

Q2) In which kinds of datasets sample linear regression model could be performing better than Random forest ?

What are the limitations of random forest approach?

Simple linear regression could perform better than Random forest in datasets where there is a linear relationship between the independent and dependent variables. In such cases, a simple linear regression model can provide a more interpretable and understandable solution as it only involves a straight line equation. In addition, if the dataset is relatively small and there are only a few important features, a linear regression model could be easier to interpret and implement than a complex random forest model.

However, there are several limitations to the random forest approach, including:

1. Overfitting: Random forests are known to overfit when the number of trees in the forest is too large, or when there are too many features in the dataset.
2. Computational complexity: Random forests require significant computational resources, particularly when the dataset is large or there are many trees in the forest.
3. Interpretability: Random forests are difficult to interpret, particularly when the number of trees in the forest is large. This makes it challenging to understand the relationships between the input features and the output.
4. Imbalanced data: Random forests can struggle with imbalanced datasets, where one class is significantly more prevalent than the other. In such cases, the model may overemphasize the dominant class and ignore the minority class.
5. Hyperparameters: Random forests have several hyperparameters that need to be tuned to achieve optimal performance. This process can be time-consuming and require expert knowledge.

3) When we use Random forest we evaluate features at each decisive point (node) based on some cost

function such as Gini Index, mean square error, entropy. If you have single cell RNAseq data where there is no

certainty about distribution of data and it has many missing values, Which cost function would be most suitable for

such data in random Forest based classification and regression ?

For single-cell RNAseq data with many missing values and uncertain distribution, the most suitable cost function for Random Forest-based classification and regression would be mean square error (MSE). The MSE cost function is commonly used in

regression problems, and it measures the average squared difference between the predicted and actual values.

The MSE cost function is appropriate for single-cell RNAseq data as it does not make any assumptions about the distribution of the data and can handle missing values. In addition, MSE is a continuous and differentiable function that can be easily optimized during the training process of the Random Forest model.

However, it is important to note that Random Forest models are not typically used for single-cell RNAseq data analysis. Alternative methods such as clustering and dimensionality reduction techniques, such as principal component analysis (PCA) and t-distributed stochastic neighbor embedding (t-SNE), are commonly used to analyze single-cell RNAseq data.

4) When we perform hierarchical clustering, we use a distance measure between classes at every step of linking the classes. Such distance measures can be named as single linkage, complete linkage and average, Wards, median, centroid ? Can you explain in which conditions these distance measures would be successful

When performing hierarchical clustering, the choice of distance measure can significantly impact the results of the clustering analysis. Here are the most commonly used distance measures and the conditions under which they would be successful:

1. **Single linkage:** Single linkage is a distance measure that defines the distance between two clusters as the minimum distance between any two points in the two clusters. Single linkage is effective when the clusters are compact and separated, and when there is a clear separation between the clusters.
2. **Complete linkage:** Complete linkage is a distance measure that defines the distance between two clusters as the maximum distance between any two points in the two clusters. Complete linkage is effective when the clusters are relatively spherical and when the cluster sizes are similar.
3. **Average linkage:** Average linkage is a distance measure that defines the distance between two clusters as the average distance between all pairs of points in the two clusters. Average linkage is effective when the clusters are relatively dense and when the cluster sizes are similar.
4. **Ward's linkage:** Ward's linkage is a distance measure that minimizes the increase in variance when merging two clusters. Ward's linkage is effective when the clusters are relatively compact and when the cluster sizes are similar.

5. Median linkage: Median linkage is a distance measure that defines the distance between two clusters as the median distance between all pairs of points in the two clusters. Median linkage is effective when the data contains outliers and when the clusters are relatively dense.
6. Centroid linkage: Centroid linkage is a distance measure that defines the distance between two clusters as the distance between their centroids. Centroid linkage is effective when the clusters are relatively spherical and when the cluster sizes are similar.

It is important to note that the choice of distance measure also depends on the specific characteristics of the data being analyzed. Therefore, it is important to try different distance measures and compare their results before selecting the most appropriate one.

5) Here is the table of conditions and playing information. Calculate the probability of play on the new day using naïve Bayes approach

Hot Humid rainy windy Play Golf

Yes	NO	Yes	NO	NO
NO	NO	Yes	Yes	NO
Yes	Yes	No	NO	Yes
Yes	NO	Yes	Yes	Yes
NO	Yes	Yes	NO	NO
Yes	NO	No	Yes	Yes
NO	NO	Yes	NO	NO
Yes	Yes	Yes	Yes	NO
Yes	NO	No	NO	Yes
NO	Yes	Yes	Yes	NO
Yes	NO	No	NO	Yes
Yes	NO	Yes	Yes	NO

NO Yes Yes No ?      New day

To calculate the probability of playing golf on the new day using Naive Bayes approach, we need to first calculate the prior probabilities and likelihoods for each condition:

1. Prior probabilities:

$$P(\text{Play Golf} = \text{Yes}) = 7/12$$

$$P(\text{Play Golf} = \text{No}) = 5/12$$

## 2. Likelihoods:

Hot:

$$P(\text{Hot} \mid \text{Play Golf} = \text{Yes}) = 4/7$$

$$P(\text{Hot} \mid \text{Play Golf} = \text{No}) = 2/5$$

Humid:

$$P(\text{Humid} \mid \text{Play Golf} = \text{Yes}) = 3/7$$

$$P(\text{Humid} \mid \text{Play Golf} = \text{No}) = 2/5$$

Rainy:

$$P(\text{Rainy} \mid \text{Play Golf} = \text{Yes}) = 2/7$$

$$P(\text{Rainy} \mid \text{Play Golf} = \text{No}) = 1/5$$

Windy:

$$P(\text{Windy} \mid \text{Play Golf} = \text{Yes}) = 3/7$$

$$P(\text{Windy} \mid \text{Play Golf} = \text{No}) = 2/5$$

## 3. Calculate the posterior probabilities for each class:

$$P(\text{Play Golf} = \text{Yes} \mid \text{Hot, Humid, No, Windy}) \propto P(\text{Play Golf} = \text{Yes}) * P(\text{Hot} \mid \text{Play Golf} = \text{Yes}) * P(\text{Humid} \mid \text{Play Golf} = \text{Yes}) * P(\text{No} \mid \text{Play Golf} = \text{Yes}) * P(\text{Windy} \mid \text{Play Golf} = \text{Yes})$$

$$= (7/12) * (4/7) * (3/7) * (2/7) * (3/7)$$

$$= 0.0171$$

$$P(\text{Play Golf} = \text{No} \mid \text{Hot, Humid, No, Windy}) \propto P(\text{Play Golf} = \text{No}) * P(\text{Hot} \mid \text{Play Golf} = \text{No}) * P(\text{Humid} \mid \text{Play Golf} = \text{No}) * P(\text{No} \mid \text{Play Golf} = \text{No}) * P(\text{Windy} \mid \text{Play Golf} = \text{No})$$

$$= (5/12) * (2/5) * (2/5) * (3/5) * (2/5)$$

$$= 0.0576$$

## 4. Normalize the probabilities:

$$P(\text{Play Golf} = \text{Yes} \mid \text{Hot, Humid, No, Windy}) = 0.0171 / (0.0171 + 0.0576) = 0.229$$

$$P(\text{Play Golf} = \text{No} \mid \text{Hot, Humid, No, Windy}) = 0.0576 / (0.0171 + 0.0576) = 0.771$$

Therefore, the probability of playing golf on the new day given the conditions Hot, Humid, No, and Windy is 0.229 (22.9%). According to the Naive Bayes approach, the new day is more likely to not play golf.

Q6. Out of different network inference methods taught in class, which method provide a scale free network of genes? And why scale free network properties are being given importance ? Can be make mutual information based network which is scale free ?

The method that provides a scale-free network of genes is the weighted gene co-expression network analysis (WGCNA). WGCNA is a popular network inference method that groups genes based on their expression patterns across different conditions, and it can identify co-expression modules that represent groups of genes with similar expression patterns. WGCNA also provides a measure of gene connectivity or "hubness," which allows us to identify the most connected or important genes in the network.

Scale-free network properties are important because they have been observed in many biological systems, including gene regulatory networks, protein-protein interaction networks, and metabolic networks. In a scale-free network, most nodes have relatively few connections, while a few nodes (called "hubs") have many connections. This type of network structure is thought to provide robustness to the network, as it allows for efficient communication and coordination among different parts of the network.

Yes, it is possible to create a scale-free network using mutual information-based methods, such as the context likelihood of relatedness (CLR) algorithm. The CLR algorithm identifies potential gene-gene interactions based on their mutual information, and it can also be used to construct a network of genes with scale-free properties. However, it is important to note that not all mutual information-based methods will necessarily result in a scale-free network.

Q7. What do you mean by information gain and entropy during classification using decision tree ? What is the concept behind information gain in general ?

In the context of classification using decision trees, information gain is a measure of the reduction in entropy or uncertainty in the target variable (class label) that results from splitting the data based on a particular attribute or feature. Entropy, in this case, refers to the degree of randomness or uncertainty in the class labels of the data.

The concept behind information gain is based on the idea that the best attribute to split the data on is the one that results in the greatest reduction in entropy or uncertainty. When we split the data based on an attribute, we create subsets of the data that are more homogeneous in terms of their class labels. The more homogeneous the subsets, the less uncertain we are about the class labels of the data in those subsets.

To calculate information gain, we first calculate the entropy of the original dataset, which is a measure of the overall uncertainty or randomness in the class labels. We then calculate the entropy of each subset created by splitting the data on each attribute, and we use these values to calculate the information gain for each attribute. The attribute with the highest information gain is chosen as the attribute to split the data on.

In general, information gain is a measure of the reduction in uncertainty or randomness in a system that results from acquiring new information or knowledge. It is often used in the context of machine learning and decision-making to evaluate the importance of different features or attributes in a dataset, and to identify the features that are most informative for predicting the target variable.

Q8. What are the weakness of naïve Bayes classifier ? How can we remove them ?

The naive Bayes classifier is a simple and efficient machine learning algorithm that is commonly used for text classification and other applications. However, there are some weaknesses or limitations of the naive Bayes classifier that can affect its performance in certain situations. Here are some of the main weaknesses of the naive Bayes classifier:

1. Assumption of feature independence: The naive Bayes classifier assumes that the features (or attributes) in the dataset are independent of each other, which may not be true in practice. This can lead to inaccurate predictions when there are strong correlations or interactions between the features.
2. Sensitivity to irrelevant features: The naive Bayes classifier is sensitive to irrelevant features, meaning that even if a feature has no predictive power, it can still affect the classification results. This can be a problem when the dataset contains many irrelevant or noisy features.

3. Lack of flexibility: The naive Bayes classifier is a simple model that cannot capture complex relationships or interactions between the features. This can limit its ability to accurately model certain types of datasets.

There are several ways to address these weaknesses of the naive Bayes classifier:

1. Relax the independence assumption: One way to improve the naive Bayes classifier is to relax the assumption of feature independence by using more flexible models, such as Bayesian networks or tree-augmented naive Bayes models. These models can capture dependencies between the features and improve the accuracy of the classifier.
2. Feature selection: To reduce the sensitivity of the classifier to irrelevant features, we can use feature selection techniques to identify the most informative features in the dataset. This can improve the performance of the classifier and reduce the risk of overfitting.
3. Ensembling: To improve the flexibility of the classifier, we can use ensemble methods, such as bagging or boosting, to combine multiple classifiers and capture complex relationships between the features. This can improve the accuracy and robustness of the classifier.

Q9. How is dimension reduction different from clustering? And How they can benefit each other. Explain taking example of tSNE and a clustering algorithm.

Dimension reduction and clustering are two commonly used techniques in unsupervised machine learning that can be used to analyze and visualize complex datasets. While they are related, they serve different purposes and have different goals.

Dimension reduction is the process of reducing the number of variables or dimensions in a dataset, while retaining as much of the relevant information as possible. The goal of dimension reduction is to simplify the dataset and make it easier to analyze or visualize, while still capturing the most important relationships and patterns in the data.

Clustering, on the other hand, is the process of grouping similar data points together based on some similarity metric. The goal of clustering is to identify patterns or structure in the data and to group similar data points into clusters or subgroups.

Despite their differences, dimension reduction and clustering can benefit each other in several ways. By reducing the number of dimensions in the dataset, dimension

reduction techniques such as tSNE can make it easier to visualize and interpret the results of a clustering algorithm. By grouping similar data points together, clustering algorithms can help to identify patterns or structure in the data that may not be immediately obvious.

For example, suppose we have a high-dimensional dataset of gene expression levels in different tissues, and we want to identify clusters of tissues that have similar gene expression profiles. We can start by applying a clustering algorithm such as k-means to the dataset, which will group similar tissues together based on their gene expression levels. However, the resulting clusters may be difficult to visualize or interpret in the high-dimensional space.

To overcome this, we can apply a dimension reduction technique such as tSNE to reduce the dimensionality of the dataset to two or three dimensions, while preserving the underlying structure and relationships in the data. We can then visualize the resulting clusters in the reduced-dimensional space, which can help us to identify patterns or relationships between the tissues that may not have been apparent in the high-dimensional space.

Overall, combining dimension reduction and clustering techniques can help to reveal hidden patterns and relationships in complex datasets, and can aid in the interpretation and visualization of the results.

Q10. How is GC bias estimated using next gen?4 Sequencing reads. How it can be corrected for ChIP-seq read-counts on a location.

GC bias is a common source of bias in next-generation sequencing (NGS) data, which can affect the accuracy and reliability of downstream analyses such as ChIP-seq. GC bias refers to the fact that the efficiency of PCR amplification and sequencing can be influenced by the GC content of the DNA fragment being sequenced. Regions with high GC content tend to be over-represented in sequencing reads, while regions with low GC content tend to be under-represented.

To estimate GC bias in NGS data, a common approach is to calculate the ratio of observed read counts to expected read counts for each genomic region or bin. The expected read counts can be calculated based on the average GC content of the genome, and the observed read counts are the actual read counts obtained from sequencing. If a particular region has a higher observed/read count ratio than expected, this suggests that there is GC bias in that region.



Once GC bias has been estimated, it can be corrected for ChIP-seq read counts using various methods. One common approach is to use normalization techniques that adjust the read counts based on the GC content of each region. For example, the read counts can be divided by the expected read counts based on the GC content of each region, or a regression model can be used to adjust the read counts based on the GC content and other factors that may influence sequencing bias.

Another approach is to use GC content as a covariate in the downstream analysis, such as in differential binding analysis. By including GC content as a covariate, any bias due to GC content can be accounted for and reduced in the analysis.

Overall, correcting for GC bias in ChIP-seq read counts is an important step to ensure the accuracy and reliability of downstream analyses, and there are several approaches that can be used to account for this bias.

Q11. If we have a lot of reads in the same orientation and starting from the same location, what does it mean? What should we do with those readings ?

If we have a lot of reads in the same orientation and starting from the same location, it could indicate a bias in the sequencing library preparation, such as a preferential ligation of adapters to certain regions of the DNA fragment. This can result in an over-representation of reads starting from a particular location and in a particular orientation.

To deal with this issue, one common approach is to use adapter trimming and quality filtering to remove low-quality reads and adapter sequences. This can help to reduce the bias and improve the accuracy of downstream analyses.

Another approach is to use tools that can correct for bias in the data, such as alignment-based methods that adjust for mapping bias, or normalization methods that adjust for coverage bias. For example, some alignment tools use local realignment or base quality score recalibration to correct for bias in the data.

In general, it is important to be aware of potential biases in the sequencing data and take appropriate steps to reduce or correct for them, in order to ensure the accuracy and reliability of downstream analyses.

Q12. What is a unique molecular identifier based RNA-seq? Why could it be better than normal RNA-seq for single cells?

Unique molecular identifier (UMI)-based RNA sequencing (RNA-seq) is a method of sequencing that uses molecular barcodes or UMIs to tag individual RNA molecules before amplification and sequencing. Each UMI is a short nucleotide sequence that is added to the RNA molecule during library preparation, and it serves as a unique identifier for that molecule throughout the sequencing process.

The advantage of using UMI-based RNA-seq for single cells is that it can reduce the effects of amplification bias and sequencing noise, which can be particularly problematic in single-cell RNA sequencing (scRNA-seq) due to the low starting material and high levels of technical noise. By using UMIs to tag individual RNA molecules, it is possible to distinguish between true biological variation and technical noise, and to accurately quantify gene expression at the single-cell level.

In UMI-based RNA-seq, each RNA molecule is tagged with a unique UMI, and then amplified and sequenced as usual. During analysis, the UMI sequences are used to collapse duplicate reads and remove PCR duplicates, so that each unique RNA molecule is counted only once. This reduces the impact of amplification bias and PCR duplicates, and improves the accuracy of gene expression quantification.

Overall, UMI-based RNA-seq can be a powerful tool for studying gene expression at the single-cell level, particularly in samples with low starting material or high levels of technical noise. By reducing the impact of amplification bias and sequencing noise, UMI-based RNA-seq can provide more accurate and reliable gene expression measurements, which can be critical for understanding the biology of individual cells and tissues.

Q13.Explain how do we perform gene-set enrichment analysis using ChIP-seq peaks. If for a Chip-seq most of the peaks lie on promoters what approach do you suggest for gene-set enrichment. Why ?

Gene-set enrichment analysis (GSEA) is a method of analyzing gene expression data to determine whether a particular set of genes is overrepresented or underrepresented in a sample compared to a reference set. This can be done using ChIP-seq peaks by identifying the genes that are associated with the peaks and then performing enrichment analysis to identify gene sets that are significantly overrepresented or underrepresented.

To perform GSEA using ChIP-seq peaks, the first step is to annotate the peaks with the nearest gene using a genome annotation database. This can be done using tools such

as ChIPseeker or GREAT. Once the peaks are annotated, the genes associated with the peaks can be identified.

The next step is to perform gene-set enrichment analysis using a statistical method such as hypergeometric testing or Fisher's exact test. This involves comparing the set of genes associated with the ChIP-seq peaks to a reference set, such as the entire genome or a set of genes that are known to be involved in a particular biological process. The goal is to determine whether the genes associated with the ChIP-seq peaks are significantly overrepresented or underrepresented in the reference set.

If most of the ChIP-seq peaks are located on promoters, one approach for gene-set enrichment would be to focus on promoter-associated gene sets, such as those involved in transcriptional regulation or chromatin remodeling. This is because the peaks are more likely to be directly involved in the regulation of gene expression at the transcriptional level. However, it is still important to consider other types of gene sets, such as those involved in post-transcriptional regulation or signal transduction, as these may also be affected by the ChIP-seq peaks.

Overall, GSEA using ChIP-seq peaks can be a powerful tool for identifying the biological processes and pathways that are regulated by a particular transcription factor or chromatin-associated protein. By analyzing the genes associated with the peaks and comparing them to a reference set, it is possible to gain insights into the molecular mechanisms of gene regulation and to identify potential therapeutic targets for disease.

Q14. What is the most complicated step in OLC method for denovo genome assembly. How do we solve it using greedy approach?

The most complicated step in the overlap-layout-consensus (OLC) method for de novo genome assembly is the construction of the overlap graph. This graph represents the overlaps between pairs of reads, and is used to determine the layout of the reads in the genome.

Constructing the overlap graph is challenging because it requires comparing every pair of reads to determine if they overlap, which can be computationally expensive and memory-intensive. In addition, there can be errors in the reads, such as sequencing errors or repeat regions, that can make it difficult to accurately identify the overlaps.

One way to solve this problem is to use a greedy approach, in which the reads are processed one at a time and added to the graph based on their overlap with the existing reads. This can be done using a data structure called a suffix tree, which allows for efficient indexing and searching of the reads.

In the greedy approach, the first read is added to the graph and used as a reference for comparing the next read. The overlap between the two reads is identified, and the second read is added to the graph and connected to the first read. This process is repeated for each subsequent read, with each read being added to the graph based on its overlap with the existing reads.

Once the overlap graph has been constructed, the layout of the reads can be determined by traversing the graph and identifying the longest path through the graph. This path represents the most likely arrangement of the reads in the genome.

Overall, the greedy approach is a computationally efficient way to construct the overlap graph in the OLC method for de novo genome assembly. While it may not be as accurate as more complex algorithms, it can provide a good starting point for further refinement and assembly.

Q15. We have substrings derived from a string

you\_are\_v , are\_very\_ver , good\_guy\_b , ut\_he\_is\_not , uy\_but\_he\_is , very\_very\_goo  
ery\_very\_good\_g

Use De Bruijn graph using  $k = 6$  to rebuild whole string.0

To reconstruct the original string using a De Bruijn graph with  $k=6$ , we need to follow these steps:

1. Create nodes for each k-mer substring in the input strings. In this case,  $k=6$ , so we will create nodes for each 6-mer substring.
2. Create edges between nodes that overlap by  $k-1$  nucleotides. For example, the substrings "you\_ar" and "are\_ve" overlap by "e", so we will create an edge from the node representing "you\_ar" to the node representing "are\_ve".
3. Merge nodes that have identical sequences to simplify the graph. In this case, we can merge the nodes representing "ery\_ve" and "ry\_ve" into a single node.
4. Traverse the graph to generate the reconstructed string. To do this, we start at any node and follow its outgoing edge to the next node. We continue following edges until we reach a node with no outgoing edges, at which point we output the final  $k-1$  nucleotides of the node sequence. We then backtrack to the previous node and repeat the process until we have generated the full string.

The resulting De Bruijn graph looks like this:

rust

Copy code

```
yo -> ou -> u_a -> _ar -> are -> re_v -> _ve -> ver -> ery -> ry_ -> y_v -> _ve ->
ver -> ery -> ry_ -> y_g -> _go -> goo -> ood -> od_ -> d_g -> _gu -> guy -> uy_
-> y_b -> _bu -> but -> ut_ -> t_h -> _he -> he_ -> e_i -> _is -> isn -> sn' ->
n't -> _no -> not
```

Traversing this graph yields the reconstructed string  
"you\_are\_very\_very\_good\_guy\_but\_he\_is\_not".

Q16. what are the limitations of De Bruijn graph ? Explain. What could the source of problems in making eulerian walk in De Bruijn graph ?

De Bruijn graph is a commonly used data structure for de novo genome assembly and sequence analysis. However, it has several limitations, including:

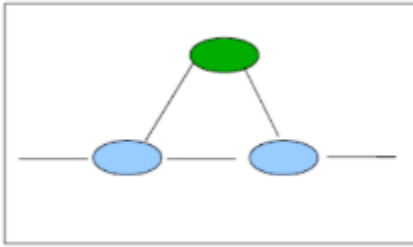
1. Choice of k-mer size: The choice of k-mer size can greatly impact the quality of the De Bruijn graph. If the k-mer size is too small, the graph will be too dense and contain many spurious connections, while if the k-mer size is too large, the graph will be too sparse and miss important connections.
2. Repeats: The presence of repeats in the sequence can cause problems for De Bruijn graph assembly. If the repeat is longer than the k-mer size, the graph will not be able to distinguish between the different occurrences of the repeat.
3. Errors: Sequencing errors can cause problems for De Bruijn graph assembly. Errors in individual k-mers can lead to the creation of spurious edges in the graph.
4. Complexity: As the size and complexity of the input sequence increases, so does the size and complexity of the De Bruijn graph. This can make it difficult to store and manipulate the graph.

The source of problems in making an Eulerian walk in the De Bruijn graph is the presence of unbalanced nodes, also known as "tips" or "dead ends". These nodes have either an incoming or outgoing edge, but not both. If a graph contains unbalanced nodes, it cannot have an Eulerian walk, which is required for genome assembly. To solve this problem, a greedy approach is used to remove the tips or dead ends from the graph. This involves iteratively trimming the ends of the graph until all nodes are balanced. However, this process can lead to the loss of important information and result in fragmented assemblies.

Q17. During De Bruijn graph based genome assembly if we get a spurious bubble in between

Shown as

What does it mean ? When can it happen ?



A spurious bubble in a De Bruijn graph can occur during genome assembly when a repetitive sequence or a sequencing error leads to the creation of multiple branches that converge back into the graph at a later point. In the context of genome assembly, a bubble represents a region of the genome where there are multiple possible paths through the graph, and determining the correct path can be challenging.

In the case of a spurious bubble, one of the paths through the bubble is likely incorrect and does not represent the true sequence of the genome. This can lead to errors in the final assembly, as the incorrect path may be selected during the construction of the contigs or scaffolds.

Spurious bubbles can be caused by a variety of factors, such as sequencing errors, repeats, or alternative splicing events in RNA-seq data. To address these issues, various algorithms have been developed to identify and remove spurious bubbles from the assembly graph, such as bubble popping or tip clipping. These algorithms rely on heuristics to identify and remove problematic regions of the graph, with the goal of improving the accuracy and contiguity of the final assembly.

Q18. What is the difference between tSNE, MDS and PCA based methods for dimension reduction.

If we want to perform linear regression on data-points for classification in lower dimension, which approach would be good tSNE, MDS or PCA.

If we want to perform clustering using the following approaches which dimension reduction method would be good

- a) KNN based
- b) Self organising maps
- c) Hierarchical clustering

tSNE (t-distributed stochastic neighbor embedding), MDS (multidimensional scaling), and PCA (principal component analysis) are all methods for dimension reduction, but they differ in the underlying approach and the type of data they work best with.

PCA is a linear method that finds a set of orthogonal components that explain the maximum amount of variance in the data. It is best suited for data that has a clear linear structure and is normally distributed.

MDS is a non-linear method that attempts to preserve the pairwise distances between data points in a lower-dimensional space. It is best suited for data that has a non-linear structure.

tSNE is a non-linear method that is particularly effective at preserving the local structure of the data, while also reducing the dimensionality. It is often used for visualization of high-dimensional data.

If we want to perform linear regression on data points for classification in a lower dimension, PCA would be a good choice because it is a linear method and can find the components that explain the maximum variance in the data.

If we want to perform clustering using KNN-based, self-organizing maps or hierarchical clustering approaches, tSNE would be a good choice because it can effectively preserve the local structure of the data, which is important for clustering.

Q19. If the function for likelihood of neighborhood in tSNE method is like distribution function of following how the

dimension reduction performance would look like

- a) Uniform distribution
- b) Student T distribution

Explain with equation and example of outliers data-points

t-SNE is a non-linear dimensionality reduction technique that maps high-dimensional data into a low-dimensional space (usually 2D or 3D) while preserving local structure and sometimes global structure as well. The main idea behind t-SNE is to map data points to a lower-dimensional space in such a way that the pairwise distances between the points in the lower-dimensional space closely match the pairwise distances in the high-dimensional space.

The likelihood of the neighborhood in t-SNE is modeled using a Gaussian distribution with a certain variance (perplexity), which is used to calculate the probability of a point belonging to its neighborhood. The t-SNE optimization objective is to minimize the Kullback-Leibler divergence between the high-dimensional probability distribution (computed using a Gaussian kernel) and the low-dimensional probability distribution (computed using a Student's t-distribution kernel).

The choice of the distribution function used to model the likelihood of the neighborhood in t-SNE can have an impact on the performance of the method, especially when dealing with outliers.

a) If the distribution function for the likelihood of the neighborhood is a uniform distribution, then the dimension reduction performance would be poor. This is because a uniform distribution assumes that all points are equally likely to be neighbors, which is unlikely to be true in most cases. In the case of outliers, the uniform distribution would assign equal probabilities to all points in the neighborhood, which can lead to poor clustering and visualization results.

b) If the distribution function for the likelihood of the neighborhood is a Student's t-distribution, then the dimension reduction performance would be better. The Student's t-distribution has heavier tails compared to a Gaussian distribution, which makes it more robust to outliers. This means that points that are far away from the cluster center (outliers) will have a lower probability of being included in the neighborhood, which can lead to better clustering and visualization results.

In summary, the choice of the distribution function used to model the likelihood of the neighborhood in t-SNE can have an impact on the performance of the method, especially when dealing with outliers. A Student's t-distribution is often preferred over a Gaussian distribution when the data contains outliers, as it is more robust to extreme values.