



Team Name – Machine Learners

Statistical Change Detection for Multi-Dimensional Data

Course Instructor: Prof P Balamurugan

Course Code: IE-506

Machine Learning Principles and Techniques

Name

Roll Number

Mayur Dhanawade

23M1512

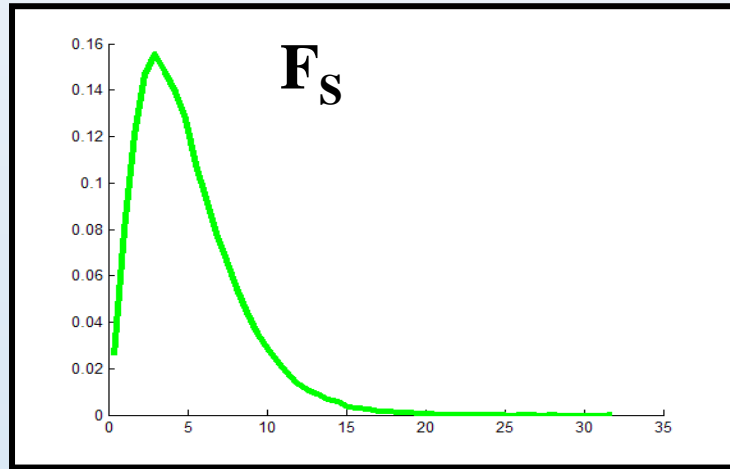
Shivam Negi

23M1508

Contents

1. Problem Explanation.
2. Motivation for the problem.
3. Introduction to Prior Work.
4. Density test high-level overview.
5. Kernel Density Estimate (KDE).
6. Optimal bandwidth.
7. Algorithms.
8. Calculate Test Statistic.
9. Derive the null distribution.
10. Decision making.
11. Two way test implementation.
12. Research paper experiments.
13. Our experiments result.
14. Data Pre-processing/Computational Framework.
15. Work done by team members.
16. Future work.

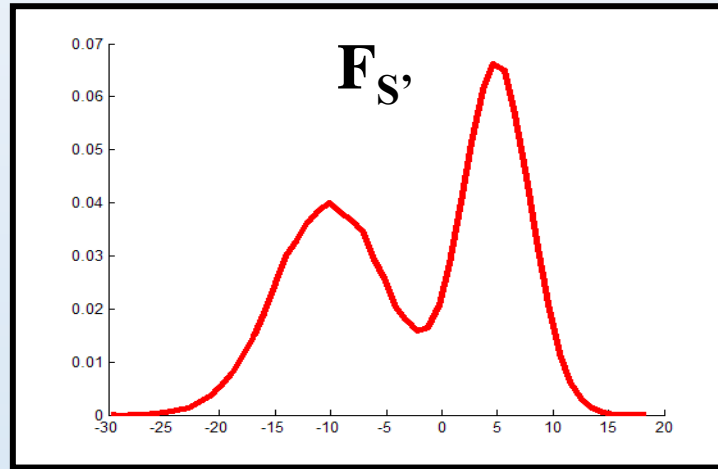
Problem Explanation.



Source: [3]

Data Set (S)

Baseline data



Source: [3]

Data Set (S')

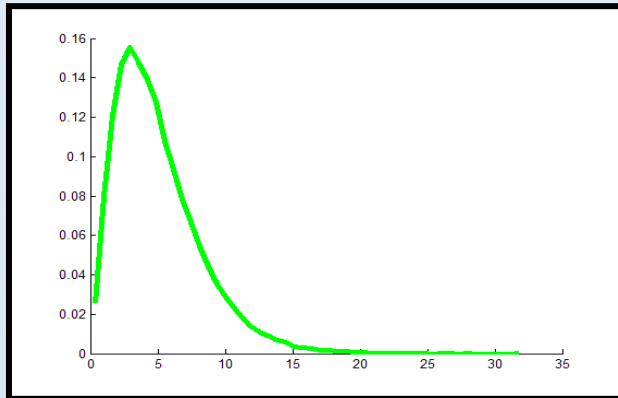
Recently observed data

Multi-dimensional space

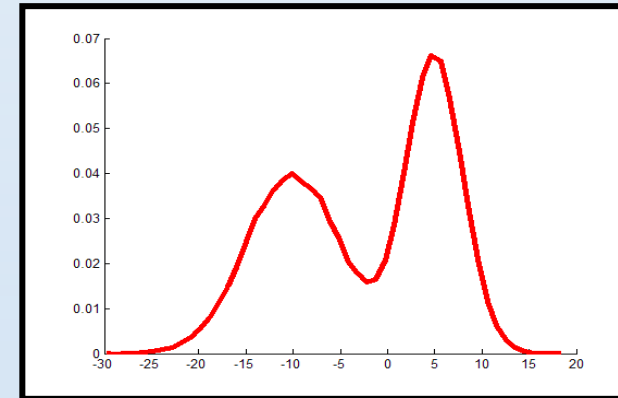
Unknown
distributions

Motivation for the Problem

- Analyzing the changes in the distribution of financial market.
- **Question: Does financial market shows different pattern recently?**
- We need a distributional change detection method to answer this question.



Source: [3]



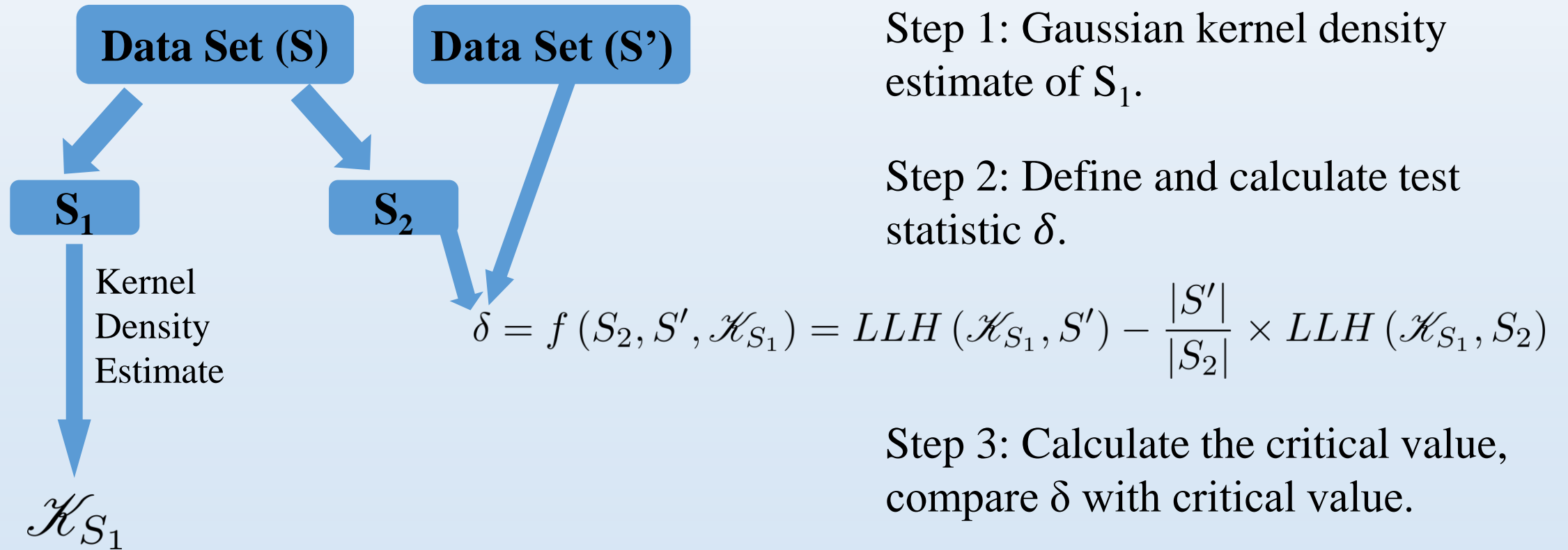
Source: [3]

- If Distribution change Observed, In this case, we will take further steps accordingly.

Introduction to Prior Work

- For uni-dimensional data, many existed tests, such as K-S test, chi-square test and many more.
- Only two tests to detect a generic distributional change in multi-dimensional space.
 - Kdq-tree test: suffer from curse of dimensionality. [1]
 - Cross-match test: computationally expensive due to maximum matching algorithm. [2]

Density test high-level overview



Step 1: Gaussian kernel density estimate of S_1 .

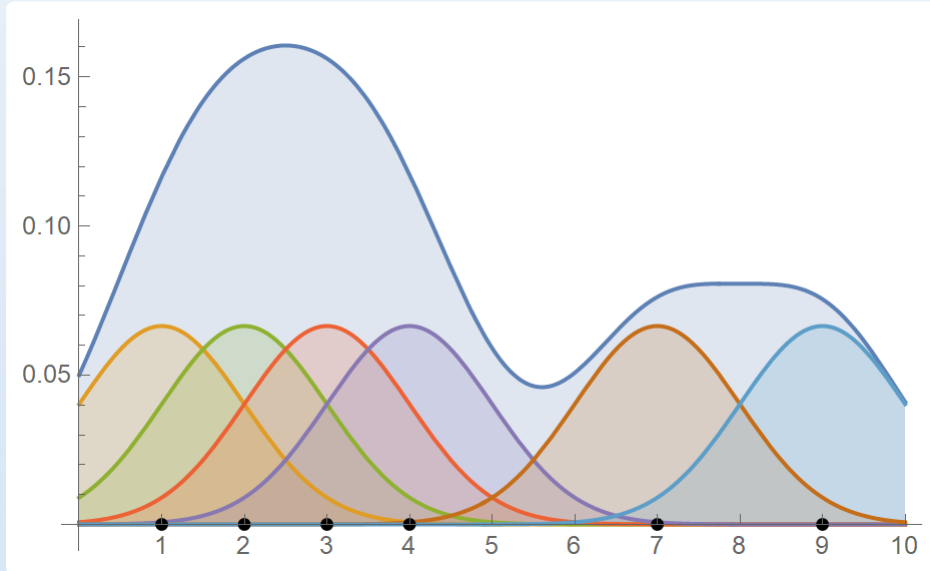
Step 2: Define and calculate test statistic δ .

Step 3: Calculate the critical value, compare δ with critical value.

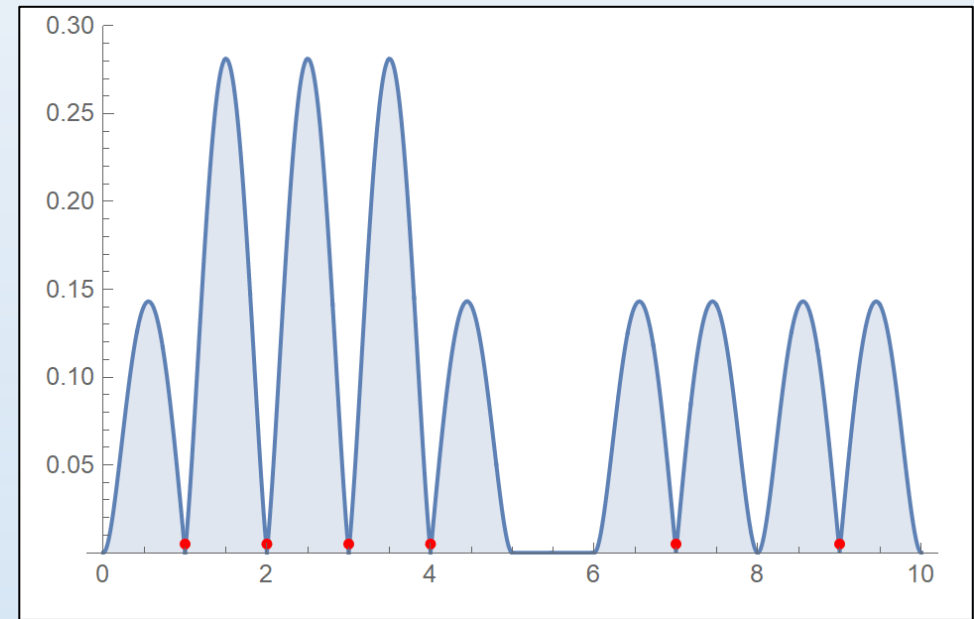
- If $\delta < \text{critical value}$, we will declare a change.
- Otherwise, it means no change.

Step 1: Kernel Density Estimate (KDE)

Bandwidth selection



[Source: ekamperi.github.io](https://ekamperi.github.io)



[Source: ekamperi.github.io](https://ekamperi.github.io)

- Data-driven bandwidth: converge better to the true distribution.
- Accuracy and power of test is increased when estimate is accurate.

Optimal bandwidth by MLE/EM

(maximum likelihood estimation / Expectation Maximization)

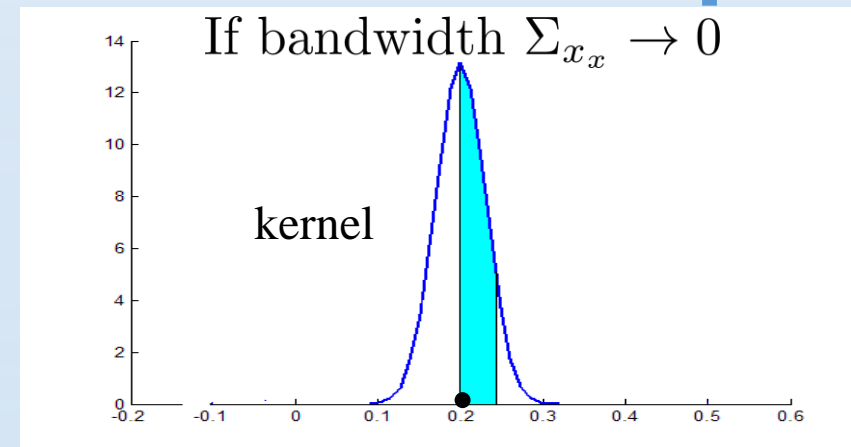
$$\text{Log Likelihood}(\mathcal{K}_{S_1}, S_1) = \sum_{x_j \in S_1} \log \left[\sum_{x_i \in S_1 \wedge i \neq j} \frac{1}{|S_1| - 1} G(\Sigma_{x_i}, x_j - x_i) \right]$$

Diagram illustrating the relationship between the Log Likelihood function and the Pseudo-LLH Function:

- A horizontal line under the Log Likelihood expression has a blue arrow pointing down to a blue box labeled "Pseudo – LLH Function".
- A blue arrow points down from the Log Likelihood expression to the symbol ∞ .
- A blue arrow points up from the expression $G(\Sigma_{x_i}, x_i - x_i) \rightarrow \infty$ to the term $G(\Sigma_{x_i}, x_j - x_i)$ in the denominator of the Log Likelihood expression.

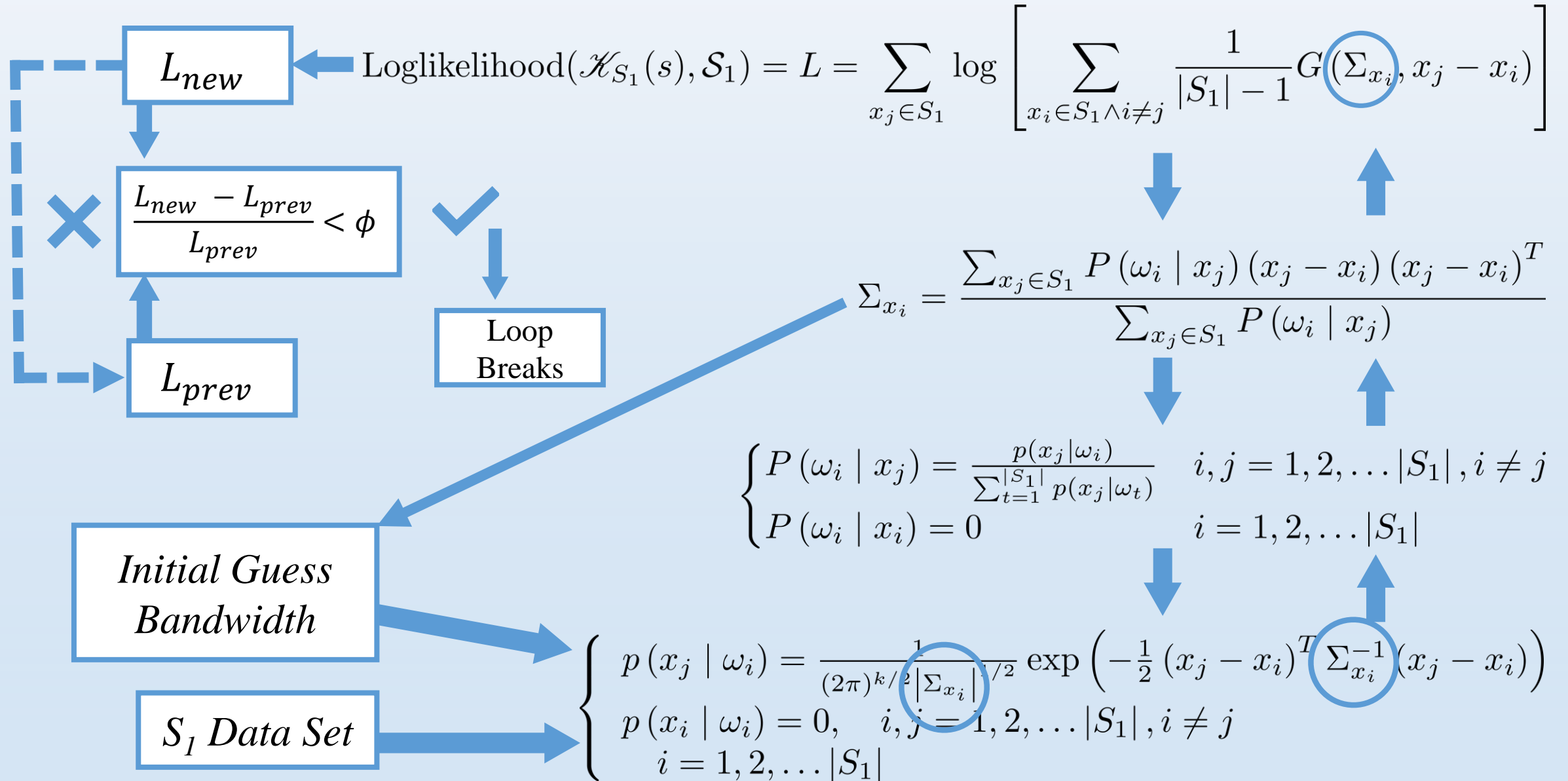
Pseudo – LLH
Function

Adding Constraints $G(\Sigma_{x_i}, x_i - x_i) = 0$
for all $x_i \in S_1, \quad i = 1, 2, \dots, |S_1|$

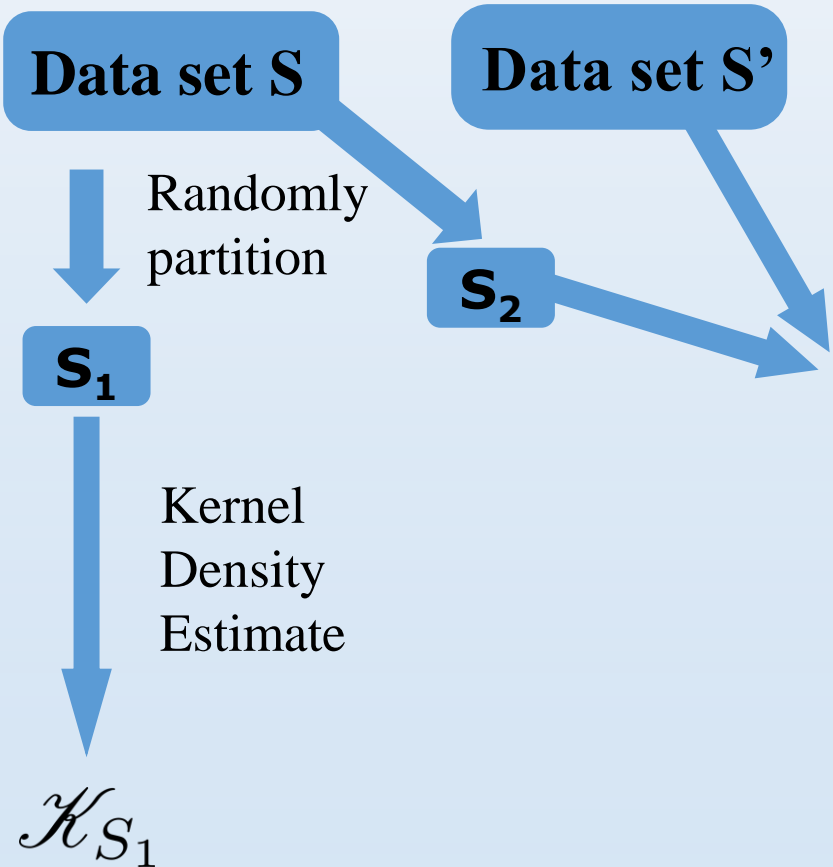


Source: [3]

Algorithm 1: Learn Bandwidth



Step 2: Calculate Test Statistic (δ)



- Small if S' different from S1.
- Large otherwise

Always Large

$$\begin{aligned}
 \delta &= \overline{\text{LLH}(\mathcal{K}_{S_1}, S')} - \frac{|S'|}{|S_2|} \times \overline{\text{LLH}(\mathcal{K}_{S_1}, S_2)} \\
 &= \log \left\{ \prod_{y \in S'} \mathcal{K}_{S_1}(y) \right\} - \frac{|S'|}{|S_2|} \times \log \left\{ \prod_{y \in S_2} \mathcal{K}_{S_1}(y) \right\} \\
 &= \sum_{y \in S'} \log \sum_{x \in S_1} \frac{1}{|S_1|} G(\Sigma_x, y - x) \\
 &\quad - \sum_{y \in S_2} \frac{|S'|}{|S_2|} \times \log \left\{ \sum_{x \in S_1} \frac{1}{|S_1|} G(\Sigma_x, y - x) \right\}
 \end{aligned}$$

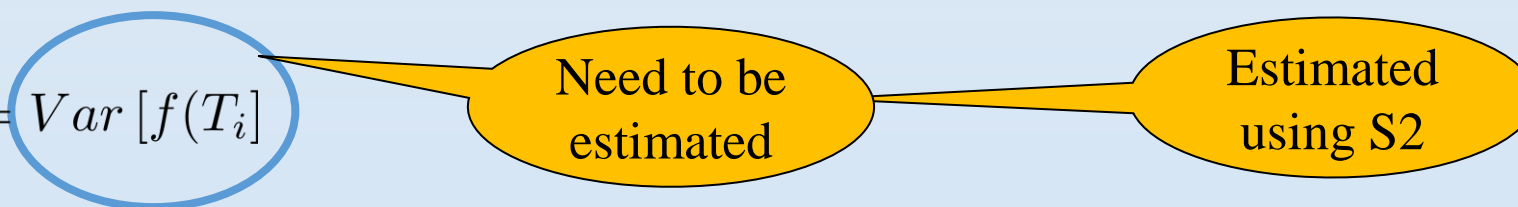
Step 3: Derive the null distribution (Algorithm 2)

$$\Delta = \underbrace{\sum_{i=1}^{|S'|} \log \sum_{x \in S_1} \frac{1}{|S_1|} G(\Sigma_x, T_i - x)}_{\Delta_1 \sim \text{normal}} - \underbrace{\sum_{i=|S'|+1}^{|S'|+|S_2|} \frac{|S'|}{|S_2|} \times \log \left\{ \sum_{x \in S_1} \frac{1}{|S_1|} G(\Sigma_x, T_i - x) \right\}}_{\Delta_2 \sim \text{normal}}$$

$\Delta \sim$ Normal Distribution by Central Limit *Theorem*.

$$E[\Delta] = 0, \text{Var}[\Delta] = \left(|S'| + \frac{|S'|^2}{|S_2|} \right) \sigma^2$$

Where, $\sigma^2 = \text{Var}[f(T_i)]$



Need to be estimated

Estimated using S2

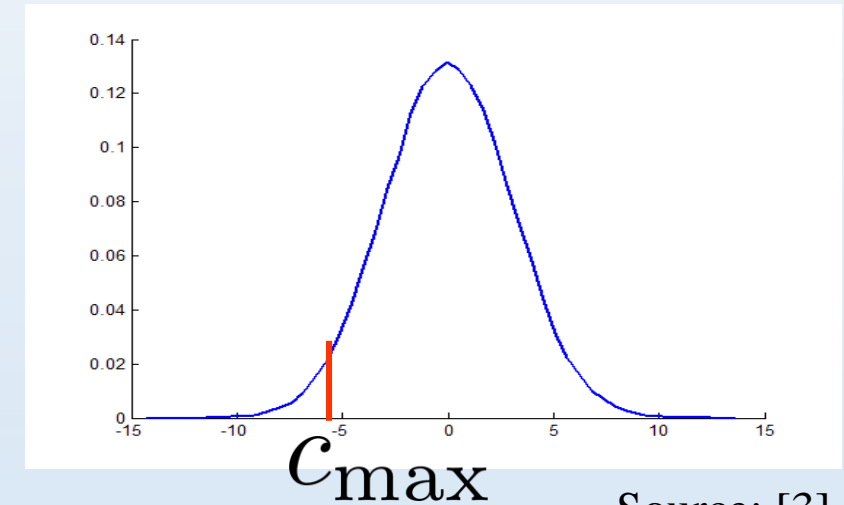
$$f(T_i) = \log \sum_{x \in S_1} \frac{1}{|S_1|} G(\Sigma_x, T_i - x)$$

All random variables T'_i s, $i = 1, 2, \dots$ follow the same distribution F_S .

Step 4: Calculate critical value and make a decision

- C_{\max} Is chosen in such a way that maximum false positive rate is less than user supplied value.
- This is done using algorithm 3.

estimated null distribution Δ



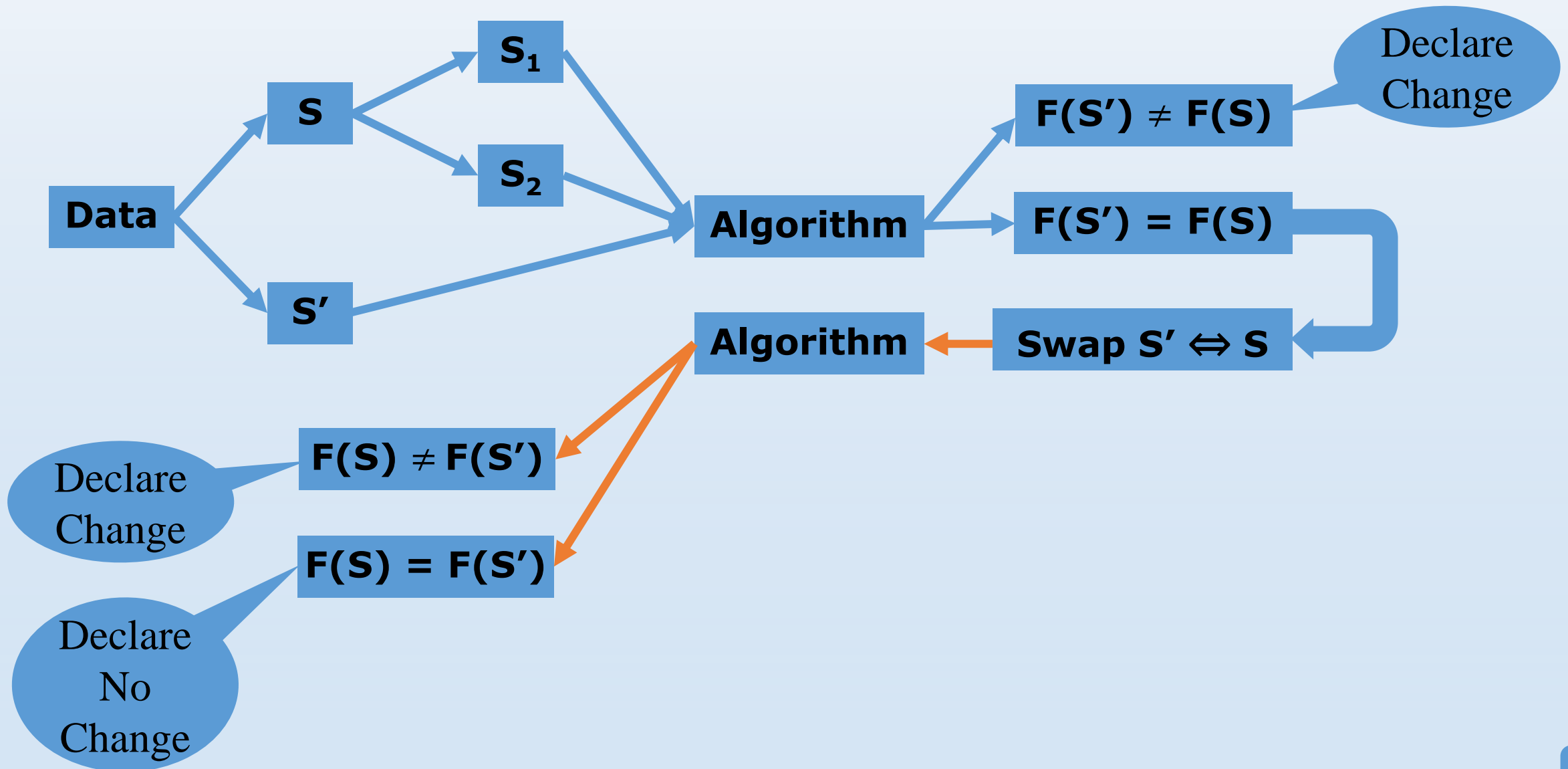
■ Inference:

- If $\delta < C_{\max}$
- Otherwise

Declare
Change

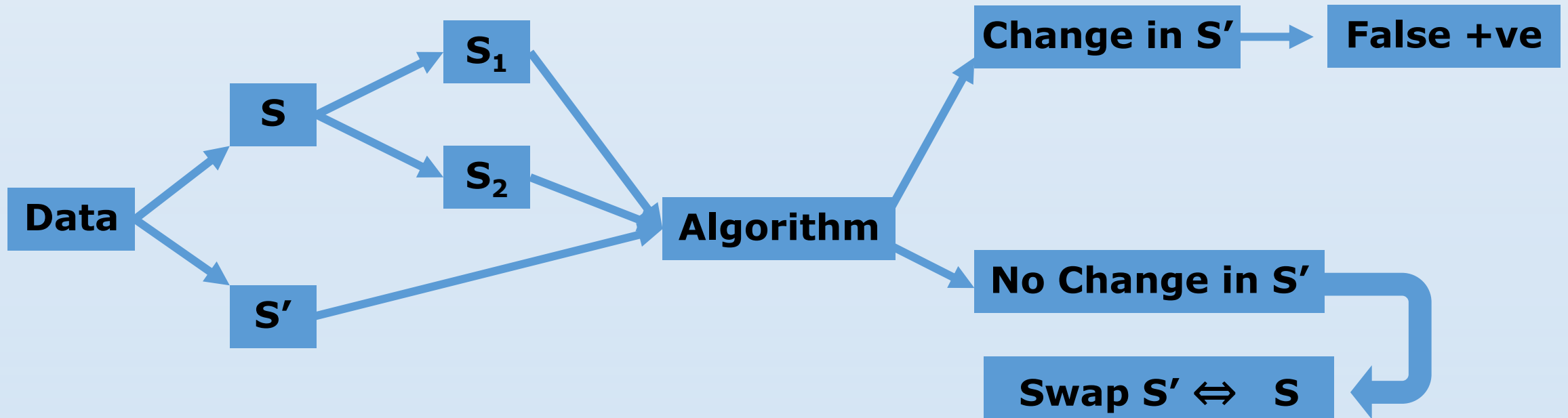
Declare No
Change

Two way density test implementation



Experiment 1

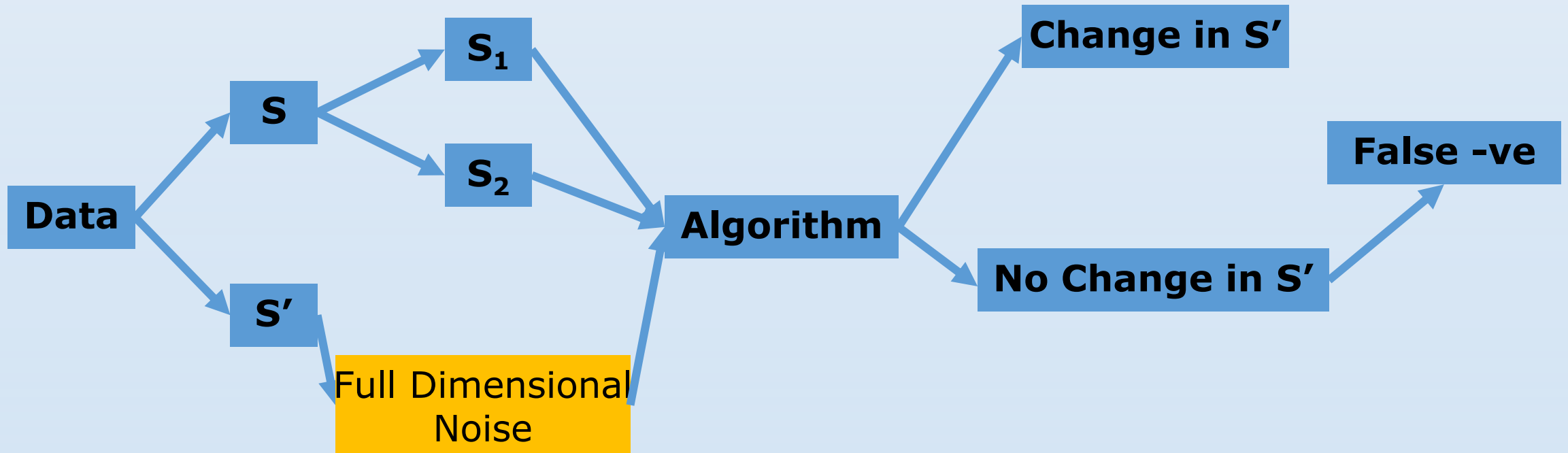
- False Positive Rate Estimation
 - Test for samples from the same distribution.
 - Run for 20 instances.
 - Test distributional change in each instance using two way test.
 - False positive rate = (Number of times it detects change when actually there is no change) / 20.



Experiment 2

Full-Dimensional Changes:

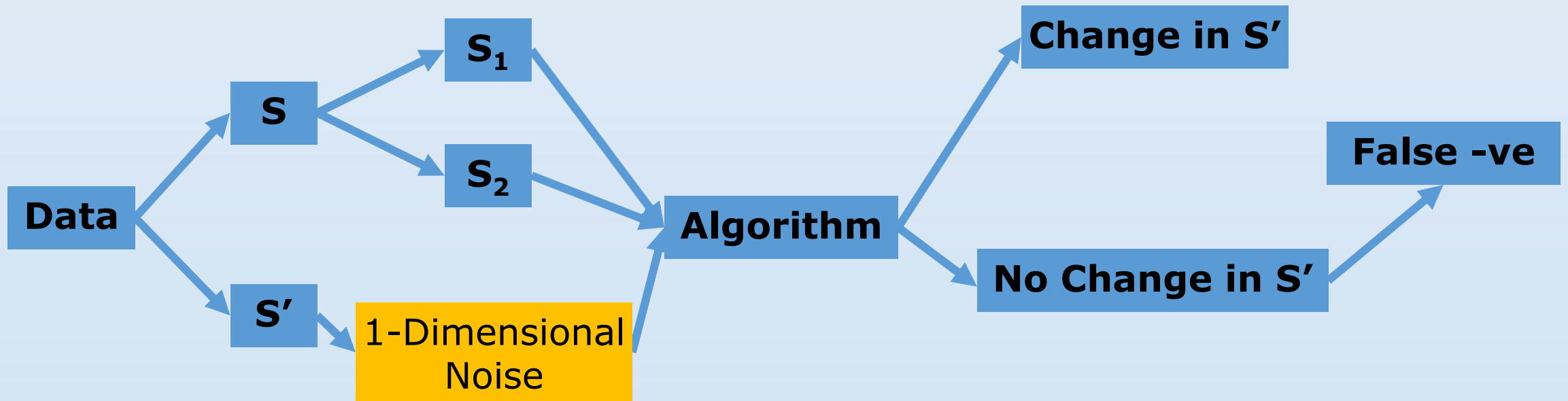
- Add some gaussian noise to all dimensions of some sample points.



Experiment 2

Single-dimensional changes:

- Add Gaussian Noise to a randomly selected dimension.
- Add Noise by multiplying by 2 to a randomly selected dimension.



Our Experiment Results

Experiment 1 Results:

- Expected false positives according to paper is 3 in 100 instances.
- False positives we got is 0 in 20 instances.

Experiment 2 Results:

- Expected false negatives according to paper for gaussian noise is 0 in 100 instances, for add 1-d noise addition is 33 in 100 instances, for scale-1d noise addition is 2 in 100 instances.
- False negatives we got for full dimensional gaussian noise is 0 in 10 instances,
- For add 1-d noise is 3 in 10 instances.
- For scale 1-d noise is 0 in 10 instances.

1) **Data Sets:**

- Worked on El-Nino dataset.

2) **Data Pre-processing:**

- Removed rows having missing values.
- Removed spatio - temporal attributes.
- Final dimensions of dataset: 93935 rows, 5 columns (5-D dataset).

3) **Computational Framework:**

- Intel(R) Xeon(R) CPU E5506, 2.13 GHz, Number of cores = 32, Ram = 128GB

4) **Programming Language:**

- Python

Work done by the team.

Shivam Negi (23M1508):

- Data preprocessing.
- Algorithm 2.
- Experiment 2.1
- Experiment 2.2

Mayur Dhanawade (23M1512):

- Algorithm 1
- Algorithm 3.
- Two way test.
- Experiment 1.

Modification/ Future work

- Currently work done on low dimensional dataset which was 'El-Nino' (5 dimensional).
- Further we intend to work on higher dimensional datasets, which are given in the research paper.

References

- [1] Dasu, Tamraparni et al. “An Information-Theoretic Approach to Detecting Changes in Multi-Dimensional Data Streams.” (2006). Accessed 30th Mar. 2024.
- [2] Rosenbaum, Paul R. “An Exact Distribution-Free Test Comparing Two Multivariate Distributions Based on Adjacency.” Journal of the Royal Statistical Society. Series B (Statistical Methodology), vol. 67, no. 4, 2005, pp. 515–30. JSTOR, <http://www.jstor.org/stable/3647642>, Accessed 30th Mar. 2024.
- [3] Image Material: <https://www.merlot.org/merlot/viewMaterial>, Accessed on: 5th Feb, 2024.
- [4] Alicia Horsch towardsdatascience.com, Accessed on: 5th Feb 2024.
- [5] Prof Manjesh Kumar Hanawal, youtube.com/playlist, Accessed on: 5th Feb 2024.
- [7] Jaroslaw Drapala: towardsdatascience.com, Accessed on: 5th Feb 2024.
- [8] Central Limit Theorem: web.stanford.edu, Accessed 2th Feb 2024.

Thank You