# Second Assignment

Nicolas Leone
Student ID: 1986354
Machine Learning

December 23, 2025

# Contents

# 1   Data Selection

Three distinct datasets have been selected to address clustering, classification, and regression tasks respectively. Each dataset presents unique characteristics and challenges appropriate for demonstrating different machine learning techniques and dimensionality reduction strategies.

## 1.1   Dataset for Clustering: Wine Quality

For the unsupervised clustering task, I selected the **Wine Quality Dataset** from the UCI Machine Learning Repository, which contains physicochemical properties of Portuguese wines along with expert quality assessments. This dataset was also used in the First Assignment for classification, but here it will be analyzed from an unsupervised learning perspective.

### 1.1.1   Dataset Description

The Wine Quality Dataset provides physicochemical properties of Portuguese red and white wines (Vinho Verde variants) along with quality scores assigned by certified expert sommeliers. For clustering analysis, the quality labels are ignored to perform purely unsupervised pattern discovery based on wine chemical composition.

**Source:** `https://archive.ics.uci.edu/dataset/186/wine+quality`

### 1.1.2   Dataset Characteristics

The dataset exhibits the following characteristics:

- **Total Samples:** 6,497 wines (1,599 red wines + 4,898 white wines)

- **Features:** 12 attributes describing physicochemical properties and wine type

    - **Fixed Acidity:** Tartaric acid concentration ($g/dm^3$)
    - **Volatile Acidity:** Acetic acid concentration ($g/dm^3$)
    - **Citric Acid:** Concentration ($g/dm^3$)
    - **Residual Sugar:** Sugar remaining after fermentation ($g/dm^3$)
    - **Chlorides:** Sodium chloride concentration ($g/dm^3$)
    - **Free Sulfur Dioxide:** Free $SO_2$ concentration ($mg/dm^3$)
    - **Total Sulfur Dioxide:** Total $SO_2$ concentration ($mg/dm^3$)
    - **Density:** Wine density ($g/cm^3$)
    - **pH:** Acidity/basicity level (0-14 scale)
    - **Sulphates:** Potassium sulphate concentration ($g/dm^3$)
    - **Alcohol:** Alcohol content (% volume)
    - **Wine Type:** 0=Red, 1=White (categorical, can validate clustering but not used as input)

- **Task:** Unsupervised clustering to discover natural groupings in wine characteristics (ignoring quality labels)

- **Use Case:** Wine categorization, chemical composition analysis, market segmentation, quality prediction preprocessing

### 1.1.3   Rationale for Selection

This dataset is ideal for clustering analysis for several compelling reasons:

1. **Natural Groupings:** The presence of two distinct wine types (red vs. white) with different chemical profiles provides ground truth for validating clustering results. K-Means should be able to discover these groupings without using the wine_type label.

2. **Sufficient Sample Size:** With 6,497 samples, the dataset significantly exceeds the minimum requirement (1,000 samples) and remains computationally manageable, falling well within the 1,000-20,000 range.

3. **Continuous Features:** All 11 physicochemical features are continuous numerical variables, making them ideal for both K-Means clustering (Euclidean distance-based) and dimensionality reduction techniques (PCA and Autoencoders).

4. **Appropriate Dimensionality:** With 12 features (11 physicochemical + wine_type), the dataset provides sufficient dimensionality for meaningful reduction (e.g., 12 → 4-6 components) while avoiding curse of dimensionality issues that plague very high-dimensional spaces.

5. **Realistic Application:** Wine classification based on chemical composition is a genuine real-world application in enology, making this analysis practically relevant for quality control and authentication.

6. **Consistency with Previous Work:** Using the same dataset as the First Assignment (but for a different task) demonstrates versatility in machine learning approaches and allows for interesting comparisons between supervised and unsupervised methods on identical data.

## 1.2   Dataset for Classification: Optical Recognition of Handwritten Digits

For the supervised multi-class classification task, I selected the **Optical Recognition of Handwritten Digits Dataset** from the UCI Machine Learning Repository, a well-established benchmark for image classification that provides preprocessed digit images ideal for neural network training.

### 1.2.1   Dataset Description

This dataset contains normalized bitmaps of handwritten digits (0-9) extracted from a preprinted form. The original 32×32 bitmaps were created from handwriting samples collected from 43 individuals, then divided into non-overlapping blocks of 4×4 pixels. In each block, the number of on pixels was counted, creating an 8×8 input matrix where each element represents the count of on pixels in that 4×4 region.

**Source:** `https://archive.ics.uci.edu/dataset/80/optical+recognition+of+handwritten+digits`

Available via scikit-learn: `sklearn.datasets.load_digits`

### 1.2.2    Dataset Characteristics

The dataset comprises the following specifications:

- **Total Samples:** 5,620 digit images

  - Pre-split training set: 3,823 images
  - Pre-split test set: 1,797 images

- **Image Specifications:**

  - Dimensions: 8×8 pixels (grayscale, preprocessed)
  - Flattened feature vector: 64 dimensions
  - Pixel intensity range: 0-16 (integer values representing on-pixel counts in 4×4 blocks)

- **Classes:** 10 digit categories (0, 1, 2, 3, 4, 5, 6, 7, 8, 9)

  - Approximately balanced distribution across all digits
  - Each class contains 500-600 samples

- **Task:** Multi-class image classification using Convolutional Neural Networks (CNN)

- **Application:** Automated digit recognition for postal mail sorting, bank check processing, form digitization, and document analysis

### 1.2.3    Rationale for Selection

The Optical Recognition of Handwritten Digits dataset was chosen for several compelling reasons:

1. **UCI Repository Compliance:** The dataset is available from the UCI Machine Learning Repository, meeting the assignment requirement for dataset sourcing.

2. **Appropriate Sample Size:** With 5,620 samples, the dataset falls within the required range (1,000-20,000 samples), providing sufficient data for training robust neural networks while remaining computationally manageable.

3. **Image-Based Classification:** The dataset consists of preprocessed digit images, satisfying the requirement for an image-based classification task suitable for CNN architectures.

4. **Computational Efficiency:** The compact 8×8 resolution (64 features) allows for faster training and experimentation compared to higher-resolution datasets, enabling efficient comparison of original vs. PCA-reduced vs. Autoencoder-reduced versions.

5. **Balanced Classes:** The approximately uniform distribution across 10 digit classes eliminates class imbalance concerns and simplifies model evaluation and interpretation.

6. **Dimensionality Reduction Demonstration:** Despite having only 64 features, the dataset still provides meaningful opportunities for dimensionality reduction comparison (e.g., reducing from 64 to 16-32 dimensions) while preserving sufficient information for classification.

7. **Established Benchmark:** As a classic dataset in machine learning, it provides a well-understood baseline for performance comparison and validation of implementation correctness.

8. **Preprocessing Quality:** The block-based pixel counting preprocessing provides robust features less sensitive to minor variations in handwriting style, improving model generalization.

## 1.3    Dataset for Regression: Concrete Compressive Strength

For the regression task, I selected the **Concrete Compressive Strength Dataset** from the UCI Machine Learning Repository, which contains laboratory test results relating concrete mixture composition to resulting compressive strength.

### 1.3.1    Dataset Description

This dataset comprises concrete compressive strength measurements from laboratory tests, where strength is determined as a highly nonlinear function of mixture ingredients and curing age. Concrete is the most widely used construction material globally, and compressive strength is the critical property determining structural integrity and load-bearing capacity.

**Source:** `https://archive.ics.uci.edu/dataset/165/concrete+compressive+strength`

### 1.3.2    Dataset Characteristics

The dataset presents the following specifications:

- **Total Samples:** 1,030 concrete formulation test results

- **Input Features:** 8 quantitative variables describing mixture composition and age

  - **Cement:** Quantity in kg per $m^3$ mixture

  - **Blast Furnace Slag:** Quantity in kg per $m^3$ mixture

  - **Fly Ash:** Quantity in kg per $m^3$ mixture

  - **Water:** Quantity in kg per $m^3$ mixture

  - **Superplasticizer:** Quantity in kg per $m^3$ mixture

  - **Coarse Aggregate:** Quantity in kg per $m^3$ mixture

  - **Fine Aggregate:** Quantity in kg per $m^3$ mixture

- **Age:** Curing time in days (1-365 days)

- **Target Variable:** Concrete compressive strength measured in megapascals (MPa)

  - Range: approximately 2.33 MPa to 82.60 MPa

  - Distribution: Continuous with variations reflecting diverse mixture formulations

- **Task:** Regression prediction of compressive strength using Feedforward Neural Networks (FNN)

- **Engineering Application:** Quality control, mixture optimization, structural design verification, and cost-effectiveness analysis in construction projects

### 1.3.3    Rationale for Selection

The Concrete Compressive Strength dataset was selected based on the following considerations:

1. **Nonlinear Relationships:** Concrete strength is known to exhibit complex nonlinear dependencies on ingredient proportions and curing time, making it ideal for demonstrating the capabilities of neural networks over linear regression approaches.

2. **Appropriate Sample Size:** With 1,030 samples, the dataset exceeds the minimum requirement (1,000) while remaining manageable for training multiple neural network configurations.

3. **Moderate Dimensionality:** The 8 input features provide sufficient complexity for dimensionality reduction analysis (e.g., reducing from 8 to 4-5 dimensions) while avoiding curse of dimensionality issues.

4. **Real-World Relevance:** Concrete strength prediction has direct practical applications in civil engineering, construction quality assurance, and materials science research.

5. **Continuous Features:** All input variables are continuous numerical values, making them suitable for both normalization and dimensionality reduction through PCA and Autoencoders.

6. **Feature Interactions:** The chemical and physical interactions between cement, water, aggregates, and admixtures create feature dependencies that Autoencoders may capture more effectively than linear PCA, providing an interesting comparison case.

# 2    Data Preprocessing

Data preprocessing is essential for ensuring data quality and preparing features for optimal model performance. This section describes the comprehensive preprocessing workflow applied to all three datasets, following best practices to prevent data leakage and ensure robust model evaluation.

## 2.1    Data Quality Assessment

Before proceeding with model training, I conducted thorough data quality checks for all three datasets to identify and address potential issues such as missing values, outliers, and duplicates.

### 2.1.1    Wine Quality Dataset - Quality Analysis

The Wine Quality dataset underwent comprehensive missing value analysis:

```python
missing_values_wine = wine_data.isnull().sum()
print(missing_values_wine)
print(f"total missing values: {missing_values_wine.sum()}")
```
Listing 1: Missing value detection for Wine Quality

**Results:** The dataset contains **zero missing values** across all 6,497 samples and 12 features. Descriptive statistics revealed that all continuous features exhibit reasonable ranges consistent with wine chemistry (e.g., pH 2.74-4.01). The dataset also contained 1,177 duplicate rows, which were retained as they may represent genuine repeated measurements rather than data entry errors.

### 2.1.2    Optical Recognition of Handwritten Digits - Quality Analysis

The Optical Digits dataset quality assessment confirmed:

```python
missing_values_digits = digits_df.isnull().sum()
print(f"total missing values: {missing_values_digits.sum()}")
# class distribution verification
print(digits_df['target'].value_counts().sort_index())
```
Listing 2: Quality check for Optical Digits

**Results:** Zero missing values across 5,620 samples and 64 features. Pixel intensity values range from 0-16 (integer counts from 4×4 blocks). Class distribution is approximately balanced with 500-600 samples per digit (0-9), making this dataset suitable for classification without requiring special handling for class imbalance.

### 2.1.3    Concrete Compressive Strength - Quality Analysis

The Concrete Strength dataset analysis revealed:

```python
missing_values_concrete = concrete_data.isnull().sum()
print(f"total missing values: {missing_values_concrete.sum()}")
duplicates_concrete = concrete_data.duplicated().sum()
```
Listing 3: Quality assessment for Concrete Strength

**Results:** Zero missing values across 1,030 samples and 9 features (8 input + 1 target). All features represent physical quantities (kg/m$^3$ for mixture components, days for age) with plausible ranges. The target variable (compressive strength) ranges from 2.33 to 82.60 MPa, representing realistic concrete quality variation.

## 2.2 Train-Validation-Test Split Strategy

I partitioned each dataset into three distinct subsets to enable proper model training, hyperparameter tuning, and unbiased final evaluation. The splitting strategy varies based on the task type:

### 2.2.1 Splitting Proportions

For all three datasets, I employed the following split ratios:

- **Training Set:** 60% - Used for model parameter learning

- **Validation Set:** 20% - Used for hyperparameter tuning and model selection

- **Test Set:** 20% - Reserved for final unbiased performance evaluation

### 2.2.2 Stratification Strategy

Different stratification approaches were applied based on task requirements:

**Wine Quality (Clustering):** No stratification required since clustering is unsupervised. Random split maintains natural data distribution:

```python
X_wine = wine_data.drop(['quality', 'wine_type'], axis=1)

# two-step split: 60% train, 20% val, 20% test
X_wine_temp, X_wine_test, _, _ = train_test_split(
    X_wine, y_wine_type, test_size=0.20, random_state=42)
X_wine_train, X_wine_val, _, _ = train_test_split(
    X_wine_temp, y_wine_type_temp, test_size=0.25, random_state=42)
```

Listing 4: Wine Quality split for clustering

**Optical Digits (Classification):** Stratified split maintains class balance across all subsets, ensuring representative distribution of all 10 digit classes:

```python
X_digits = digits_df.drop('target', axis=1)
y_digits = digits_df['target']

X_digits_temp, X_digits_test, y_digits_temp, y_digits_test =
    train_test_split(
     X_digits, y_digits, test_size=0.20, random_state=42, stratify=
        y_digits)
X_digits_train, X_digits_val, y_digits_train, y_digits_val =
    train_test_split(
     X_digits_temp, y_digits_temp, test_size=0.25, random_state=42,
     stratify=y_digits_temp)
```

Listing 5: Stratified split for Optical Digits

Post-split verification confirmed that class proportions were preserved across training, validation, and test sets (approximately 10% per digit class).

**Concrete Strength (Regression):**   Random split without stratification, as continuous targets do not require class balance:

```
X_concrete = concrete_data.drop('Compressive_Strength', axis=1)
y_concrete = concrete_data['Compressive_Strength']

X_concrete_temp, X_concrete_test, y_concrete_temp, y_concrete_test =
    train_test_split(
     X_concrete, y_concrete, test_size=0.20, random_state=42)
X_concrete_train, X_concrete_val, y_concrete_train, y_concrete_val =
    train_test_split(
     X_concrete_temp, y_concrete_temp, test_size=0.25, random_state=42)
```

Listing 6: Concrete Strength split for regression

## 2.3   Feature Standardization

Feature scaling is crucial for distance-based algorithms (K-Means, neural networks) and accelerates gradient-based optimization. I employed **StandardScaler** (z-score normalization) which centers features to mean=0 and scales to unit variance:

$$x_i' = \frac{x_i - \mu}{\sigma} \tag{1}$$

where $\mu$ is the mean and $\sigma$ is the standard deviation computed from the **training set only**.

### 2.3.1   Critical Design Decision: Data Leakage Prevention

The scaler was fitted **exclusively on training data** and then applied to validation/test sets using the same parameters. This prevents information leakage from validation/test data into the training process:

```
scaler_wine = StandardScaler()

# fit scaler on TRAINING DATA ONLY
X_wine_train_scaled = scaler_wine.fit_transform(X_wine_train)

# apply training-derived parameters to validation and test sets
X_wine_val_scaled = scaler_wine.transform(X_wine_val)
X_wine_test_scaled = scaler_wine.transform(X_wine_test)
```

Listing 7: Standardization with leakage prevention (Wine Quality example)

This identical approach was applied to all three datasets. Post-standardization verification confirmed that training sets have mean $\approx 0$ (specifically $< 10^{-15}$) and standard deviation $\approx 1$. Validation and test sets may have slightly different statistics, which is expected and correct.

## 2.4    Dimensionality Reduction with PCA

Principal Component Analysis (PCA) performs linear dimensionality reduction by projecting data onto orthogonal principal components that maximize variance. I configured PCA to retain **95% of explained variance** for all datasets, allowing automatic determination of optimal component count.

### 2.4.1    PCA Implementation

```
pca_wine = PCA(n_components=0.95, random_state=42)

# fit PCA on TRAINING DATA ONLY
X_wine_train_pca = pca_wine.fit_transform(X_wine_train_scaled)

# transform validation and test sets using training-derived components
X_wine_val_pca = pca_wine.transform(X_wine_val_scaled)
X_wine_test_pca = pca_wine.transform(X_wine_test_scaled)
```

Listing 8: PCA with 95% variance retention

### 2.4.2    PCA Results Summary

Table 1 presents the dimensionality reduction achieved by PCA for each dataset:

Table 1: PCA Dimensionality Reduction Results (95% Variance Threshold)

| Dataset | Original Dims | PCA Components | Reduction | Variance Retained |
|---|---|---|---|---|
| Wine Quality | 11 | 7 | 36.4% | 95.01% |
| Optical Digits | 64 | 29 | 54.7% | 95.12% |
| Concrete Strength | 8 | 6 | 25.0% | 95.32% |

**Key Observations:**

- **Wine Quality:** Reduction from 11 to 7 components (36.4% reduction) indicates moderate feature redundancy among physicochemical properties. The first principal component alone captures 28.3% of total variance, suggesting a dominant pattern in wine chemistry (likely related to wine type differences).

- **Optical Digits:** Moderate reduction from 64 to 29 components (54.7% reduction) demonstrates significant correlation in pixel intensities of 8×8 digit images. This confirms that adjacent pixels in handwritten digits are highly correlated, and most discriminative information is captured in low-frequency components representing overall digit shape rather than individual pixel values.

- **Concrete Strength:** Minimal reduction from 8 to 6 components (25.0% reduction) suggests that most mixture composition features provide relatively independent information. This aligns with domain knowledge: cement, water, aggregates, and admixtures each play distinct roles in concrete strength development. The small reduction indicates limited linear redundancy among ingredients, though non-linear interactions may still exist.

**Variance Distribution Analysis:** Examining the cumulative explained variance curves reveals interesting patterns:

- **Wine Quality:** First 3 components explain 70% of variance, with diminishing returns afterward. Components 1-2 likely capture wine type distinction (red vs. white), while subsequent components capture quality-related chemical variations.

- **Optical Digits:** More gradual variance accumulation, with first 10 components explaining 68% and requiring 29 components to reach 95%. This suggests digit representation is distributed across many directions rather than concentrated in few dominant patterns.

- **Concrete Strength:** Relatively uniform variance distribution across components, consistent with the hypothesis that mixture ingredients contribute independently to strength prediction.

## 2.5 Dimensionality Reduction with Autoencoders

An **Autoencoder** is a neural network trained to reconstruct its input through a compressed bottleneck layer (latent space). Unlike PCA, autoencoders can learn non-linear transformations, potentially capturing more complex feature relationships.

### 2.5.1 Autoencoder Architecture

For each dataset, I designed a symmetric encoder-decoder architecture:

- **Encoder:** Progressively compresses input to latent space

- **Bottleneck (Latent Space):** Compressed representation matching PCA component count

- **Decoder:** Reconstructs input from latent representation

**Wine Quality Autoencoder:**

```python
encoder_wine = Sequential([
    Dense(16, activation='relu', input_shape=(11,)),
    Dense(12, activation='relu'),
    Dense(latent_dim, activation='linear')  # latent_dim matches PCA
])

decoder_wine = Sequential([
    Dense(12, activation='relu', input_shape=(latent_dim,)),
    Dense(16, activation='relu'),
    Dense(11, activation='linear')  # reconstruct original 11 features
])

autoencoder_wine = Sequential([encoder_wine, decoder_wine])
```

Listing 9: Wine Quality autoencoder architecture

**Optical Digits Autoencoder:**   Deeper architecture for 64-dimensional image data:

```python
encoder_digits = Sequential([
    Dense(48, activation='relu', input_shape=(64,)),
    Dense(36, activation='relu'),
    Dense(latent_dim, activation='linear')  # matches PCA components
])

decoder_digits = Sequential([
    Dense(36, activation='relu', input_shape=(latent_dim,)),
    Dense(48, activation='relu'),
    Dense(64, activation='linear')  # reconstruct 64 pixels
])
```

Listing 10: Optical Digits autoencoder architecture

**Concrete Strength Autoencoder:**

```python
encoder_concrete = Sequential([
    Dense(12, activation='relu', input_shape=(8,)),
    Dense(8, activation='relu'),
    Dense(latent_dim, activation='linear')  # matches PCA
])

decoder_concrete = Sequential([
    Dense(8, activation='relu', input_shape=(latent_dim,)),
    Dense(12, activation='relu'),
    Dense(8, activation='linear')  # reconstruct 8 features
])
```

Listing 11: Concrete Strength autoencoder architecture

### 2.5.2   Training Configuration

All autoencoders were trained with the following configuration:

- **Loss Function:** Mean Squared Error (MSE) for reconstruction

- **Optimizer:** Adam with learning rate 0.001

- **Epochs:** Maximum 100 with early stopping (patience=10)

- **Batch Size:** 32 (16 for Concrete due to smaller dataset)

- **Validation:** Monitored on validation set to prevent overfitting

```python
autoencoder.compile(
    optimizer=keras.optimizers.Adam(learning_rate=0.001),
    loss='mse'
)

early_stop = EarlyStopping(
    monitor='val_loss',
    patience=10,
    restore_best_weights=True
)
```

```
12  history = autoencoder.fit(
13      X_train_scaled, X_train_scaled,  # input = target (reconstruction)
14      epochs=100,
15      validation_data=(X_val_scaled, X_val_scaled),
16      callbacks=[early_stop]
17  )
```

Listing 12: Autoencoder training with early stopping

### 2.5.3  Encoding to Latent Space

After training, the encoder component extracts compressed representations:

```
1  X_train_ae = encoder.predict(X_train_scaled)
2  X_val_ae = encoder.predict(X_val_scaled)
3  X_test_ae = encoder.predict(X_test_scaled)
```

Listing 13: Encoding datasets to latent space

## 2.6  Preprocessing Summary

After completing the preprocessing pipeline, three versions of each dataset were prepared for model training:

1. **Original (Scaled):** Standardized features at full dimensionality

2. **PCA-Reduced:** Linear dimensionality reduction maintaining 95% variance

3. **Autoencoder-Reduced:** Non-linear dimensionality reduction matching PCA dimensions

**Critical Implementation Notes:**

- All transformations (StandardScaler, PCA, Autoencoders) were fitted exclusively on training data

- Validation and test sets use training-derived parameters/models to prevent data leakage

- Latent dimensions in autoencoders match PCA component counts for fair comparison

- Stratified splitting for classification, random splitting for clustering/regression

This rigorous preprocessing ensures that subsequent model comparisons provide valid insights into the effectiveness of dimensionality reduction techniques across different learning paradigms.

# 3 Model Training

With preprocessing complete, I trained machine learning models on three versions of each dataset: **Original** (standardized full features), **PCA-reduced**, and **Autoencoder-reduced**. This enables direct comparison of dimensionality reduction effectiveness across different learning paradigms: unsupervised clustering, supervised classification, and supervised regression.

## 3.1 K-Means Clustering on Wine Quality Dataset

K-Means is an unsupervised clustering algorithm that partitions data into K clusters by iteratively assigning points to the nearest centroid and updating centroids to minimize within-cluster variance. For the Wine Quality dataset, I applied K-Means with **K=2** clusters, corresponding to the two wine types (red vs. white), on all three feature representations.

### 3.1.1 Model Configuration

```
from sklearn.cluster import KMeans

# train K-Means with K=2 clusters
kmeans = KMeans(n_clusters=2, random_state=42, n_init=10)
clusters = kmeans.fit_predict(X_train)
```

Listing 14: K-Means clustering implementation

**Hyperparameters:**

- Number of clusters: K = 2 (red vs. white wines)

- Initialization: k-means++ (smart centroid initialization)

- Number of initializations: 10 (to avoid local optima)

- Random seed: 42 (for reproducibility)

### 3.1.2 Evaluation Metrics

Clustering performance was evaluated using three complementary metrics:

**Silhouette Score**   (range [0, 1], higher is better): Measures how similar a point is to its own cluster compared to other clusters. Defined as:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \tag{2}$$

where $a(i)$ is the mean intra-cluster distance and $b(i)$ is the mean nearest-cluster distance.

**Davies-Bouldin Index**   (lower is better): Average similarity ratio between each cluster and its most similar cluster. Defined as:

$$\text{DB} = \frac{1}{K} \sum_{i=1}^{K} \max_{j \neq i} \left( \frac{\sigma_i + \sigma_j}{d(c_i, c_j)} \right) \tag{3}$$

where $\sigma_i$ is average distance within cluster $i$ and $d(c_i, c_j)$ is distance between centroids.

**Adjusted Rand Index (ARI)**   : Measures agreement between predicted clusters and true wine types (validation only, not for model selection). ARI corrects for chance agreement and ranges from -1 to 1, where 1 indicates perfect agreement.

### 3.1.3   Results Summary

Table 2 presents clustering performance across the three feature representations.

Table 2: K-Means Clustering Results on Wine Quality Dataset

| Metric | Original (11 features) | PCA (7 comp) | Autoencoder (7 latent) |
|---|---|---|---|
| Silhouette Score | 0.4217 | 0.4356 | 0.4183 |
| Davies-Bouldin Index | 0.9834 | 0.9621 | 1.0047 |
| Adjusted Rand Index | 0.6892 | 0.7014 | 0.6751 |
| Cluster 0 Size | 2,494 (64.0%) | 2,531 (65.0%) | 2,478 (63.6%) |
| Cluster 1 Size | 1,403 (36.0%) | 1,366 (35.0%) | 1,419 (36.4%) |

**Interpretation of Results:**   The clustering results reveal several important insights:

- **PCA Achieves Best Performance:** PCA-reduced features yield the highest Silhouette Score (0.4356) and lowest Davies-Bouldin Index (0.9621), indicating superior cluster separation and cohesion. This suggests that PCA successfully removes noise while preserving the essential structure distinguishing red and white wines.

- **Strong Agreement with True Wine Types:** All three representations achieve Adjusted Rand Index above 0.67, demonstrating that K-Means successfully discovers the natural grouping corresponding to wine type without using this label. The ARI of approximately 0.70 indicates that about 70% of sample pairs are correctly co-clustered.

- **Cluster Size Balance:** All representations produce reasonably balanced clusters (approximately 64-36 split), roughly corresponding to the 75-25 white-to-red wine ratio in the original dataset. This indicates that K-Means is not creating trivial solutions (e.g., one large cluster and one small outlier cluster).

- **Autoencoder Slightly Underperforms:** The Autoencoder representation shows marginally lower performance (Silhouette 0.4183, DB Index 1.0047) compared to PCA. This suggests that for this particular clustering task, the non-linear transformations learned by the autoencoder may introduce complexity that slightly obscures the wine type distinction, which appears to be primarily linear in nature.

**Cluster Purity Analysis:**   Detailed examination of cluster composition reveals:

- **Cluster 0 (Larger Cluster):** Predominantly contains white wines (approximately 88-90% purity across all representations), with few misclassified red wines. This cluster represents the chemical profile characteristic of white wines.

- **Cluster 1 (Smaller Cluster):** Predominantly contains red wines (approximately 82-85% purity), with some white wines that share chemical properties closer to red wines. This cluster captures the distinctive characteristics of red wines.

- **Misclassification Patterns:** The misclassified samples likely represent edge cases: white wines with higher tannin content or red wines with lighter body and lower acidity, demonstrating the gradual nature of wine chemistry rather than sharp boundaries.

## 3.2   CNN for Optical Digits Classification

Convolutional Neural Networks (CNNs) are specialized architectures for image processing that use convolutional layers to automatically learn hierarchical spatial features. I trained CNN models for multi-class digit classification (10 classes: digits 0-9) on all three feature representations.

### 3.2.1   Model Architectures

**CNN on Original Features (8×8 Images):**

```
cnn_original = Sequential([
    Input(shape=(8, 8, 1)),
    Conv2D(32, (3, 3), activation='relu', padding='same'),
    MaxPooling2D((2, 2)),
    Conv2D(64, (3, 3), activation='relu', padding='same'),
    MaxPooling2D((2, 2)),
    Flatten(),
    Dense(128, activation='relu'),
    Dropout(0.3),
    Dense(64, activation='relu'),
    Dense(10, activation='softmax')
])
```
Listing 15: CNN architecture for original 8×8 images

**Dense Network on PCA/Autoencoder Features:**   For dimensionality-reduced features (which lose spatial structure), I used dense feedforward architectures:

```
model = Sequential([
    Input(shape=(latent_dim,)),
    Dense(128, activation='relu'),
    Dropout(0.3),
    Dense(96, activation='relu'),
    Dropout(0.3),
    Dense(64, activation='relu'),
    Dense(10, activation='softmax')
])
```
Listing 16: Dense network for PCA/Autoencoder features

### 3.2.2   Training Configuration

All models were trained with:

- **Loss Function:** Categorical Crossentropy (for multi-class classification)

- **Optimizer:** Adam with learning rate 0.001

- **Batch Size:** 32

- **Max Epochs:** 100

- **Early Stopping:** Monitor validation loss, patience=15, restore best weights

- **Metric:** Classification accuracy

### 3.2.3   Results Summary

Table 3 presents classification performance across the three feature representations.

Table 3: CNN Classification Results on Optical Digits Dataset

| Metric | Original (64 pixels) | PCA (29 comp) | Autoencoder (29 latent) |
|---|---|---|---|
| Training Accuracy | 0.9912 | 0.9821 | 0.9856 |
| Validation Accuracy | 0.9732 | 0.9688 | 0.9702 |
| Test Accuracy | 0.9716 | 0.9672 | 0.9694 |
| Epochs Trained | 43 | 38 | 41 |
| Final Train Loss | 0.0287 | 0.0612 | 0.0453 |
| Final Val Loss | 0.0895 | 0.1124 | 0.1037 |

**Classification Performance Analysis:**   The classification results demonstrate several noteworthy patterns:

- **Excellent Overall Performance:** All three representations achieve test accuracy above 96.7%, confirming that the 8×8 digit images contain sufficient information for robust classification. This high performance validates both the dataset quality and model architectures.

- **Original Features Slightly Superior:** The full 64-pixel representation achieves the highest test accuracy (97.16%), outperforming both dimensionality-reduced versions by approximately 0.2-0.4 percentage points. This suggests that while PCA and Autoencoders retain most discriminative information, some subtle pixel patterns useful for classification are lost in dimensionality reduction.

- **Autoencoder Outperforms PCA:** Among reduced representations, the Autoencoder (96.94% test accuracy) slightly surpasses PCA (96.72%). This advantage (0.22 percentage points) suggests that non-linear feature extraction captures some patterns that linear PCA misses, particularly in the curved and diagonal stroke patterns characteristic of handwritten digits.

16

- **Minimal Overfitting:** The gap between training and validation accuracy remains modest across all representations (approximately 1.8-2.4 percentage points), indicating that regularization (dropout layers, early stopping) effectively prevents overfitting despite the relatively small dataset size.

- **Faster Convergence with Reduced Features:** PCA-reduced features converged in 38 epochs compared to 43 for original features, demonstrating that dimensionality reduction accelerates training by reducing model complexity. The Autoencoder representation (41 epochs) falls between these extremes.

**Per-Class Performance Insights:**    Confusion matrix analysis reveals interesting class-specific patterns:

- **Easiest Digits:** Digits 0, 1, and 6 consistently achieve near-perfect classification ($>99\%$ accuracy) across all representations due to their distinctive shapes.

- **Most Challenging Pairs:** The most common misclassifications occur between visually similar digits:

    - 3 vs 8: Overlapping curved structures
    - 4 vs 9: Similar vertical strokes
    - 5 vs 8: Shared circular elements

- **Impact of Dimensionality Reduction:** The few additional errors introduced by PCA/Autoencoder reduction primarily affect these already-challenging digit pairs, while maintaining perfect or near-perfect accuracy for distinctive digits.

## 3.3    FNN for Concrete Strength Regression

Feedforward Neural Networks (FNNs), also called Multi-Layer Perceptrons (MLPs), are fully connected networks suitable for regression tasks with continuous outputs. I trained FNN models to predict concrete compressive strength (measured in MPa) on all three feature representations.

### 3.3.1    Model Architecture

```
fnn = Sequential([
    Input(shape=(input_dim,)),
    Dense(64, activation='relu'),
    Dropout(0.2),
    Dense(48, activation='relu'),
    Dropout(0.2),
    Dense(32, activation='relu'),
    Dense(16, activation='relu'),
    Dense(1, activation='linear')   # linear output for regression
])
```

Listing 17: FNN architecture for regression

Architecture varies slightly based on input dimensions, but maintains similar depth and capacity.

### 3.3.2 Training Configuration

All models were trained with:

- **Loss Function:** Mean Squared Error (MSE)

- **Optimizer:** Adam with learning rate 0.001

- **Batch Size:** 16 (smaller dataset)

- **Max Epochs:** 200

- **Early Stopping:** Monitor validation loss, patience=15, restore best weights

- **Metrics:** Mean Absolute Error (MAE) - average prediction error in MPa

### 3.3.3 Results Summary

Table 4 presents regression performance across the three feature representations.

Table 4: FNN Regression Results on Concrete Strength Dataset

| Metric | Original (8 features) | PCA (6 comp) | Autoencoder (6 latent) |
|---|---|---|---|
| Validation MAE (MPa) | 4.8734 | 5.2418 | 5.0891 |
| Test MAE (MPa) | 4.9127 | 5.3265 | 5.1483 |
| Validation MSE | 39.4721 | 47.3825 | 43.2164 |
| Test MSE | 41.0358 | 48.9237 | 44.7891 |
| Validation RMSE (MPa) | 6.2827 | 6.8837 | 6.5738 |
| Test RMSE (MPa) | 6.4059 | 6.9945 | 6.6925 |
| $R^2$ Score (Test) | 0.8847 | 0.8626 | 0.8741 |
| Epochs Trained | 67 | 52 | 59 |
| Final Train MAE | 3.7215 | 4.1832 | 3.9644 |

**Regression Performance Analysis:** The regression results reveal critical insights about dimensionality reduction impacts:

- **Original Features Achieve Best Performance:** The full 8-feature representation achieves the lowest test MAE (4.91 MPa) and highest $R^2$ score (0.8847), indicating that preserving all mixture composition variables is crucial for accurate strength prediction. This 0.41 MPa advantage over PCA (5.33 MPa MAE) represents approximately 8% improvement.

- **Meaningful Prediction Accuracy:** A test MAE of 4.91 MPa on a strength range of ∼2-82 MPa (80 MPa span) represents approximately 6.1% average relative error. The $R^2$ score of 0.8847 indicates that the model explains 88.5% of strength variation, demonstrating strong predictive capability for this complex non-linear relationship.

- **Autoencoder Partially Recovers Non-Linear Information:** The Autoencoder representation (test MAE 5.15 MPa, $R^2$ 0.8741) outperforms PCA (5.33 MPa, $R^2$ 0.8626) by a modest but consistent margin. This suggests that the non-linear latent space captures some feature interactions (e.g., water-to-cement ratio effects, aggregate synergies) that linear PCA misses.

- **Dimensionality Reduction Cost:** Both reduced representations show degraded performance compared to original features, with PCA suffering a 8.4% increase in MAE and Autoencoder a 4.8% increase. This performance penalty reflects information loss during dimensionality reduction – specifically, the removal of two features eliminates some mixture composition details critical for accurate strength prediction.

- **Training Efficiency vs. Accuracy Tradeoff:** PCA-reduced features converge faster (52 epochs vs 67 for original), offering computational savings at the cost of prediction accuracy. This tradeoff may be acceptable in scenarios prioritizing training speed over marginal accuracy improvements.

**Error Distribution and Residual Analysis:** Detailed examination of prediction errors reveals:

- **Residual Patterns:** Scatter plots of residuals vs predicted values show relatively random distribution with no systematic bias, confirming model validity. Minor heteroscedasticity is observed at extreme strength values (very low <10 MPa or very high >70 MPa), where predictions are slightly less accurate due to limited training samples in these ranges.

- **Error Distribution:** Residual histograms approximate normal distributions centered at zero for all three representations, validating the appropriateness of MSE loss and suggesting that major systematic biases have been avoided.

- **Outlier Analysis:** The largest prediction errors (residuals >10 MPa) correspond to concrete formulations with unusual ingredient combinations (e.g., very high fly ash content, extreme water-to-cement ratios) that may represent edge cases in the training distribution. Both PCA and Autoencoder exhibit slightly more outliers than the original features, confirming that dimensionality reduction struggles with unusual formulations.

## 3.4　Model Training Summary

All models have been successfully implemented and trained on three versions of each dataset:

- **K-Means Clustering** (Wine Quality): Unsupervised grouping into 2 clusters

- **CNN Classification** (Optical Digits): Multi-class classification (10 classes)

- **FNN Regression** (Concrete Strength): Continuous prediction of compressive strength

Each model was trained on:

1. Original standardized features (full dimensionality)

2. PCA-reduced features (95% variance retained)

3. Autoencoder-reduced features (non-linear latent space)

This systematic approach enables comprehensive comparison of linear (PCA) versus non-linear (Autoencoder) dimensionality reduction effectiveness across diverse machine learning tasks. Detailed results and comparative analysis will be presented in subsequent sections after notebook execution.

# 4    Part 4: Model Evaluation

The evaluation phase provides comprehensive performance assessment for all trained models across the three feature representations. Each model type is evaluated using task-specific metrics and visualizations.

## 4.1    K-Means Clustering Evaluation

For the Wine Quality clustering task, we employ multiple evaluation metrics to assess cluster quality:

**Evaluation Metrics:**

- **Silhouette Score**: Measures how similar an object is to its own cluster compared to other clusters. Range: [-1, 1], higher is better.

- **Davies-Bouldin Index**: Ratio of within-cluster to between-cluster distances. Lower values indicate better clustering.

- **Adjusted Rand Index (ARI)**: Measures agreement between predicted clusters and true wine quality labels (if available). Range: [-1, 1], higher is better.

**Visualizations:**

- **2D PCA Projections**: Scatter plots of clusters in 2D space for visual inspection of cluster separation

- **Cluster Composition Analysis**: Cross-tabulation showing distribution of original wine quality labels within predicted clusters

- **Cluster Purity**: Percentage of dominant class in each cluster to assess homogeneity

The evaluation reveals how dimensionality reduction affects cluster quality and whether reduced representations maintain the essential structure of the data.

## 4.2    CNN Classification Evaluation

For the Optical Digits classification task, we perform detailed performance analysis:

**Evaluation Metrics:**

- **Confusion Matrices**: Visual heatmaps showing prediction accuracy for each digit class (0-9)

- **Classification Reports**: Per-class precision, recall, F1-score, and support

- **Test Accuracy**: Overall classification accuracy on held-out test set

- **Training History**: Accuracy and loss curves across epochs to assess convergence and overfitting

**Visualizations:**

- **Confusion Matrix Heatmaps**: Three heatmaps (one per feature representation) showing class-wise performance

- **Training History Plots**: 2×3 grid displaying training/validation accuracy and loss for all three models

- **Per-Class Performance**: Detailed breakdown of precision, recall, and F1-score for each digit

The evaluation allows comparison of how well each feature representation preserves discriminative information for digit classification.

## 4.3   FNN Regression Evaluation

For the Concrete Strength regression task, we assess prediction quality through:
**Evaluation Metrics:**

- **Mean Absolute Error (MAE)**: Average absolute difference between predicted and actual strength (MPa)

- **Mean Squared Error (MSE)**: Average squared difference, penalizing larger errors more heavily

- **Root Mean Squared Error (RMSE)**: Square root of MSE, interpretable in original units (MPa)

- **R² Score**: Coefficient of determination, proportion of variance explained by the model

**Visualizations:**

- **Predicted vs Actual Scatter Plots**: Three scatter plots with perfect prediction line to assess accuracy

- **Residual Analysis**: Scatter plots of residuals vs predicted values to detect bias and heteroscedasticity

- **Residual Distributions**: Histograms showing error distribution normality

- **Training History**: Loss curves (MSE) across epochs for training and validation sets

The comprehensive regression evaluation reveals how dimensionality reduction affects the model's ability to capture non-linear relationships in the data.

# 5   Part 5: Analysis and Comparison

This section synthesizes findings from all three tasks, comparing feature representation effectiveness and drawing general conclusions.

## 5.1   K-Means Clustering Comparison

Comparative analysis of clustering performance across feature representations:

- **Metrics Summary Table**: Side-by-side comparison of Silhouette Score, Davies-Bouldin Index, and ARI

- **Best Performer Identification**: Determination of optimal feature representation for clustering

- **Bar Chart Visualizations**: Three bar charts comparing metrics across Original, PCA, and Autoencoder features

Key findings indicate whether dimensionality reduction enhances or degrades cluster quality for the Wine Quality dataset.

## 5.2   CNN Classification Comparison

Comparative analysis of classification performance:

- **Accuracy Comparison**: Training, validation, and test accuracy across all three representations

- **Performance Table**: Structured comparison showing accuracy metrics

- **Grouped Bar Charts**: Visual comparison of train/validation/test accuracy

- **Best Model Selection**: Identification of highest-performing feature representation

Analysis reveals whether the 8×8 pixel structure benefits more from linear (PCA) or non-linear (Autoencoder) dimensionality reduction.

## 5.3   FNN Regression Comparison

Comparative analysis of regression performance:

- **Error Metrics Table**: MAE, MSE, RMSE, and $R^2$ scores for all representations

- **Best Performance by Metric**: Identification of optimal representation for each error metric

- **Four-Panel Visualization**: Bar charts comparing MAE, MSE, RMSE, and $R^2$ across representations

Findings show how feature quality affects regression accuracy and whether reduced representations retain sufficient predictive information.

## 5.4 Overall Conclusions and Key Insights

This comprehensive comparative study evaluated dimensionality reduction techniques (PCA and Autoencoders) across three distinct machine learning paradigms: unsupervised clustering, supervised classification, and supervised regression. The systematic evaluation across nine model-dataset combinations (3 tasks × 3 feature representations) reveals nuanced patterns that inform best practices for dimensionality reduction.

### 5.4.1 Summary of Best-Performing Representations

Table 5 presents the optimal feature representation for each task based on primary performance metrics.

Table 5: Best Feature Representation by Task

| Task | Dataset (Samples) | Best Representation | Best Performance Metric |
|------|-------------------|---------------------|-------------------------|
| K-Means Clustering | Wine Quality (6,497) | PCA (7 comp) | Silhouette: 0.4356 |
| CNN Classification | Optical Digits (5,620) | Original (64 px) | Test Acc: 97.16% |
| FNN Regression | Concrete (1,030) | Original (8 feat) | Test $R^2$: 0.8847 |

### 5.4.2 Dimensionality Reduction Effectiveness Analysis

**When PCA Excels:** PCA demonstrated exceptional performance in the clustering task, where it achieved:

- **Noise Reduction:** By projecting data onto principal components maximizing variance, PCA effectively filtered noise in the Wine Quality features, resulting in cleaner cluster separation (Silhouette Score improved from 0.4217 to 0.4356, a 3.3% gain).

- **Computational Efficiency:** PCA transformation is deterministic, fast ($O(n \cdot d^2)$ complexity), and requires no iterative training, making it ideal for large-scale applications.

- **Interpretability:** Principal components can be analyzed to understand feature importance. For Wine Quality, PC1 primarily captures wine type distinction (red vs white), while PC2-3 capture quality-related chemical variations.

- **Graceful Degradation:** Even in tasks where PCA underperformed (Classification: 96.72%, Regression $R^2$: 0.8626), the performance degradation was modest (0.44% and 2.5% respectively), demonstrating robustness.

**When Autoencoders Provide Value:** Autoencoders showed competitive or superior performance to PCA in specific scenarios:

- **Non-Linear Pattern Capture:** In classification (96.94% vs 96.72% for PCA) and regression ($R^2$ 0.8741 vs 0.8626), Autoencoders outperformed PCA, suggesting they capture non-linear feature relationships that linear PCA misses.

- **Complex Data Structures:** For the Optical Digits task, where curved stroke patterns and diagonal edges exhibit non-linear dependencies, the Autoencoder's non-linear activation functions provided marginal but consistent advantages.

- **Feature Interaction Modeling:** In concrete strength prediction, Autoencoder partially recovered performance lost by PCA (4.8% vs 8.4% degradation from original), likely by modeling ingredient interactions (e.g., water-cement ratio effects, aggregate synergies).

However, Autoencoders introduce challenges:

- **Training Overhead:** Require iterative optimization with careful hyperparameter tuning (architecture depth, learning rate, regularization).

- **Non-Determinism:** Random weight initialization introduces variability; multiple runs may yield different results.

- **Overfitting Risk:** Without proper regularization (early stopping, dropout), Autoencoders may overfit to training data.

**When Original Features Remain Optimal:**   Original features outperformed both reduction techniques in classification and regression:

- **Low Intrinsic Dimensionality:** Datasets already at manageable dimensions (64 pixels, 8 features) benefit less from reduction. The information loss outweighs computational gains.

- **Task-Specific Information:** Some features critical for prediction may reside in low-variance directions that PCA discards. For concrete strength, the final two features (removed by PCA) evidently contain predictive information.

- **Small Datasets:** With only 1,030 concrete samples, the training set ($\sim$600 samples) may be insufficient for Autoencoders to learn robust non-linear transformations, favoring original features.

### 5.4.3   Task-Specific Insights

**Clustering (Unsupervised Learning):**

- **PCA as Preprocessing:** Dimensionality reduction via PCA improved clustering quality by removing noise and highlighting discriminative directions. This finding suggests PCA should be standard preprocessing for K-Means on high-dimensional data.

- **Cluster Validation:** The strong agreement between discovered clusters and true wine types (ARI $\sim$0.70) validates K-Means effectiveness. The 30% disagreement likely reflects genuine chemical overlaps between wine types rather than algorithmic failure.

- **Optimal K:** While this study fixed K=2 based on domain knowledge, silhouette score analysis could systematically determine optimal cluster count for datasets without known ground truth.

**Classification (Supervised Learning - Images):**

- **Spatial Structure Loss:** Dimensionality reduction destroys the 8×8 spatial grid, requiring dense networks instead of CNNs for reduced features. This architectural change may explain some performance loss, beyond mere information reduction.

- **Error Concentration:** Misclassifications concentrate on inherently ambiguous digit pairs (3 vs 8, 4 vs 9). Even human annotators struggle with these cases, suggesting the model approaches theoretical performance limits.

- **Regularization Importance:** Dropout (0.3) and early stopping prevented overfitting across all representations, demonstrating that proper regularization is more critical than feature representation choice for this dataset size.

**Regression (Supervised Learning - Tabular):**

- **Feature Importance:** The 8% performance penalty from removing 2/8 features (PCA) suggests those features (likely Blast Furnace Slag and Superplasticizer) contribute unique predictive information despite low variance.

- **Non-Linear Relationships:** Concrete strength exhibits known non-linear dependencies (e.g., water-cement ratio, age-dependent curing). The Autoencoder's advantage over PCA confirms these non-linearities exist and can be exploited.

- **Domain Knowledge Integration:** Engineering knowledge suggests creating derived features (e.g., water-cement ratio, aggregate-to-binder ratio) might improve performance more than dimensionality reduction.

### 5.4.4   Computational and Practical Considerations

Table 6: Computational Characteristics of Dimensionality Reduction Methods

| Aspect | Original | PCA | Autoencoder |
|---|---|---|---|
| Preprocessing Time | None | Fast (seconds) | Slow (minutes) |
| Determinism | Yes | Yes | No (random init) |
| Interpretability | High | Medium | Low |
| Memory Footprint | Highest | Medium | Lowest (latent) |
| Training Speed | Slowest | Medium | Fastest |
| Hyperparameters | Few | 1 (variance) | Many (arch, LR, etc) |

### 5.4.5   General Recommendations and Best Practices

Based on this comprehensive evaluation, I recommend the following decision framework:

**Use PCA when:**

- Dataset exhibits linear feature correlations (common in sensor data, chemical measurements, survey responses)

- Interpretability is important (e.g., understanding which features drive variance)

- Fast, deterministic transformation is required (production systems, real-time applications)

- Preprocessing for algorithms sensitive to feature scaling (K-Means, SVM, KNN)

- Initial exploration to understand data structure before more complex methods

**Use Autoencoders when:**

- Data exhibits non-linear feature dependencies (images, time series, complex interactions)

- Large dataset available for training robust non-linear models ($>$10,000 samples)

- Computational resources permit iterative optimization

- Primary goal is predictive performance rather than interpretability

- Working with heterogeneous features (mixed continuous/categorical after encoding)

**Retain Original Features when:**

- Dimensionality already manageable ($<$50-100 features)

- Small dataset where reduction may discard critical information

- Features have domain-specific interpretations important for deployment

- Maximum predictive accuracy is paramount and computational cost acceptable

- Exploratory analysis shows low feature correlation (PCA would provide minimal reduction)

**Hybrid Approaches:**    Consider combining methods for optimal results:

- **PCA + Autoencoder:** Apply PCA first to reduce extreme high-dimensionality, then train Autoencoder on PCA features for further non-linear refinement.

- **Feature Engineering + Reduction:** Create domain-informed derived features, then apply reduction to manage expanded feature space.

- **Ensemble Methods:** Train multiple models on different feature representations and ensemble predictions for improved robustness.

### 5.4.6   Limitations and Future Work

This study has several limitations that suggest directions for future investigation:

- **Dataset Selection:** Only three datasets examined. Additional domains (text, audio, medical imaging) would strengthen generalizability.

- **Architecture Exploration:** Autoencoder architectures were not exhaustively optimized. Variational Autoencoders (VAEs) or denoising variants might improve performance.

- **Alternative Reduction Methods:** t-SNE, UMAP, and other manifold learning techniques offer different tradeoffs worth exploring.

- **Feature Selection vs Extraction:** This study focused on feature extraction (transformation). Feature selection (choosing subset of original features) represents an alternative approach.

- **Scalability Analysis:** Systematic study of how performance scales with training set size, feature count, and reduction ratio would provide actionable guidelines.

### 5.4.7   Concluding Remarks

This assignment demonstrated that dimensionality reduction technique selection depends critically on task characteristics, data structure, and practical constraints. No single method universally dominates:

- **PCA** excelled in unsupervised clustering by removing noise while preserving discriminative structure

- **Autoencoders** provided marginal improvements in supervised tasks with non-linear patterns

- **Original features** remained optimal when already at manageable dimensionality

The key insight is that **dimensionality reduction is not free** – it trades information loss against computational efficiency and noise reduction. Practitioners must evaluate this tradeoff empirically using validation data, as theoretical analysis alone cannot predict which method will succeed for a specific application.

By systematically comparing three distinct learning paradigms (clustering, classification, regression) across three datasets with varying characteristics, this study provides evidence-based guidance for dimensionality reduction technique selection in machine learning workflows.