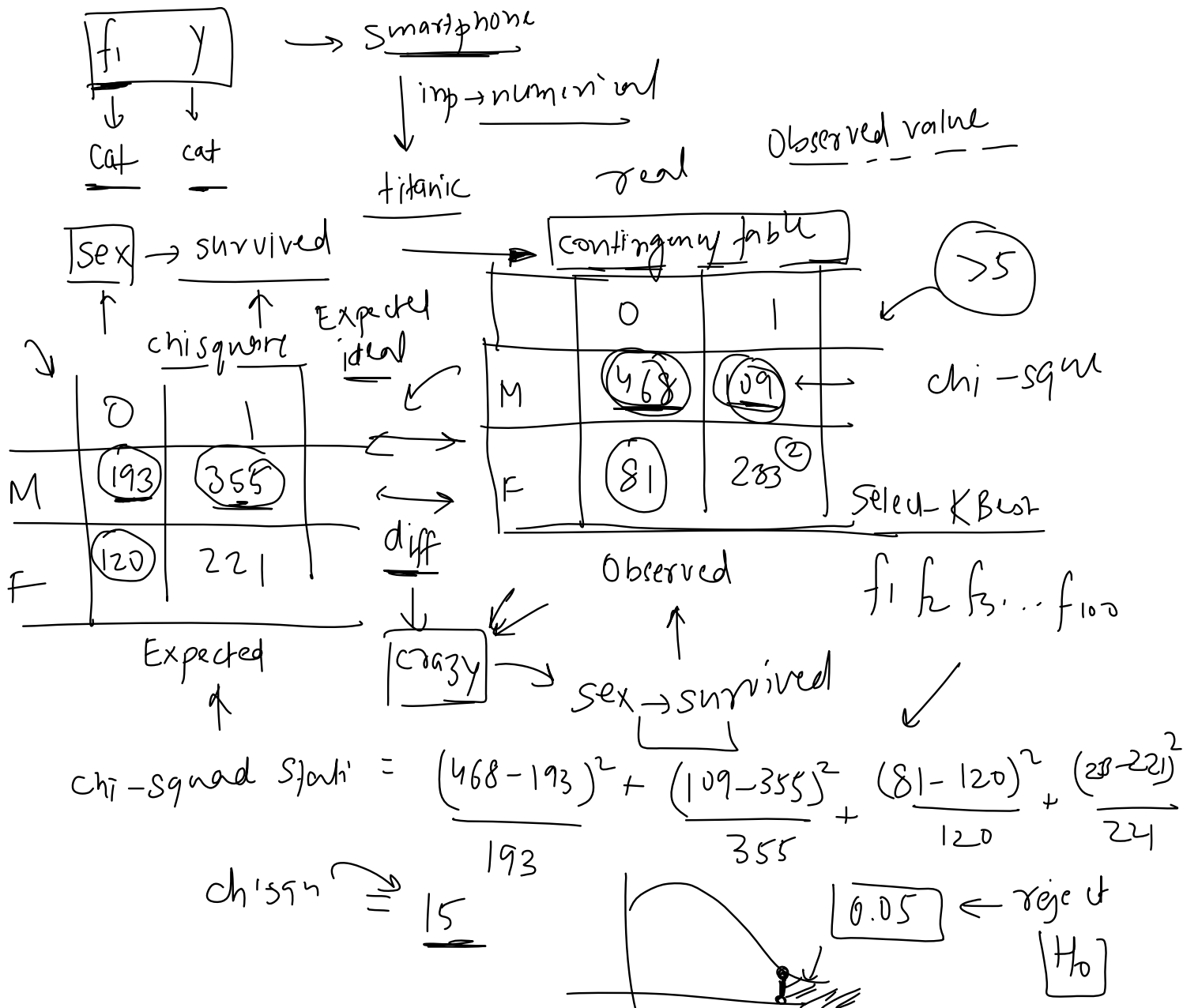
Disadvantages

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 distributions.
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 (e.g., time series data, 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 between features.



5. Chi-Square

10 May 2023 14:48



Disadvantages

- 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.
- 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 (e.g., time series data, nested data).
- Sufficient Sample Size:** Chi-square test requires a sufficiently large sample size. The results 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).
- No Variable Interactions:** Chi-square test, like other univariate feature selection methods, does not consider interactions between features. It might miss out on identifying important features that are significant in combination with other features.

Age

0-10
10-20
20-30

f_1, f_2

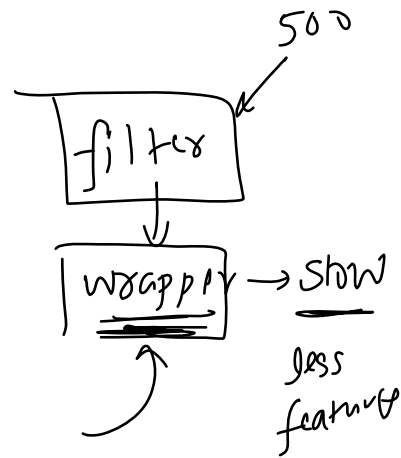
$$\underline{f_1 \rightarrow y}$$

Advantages and Disadvantages

11 May 2023 16:07

Advantages

1. **Simplicity**: Filter methods are generally straightforward and easy to understand. They involve calculating a statistic that measures the relevance of each feature, and selecting the top features based on this statistic.
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.



Disadvantages

f-h

1. **Lack of Feature Interaction**: Filter methods treat each feature individually and hence do not consider the interactions between features. They might miss out on identifying important features that don't appear significant individually but are significant in combination with other features.
2. **Model Agnostic**: Filter methods are agnostic to the machine learning model that will be used for the prediction. This means that the selected features might not necessarily contribute to the accuracy of the specific model you want to use.
3. **Statistical Measures Limitation**: The statistical measures used in these methods have their own limitations. For example, correlation is a measure of linear relationship and might not capture non-linear relationships effectively. Similarly, variance-based methods might keep features with high variance but low predictive power.
4. **Threshold Determination**: For some methods, determining the threshold to select features can be a bit subjective. For example, what constitutes "low" variance or "high" correlation might differ depending on the context or the specific dataset.

Vari
ANOVA
chi-sqr