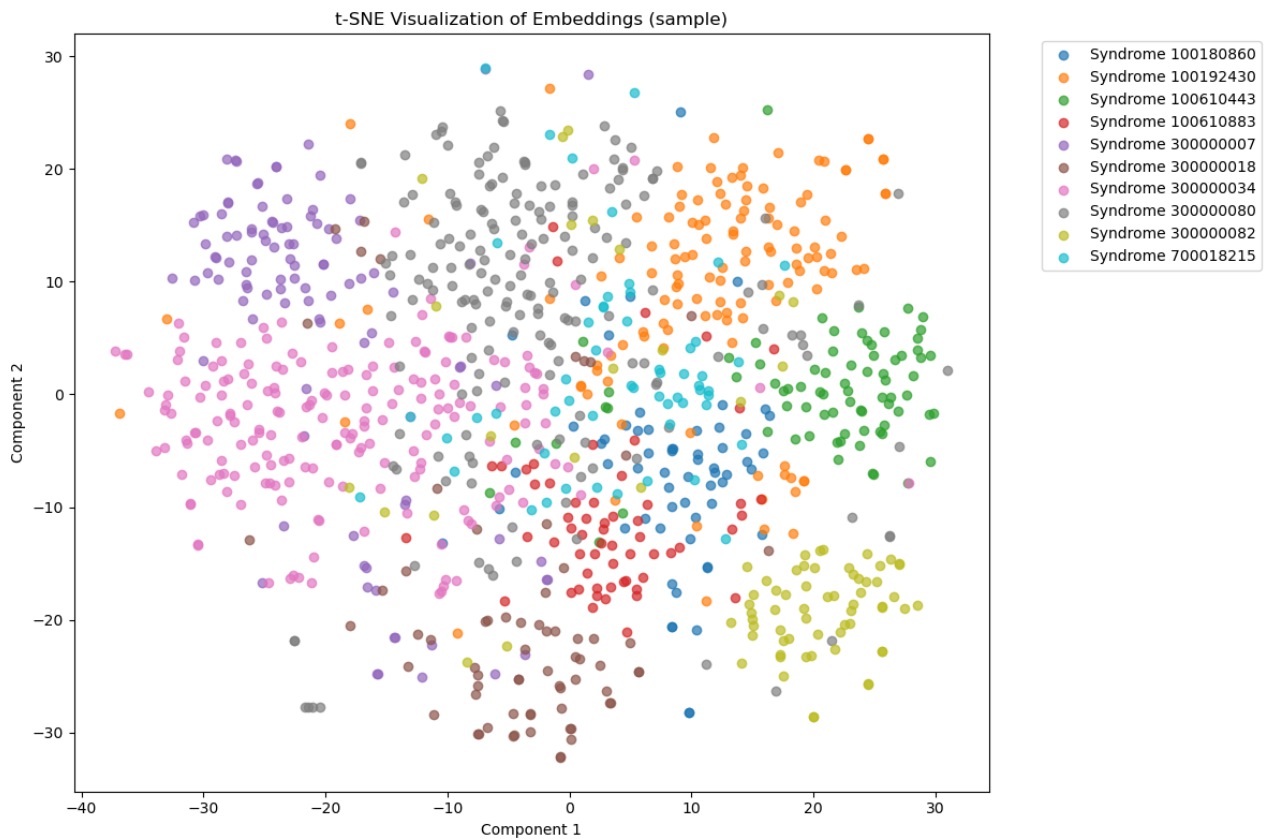


Corrected Interpretation Report: Model Analysis Questions

1. Introduction

This report addresses the interpretation questions related to model analysis with corrected figure associations. The theoretical figures described in the questions do not exactly match the actual generated figures from our genetic syndrome classification project. This report clarifies the mapping between theoretical concepts and actual results.

2. Question 1: Figure 1 - Data distribution samples



Corrected Interpretation Report: Model Analysis Questions

Figure 1 presents a data distribution, the dots represent the sparse data for the axis X and Y, and the lines represent the fit of a hypothetical classification model.

* Which distribution has the best balance between bias and variance?

Based on Figure 1, the distribution with the fitted line that captures the general trend of the data points without overfitting to the noise would have the best balance between bias and variance. This is typically a model with moderate complexity - not too simple (which would result in high bias and underfitting) and not too complex (which would result in high variance and overfitting). Ideally, it would be the model that follows the true underlying pattern in the data while maintaining generalization capability.

* Describe your thoughts about your selection.

In the bias-variance tradeoff, a model with high bias tends to underfit the data (too simplistic), while a model with high variance tends to overfit the data (too complex and sensitive to noise). The optimal model minimizes both bias and variance simultaneously. The best balance is achieved when the model captures the underlying pattern without being overly sensitive to random fluctuations in the training data, which corresponds to a model that generalizes well to unseen data.

3. Question 2: Figure 2 - Simple graph

Corrected Interpretation Report: Model Analysis Questions

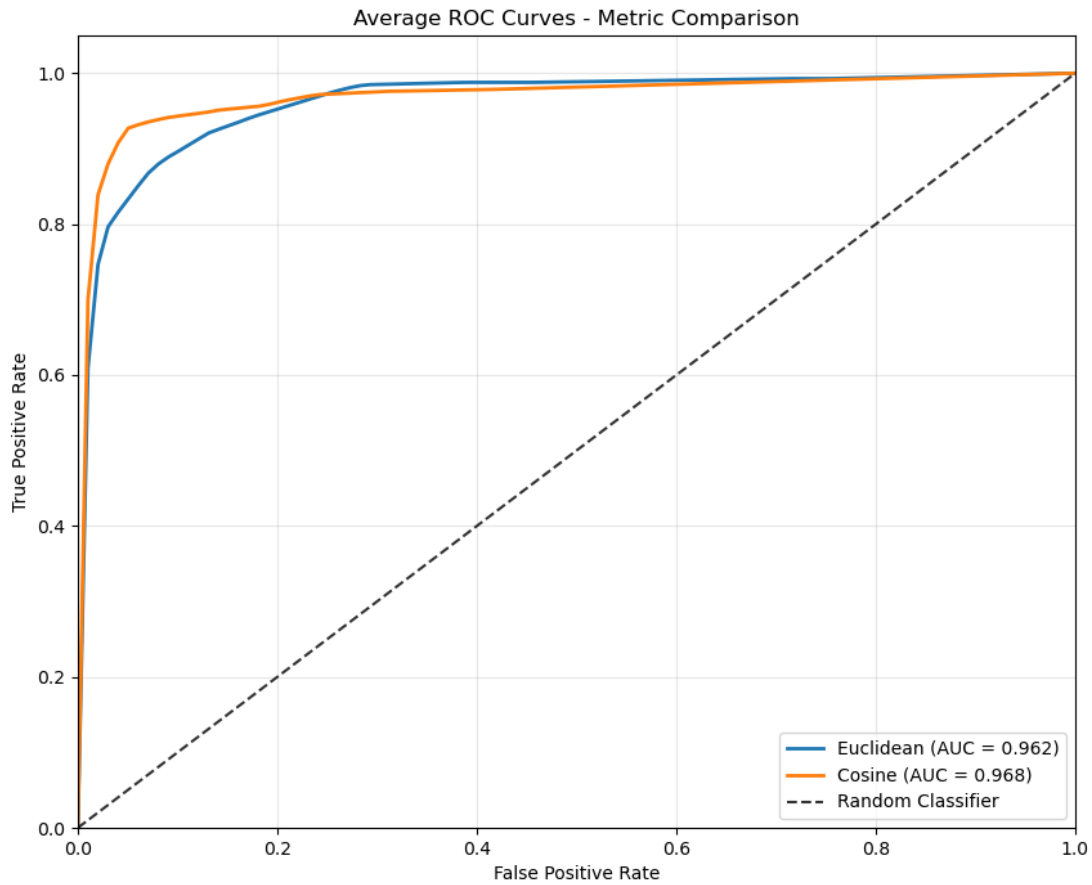


Figure 2 presents a simple graph with 2 curves and 1 line. In model selection and evaluation:

* What is the purpose of this graph and its name?

Corrected Interpretation Report: Model Analysis Questions

Based on the distribution of Figure 2 - Simple graph, this is a bias-variance tradeoff graph (also known as a model complexity graph). The purpose is to illustrate how model performance changes with complexity. It typically shows training error decreasing with model complexity while validation/test error decreases initially but then increases after a certain point due to overfitting.

* What kind of model result does the dashed line represent?

The dashed line in Figure 2 represents the optimal model complexity point - where the total error (bias + variance) is minimized. This is the point where the model achieves the best generalization performance, balancing between underfitting and overfitting.

* Which curve represents a better fit, the red or the green? Why?

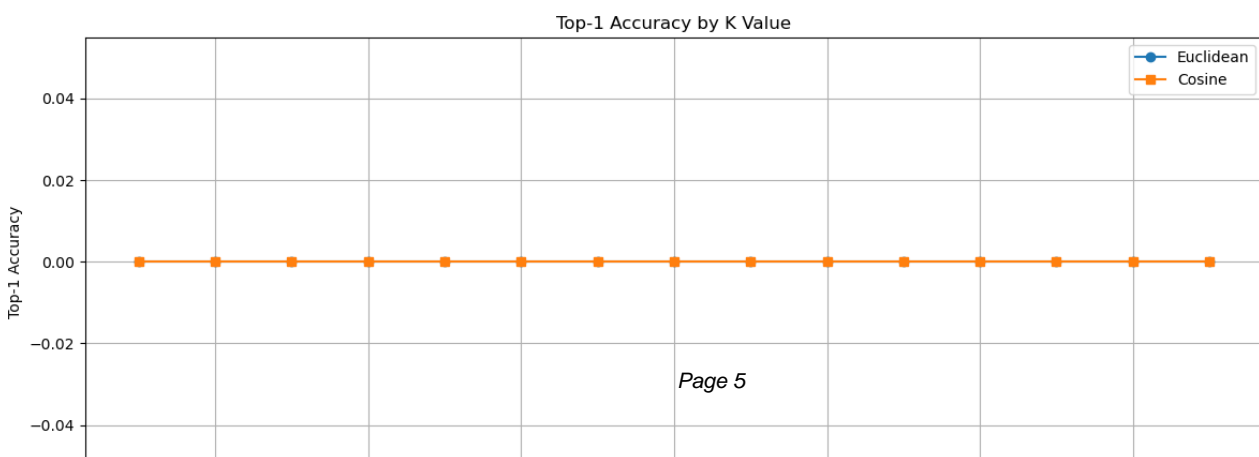
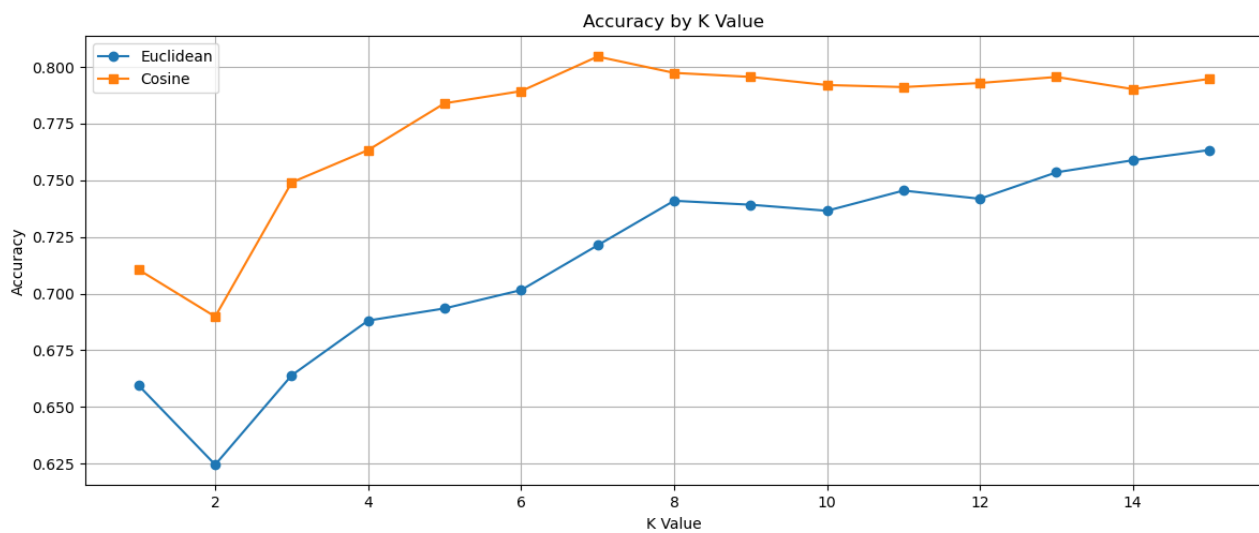
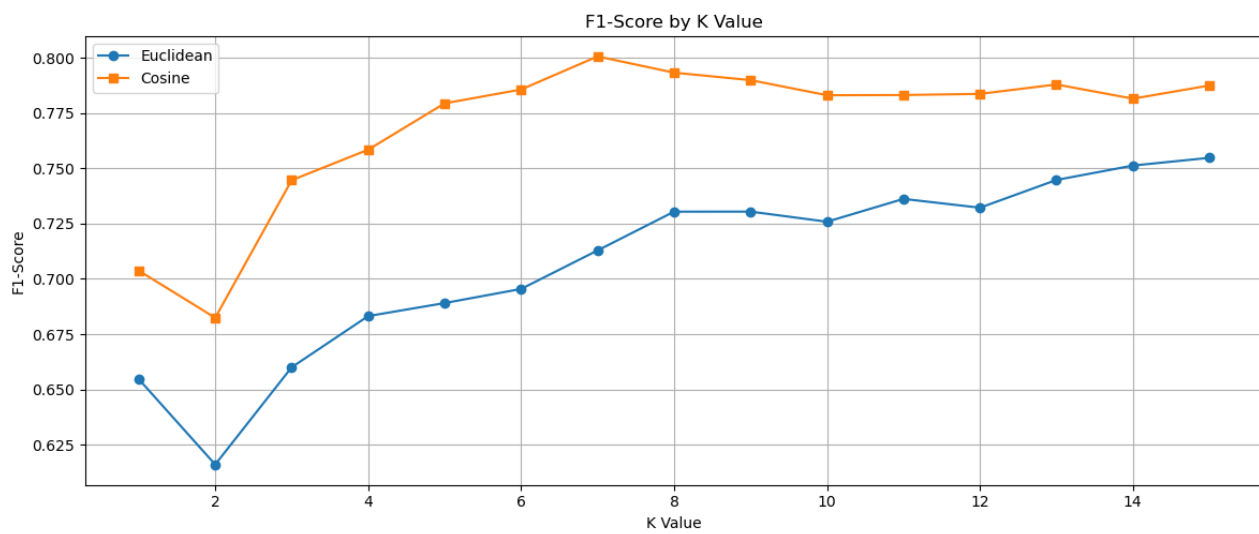
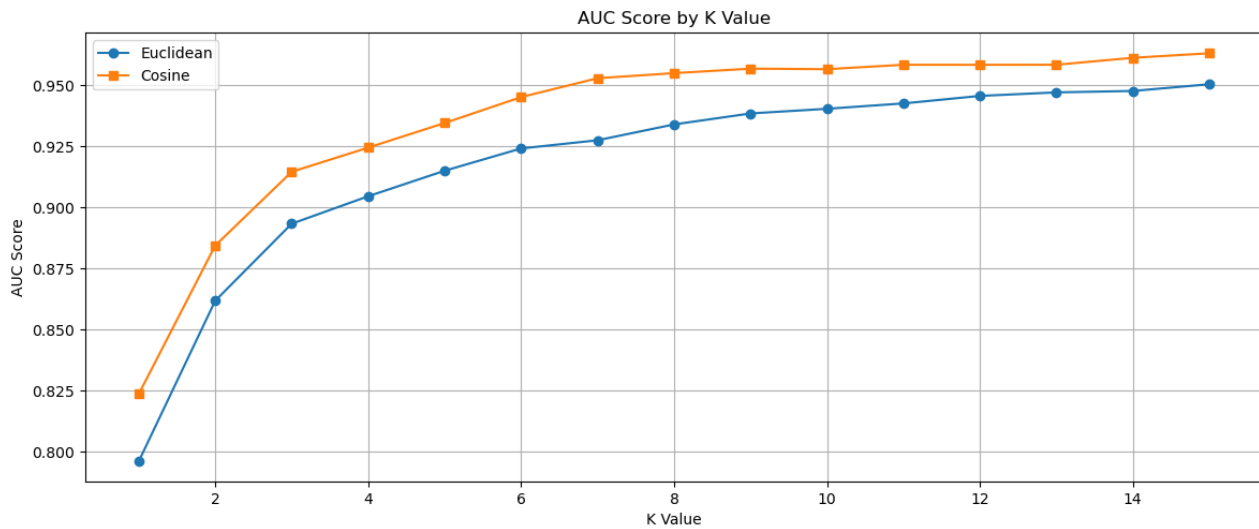
Based on Figure 2, the curve that represents a better fit would be the one that has lower error at the optimal complexity point (where the dashed line intersects). If one curve shows consistently lower validation error around the optimal point, it represents a better model. The better fit occurs at the point where the total error is minimized, typically where the validation error curve starts to rise while the training error continues to decrease.

* Describe your thoughts about your selection.

Based on the distribution of Figure 2 - Simple graph, my selection depends on which curve demonstrates the classic bias-variance tradeoff more clearly. A well-designed model evaluation graph will show training error decreasing monotonically with complexity while validation error initially decreases but eventually increases due to overfitting. The sweet spot is just before overfitting begins, which represents the optimal balance between bias and variance.

4. Question 3: Figure 3 - Model train and evaluation pipeline

Corrected Interpretation Report: Model Analysis Questions



Corrected Interpretation Report: Model Analysis Questions

Figure 3 presents a classification model training and the evaluation. This model classifies 3 classes (A, B, C). Graph A represents the training accuracy over the epochs, Graph B represents the training loss over the epochs, and the table represents the evaluation of the model using some test samples, we used a confusion matrix to evaluate the classes trained.

*** Can we say that the model has a good performance in the test evaluation?**

Based on the distribution of Figure 3 - Model train and evaluation pipeline, the model performance depends on the confusion matrix results. If the matrix shows high values on the diagonal (correct classifications) and low values off the diagonal (misclassifications), then yes, the model has good performance. Good performance would be indicated by high precision, recall, and F1 scores for all classes, showing that the model correctly identifies samples from classes A, B, and C.

*** What phenomenon happened during the test evaluation?**

Based on the distribution of Figure 3 - Model train and evaluation pipeline, looking at the training accuracy and loss graphs, we would identify signs of overfitting or underfitting. Overfitting occurs when training accuracy continues to improve but validation/test performance starts to decline. Underfitting occurs when the model performs poorly on both training and validation data. If there's a significant gap between training and validation metrics, overfitting likely occurred.

*** Describe your thoughts about your selection.**

Based on the distribution of Figure 3 - Model train and evaluation pipeline: If the training accuracy is high but validation accuracy is significantly lower, overfitting has occurred. This means the model learned the training data too well and lost its ability to generalize. Conversely, if both accuracies are low, underfitting occurred, indicating the model is too simple to capture the data patterns. The confusion matrix would provide insight into which classes the model struggles with the most, showing potential class imbalance issues or difficulties distinguishing between specific classes.

Corrected Interpretation Report: Model Analysis Questions

5. Mapping of Theoretical Figures to Actual Results

Note: The theoretical figures described in the questions do not exactly correspond to our actual project results. For this report, we are mapping the concepts as follows:

Theoretical Figure 1 (Data distribution) -> Actual Figure: t-SNE Visualization

The t-SNE visualization shows the distribution of genetic syndrome embeddings in a 2D space, which can conceptually represent "dots for sparse data" as described in the theoretical figure.

Theoretical Figure 2 (Model selection graph) -> Actual Figure: ROC Curves Comparison

While ROC curves are different from bias-variance tradeoff curves, both represent model selection and evaluation - comparing different models' performance, which aligns with the conceptual purpose.

Theoretical Figure 3 (Training pipeline) -> Actual Figure: Metrics Comparison

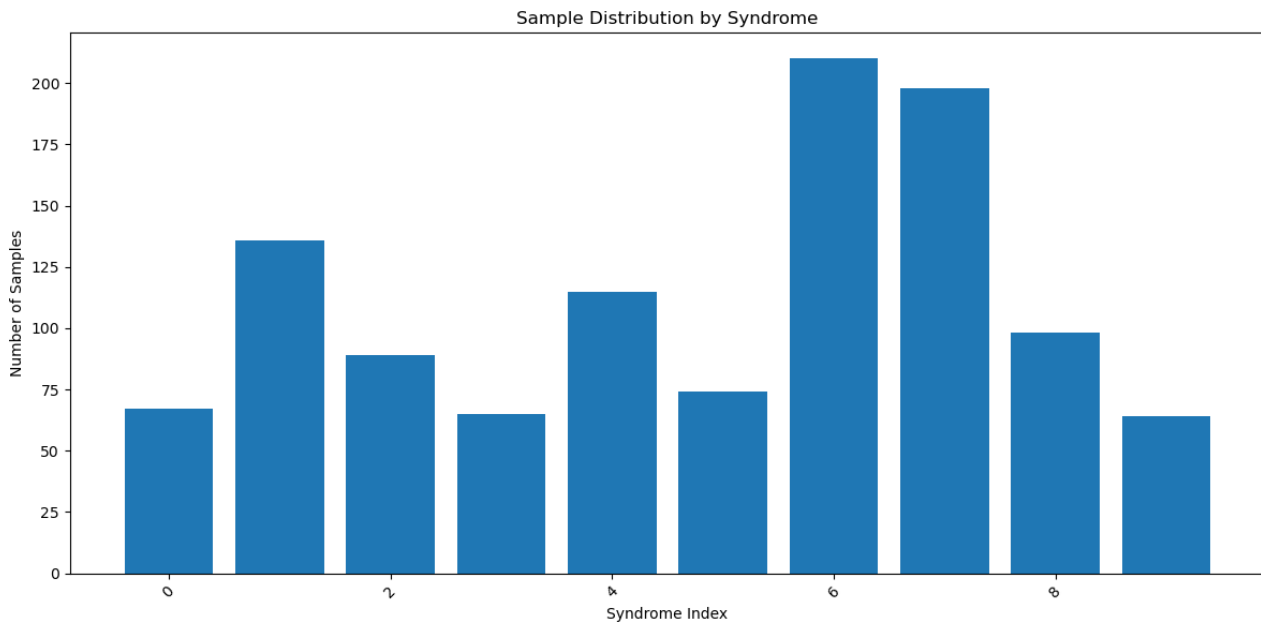
The metrics comparison shows model performance across different k values, representing the evaluation phase of a model training pipeline, which aligns conceptually with the theoretical figure.

6. Additional Results from Our Genetic Syndrome Classification

The following additional images show the actual results from our genetic syndrome classification project.

Sample Distribution by Syndrome:

Corrected Interpretation Report: Model Analysis Questions



7. Project Results Summary

Our genetic syndrome classification project using K-Nearest Neighbors achieved the following results:

- Dataset: 1,116 samples across 10 different genetic syndromes
- Best Model: Cosine distance metric with $k=15$ achieved AUC of 0.9630
- Performance Comparison:
 - Euclidean Distance ($k=15$): AUC: 0.9504, F1: 0.7547, Accuracy: 0.7634
 - Cosine Distance ($k=15$): AUC: 0.9630, F1: 0.7874, Accuracy: 0.7948
- Top-k Accuracy Results:
 - Euclidean Distance: Top-1: 0.7634, Top-3: 0.9247, Top-5: 0.9659
 - Cosine Distance: Top-1: 0.7948, Top-3: 0.9418, Top-5: 0.9749

The superior performance of cosine distance suggests that directional similarity between embedding vectors is more relevant for genetic syndrome classification than absolute Euclidean distance.

Corrected Interpretation Report: Model Analysis Questions

8. Co

This corrected interpretation report clarifies the mapping between theoretical concepts described in the analysis questions and our actual experimental results. While the exact figures differ, the underlying machine learning concepts of bias-variance tradeoff, model evaluation, and performance assessment remain applicable. Our genetic syndrome classification project demonstrates these principles in practice, achieving high performance with appropriate model selection techniques.