

# Pattern Recognition and Machine Learning - Lab 01

## Member Contribution

This assignment is generated based on the list of requirements in the assignment file.

Tasks	Member
Data Loading & Preprocessing	Đào Xuân Tân
Linear Model - Logistic Regression	Lý Quang Thắng
Support Vector Machine (SVM)	Trần Lê Hữu Vinh
Curse of Dimensionality & Dimensionality Reduction	Lê Phú Trường
Model Evaluation & Comparison	Hồ Quốc Nhân Hòa
Compile and write the report	Lê Phú Trường, Trần Lê Hữu Vinh

## Table of contents

- [Pattern Recognition and Machine Learning - Lab 01](#)
  - [Member Contribution](#)
  - [Table of contents](#)
- [1. Introduction](#)
  - [The goal of this assignment](#)
  - [Workflow of this assignment](#)
    - [i. Data Loading & Preprocessing](#)
      - [Step 1: Download Dataset](#)
      - [Step 2: Data Splitting](#)
      - [Step 3: Data Transformation](#)
      - [Step 4: Exploratory Data Analysis \(EDA\)](#)
    - [ii. Linear Model - Logistic Regression](#)
      - [Step 1: Train the Model](#)
      - [Step 2: Model Evaluation](#)

- *iii.* Support Vector Machine (SVM)
    - Step 1: Train SVM with Linear Kernel
    - Step 2: Train SVM with RBF Kernel
    - Step 3: Performance Comparison
  - *iv.* Curse of Dimensionality & Dimensionality Reduction
    - Step 1: Curse of Dimensionality Discussion
    - Step 2: Apply Dimensionality Reduction
    - Step 3: Retrain Models on Reduced Data
    - Step 4: Analyze Trade-offs
  - *v.* Model Evaluation & Comparison
    - Step 1: Model Comparison on Full vs. Reduced Data
    - Step 2: Visualization
    - Step 3: Discussion
- 2. Discussion of model performance on full and reduced data
  - Model Performance Comparison
    - Logistic Regression
      - Before Dimensionality Reduction
      - After Dimensionality Reduction
      - Conclusion
    - Support Vector Machine (SVM)
      - Before Dimensionality Reduction
      - After Dimensionality Reduction
      - Conclusion
  - Model Performance In Full vs. Reduced Data
    - Full-Dimensional Data:
      - Observations for Full-Dimensional Data:
      - Summary
    - Reduced-Dimensional Data:
      - Observations for Reduced-Dimensional Data:
      - Summary
    - Overall Summary
    - Conclusion
    - Extending the analysis
      - RBF Kernel SVC
      - Linear Kernel SVC
      - Logistic Regression
      - Summary:
- 3. The Curse of Multidimensionality and Its Implications for Model Performance

- Introduction to the Curse of Dimensionality
- Impact on Machine Learning Models
  - Reduced Model Performance and Generalization
  - Difficulty in Finding Meaningful Patterns
  - Increased Computational Complexity
  - Increased Risk of Overfitting
  - Challenges with Distance Metrics
- Reflect on the trade-offs between dimensionality reduction and model performance
  - Model Efficiency and Computational Cost
  - Impact on Accuracy and Predictive Performance
  - Summary of Trade-offs
  - Conclusion
- 4. Option for choosing the best model
  - Logistic Regression
  - Linear SVM (Support Vector Machine)
  - RBF SVM (Radial Basis Function Kernel)
  - Summary

# 1. Introduction

In this assignment, we explore the application of Logistic Regression and Support Vector Machines (SVM) for classifying fashion attires from **Fashion-MNIST** dataset which was provided by Zalando Research. This dataset consists of 70.000 grayscale images, each of size  $28 \times 28$  pixels, include 784 features. Fashion-MNIST is categorized into 10 distinct classes represented by 10 fashion items:

- T-shirts/top
- Trouser
- Pullover sweater
- Dress
- Coat
- Sandal
- Shirt
- Sneaker
- Bag
- Ankle boot



Figure 1: Sample images from the Fashion-MNIST dataset, showcasing the diversity of fashion items in the dataset

This diversity in fashion categories presents unique challenges for classification algorithms, as some items may share similar features, leading to confusion in model predictions.

# The goal of this assignment

- Applying and evaluating the performance of Logistic Regression and SVC in the context of image classification along with handling training time in the large dataset problem.
- Understanding how high-dimensional data affects these machine learning models. Specifically, we aim to examine how dimensionality reduction techniques like PCA can mitigate the "**curse of dimensionality**" problem which often stifles model's performance in such large feature spaces. We can analyze how these approaches cope with the complexities of high-dimensional feature spaces. To sum up, the ultimate goal is to enhance model performance while understanding the trade-offs involved in dimensionality reduction and the inherent challenges posed by the curse of dimensionality.

## Workflow of this assignment

This workflow includes steps for data loading, preprocessing, model training with logistic regression and SVM, dimensionality reduction, and performance evaluation using the Fashion-MNIST dataset. The process leverages key machine learning concepts and algorithms to provide insights into model performance and trade-offs.

### i. Data Loading & Preprocessing

#### Step 1: Download Dataset

- **Source:** [Fashion-MNIST GitHub Repository](#).
- Download the dataset and load it into the workspace.

#### Step 2: Data Splitting

- **Purpose:** Split the data into training and test sets.
- Use an appropriate ratio for training and testing (e.g., 80/20 or 70/30 split).

#### Step 3: Data Transformation

- **Normalization:** Scale pixel values to the range  $[0, 1]$ .
- **Standardization:** Standardize feature values if needed.
- **Rotation, Scaling, and Cropping:** Apply as necessary to augment the dataset.

## Step 4: Exploratory Data Analysis (EDA)

- Understand the class distribution.
- Visualize some sample images from each class.
- Plot the distribution of pixel intensities to assess variance.

## *ii.* Linear Model - Logistic Regression

### Step 1: Train the Model

Train a **Logistic Regression** classifier on all 784 features of the Fashion-MNIST dataset.

### Step 2: Model Evaluation

- **Metrics:**
  - Accuracy;
  - Precision;
  - Recall;
  - F1 Score;
  - Confusion Matrix.
- **Analysis:**
  - Identify patterns in model performance.
  - Determine classes that are often misclassified and discuss possible reasons.

## *iii.* Support Vector Machine (SVM)

### Step 1: Train SVM with Linear Kernel

Train an SVM model using a **linear kernel** on the same dataset.

### Step 2: Train SVM with RBF Kernel

- Train an SVM model using an **RBF (Radial Basis Function) kernel**.
- Compare results with the **linear kernel SVM** and **logistic regression model**.

### Step 3: Performance Comparison

Evaluate and compare computation time for training and prediction across all models, with particular focus on **linear vs. RBF kernel SVMs**.

## *iv.* **Curse of Dimensionality & Dimensionality Reduction**

### **Step 1: Curse of Dimensionality Discussion**

Explain the concept of the curse of dimensionality and its implications for model performance in high-dimensional spaces.

### **Step 2: Apply Dimensionality Reduction**

Utilize **Principal Component Analysis (PCA)** to reduce the dimensionality from **784 features to a lower number** (100 for example).

### **Step 3: Retrain Models on Reduced Data**

- Retrain **Logistic Regression** and **SVM (linear or RBF kernels)** on the reduced-dimensional dataset.
- Compare the performance of these models with those trained on the full dataset.

### **Step 4: Analyze Trade-offs**

- Reflect on the trade-offs between dimensionality reduction and model performance.
- Discuss the impact on accuracy, computation time, and memory efficiency.

## *v.* **Model Evaluation & Comparison**

### **Step 1: Model Comparison on Full vs. Reduced Data**

- Accuracy;
- Confusion Matrix;
- Precision, Recall, F1 Score for each model and setting.

### **Step 2: Visualization**

- Accuracy comparison bar chart for each model (logistic regression, linear SVM, RBF SVM).
- Confusion matrices for each model and dataset configuration.
- Precision, Recall, and F1 score for each model.

### **Step 3: Discussion**

- **Pros and Cons** of Each Model: Considerations such as accuracy, computation time, and ability to handle high-dimensional data.
- Reflect on optimal models for different use cases based on observed results.

# 2. Discussion of model performance on full and reduced data

## Model Performance Comparison

### Logistic Regression

#### Before Dimensionality Reduction

The logistic regression model applied to the dataset before PCA achieved:

- Accuracy, Precision, Recall, and F1 Score of about 0.85.
- The confusion matrix shows that certain classes, such as Shirt and Coat, are frequently misclassified, indicating that these categories might have similar features that the model struggles to distinguish without dimensionality reduction.

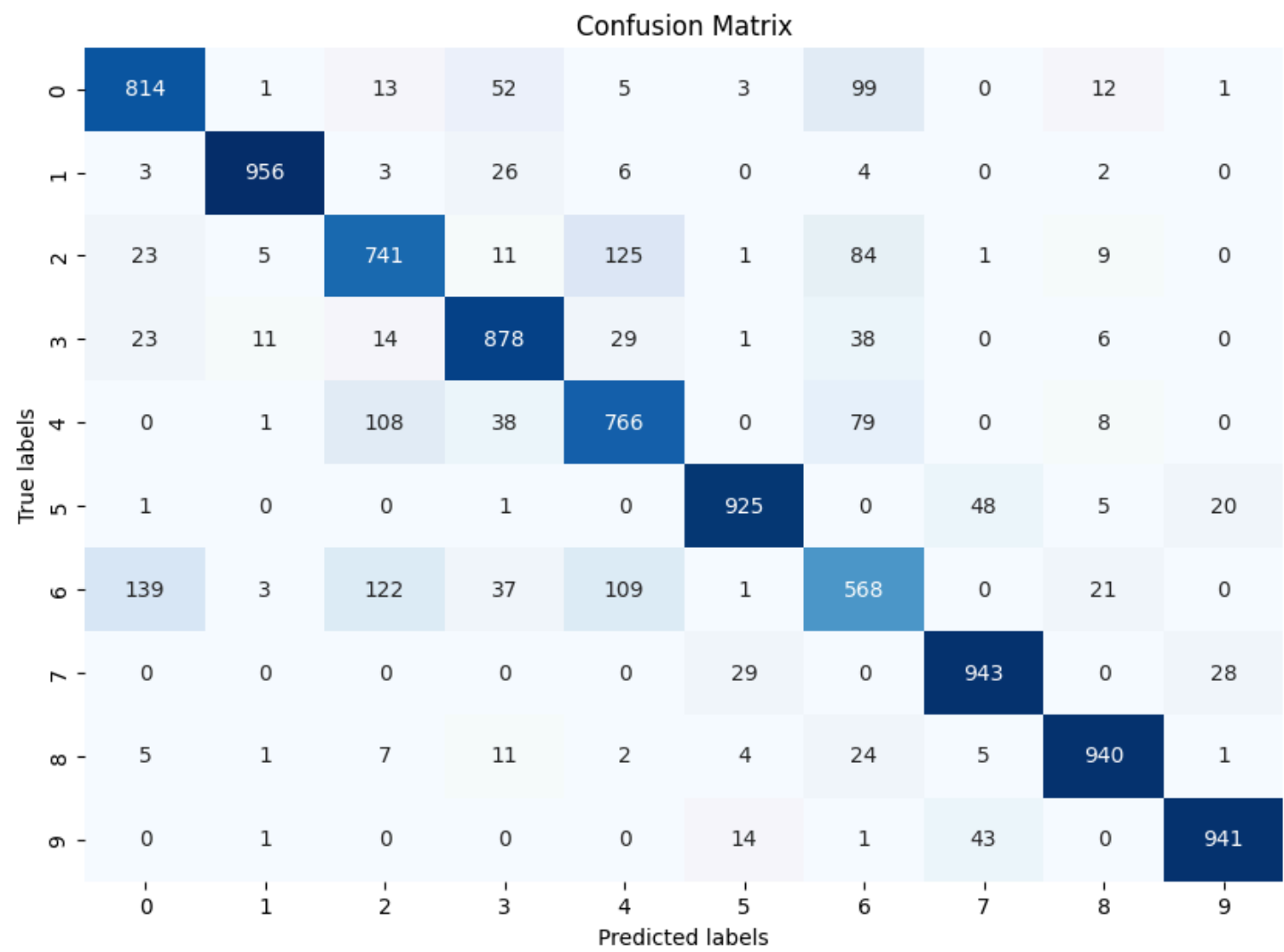




Figure 2: Confusion matrix for logistic regression model before dimensionality reduction

- The misclassification images also highlight that the model sometimes confuses classes with visually similar outlines or textures (e.g., Coat as Shirt or Pullover), likely due to overlapping feature spaces.

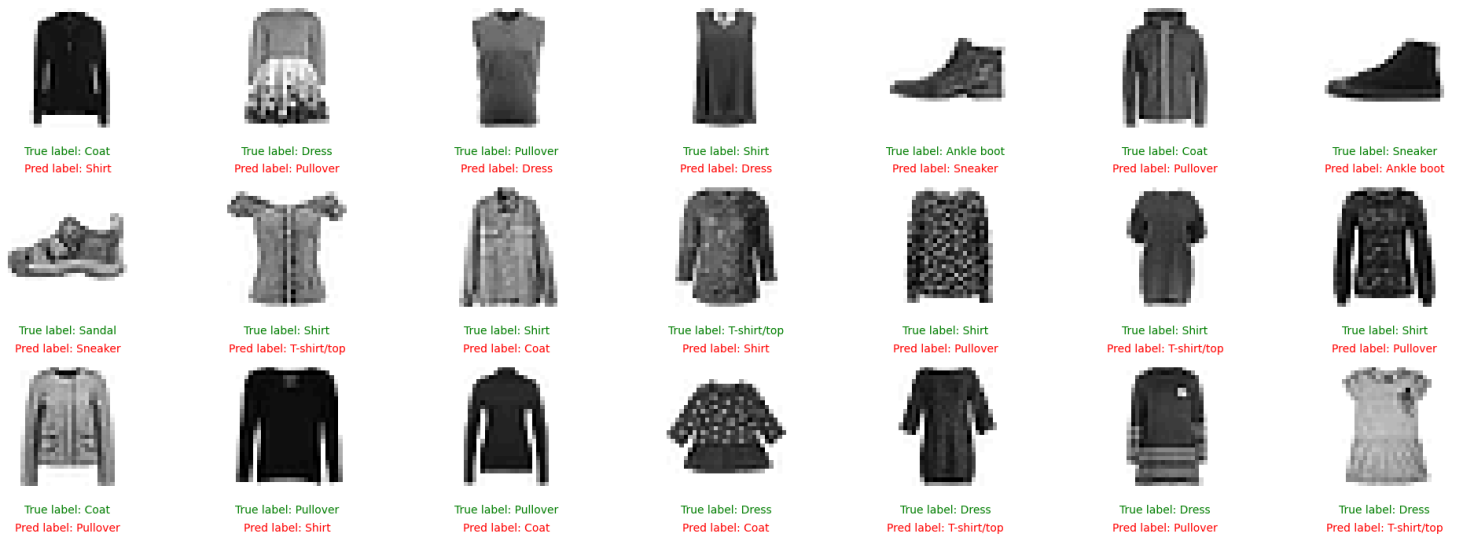


Figure 3: Misclassified images by the logistic regression model before dimensionality reduction

## After Dimensionality Reduction

With PCA applied, the logistic regression model shows:

- Slightly lower accuracy ( $\approx 0.839$ ), along with minor drops in precision and F1 score. However, the scores remain consistent, showing that dimensionality reduction did not severely impact model performance.
- The confusion matrix still displays misclassifications, but PCA has preserved the essential features for most classes, allowing the model to make similar predictions with reduced data dimensionality.

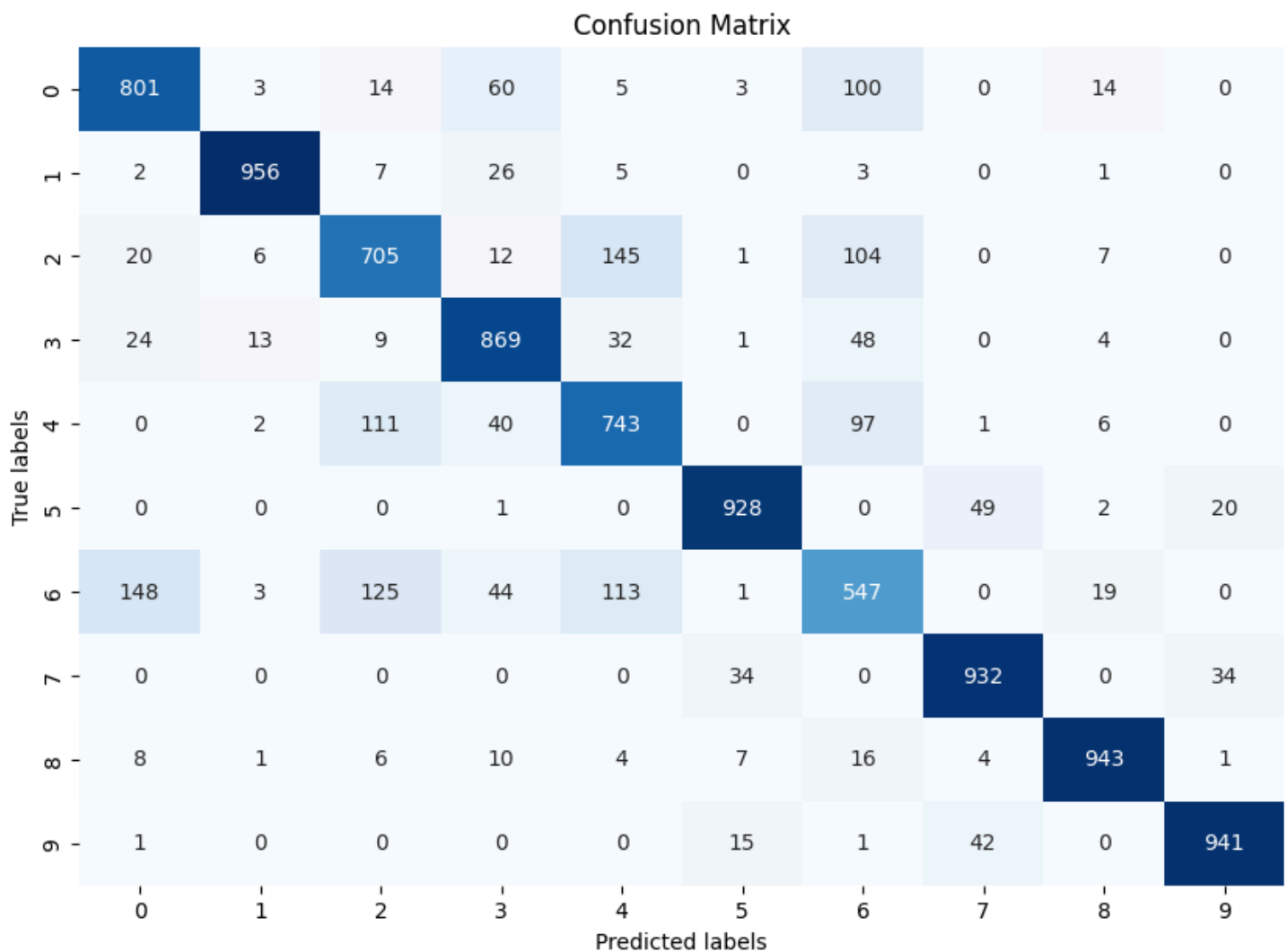


Figure 4: Confusion matrix for logistic regression model after dimensionality reduction

- The misclassified images after PCA reflect similar patterns of confusion as before PCA, but the model's performance is more efficient due to lower complexity.

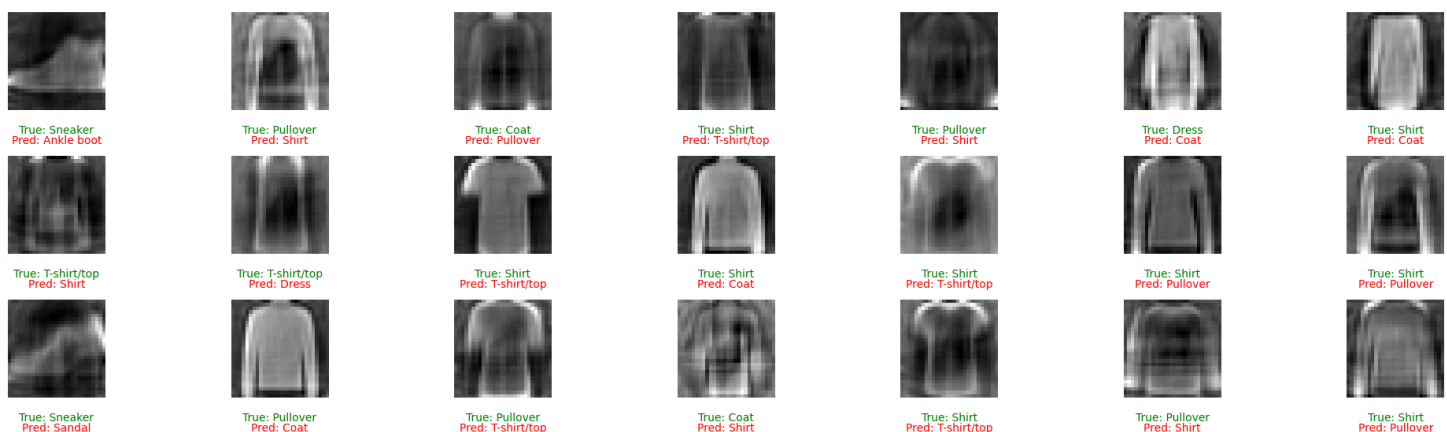


Figure 5: Misclassified images by the logistic regression model after dimensionality reduction

## Conclusion

- **Performance Consistency:** The model's accuracy, precision, recall, and F1 score remain close to the original values after applying PCA with 100 components, indicating that PCA successfully retains the essential features for classification.
- **Efficiency:** The model with PCA is computationally more efficient, processing fewer dimensions without significant loss in classification performance. This makes PCA a valuable step, especially for large datasets where dimensionality reduction can speed up training and inference times.
- **Error Patterns:** The same classes (like Shirt and Coat) are often misclassified both before and after PCA, suggesting that these classes have overlapping features that logistic regression alone struggles to separate. Addressing this might require feature engineering or a more advanced model to better capture subtle differences.

In summary, applying PCA with 100 components allows the model to work more efficiently while maintaining similar accuracy, making it a practical choice for reducing dimensionality in this context.

## Support Vector Machine (SVM)

### Before Dimensionality Reduction

- Best Validation Accuracy: 0.90275.
- Overall Metrics: Accuracy of 0.8995, Precision of 0.8993, Recall of 0.8995, and F1 Score of 0.8993.
- Confusion Matrix Insights: The performance across classes is more consistent, with fewer misclassifications compared to the PCA-reduced model. The lack of PCA allows the model to utilize the full feature set, which may help it capture more detailed variations within each class. This results in higher accuracy, precision, recall, and F1 scores, indicating better classification ability.

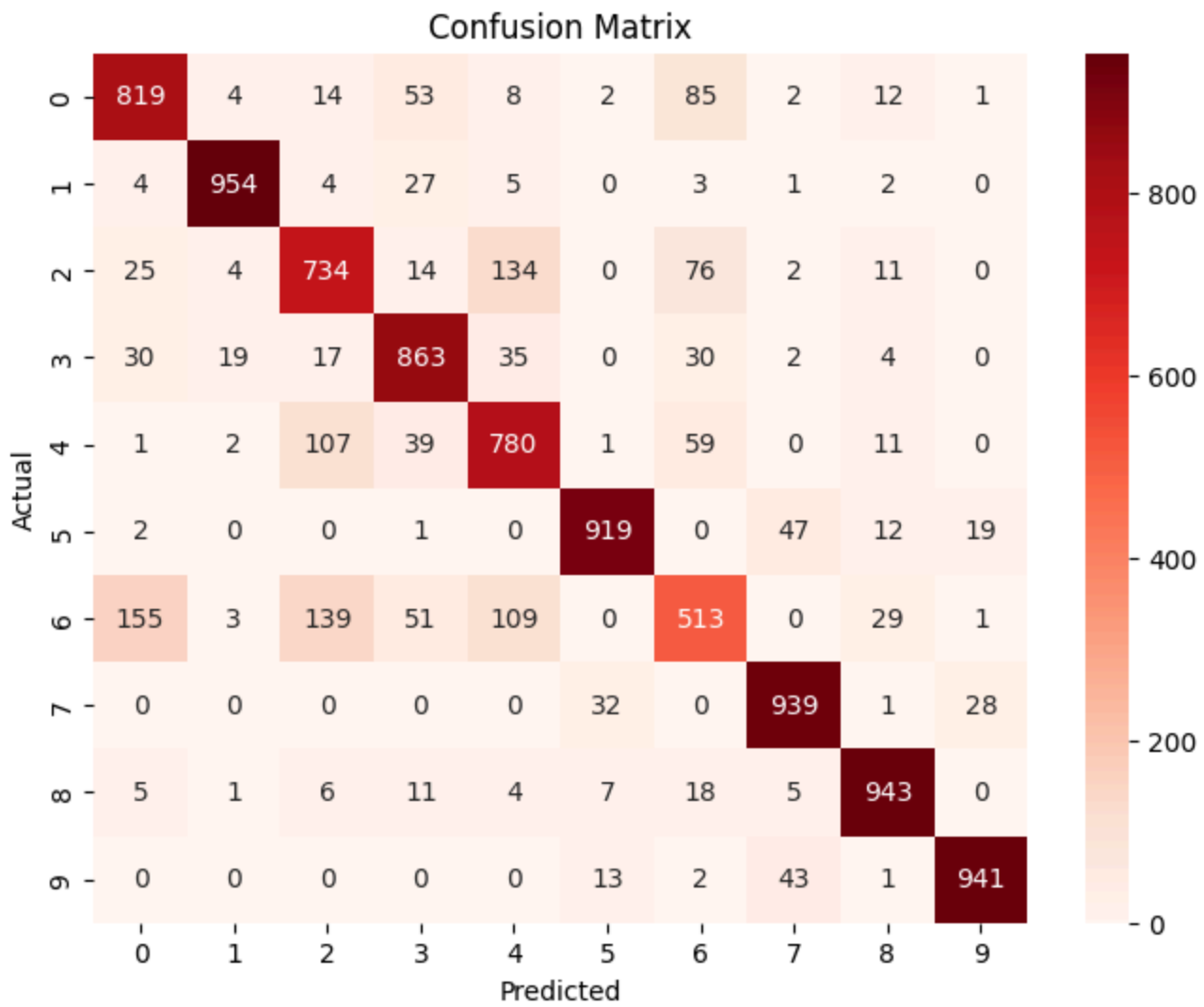


Figure 6: Confusion matrix for Linear Kernel SVC model before dimensionality reduction

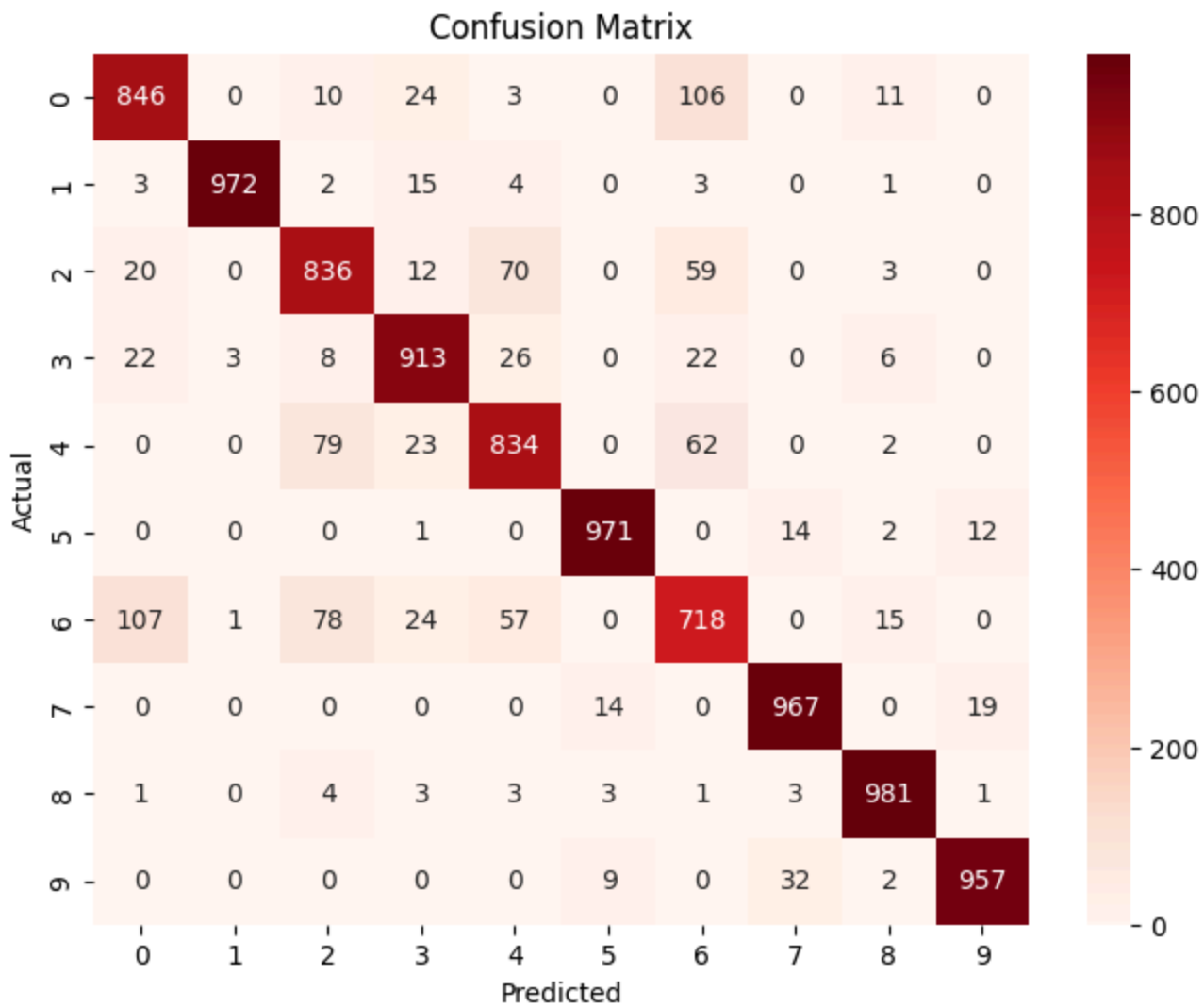


Figure 7: Misclassified images by the RBF SVC model before dimensionality reduction

## After Dimensionality Reduction

- Best Validation Accuracy: 0.8965.
- Test Accuracy: 0.8923.
- Overall Metrics: Accuracy of 0.8923, Precision of 0.8919, Recall of 0.8923, and F1 Score of 0.8920.
- Confusion Matrix Insights: The model performs reasonably well across most classes. However, certain classes still show misclassifications, especially where visual similarity between classes might be high. For example, items like Shirt and Coat are frequently confused, which might indicate that the reduced feature set with PCA still captures relevant details but may lose some finer distinctions.

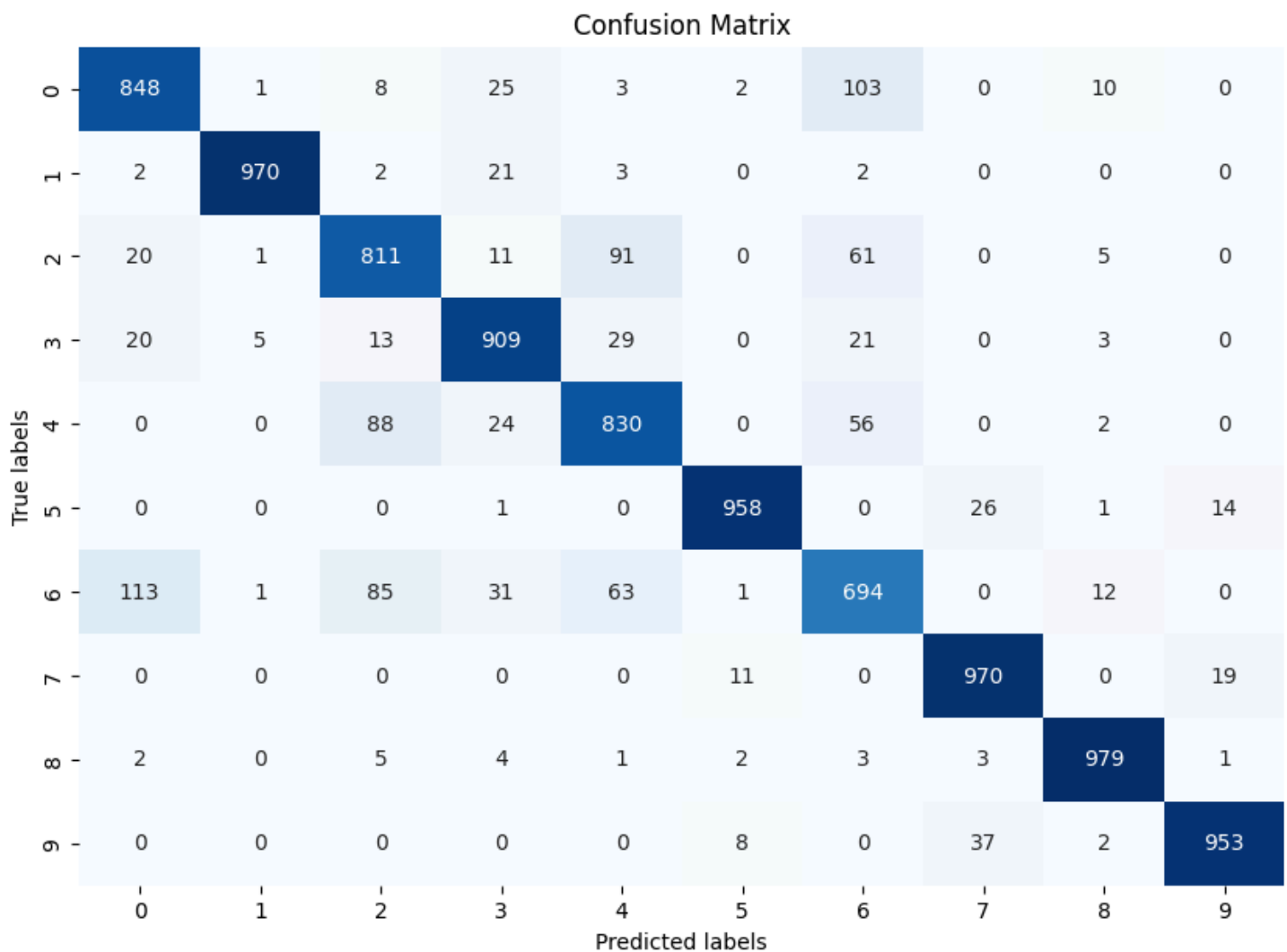


Figure 8: Confusion matrix for the best SVM model after dimensionality reduction

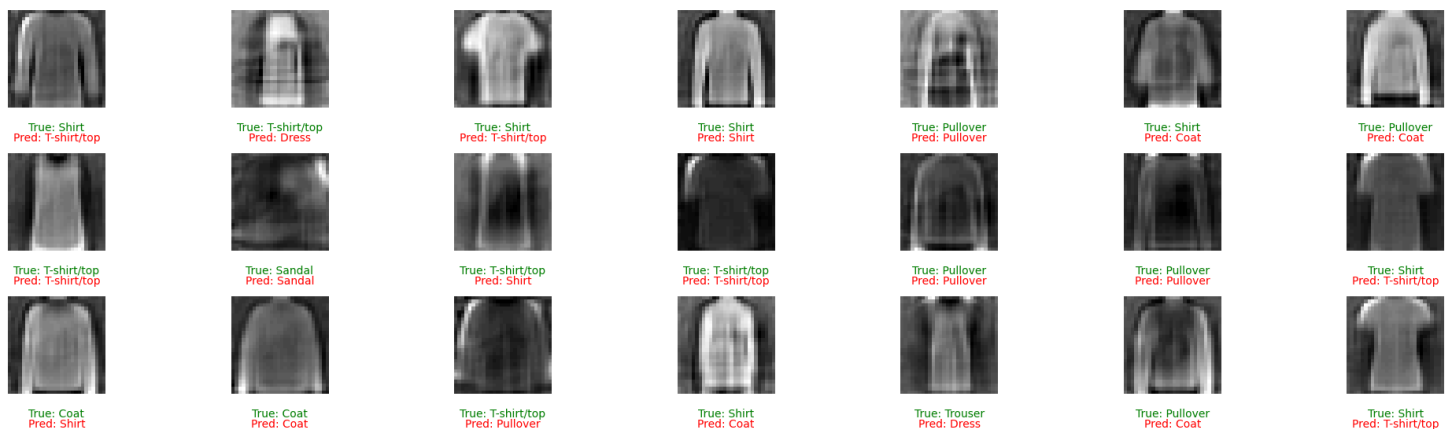


Figure 9: Misclassified images by the best SVM model after dimensionality reduction

## Conclusion

- **Performance Difference:** The SVM model without PCA has slightly better performance across all metrics. This suggests that while PCA effectively reduces dimensionality, the full feature set

provides additional details, leading to marginally higher classification accuracy.

- **Efficiency vs. Accuracy:** Using PCA reduces computational load, which is advantageous for large datasets. However, there’s a slight trade-off in accuracy. For tasks requiring high precision, the non-PCA model is preferable, despite the higher computational demand.
- **Class-Specific Misclassifications:** Both models show similar patterns in misclassifying visually alike classes, indicating that additional feature engineering or a more complex model may help improve separation between these classes.

In summary, the PCA-reduced SVM model offers computational efficiency with strong accuracy, while the non-PCA model provides superior accuracy and might be better suited for applications where computational resources aren’t a constraint.

## Model Performance In Full vs. Reduced Data

### Full-Dimensional Data:

Model	Training time	Evaluation time	Accuracy	Precision	Recall	F1 Score
Logistic Regression	53.04 s	0.052 s	0.8467	0.8455	0.8467	0.8458
SVM (Linear Kernel)	6.439 s	0.0417 s	0.8405	0.8381	0.8405	0.8382
SVM (RBF Kernel)	24.44 s	8.4204 s	0.8995	0.8993	0.8995	0.8993

### Observations for Full-Dimensional Data:

- **Logistic Regression:** Logistic regression on full data offers good accuracy and balanced precision-recall, indicating reliable performance. However, it has the highest training time (53.04 seconds), making it less efficient than SVM models on full data.
- **Linear Kernel SVM:** The linear SVM model is faster than logistic regression, with a much lower training time (6.439 seconds) and a similar evaluation time (0.0417 seconds). However, its accuracy and F1 score (0.8405 and 0.8382, respectively) are slightly lower than logistic regression’s, making it a slightly less effective model for full data.
- **RBF Kernel SVM:** The RBF SVM outperforms both logistic regression and linear SVM in accuracy and F1 score on full data (0.8995 and 0.8993, respectively), demonstrating its

effectiveness for complex, non-linear data relationships. However, this comes at the cost of significantly longer evaluation time (8.42 seconds), though its training time (24.44 seconds) is still shorter than logistic regression's.

### Summary

For full data, **RBF SVM** is the best performer in terms of accuracy and overall metrics, but it has the highest evaluation time. **Linear SVM** offers a good balance between computational efficiency and performance, while **Logistic Regression** performs well in terms of accuracy but has the longest training time.

### Reduced-Dimensional Data:

Model	Training Time	Evaluation Time	Accuracy	Precision	Recall	F1 Score
Logistic Regression	18.81 s	0.109 s	0.8365	0.8352	0.8365	0.8356
SVM (RBF Kernel)	21.78 s	0.351 s	0.8922	0.8918	0.8922	0.8918

### Observations for Reduced-Dimensional Data:

- **Logistic Regression:** With reduced data, logistic regression becomes significantly faster, with a training time of 18.81 seconds, and maintains a high level of accuracy (0.8365). The metrics are slightly lower than those on full data but indicate that logistic regression handles dimensionality reduction well, with minimal loss in performance.
- **RBF Kernel SVM:** For SVM with RBF kernel, reduced-dimensional data improves accuracy significantly to 0.8922, with a training time of 21.78 seconds. Although the evaluation time (0.351 seconds) is longer than logistic regression, the improved accuracy and balanced precision-recall metrics make RBF SVM highly effective when working with reduced-dimensional data.

### Summary

With reduced data, **RBF SVM** achieves the best accuracy and performance metrics, demonstrating that it benefits substantially from dimensionality reduction. **Logistic Regression** also performs well but remains slightly less effective than **RBF SVM** in reduced dimensions.



# Overall Summary

- **Logistic Regression:** Suitable for simpler, low-dimensional data, providing consistent accuracy when dimensions are reduced. However, it has a longer training time on full data.
- **Linear Kernel SVM:** Offers a good balance of speed and accuracy, especially with reduced-dimensional data, where it achieves a notable accuracy. The evaluation time is longer than that of Logistic Regression, but it is effective for datasets with many features that can be reduced without losing significant information.
- **RBF Kernel SVM:** Achieves the highest accuracy and F1 score on full data and continues to perform well with reduced-dimensional data. It is ideal for complex datasets where capturing non-linear relationships is essential. However, it is computationally intensive, particularly during evaluation.

# Conclusion

- Use **Logistic Regression** for simple, interpretable tasks where time is not a critical factor.
- Choose **Linear Kernel SVM** when working with high-dimensional data that requires good accuracy without excessive computation time.
- Opt for **RBF Kernel SVM** for the best accuracy on complex, non-linear data, prioritizing model performance over computational time, especially after dimensionality reduction.

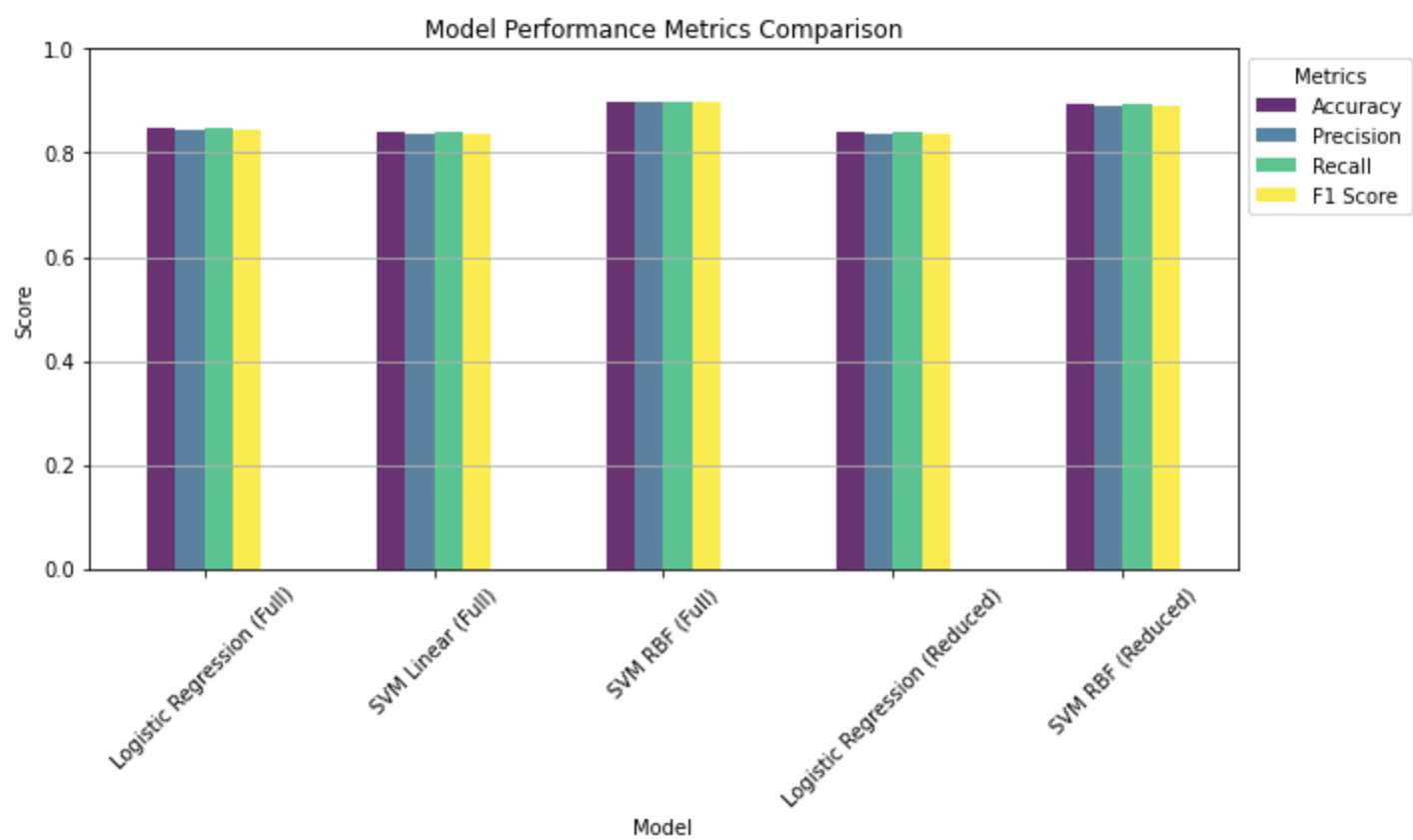


Figure 10: Comparison of model performance on full and reduced data

# Extending the analysis

We can observe that in terms of both Training time and Evaluation time across all three algorithms, **RBF SVC** takes the most time, followed by **Logistic Regression**, while **LinearSVC** leads in speed.

So, what is the reason behind this? Let's delve deeper into how these algorithms operate.

## RBF Kernel SVC

- **Mathematical Formula of the RBF Kernel:**

$$K(x, x') = \exp(-\gamma ||x - x'||^2)$$

Here,  $\gamma$  is a parameter that adjusts the influence of a data point. If  $x$  and  $x'$  are close,  $K(x, x')$  will be near 1; conversely, if the distance is large, the kernel value will be close to 0.

- **Reasons for High Computation Time:**
  - **Space Transformation:** The RBF Kernel transforms data from the input space to a multi-dimensional space where data points become more linearly separable. This increases complexity, as each data point must compute relationships with all other points using the above formula.
  - **Computational Complexity:** With the RBF kernel, the algorithm calculates every data point with the kernel formula, so the number of calculations increases quadratically with the number of points ( $O(n^2)$ ), significantly increasing training time.
  - **Iteration Count for Optimization:** Since the model aims to optimize the boundary for accurate classification in a high-dimensional space, the iteration count for convergence also increases, extending training time.

## Linear Kernel SVC

- **Mathematical Formula of the Linear Kernel:**

$$K(x, x') = x \cdot x'$$

Here,  $x \cdot x'$  is the dot product between two data points. For the linear kernel, no space transformation is needed, allowing the model to train directly on the original data.

- **Reasons for Fast Computation Time:**
  - **Linear Space:** Linear SVC performs well on linearly separable datasets and does not require space transformation, which significantly reduces computation time.

- **Simpler Optimization:** Without the complex calculations of the RBF kernel, Linear SVC only needs to optimize a linear separation boundary, reducing the number of calculations substantially to  $O(n)$ , where  $n$  is the number of data points.
- **Direct Classification:** During evaluation, classification only requires calculating the dot product of the new point with the learned weight vector, making evaluation time extremely fast.

## Logistic Regression

- **Mathematical Formula of Logistic Regression:**

$$P(y = 1|x) = \frac{1}{1 + e^{-(w \cdot x + b)}}$$

Here,  $w$  is the weight vector and  $b$  is the intercept. Logistic Regression is a linear probability model that seeks to optimize weights to determine the probability of a point belonging to a specific class.

- **Reasons for Fast Training and Evaluation Time:**
  - **Linear Space:** Like Linear SVC, Logistic Regression operates in linear space and does not require space transformation, reducing the number of computations.
  - **Loss Function:** Logistic Regression uses a logarithmic loss function, which is easily optimized using gradient descent or other optimization methods, with a complexity of  $O(n)$ , where  $n$  is the number of data points.
  - **Fast Evaluation:** During evaluation, Logistic Regression only requires calculating the dot product of the weights and new data point, followed by applying the sigmoid function. This is simpler and faster compared to the RBF kernel.

## Summary:

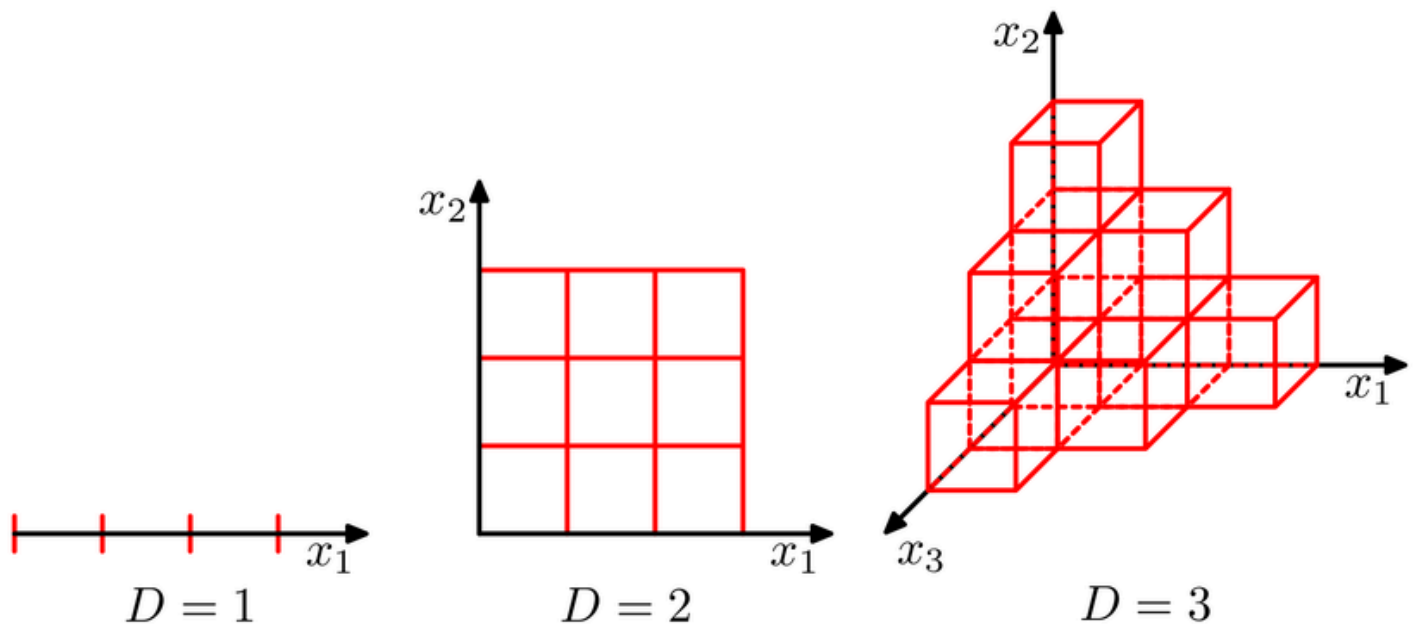
- **RBF Kernel SVC** is complex due to space transformation and requires calculations for each data point in a non-linear space, leading to significantly higher training and evaluation times.
- **Linear Kernel SVC** and **Logistic Regression** are simpler and faster as they operate in a linear space, with calculations mainly involving dot products and linear optimization, making them suitable for tasks requiring fast computation times.

In general, if the data can be linearly separated, Linear SVC or Logistic Regression would be optimal choices for speed. On the other hand, if the data is highly non-linear, RBF Kernel SVC may provide better accuracy at the cost of higher computation time.

# 3. The Curse of Multidimensionality and Its Implications for Model Performance

## Introduction to the Curse of Dimensionality

The Curse of Dimensionality refers to problems that arise when working with data in high-dimensional spaces. *Figure 11* illustrates this problem in case of the incremental number of dimension:



*Figure 11: Illustration of the curse of dimensionality, showing how the number of regions of a regular grid grows exponentially with the dimensionality  $D$  of the space. For clarity, only a subset of the cubical regions are shown for  $D = 3$*

Particularly, as the number of dimensions  $D$  increases, the space becomes exponentially larger. Imagine dividing the space into small cells, where each cell could contain some data points. In a one-dimensional space ( $D = 1$ ), dividing it into cells is straightforward, like cutting a line into segments. In two dimensions ( $D = 2$ ), it's like dividing a plane into squares. But as the number of dimensions increases, the number of cells increases exponentially. For example, if you have 10 divisions per dimension, a 1D space has 10 cells, a 2D space has 100 cells, and a 3D space has 1000 cells.

Now, if the data is sparse (spread out) in these cells, many cells might not contain any data at all, or contain only a few data points. This creates a major problem for machine learning models because they rely on having enough data in each region of the space to make accurate predictions. When the data is spread too thin across many dimensions, the models struggle to generalize well because they

don't have enough data in each part of the space to learn meaningful patterns. For example, in a 10-dimensional space, data points are scattered throughout a vast volume, making it unlikely that any two points are close to each other. This sparsity weakens the reliability of nearest-neighbor approaches and increases the risk of overfitting as models struggle to generalize from few, isolated samples.

This sparsity is the core of the Curse of Dimensionality, making it harder for machine learning models to work effectively as dimensions increase.

We will practice the following example with noting that the volume of a sphere of radius  $r$  in  $D$  dimensions must scale as  $r^D$ , and so we write  $V_D(r) = K_D r^D$  where  $K_D$  is a constant that depends on the dimension  $D$ . Thus the fraction of the volume of the sphere that lies between radius  $r = 1 - \epsilon$  and  $r = 1$  based on relative error as follows:

$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D$$

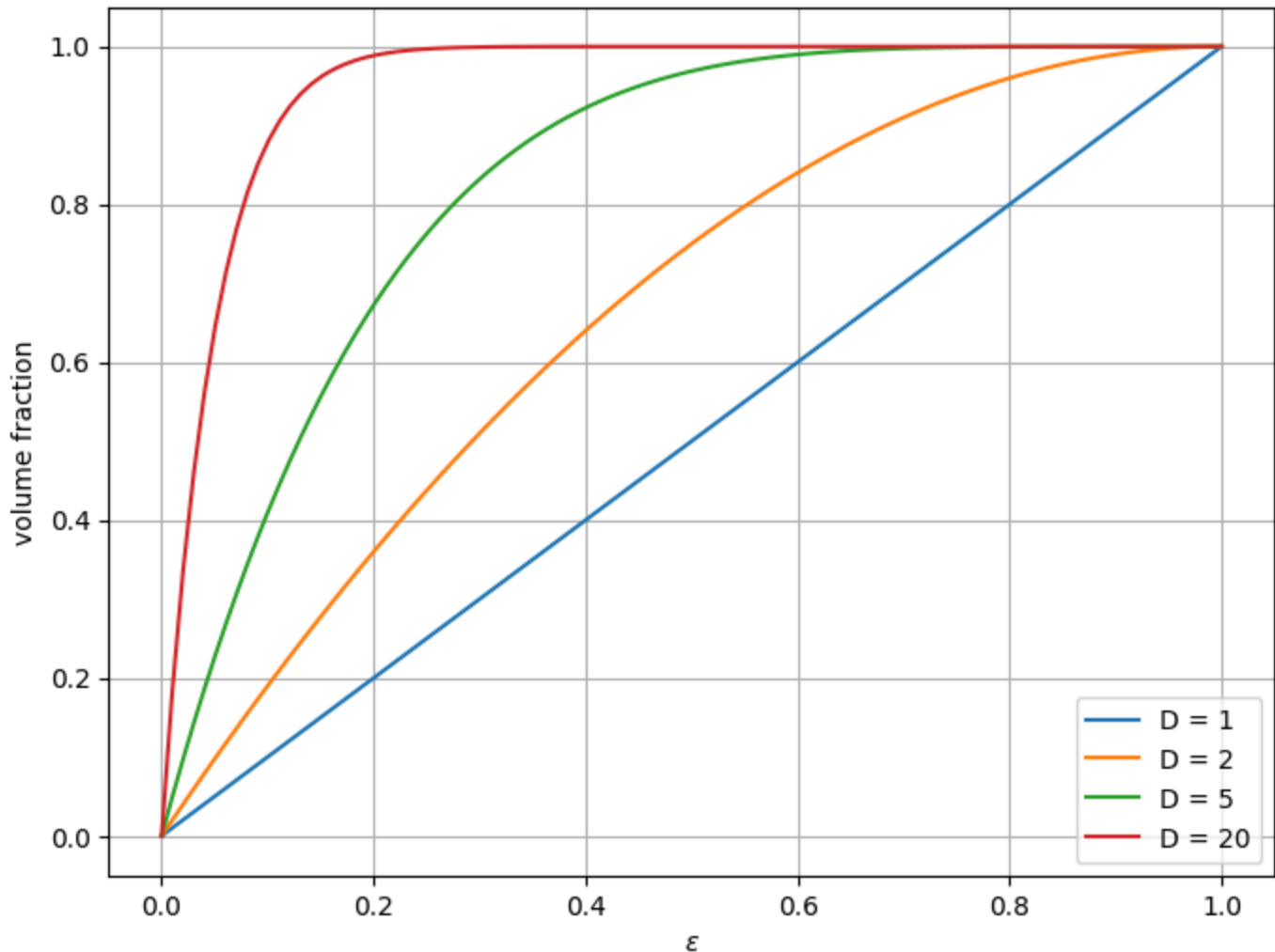


Figure 12: Visualization of the volume of a sphere in high-dimensional space, showing the concentration of volume near the surface as dimensionality increases

We can see that, in *Figure 12*, the dimension  $D$  increases, most of the volume of the sphere is located near the surface. This suggests that the data in high-dimensional space is not uniformly distributed, but rather concentrated at the edges of the space, which reduces the significance of distance measures such as Euclidean distance. As data points become sparse in large spaces, it becomes more difficult to identify patterns or groups in the data. Many machine learning algorithms rely on distance metrics (like Euclidean distance) that lose effectiveness in high-dimensional spaces. In low-dimensional spaces, distance measures are meaningful, but in high dimensions, all points can become almost equidistant from each other, reducing the effectiveness of algorithms that depend on nearest-neighbor calculations.

With each additional dimension, the complexity of fitting models increases. Polynomial models, for example, require an increasing number of coefficients, and capturing interactions between features demands even more data, making models computationally intensive and harder to train effectively.

## **Impact on Machine Learning Models**

### **Reduced Model Performance and Generalization**

In a classification task using the k-Nearest Neighbors (k-NN) algorithm, the Euclidean distance metric becomes less meaningful in high-dimensional spaces, as most points are equidistant from each other. This undermines the effectiveness of k-NN, resulting in incorrect or inconsistent predictions.

### **Difficulty in Finding Meaningful Patterns**

Regression algorithms, like Linear or Polynomial Regression, require substantial data to capture relationships among input features in high-dimensional spaces. Without enough data, models are prone to noise, learning incorrect relationships that don't represent the true underlying patterns.

### **Increased Computational Complexity**

The number of parameters a model needs to learn grows rapidly as dimensionality increases, leading to higher computational complexity, memory requirements, and slower training and prediction times. For instance, a model operating in a 1000-dimensional space requires considerably more computational power and time than one in a 10-dimensional space.

### **Increased Risk of Overfitting**

The curse of dimensionality increases the likelihood of overfitting because, in high-dimensional spaces, models can easily find boundaries or relationships to classify training data with high accuracy.

However, these boundaries may not reflect the true data structure, leading to errors on new data. For instance, in an image classification task with thousands of dimensions (pixels), a neural network might learn features specific to each training image rather than general features of the image class.

## Challenges with Distance Metrics

Many machine learning algorithms, such as k-NN, SVM, and clustering methods, rely on distance metrics. In high-dimensional spaces, however, distances between points become almost uniform, rendering distance metrics ineffective. For clustering tasks, the curse of dimensionality makes identifying clusters more challenging. Algorithms like k-means may lose effectiveness as data points become equidistant from each other, making cluster centroids less representative of any particular group of points.

## Reflect on the trade-offs between dimensionality reduction and model performance

Dimensionality reduction techniques, like PCA, are often employed to simplify data, especially in high-dimensional spaces. However, reducing dimensions can lead to both advantages and trade-offs that affect model performance. Here's a reflection on these trade-offs as seen with logistic regression and SVM.

## Model Efficiency and Computational Cost

- **With Dimensionality Reduction (PCA):** By reducing the number of features, models require less computational power and memory, leading to faster training and inference times. This is especially beneficial for computationally intensive algorithms like SVM, where the time complexity increases with the feature count.
- **Without Dimensionality Reduction:** Although more computational resources are required, models can fully utilize all available information, which can lead to more nuanced and accurate representations of the data. For logistic regression, the impact may be less significant since it's inherently faster, but for SVM, avoiding dimensionality reduction can noticeably increase processing time.

## Impact on Accuracy and Predictive Performance

- **With Dimensionality Reduction:** PCA helps retain the most informative features, allowing models to perform well with less data "noise" and fewer irrelevant features. However, it may lose finer distinctions within the data, particularly if complex relationships between features are not fully

- preserved in lower-dimensional space. This can result in a slight drop in accuracy, as seen with both logistic regression and SVM in the PCA-reduced versions.
- **Without Dimensionality Reduction:** Models generally achieve better accuracy and precision without PCA since they can leverage all the data's intricacies. In both logistic regression and SVM, using the full feature set led to marginally higher accuracy, indicating that some critical information was lost during PCA. This may be particularly important for complex classes with subtle differences, where maintaining high dimensionality helps the model differentiate effectively.

## Summary of Trade-offs

Aspect	With Dimensionality Reduction (PCA)	Without Dimensionality Reduction
Efficiency	Faster training and inference, especially for SVM	Higher computational cost, slower with SVM
Accuracy	Slightly reduced accuracy due to loss of information	Generally higher accuracy, retains all feature nuances

## Conclusion

For both logistic regression and SVM, dimensionality reduction with PCA offers computational efficiency, making it valuable for high-dimensional data. However, there is a trade-off in accuracy, as some critical information may be lost, impacting the model's ability to capture subtle class distinctions. Without dimensionality reduction, both models perform better in terms of accuracy but at the cost of increased computation and potential overfitting.

In practice, the choice depends on the specific requirements: if computational resources and time are limited, or if the data is highly dimensional with potential redundancy, PCA can offer a balanced approach. However, when high accuracy and interpretability are prioritized, and resources permit, avoiding dimensionality reduction may be beneficial.

# 4. Option for choosing the best model

## Logistic Regression

Advantages:



- **Simplicity and Interpretability:** Logistic regression is straightforward to implement and interpret, making it suitable for simpler tasks. It provides clear insights into the influence of each feature on the prediction.
- **Consistent Performance with Reduced Data:** Logistic regression retains relatively stable accuracy even with dimensionality reduction, indicating robustness against overfitting in high-dimensional spaces.
- **Efficiency:** It is computationally efficient for training and evaluation in low to moderate-dimensional datasets.

Disadvantages:

- **Long Training Time on Full Data:** In the case of full-dimensional data, logistic regression has a significantly longer training time (53.04 seconds), making it less suitable for large datasets where computational resources or time are constrained.
- **Limited Performance on Complex Data:** Logistic regression may struggle to capture complex, non-linear relationships in data, resulting in lower accuracy and F1 scores compared to SVM models.

## Linear SVM (Support Vector Machine)

Advantages:

- **Efficiency with High-Dimensional Data:** Linear SVM is effective for high-dimensional data, especially when the data is linearly separable, maintaining reasonable performance without excessive computational burden.
- **Good Balance of Speed and Accuracy:** The linear SVM offers a favorable trade-off, achieving competitive accuracy with significantly shorter training time (6.439 seconds) compared to logistic regression on full-dimensional data.

Disadvantages:

- **Lower Performance on Complex Relationships:** While the linear SVM performs well on linearly separable data, it may not adequately capture more complex, non-linear relationships, leading to lower accuracy and F1 scores compared to the RBF kernel SVM.
- **Interpretability:** The coefficients can be less interpretable than those in logistic regression, making it harder to derive insights about feature importance.

# RBF SVM (Radial Basis Function Kernel)

Advantages:

- **Best Accuracy and Performance Metrics:** RBF SVM consistently achieves the highest accuracy and F1 score across both full and reduced-dimensional datasets, demonstrating its effectiveness in handling complex, non-linear relationships in the data.
- **Adaptability to Reduced Dimensions:** The model performs exceptionally well with reduced-dimensional data, benefiting from the elimination of irrelevant features and noise, leading to improved accuracy (0.8922) and precision-recall metrics.

Disadvantages:

- **Computationally Intensive:** The training and evaluation times for RBF SVM are higher (24.44 seconds training, 8.4204 seconds evaluation), making it less efficient for very large datasets or real-time applications where speed is critical.
- **Sensitivity to Hyperparameters:** RBF kernel SVMs often require careful tuning of hyperparameters (e.g., kernel width and regularization parameters) to achieve optimal performance, adding complexity to the model training process.

## Summary

Overall, the choice of model should be driven by the specific requirements of the task, including data complexity, computational resources, and the need for interpretability:

- **Logistic Regression** is best for simple tasks and scenarios with low-dimensional data where interpretability and computational efficiency are prioritized.
- **Linear SVM** provides a balanced option for high-dimensional data that can be effectively modeled with linear separation, offering good accuracy with reasonable computational demands.
- **RBF SVM** is the preferred choice for tasks requiring high accuracy on complex, non-linear data, though it comes with increased computational costs and potential complexity in tuning hyperparameters.

In conclusion, the choice of model ultimately depends on the specific problem at hand, the nature of the data, and the importance of accuracy versus computation time.