

AI-Capstone: Project1 Report

111550093 朱驛庭

March 17, 2025

1 Introduction

In this project, our objective is to classify images of clothing items into predefined categories using various machine learning techniques. The task involves implementing and comparing the performance of three different models: ResNet, a deep learning-based approach; K-means, an unsupervised learning method; and SVM, a traditional supervised learning algorithm. The goal is to evaluate the effectiveness of these models in accurately classifying clothing items based on their visual characteristics.

2 Dataset

The dataset used in this project comprises clothing item images collected from Google Search and Unsplash. It is shared with a classmate (student ID: 111550160). The dataset can be accessed through the [link \(click\)](#).

2.1 Source

To assess the generalization capabilities of different models, the dataset includes images with varying levels of similarity within each category. For example, some images feature models wearing clothing items, while others showcase only clothing. The dataset comprises ten clothing categories: T-shirt, bag, dress, glasses, hat, hoodie, jacket, pants, shoes, and socks, with 50 images per category, totaling 500 images.

2.2 Train-test Split

The dataset is divided into training and testing sets, with 80% of the data used for training and 20% for testing. This split ensures that the models are trained on a substantial amount of data while retaining a separate set for evaluating their performance.

2.3 Preprocessing

I organized the images into separate folders based on their respective categories, naming each folder after the corresponding clothing type. Additionally, all images were converted to JPG format and sequentially renamed from 1 to 50.

3 Methods

3.1 Supervised Method1: Support Vector Machine

Support Vector Machine (SVM) is a supervised learning algorithm designed to find the optimal hyperplane that best separates data into different classes while maximizing the margin between them. The margin is defined as the distance between the hyperplane and the closest data points, known as **support vectors**. A larger margin generally leads to better generalization and robustness against overfitting. The illustration is shown in Fig. 1. The SVM classifier seeks to minimize the following objective function:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \quad \text{s.t.} \quad y_i(w^T x_i + b) \geq 1, \forall i$$

where w represents the weight vector, b is the bias term, x_i denotes the input feature vector, and y_i is the corresponding class label (+1 or -1).

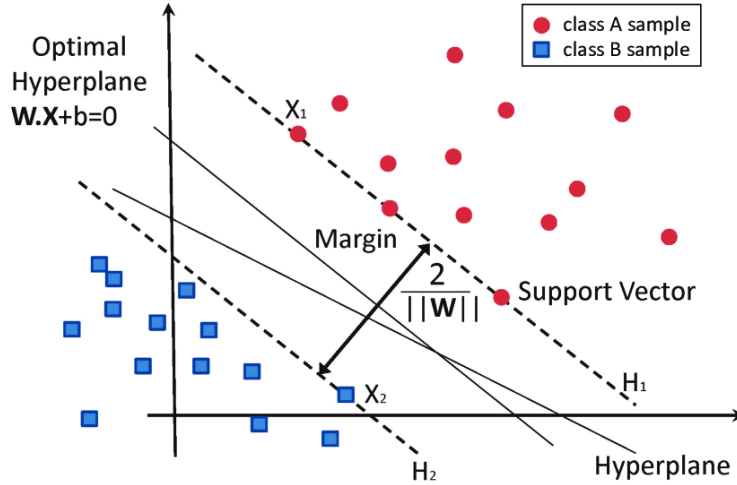


Figure 1: SVM Illustration [1].

3.1.1 Hinge Loss and Soft Margin SVM

To handle cases where data points may not be perfectly separable, SVM introduces **soft margin classification**, allowing for certain misclassifications by incorporating **slack variables** ξ_i . This is achieved by optimizing the **hinge loss function**, which penalizes misclassified points:

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

subject to:

$$y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \forall i$$

where C is a regularization parameter that balances the trade-off between maximizing the margin and minimizing classification errors.

3.1.2 Kernel Trick for Non-Linearly Separable Data

When the data is not linearly separable in the original feature space, SVM leverages the **kernel trick** to project the data into a higher-dimensional space where it becomes linearly separable. The kernel function replaces the direct computation of dot products, effectively mapping data into a transformed feature space:

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

where $\phi(x)$ is a mapping function that transforms the input space into a higher-dimensional feature space.

3.2 Supervised Method2: ResNet18

ResNet18 is a deep convolutional neural network (CNN) that belongs to the ResNet (Residual Network) family, introduced by He et al. in 2015. It is designed to address the vanishing gradient problem in deep networks by incorporating residual connections, allowing for deeper architectures while maintaining efficient training. ResNet18 consists of 18 layers, including convolutional, batch normalization, ReLU activation, and fully connected layers.

3.2.1 Architecture and Modification

ResNet18 follows a hierarchical structure with an initial convolutional layer, followed by 8 residual blocks grouped into 4 stages, and ending with a global average pooling and a fully connected layer. The modified version in this project replaces the original 1000-class output layer with a fully connected layer of size 10, corresponding to our clothing categories. The softmax activation is removed to output raw logits, which are processed by a cross-entropy loss function. The modified final layer is defined as:

$$\text{FC Layer: } y = Wx + b, \quad W \in \mathbb{R}^{10 \times 512}, \quad b \in \mathbb{R}^{10}$$

The overall structure of ResNet18 is illustrated in Fig. 2.

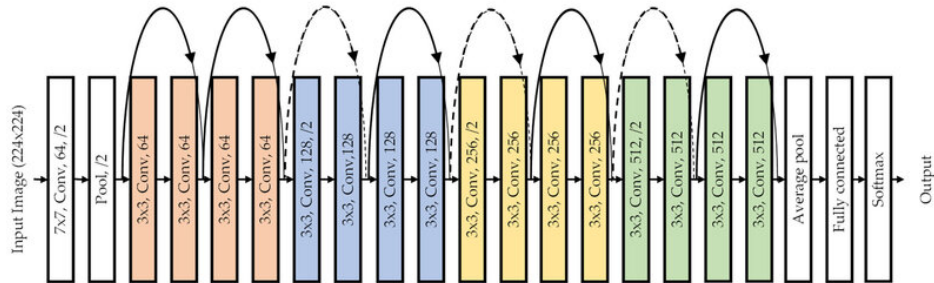


Figure 2: Structure of the ResNet18 Model [2].

3.2.2 Loss Function

The model is trained using the **cross-entropy loss function**, defined as:

$$L_{CE} = - \sum_{i=1}^N [y_i \log p_i + (1 - y_i) \log(1 - p_i)]$$

where y_i is the true class label, p_i is the predicted probability, and N is the total number of samples.

3.3 K-Means Clustering

K-Means is a widely used unsupervised learning algorithm for clustering, aiming to partition a dataset into K clusters based on feature similarity. It operates iteratively to minimize intra-cluster variance, ensuring that data points within a cluster are as close as possible to the cluster centroid. The objective function, known as the **within-cluster sum of squares (WCSS)**, is given by:

$$\sum_{j=1}^K \sum_{i \in C_j} \|x_i - c_j\|^2$$

where C_j represents the set of data points assigned to cluster j , and c_j is the centroid of that cluster.

4 Implementation Details

4.1 Evaluation Metrics

4.1.1 Supervised Methods

For the supervised methods, namely ResNet and SVM, I employ several evaluation metrics to assess the performance of the models:

- **Confusion Matrix:** A confusion matrix is used to visualize the classification model's performance by organizing predictions into four categories:

$$\begin{bmatrix} TP & FP \\ FN & TN \end{bmatrix}$$

where:

- **True Positives (TP):** Correctly predicted positive instances.
- **False Positives (FP):** Incorrectly predicted positive instances.
- **False Negatives (FN):** Incorrectly predicted negative instances.
- **True Negatives (TN):** Correctly predicted negative instances.

- **Accuracy:** This metric measures the proportion of correctly classified instances over the total number of instances, providing a general sense of model performance. It is defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Precision:** Precision measures the proportion of correctly predicted positive instances among all instances classified as positive. It is given by:

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall (Sensitivity):** Recall measures the model's ability to correctly identify all actual positive instances. It is calculated as:

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1-Score:** The F1-score is the harmonic mean of precision and recall, balancing the trade-off between them. It is defined as:

$$\text{F1-Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

These metrics are computed using the `classification_report` function from scikit-learn, providing a comprehensive overview of the model's performance across different classes.

4.1.2 Unsupervised Method

For the unsupervised method, K-means, I use different metrics to evaluate clustering performance:

- **Silhouette Score:** This metric measures how similar an object is to its own cluster compared to other clusters. It ranges from -1 to 1, where a higher score indicates better-defined clusters.
- **Adjusted Rand Index (ARI):** ARI measures the similarity between the clustering results and the true labels, adjusted for chance. It ranges from -1 to 1, with 1 indicating perfect agreement.
- **Calinski-Harabasz Score:** This metric evaluates the ratio of between-cluster dispersion to within-cluster dispersion. A higher score indicates well-separated and compact clusters.
- **Davies-Bouldin Score:** This metric assesses the average similarity between each cluster and its most similar cluster. A lower score indicates better clustering quality, as it suggests that clusters are well-separated.

4.2 Hyperparameters

- **ResNet:** The ResNet model’s hyperparameters are carefully chosen to ensure effective training and convergence. The learning rate is initialized at 0.001, and a **StepLR** scheduler is used to dynamically adjust the learning rate during training. This scheduler reduces the learning rate by a factor of 0.1 every 10 epochs, helping the model converge more smoothly. The random seed is set to 40 to ensure reproducibility of results. The model is initialized with weights pre-trained on the ImageNet dataset, which provides a strong starting point for feature extraction. The batch size is set to 32, and the model is trained for 20 epochs, balancing computational efficiency with model performance.
- **SVM:** For the SVM model, I perform a grid search over a range of hyperparameters, including the regularization parameter C , kernel type (e.g., **rbf**, **poly**), and kernel-specific parameters like gamma and degree. This search is conducted using **GridSearchCV** from scikit-learn, which performs cross-validation to find the optimal hyperparameter combination.
- **K-means:** The primary hyperparameter for K-means is the number of clusters K . I experiment with different values of K and use the ARI score to select the optimal number of clusters. Additionally, I use UMAP for dimensionality reduction before clustering, which involves tuning the number of components and the minimum distance between points.

5 Experiments

5.1 Data Augmentation

Data augmentation is a technique commonly used in machine learning to artificially expand the training dataset by applying transformations to the input images. The key idea is that by introducing non-identical inputs to the model in each epoch, data augmentation helps prevent overfitting and improves the generalization ability of the model. This is particularly beneficial for deep learning models, which tend to overfit when trained on small datasets. The comparison is shown below:

Methods	Aug.	Accuracy(%)	Precision(%)	Recall(%)	F1-score(%)
SVM	✗	73.74	76.30	73.80	73.80
SVM	✓	69.70	70.70	69.70	69.30
ResNet18	✗	86.87	89.60	86.80	86.80
ResNet18	✓	91.92	92.50	91.80	91.70

Table 1: Performance comparison of the effect of data augmentation on the two supervised methods.

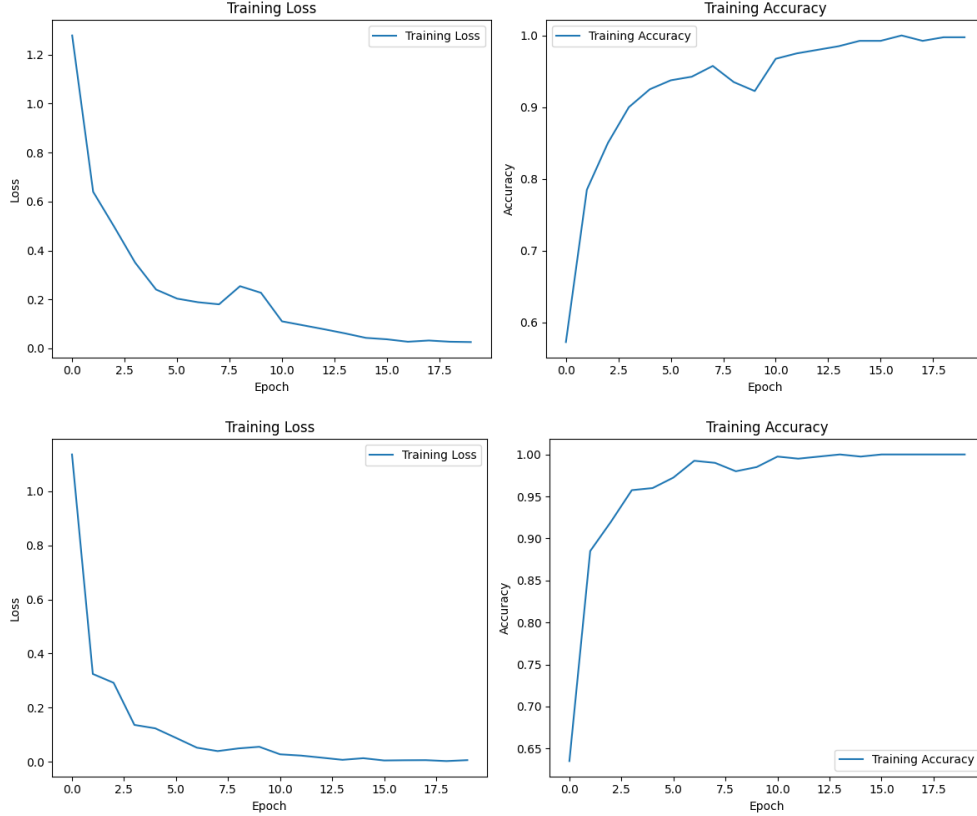


Figure 3: Training history of ResNet18 with and without data augmentation. The top row represents training with data augmentation, while the bottom row represents training without it.

From the results presented in Table 1, it is evident that data augmentation enhances the performance of ResNet18. This improvement suggests that ResNet18 benefits from the increased diversity introduced by augmentation, allowing it to generalize better to unseen data. However, the same approach does not yield positive results for SVM, as it is not a deep learning-based method. Instead of improving performance, the additional variations introduced by data augmentation negatively impact SVM, likely due to its reliance on hand-crafted features rather than learned representations.

Furthermore, as shown in Figure 3, data augmentation slows down the training process compared to the original setup. The model takes longer to converge since it needs to learn from a more diverse set of augmented samples. In typical cases, data augmentation may require additional training epochs to reach the same level of training loss as a model trained without augmentation.

5.2 t-SNE Visualization on ResNet18

t-SNE (t-distributed stochastic neighbor embedding) is a dimensionality reduction technique commonly used for visualization. By mapping high-dimensional data into a lower-dimensional space, t-SNE helps reveal underlying patterns and relationships within the data.

In Figure 4, we observe that as the data progresses through each layer of ResNet, the separation between different classes becomes more distinct. Notably, in the fully connected layer, the *hoodie*, *jacket*, and *t-shirt* classes are positioned relatively close to one another. This alignment is reasonable, as these three types of clothing share similar visual characteristics in the real world.

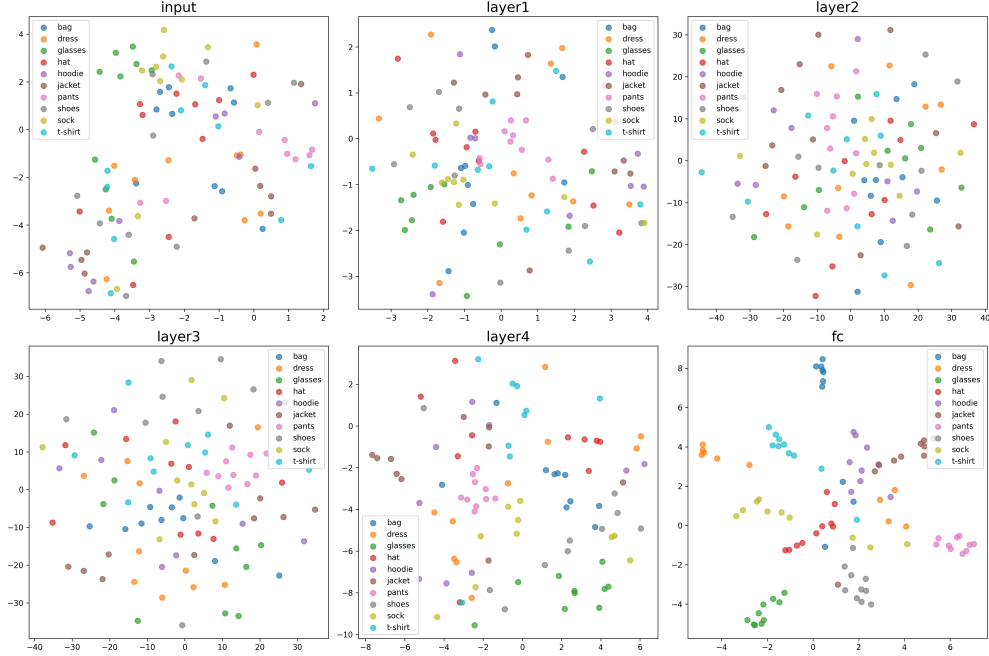


Figure 4: t-SNE visualization of different layers on ResNet18.

5.3 Multi-criteria Cluster Validation of K-means

Based on the result shown on Figure. 5. The analysis of K-means clustering metrics suggests that K=5 or K=6 provides the best balance across multiple validation metrics. The **Silhouette Score** and **Calinski-Harabasz Score** indicate that K=5 achieves strong separation, while K=6 has the highest **Adjusted Rand Index (ARI)**, suggesting better alignment with true labels. The **Davies-Bouldin Score** further supports K=5 or K=6 as optimal choices. At K=6, some clusters achieve high purity, such as bags (77.78%) and glasses (40.74%), while upper-body garments (t-shirts, jackets, and hoodies) tend to mix. The results indicate that while K-means effectively groups certain categories, visual similarities between some clothing items limit perfect separation. For general clothing classification, K=5 or K=6 is recommended, while K=2 or K=3 may suffice for broader category distinctions.

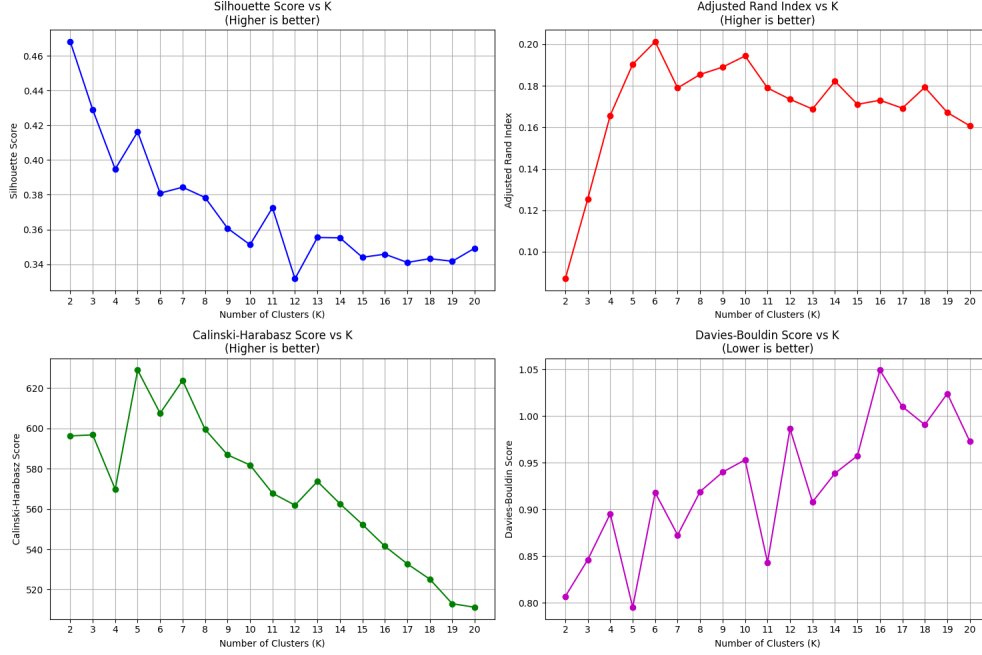


Figure 5: K-means clustering evaluation metrics for different numbers of clusters (K).

6 Discussion

6.1 Performance Comparison and Analysis

In this project, I have implemented and evaluated three different machine learning approaches for clothing image classification: ResNet18, SVM, and K-means. The experiments yielded several interesting observations:

- **ResNet18** demonstrated superior performance with a test accuracy of 91.92% when using data augmentation, outperforming the other methods. This result aligns with expectations, as deep learning methods typically excel at image classification tasks due to their ability to learn hierarchical representations directly from raw pixel data.
- **SVM** achieved a respectable accuracy of 73.74% without data augmentation. Interestingly, data augmentation decreased its performance, contradicting the common assumption that augmentation universally improves classification results. This suggests that traditional machine learning algorithms like SVM may struggle with the increased variance introduced by augmentation when using handcrafted features.
- **K-means**, as an unsupervised approach, showed limitations in perfectly separating the clothing categories, achieving a maximum ARI of 0.20 at K=6. This was expected, as unsupervised methods lack the guidance of class labels during training. However, the clustering revealed meaningful patterns, particularly in separating accessories from clothing items.

6.2 Factors Affecting Performance

Several factors influenced the performance of the different models:

- **Feature Representation:** ResNet18's superior performance can be attributed to its ability to learn complex, hierarchical features directly from raw images. In contrast, SVM relied on manually engineered features, which, while effective, may not capture all the nuances present in the data.
- **Dataset Characteristics:** The dataset's inherent challenges, such as varying backgrounds, lighting conditions, and item orientations, affected all models. Items with distinctive shapes (like bags) were generally easier to classify or cluster correctly compared to items with similar appearances (like t-shirts, hoodies, and jackets).
- **Visual Similarity:** The t-SNE visualization and K-means clustering results both revealed that visually similar categories (such as hoodie, jacket, and t-shirt) tend to overlap in the feature space. This suggests that even with sophisticated feature extraction, some ambiguity between similar clothing categories is inevitable.
- **Dimensionality Reduction:** For K-means, the UMAP dimensionality reduction played a crucial role in clustering performance. The optimal number of dimensions balanced between preserving data structure and reducing noise.

6.3 Lessons Learned

Overall, this project demonstrates the complementary nature of different machine learning approaches for image classification tasks. Each method offers unique insights and advantages, with deep learning-based approaches like ResNet18 providing state-of-the-art performance, traditional methods like SVM offering interpretability and computational efficiency, and unsupervised methods like K-means revealing natural data patterns without requiring labeled data.

References

- [1] D. Man, *Support vector machine in python*. [Online]. Available: <https://medium.com/@dattanaman213/support-vector-machine-in-python-576eaac337ae>.
- [2] ResearchGate, *Structure of the resnet-18 model*. [Online]. Available: https://www.researchgate.net/figure/Structure-of-the-Resnet-18-Model_fig1_366608244.

Appendix

Github link: <https://github.com/ChuEating1005/AI-Capstone/tree/main/Project1>

Listing 1: prepare_dataset.py

```
1 import os
2 import shutil
3 import random
4 from pathlib import Path
5 import numpy as np
6
7 def main():
8     # Define paths
9     SOURCE_DIR = Path('images')
10    OUTPUT_DIR = Path('data')
11
12    # Define split ratios
13    TRAIN_RATIO = 0.8
14    VAL_RATIO = 0
15    TEST_RATIO = 0.2
16
17    # Create output directories
18    train_dir = OUTPUT_DIR / 'train'
19    test_dir = OUTPUT_DIR / 'test'
20
21    # Create directories if they don't exist
22    for dir_path in [train_dir, test_dir]:
23        for category in os.listdir(SOURCE_DIR):
24            if category.startswith('.'): # Skip hidden files like .DS_Store
25                continue
26            os.makedirs(dir_path / category, exist_ok=True)
27
28    # Set random seed for reproducibility
29    random.seed(42)
30    np.random.seed(42)
31
32    # Process each category
33    for category in os.listdir(SOURCE_DIR):
34        if category.startswith('.'): # Skip hidden files like .DS_Store
35            continue
36
37        category_dir = SOURCE_DIR / category
38        if not os.path.isdir(category_dir):
39            continue
40
41        # Get all image files
42        image_files = [f for f in os.listdir(category_dir) if f.endswith(('.'
43        jpg', '.jpeg', '.png'))]
44
45        # Shuffle the files
46        random.shuffle(image_files)
47
48        # Calculate split indices
```

```

48     n_images = len(image_files)
49     n_train = 40
50
51     # Split the files
52     train_files = image_files[:n_train]
53     test_files = image_files[n_train:]
54
55     # Copy files to respective directories
56     for files, target_dir in zip([train_files, test_files], [train_dir,
test_dir]):
57         for file in files:
58             src_path = category_dir / file
59             dst_path = target_dir / category / file
60             shutil.copy2(src_path, dst_path)
61
62             print(f"Category {category}: {len(train_files)} train, {len(test_files
)} test")
63
64             print("Dataset preparation completed!")
65
66 if __name__ == "__main__":
67     main()

```

Listing 2: dataset.py

```

1 from torch.utils.data import Dataset
2 from torchvision import transforms
3 from PIL import Image
4 import os
5 import numpy as np
6 import cv2
7 from skimage.feature import local_binary_pattern
8 from skimage.feature import hog
9 import umap
10
11 data_aug = True
12
13 train_tfm = transforms.Compose([
14     transforms.Resize((224, 224)),
15     transforms.RandomHorizontalFlip(p=0.5),
16     transforms.RandomRotation(15),
17     transforms.ColorJitter(brightness=0.3, contrast=0.3, saturation=0.3, hue
=0.1),
18     transforms.RandomAffine(degrees=0, translate=(0.1, 0.1)),
19     transforms.RandomGrayscale(p=0.1),
20     transforms.ToTensor(),
21     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
22 ]) if data_aug else transforms.Compose([
23     transforms.Resize((224, 224)),
24     transforms.ToTensor(),
25     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
26 ])
27
28 valid_tfm = transforms.Compose([

```

```

29     transforms.Resize((224, 224)),
30     transforms.ToTensor(),
31     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
32 ])
33
34 test_tfm = transforms.Compose([
35     transforms.Resize((224, 224)),
36     transforms.ToTensor(),
37     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
38 ])
39
40 svm_tfm = transforms.Compose([
41     transforms.RandomHorizontalFlip(p=0.5),
42     transforms.RandomRotation(15),
43     transforms.ColorJitter(brightness=0.3, contrast=0.3, saturation=0.3, hue
44                             =0.1),
45     transforms.RandomAffine(degrees=0, translate=(0.1, 0.1)),
46     transforms.RandomGrayscale(p=0.1),
47 ])
48
49 class ResnetDataset(Dataset):
50     def __init__(self, data_dir, transform=None):
51         self.data_dir = data_dir
52         self.transform = transform
53         self.samples = []
54         self.classes = []
55         self.loader = Image.open
56
57         # Populate samples and classes
58         for class_name in os.listdir(data_dir):
59             class_path = os.path.join(data_dir, class_name)
60             if os.path.isdir(class_path):
61                 self.classes.append(class_name)
62                 for file_name in os.listdir(class_path):
63                     if file_name.lower().endswith(('.png', '.jpg', '.jpeg')):
64                         file_path = os.path.join(class_path, file_name)
65                         self.samples.append((file_path, self.classes.index(
66 class_name)))
67
68     def __len__(self):
69         return len(self.samples)
70
71     def __getitem__(self, index):
72         """
73         Override the __getitem__ method to skip corrupted images
74         """
75         path, label = self.samples[index]
76         try:
77             sample = self.loader(path)
78
79             if self.transform is not None:
80                 sample = self.transform(sample)

```

```

81         return sample, label
82
83     except Exception as e:
84         print(f"Error loading image {path}: {e}")
85         # Return a random valid image instead
86         valid_idx = (index + 1) % len(self)
87         return self.__getitem__(valid_idx)
88
89 class SVMDataset:
90     def __init__(self, data_dir, augmentations=3):
91         self.data_dir = data_dir
92         self.augmentations = augmentations
93         self.train_transform = svm_tfm
94         self.test_transform = test_tfm
95         self.class_names = []
96         self.class_to_idx = {}
97
98     def load_data(self):
99         X_train, y_train = [], []
100         X_test, y_test = [], []
101
102         # Get class names from train directory
103         self.class_names = [d for d in os.listdir(os.path.join(self.data_dir,
104 'train'))
105                             if os.path.isdir(os.path.join(self.data_dir, '
106 train', d))]
107
108         # Create class to index mapping
109         self.class_to_idx = {cls_name: i for i, cls_name in enumerate(self.
110 class_names)}
111
112         # Process training data with augmentation
113         print("Loading training data with augmentation...")
114         for class_name in self.class_names:
115             class_dir = os.path.join(self.data_dir, 'train', class_name)
116             for img_file in os.listdir(class_dir):
117                 if img_file.lower().endswith(('.png', '.jpg', '.jpeg')):
118                     img_path = os.path.join(class_dir, img_file)
119                     img = Image.open(img_path).convert('RGB')
120
121                     # Original image
122                     features = self.extract_features(img, augment=False)
123                     if features is not None:
124                         X_train.append(features)
125                         y_train.append(self.class_to_idx[class_name])
126
127                     # Augmented versions
128                     for _ in range(self.augmentations):
129                         aug_features = self.extract_features(img, augment=True)
130
131                         if aug_features is not None:
132                             X_train.append(aug_features)
133                             y_train.append(self.class_to_idx[class_name])

```

```

131     # Process test data (no augmentation)
132     print("Loading test data...")
133     for class_name in self.class_names:
134         class_dir = os.path.join(self.data_dir, 'test', class_name)
135         for img_file in os.listdir(class_dir):
136             if img_file.lower().endswith(('.png', '.jpg', '.jpeg')):
137                 img_path = os.path.join(class_dir, img_file)
138                 img = Image.open(img_path).convert('RGB')
139                 features = self.extract_features(img, augment=False)
140                 if features is not None:
141                     X_test.append(features)
142                     y_test.append(self.class_to_idx[class_name])
143
144     return np.array(X_train), np.array(y_train), np.array(X_test), np.
array(y_test)
145
146     def extract_features(self, img, augment=False):
147         if augment:
148             img = self.train_transform(img)
149
150         # Resize to 128x128
151         img = img.resize((128, 128))
152         img_array = np.array(img)
153
154         # 1. Color Features (HSV space)
155         img_hsv = cv2.cvtColor(img_array, cv2.COLOR_RGB2HSV)
156
157         # HSV histograms
158         hist_h = np.histogram(img_hsv[:, :, 0], bins=16)[0]
159         hist_s = np.histogram(img_hsv[:, :, 1], bins=16)[0]
160         hist_v = np.histogram(img_hsv[:, :, 2], bins=16)[0]
161
162         # Color statistics
163         color_stats = np.concatenate([
164             np.mean(img_hsv, axis=(0,1)), # mean of each channel
165             np.std(img_hsv, axis=(0,1)), # std of each channel
166             np.percentile(img_hsv, [25,75], axis=(0,1)).flatten() # quartiles
167         ])
168
169         # 2. Texture Features
170         gray = cv2.cvtColor(img_array, cv2.COLOR_RGB2GRAY)
171
172         # Convert to 8-bit unsigned integer
173         gray = (gray * 255).astype(np.uint8)
174
175         # Multi-scale LBP
176         lbp_features = []
177         for radius in [1, 2, 3]:
178             lbp = local_binary_pattern(gray, P=8*radius, R=radius, method='
uniform')
179             lbp_hist = np.histogram(lbp, bins=10)[0]
180             lbp_features.extend(lbp_hist/lbp_hist.sum()) # 正規化
181
182         # 3. Shape Features (HOG)

```

```

183     hog_features = hog(gray,
184                        orientations=9,
185                        pixels_per_cell=(16, 16),
186                        cells_per_block=(2, 2),
187                        visualize=False)
188
189     # 4. Edge Features
190     edges = cv2.Canny(gray, 100, 200)
191     edge_hist = np.histogram(edges, bins=16)[0]
192
193     # Gradient direction histograms
194     grad_x = cv2.Sobel(gray, cv2.CV_64F, 1, 0, ksize=3)
195     grad_y = cv2.Sobel(gray, cv2.CV_64F, 0, 1, ksize=3)
196     grad_mag = np.sqrt(grad_x**2 + grad_y**2)
197     grad_dir = np.arctan2(grad_y, grad_x)
198     grad_hist = np.histogram(grad_dir, bins=18, range=(-np.pi, np.pi))[0]
199
200     # Combine all features and normalize
201     features = np.concatenate([
202         hist_h/hist_h.sum(),           # 16 features
203         hist_s/hist_s.sum(),           # 16 features
204         hist_v/hist_v.sum(),           # 16 features
205         color_stats,                   # 12 features
206         np.array(lbp_features),         # 30 features
207         hog_features/np.linalg.norm(hog_features), # normalized HOG
208         edge_hist/edge_hist.sum(),     # 16 features
209         grad_hist/grad_hist.sum()      # 18 features
210     ])
211
212
213     return features

```

Listing 3: main.py

```

1 import os
2 import argparse
3 import time
4 from pathlib import Path
5
6
7 def prepare_dataset():
8     """Prepare the dataset by splitting it into train, validation, and test
9     sets"""
10    from prepare_dataset import main as prepare_main
11    prepare_main()
12
13 def train_resnet():
14     """Train the ResNet model (supervised learning with DL)"""
15     from models.ResNet import ResNetTrainer
16
17     print("\n" + "="*60)
18     print("Training ResNet Model (Supervised Learning - Deep Learning)")
19     print("="*60)

```



```

20     start_time = time.time()
21     trainer = ResNetTrainer(data_dir='data', batch_size=32, num_epochs=20,
22                             seed=40)
23     trainer.train()
24     accuracy, _ = trainer.evaluate()
25     trainer.plot_training_history()
26     trainer.visualize_tsne_comparison()
27     end_time = time.time()
28     print(f"\nResNet training completed in {end_time - start_time:.2f} seconds
29         ")
30     # print(f"Test accuracy: {accuracy:.4f}")
31     # return accuracy
32
33 def train_svm():
34     """Train the SVM model (supervised learning without DL)"""
35     from models.SVM import SVMClassifier
36
37     print("\n" + "="*60)
38     print("Training SVM Model (Supervised Learning- Traditional ML)")
39     print("="*60)
40
41     start_time = time.time()
42     classifier = SVMClassifier(data_dir='data')
43     accuracy, _ = classifier.train()
44     end_time = time.time()
45
46     print(f"\nSVM training completed in {end_time - start_time:.2f} seconds")
47     print(f"Test accuracy: {accuracy:.4f}")
48
49     return accuracy
50
51 def train_kmeans():
52     """Train the K-means model (unsupervised learning)"""
53     from models.K_means import KMeansClusterer
54
55     print("\n" + "="*60)
56     print("Training K-means Model (Unsupervised Learning)")
57     print("="*60)
58
59     start_time = time.time()
60     clusterer = KMeansClusterer(data_dir='data')
61     clusterer.train()
62     end_time = time.time()
63
64     print(f"\nK-means training completed in {end_time - start_time:.2f}
65         seconds")
66
67     return None # No standard accuracy for unsupervised learning
68
69 def main():
70     parser = argparse.ArgumentParser(description='Train and evaluate models
71         for clothing classification')

```

```

70     parser.add_argument('--prepare', action='store_true', help='Prepare the
dataset')
71     parser.add_argument('--resnet', action='store_true', help='Train ResNet
model')
72     parser.add_argument('--svm', action='store_true', help='Train SVM model')
73     parser.add_argument('--kmeans', action='store_true', help='Train K-means
model')
74     parser.add_argument('--all', action='store_true', help='Run all models')
75
76     args = parser.parse_args()
77
78     # Create results directory
79     os.makedirs('results', exist_ok=True)
80
81     # If no specific arguments are provided, run all
82     if not (args.prepare or args.resnet or args.svm or args.kmeans or args.all
):
83         args.all = True
84
85     # Prepare dataset if requested
86     if args.prepare or args.all:
87         print("\nPreparing dataset...")
88         prepare_dataset()
89
90     # Check if dataset exists
91     if not Path('data').exists():
92         print("\nDataset not found. Please run with --prepare first.")
93         return
94
95     # Ensure models directory exists
96     if not Path('models').exists() and (args.resnet or args.svm or args.kmeans
or args.all):
97         print("\nModels directory not found. Please make sure the 'models'
folder exists with model files.")
98         return
99
100     results = {}
101
102     # Train models as requested
103     if args.resnet or args.all:
104         results['resnet'] = train_resnet()
105
106     if args.svm or args.all:
107         results['svm'] = train_svm()
108
109     if args.kmeans or args.all:
110         results['kmeans'] = train_kmeans()
111
112     # Print summary of results
113     print("\n" + "="*50)
114     print("Summary of Results")
115     print("="*50)
116
117     for model, accuracy in results.items():

```

```

118         if accuracy is not None:
119             print(f"{model.upper()} Test Accuracy: {accuracy:.4f}")
120
121         print("\nAll models have been trained and evaluated successfully!")
122
123 if __name__ == "__main__":
124     main()

```

Listing 4: models/ResNet.py

```

1 import os
2 import torch
3 import torch.nn as nn
4 import torch.optim as optim
5 from torch.utils.data import DataLoader
6 import numpy as np
7 import random
8 from torchvision import models
9 import matplotlib.pyplot as plt
10 from sklearn.metrics import confusion_matrix, classification_report
11 import seaborn as sns
12 from tqdm import tqdm
13 from sklearn.manifold import TSNE
14 from dataset import ResnetDataset, train_tfm, valid_tfm, test_tfm
15
16 # Set random seeds for reproducibility
17 def set_seed(seed=40):
18     """Set all random seeds for reproducibility"""
19     random.seed(seed)
20     np.random.seed(seed)
21     torch.manual_seed(seed)
22     torch.cuda.manual_seed(seed)
23     torch.cuda.manual_seed_all(seed) # For multi-GPU setups
24     torch.backends.cudnn.deterministic = True
25     torch.backends.cudnn.benchmark = False
26     os.environ['PYTHONHASHSEED'] = str(seed)
27
28
29 class ResNetTrainer:
30     def __init__(self, data_dir='data', batch_size=32, num_epochs=40,
31                 learning_rate=0.001, seed=40):
32         # Set seed for reproducibility
33         set_seed(seed)
34
35         self.data_dir = data_dir
36         self.batch_size = batch_size
37         self.num_epochs = num_epochs
38         self.learning_rate = learning_rate
39         self.device = torch.device("cuda:0" if torch.cuda.is_available() else
40                                   "cpu")
41
42         # Load datasets with custom class that handles corrupted images
43         self.datasets = {
44             'train': ResnetDataset(os.path.join(data_dir, 'train'), transform=

```

```

train_tfm),
43     'test': ResnetDataset(os.path.join(data_dir, 'test'), transform=
test_tfm)
44 }
45
46     # Create dataloaders
47     self.dataloaders = {
48         'train': DataLoader(self.datasets['train'], batch_size=batch_size,
shuffle=True, num_workers=6),
49         'test': DataLoader(self.datasets['test'], batch_size=batch_size,
shuffle=False, num_workers=2)
50     }
51
52     self.dataset_sizes = {x: len(self.datasets[x]) for x in ['train', '
test']}
53     self.class_names = self.datasets['train'].classes
54
55     print(f"Classes: {self.class_names}")
56     print(f"Dataset sizes: {self.dataset_sizes}")
57
58     # Initialize model
59     self.model = models.resnet18(weights='IMAGENET1K_V1')
60     self.model.fc = nn.Sequential(
61         nn.Dropout(0.5),
62         nn.Linear(self.model.fc.in_features, len(self.class_names))
63     )
64     self.model = self.model.to(self.device)
65
66     # Loss function and optimizer
67     self.criterion = nn.CrossEntropyLoss()
68     self.optimizer = optim.Adam(self.model.parameters(), lr=learning_rate)
69     self.scheduler = optim.lr_scheduler.StepLR(self.optimizer, step_size
=10, gamma=0.1)
70
71     # For tracking metrics
72     self.train_losses = []
73     self.train_accs = []
74
75     def train(self):
76         best_test_acc = 0.0
77
78         for epoch in range(self.num_epochs):
79             print(f'Epoch {epoch+1}/{self.num_epochs}')
80             print('-' * 10)
81
82             # Each epoch has a training and testing phase
83             for phase in ['train', 'test']:
84                 if phase == 'train':
85                     self.model.train() # Set model to training mode
86                 else:
87                     self.model.eval() # Set model to evaluate mode
88
89                 running_loss = 0.0
90                 running_corrects = 0

```

```

91         # Wrap dataloader with tqdm for progress bar
92         pbar = tqdm(self.dataloaders[phase], desc=f'{phase} Epoch {
93 epoch+1}/{self.num_epochs}')
94
95         # Iterate over data
96         for inputs, labels in pbar:
97             inputs = inputs.to(self.device)
98             labels = labels.to(self.device)
99
100        # Zero the parameter gradients
101        self.optimizer.zero_grad()
102
103        # Forward pass
104        with torch.set_grad_enabled(phase == 'train'):
105            outputs = self.model(inputs)
106            _, preds = torch.max(outputs, 1)
107            loss = self.criterion(outputs, labels)
108
109        # Backward + optimize only if in training phase
110        if phase == 'train':
111            loss.backward()
112            self.optimizer.step()
113
114        # Statistics
115        running_loss += loss.item() * inputs.size(0)
116        running_corrects += torch.sum(preds == labels.data)
117
118        # Update progress bar
119        batch_loss = loss.item()
120        batch_acc = torch.sum(preds == labels.data).double() /
121 inputs.size(0)
122        pbar.set_postfix({'loss': f'{batch_loss:.4f}', 'acc': f'{
123 batch_acc:.4f}'})
124
125        epoch_loss = running_loss / self.dataset_sizes[phase]
126        epoch_acc = running_corrects.double() / self.dataset_sizes[
127 phase]
128
129        if phase == 'train':
130            self.train_losses.append(epoch_loss)
131            self.train_accs.append(epoch_acc.item())
132            self.scheduler.step() # Step the scheduler
133        else:
134            # Save the best model
135            if epoch_acc > best_test_acc:
136                best_test_acc = epoch_acc
137                torch.save(self.model.state_dict(), 'results/
138 resnet_best_model.pth')
139
140        print(f'{phase} Loss: {epoch_loss:.4f} Acc: {epoch_acc:.4f}')
141
142    print()

```

```

140     print(f'Best test Acc: {best_test_acc:.4f}')
141
142     # Load best model weights
143     self.model.load_state_dict(torch.load('results/resnet_best_model.pth')
144 )
145     return self.model
146
147 def evaluate(self):
148     # Set model to evaluate mode
149     self.model.eval()
150
151     # Initialize lists to store predictions and ground truth
152     all_preds = []
153     all_labels = []
154
155     # No gradient calculation needed
156     with torch.no_grad():
157         # Wrap dataloader with tqdm for progress bar
158         pbar = tqdm(self.dataloaders['test'], desc='Evaluating')
159
160         for inputs, labels in pbar:
161             inputs = inputs.to(self.device)
162             labels = labels.to(self.device)
163
164             # Forward pass
165             outputs = self.model(inputs)
166             _, preds = torch.max(outputs, 1)
167
168             # Store predictions and labels
169             all_preds.extend(preds.cpu().numpy())
170             all_labels.extend(labels.cpu().numpy())
171
172             # Update progress bar
173             batch_acc = torch.sum(preds == labels.data).double() / inputs.
174             size(0)
175             pbar.set_postfix({'batch_acc': f'{batch_acc:.4f}'})
176
177     # Calculate accuracy
178     accuracy = np.mean(np.array(all_preds) == np.array(all_labels))
179     print(f'Test Accuracy: {accuracy:.4f}')
180
181     # Generate confusion matrix
182     cm = confusion_matrix(all_labels, all_preds)
183
184     # Plot confusion matrix
185     plt.figure(figsize=(10, 8))
186     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
187                 xticklabels=self.class_names,
188                 yticklabels=self.class_names)
189     plt.xlabel('Predicted')
190     plt.ylabel('True')
191     plt.title('Confusion Matrix')
192     plt.savefig('plots/resnet_confusion_matrix.png')

```

```

192     # Print classification report
193     print("\nClassification Report:")
194     print(classification_report(all_labels, all_preds, target_names=self.
class_names))
195
196     return accuracy, cm
197
198     def plot_training_history(self):
199         plt.figure(figsize=(12, 5))
200
201         # Plot loss
202         plt.subplot(1, 2, 1)
203         plt.plot(self.train_losses, label='Training Loss')
204         plt.xlabel('Epoch')
205         plt.ylabel('Loss')
206         plt.title('Training Loss')
207         plt.legend()
208
209         # Plot accuracy
210         plt.subplot(1, 2, 2)
211         plt.plot(self.train_accs, label='Training Accuracy')
212         plt.xlabel('Epoch')
213         plt.ylabel('Accuracy')
214         plt.title('Training Accuracy')
215         plt.legend()
216
217         plt.tight_layout()
218         plt.savefig('plots/resnet_training_history.png')
219         plt.show()
220
221
222     def visualize_tsne_comparison(self, layers=['input', 'layer1', 'layer2', '
layer3', 'layer4', 'fc'], perplexity=30, n_iter=1000):
223         """
224         Visualize and compare feature representations from multiple layers
225         using t-SNE
226
227         Args:
228             layers (list): List of layers to extract features from
229             perplexity (int): Perplexity parameter for t-SNE
230             n_iter (int): Number of iterations for t-SNE
231         """
232         self.model.eval()
233
234         layer_indices = {
235             'layer1': 4,
236             'layer2': 5,
237             'layer3': 6,
238             'layer4': 7
239         }
240
241         fig, axes = plt.subplots(2, 3, figsize=(18, 12))
242         axes = axes.flatten()

```

```

243     for i, layer in enumerate(layers):
244         if layer == 'input':
245             feature_extractor = torch.nn.Identity()
246         elif layer in layer_indices:
247             feature_extractor = torch.nn.Sequential(*list(self.model.
children()[layer_indices[layer]]))
248         elif layer == 'fc':
249             feature_extractor = torch.nn.Sequential(*list(self.model.
children()[::-1]))
250         else:
251             raise ValueError(f"Layer {layer} not supported for feature
extraction")
252
253         features = []
254         labels = []
255
256         with torch.no_grad():
257             for inputs, targets in tqdm(self.dataloaders['test'], desc=f'
Extracting features from {layer}'):
258                 inputs = inputs.to(self.device)
259                 feat = feature_extractor(inputs)
260                 feat = feat.view(feat.size(0), -1)
261                 features.append(feat.cpu().numpy())
262                 labels.append(targets.numpy())
263
264         features = np.concatenate(features, axis=0)
265         labels = np.concatenate(labels, axis=0)
266
267         print(f"Applying t-SNE on {features.shape[0]} samples with {
features.shape[1]} features from {layer}...")
268         tsne = TSNE(n_components=2, perplexity=perplexity, n_iter=n_iter,
random_state=40)
269         features_tsne = tsne.fit_transform(features)
270
271         ax = axes[i]
272         for _, label in enumerate(np.unique(labels)):
273             idx = labels == label
274             ax.scatter(features_tsne[idx, 0], features_tsne[idx, 1], label
=self.class_names[label], alpha=0.7, s=50)
275
276         ax.set_title(f'{layer}', fontsize=18)
277         ax.legend()
278
279     plt.tight_layout()
280     plt.savefig('plots/resnet_tsne_comparison.png', dpi=300)
281     plt.show()
282
283
284 if __name__ == "__main__":
285     trainer = ResNetTrainer(data_dir='data', batch_size=32, num_epochs=40,
seed=40)
286     trainer.train()
287     trainer.evaluate()
288     trainer.plot_training_history()

```


Listing 5: models/SVM.py

```

1 import os
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score, confusion_matrix,
    classification_report
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.pipeline import Pipeline
8 from sklearn.model_selection import GridSearchCV
9 import seaborn as sns
10 from PIL import Image
11 import joblib
12 from pathlib import Path
13 from dataset import SVMDataset
14
15 class SVMClassifier:
16     def __init__(self, data_dir='dataset'):
17         self.data_dir = Path(data_dir)
18         self.class_names = None
19         self.model = None
20         self.scaler = None
21         self.dataset = SVMDataset(data_dir, 0)
22
23     def train(self):
24         """Train the SVM model with hyperparameter tuning"""
25         # Load data
26         X_train, y_train, X_test, y_test = self.dataset.load_data()
27
28         print(f"Training data shape: {X_train.shape}")
29         print(f"Test data shape: {X_test.shape}")
30
31         # Create a pipeline with scaling and SVM
32         pipeline = Pipeline([
33             ('scaler', StandardScaler()),
34             ('svm', SVC(probability=True))
35         ])
36
37         # Define parameter grid for grid search
38         param_grid = {
39             'svm__C': [0.1, 1, 10, 100, 1000],
40             'svm__gamma': ['scale', 'auto', 0.001, 0.01, 0.1],
41             'svm__kernel': ['rbf', 'poly'],
42             'svm__degree': [2, 3] # for poly kernel
43         }
44
45         # Perform grid search with cross-validation
46         print("Performing grid search for hyperparameter tuning...")
47         grid_search = GridSearchCV(pipeline, param_grid, cv=3, n_jobs=-1,
48 verbose=2)
49         grid_search.fit(X_train, y_train)

```

```

49
50     # Get the best model
51     self.model = grid_search.best_estimator_
52     print(f"Best parameters: {grid_search.best_params_}")
53
54     # Save the model
55     joblib.dump(self.model, 'results/svm_model.pkl')
56
57     # Evaluate on test set
58     test_preds = self.model.predict(X_test)
59     test_accuracy = accuracy_score(y_test, test_preds)
60     print(f"Test accuracy: {test_accuracy:.4f}")
61
62     # Generate confusion matrix
63     cm = confusion_matrix(y_test, test_preds)
64
65     # Check if class_names is populated
66     if not self.class_names:
67         self.class_names = [str(i) for i in range(len(set(y_test)))]
68
69     # Plot confusion matrix
70     plt.figure(figsize=(10, 8))
71     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
72                 xticklabels=self.class_names,
73                 yticklabels=self.class_names)
74     plt.xlabel('Predicted')
75     plt.ylabel('True')
76     plt.title('SVM Confusion Matrix')
77     plt.savefig('plots/svm_confusion_matrix.png')
78
79     # Print classification report
80     print("\nClassification Report:")
81     print(classification_report(y_test, test_preds, target_names=self.
82                                class_names))
83
84     return test_accuracy, cm
85
86 if __name__ == "__main__":
87     classifier = SVMClassifier(data_dir='dataset')
88     classifier.train()

```

Listing 6: models/K_means.py

```

1 import os
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.cluster import KMeans
5 from sklearn.decomposition import PCA
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.metrics import silhouette_score, adjusted_rand_score,
8     calinski_harabasz_score, davies_bouldin_score
9 import seaborn as sns
10 from PIL import Image
11 import joblib

```

```

11 from pathlib import Path
12 from collections import Counter
13 import pandas as pd
14 from skimage.feature import local_binary_pattern
15 from scipy.ndimage import sobel
16 from skimage.feature import hog
17 import cv2
18 from sklearn.feature_selection import SelectKBest
19 from sklearn.feature_selection import mutual_info_classif
20 import umap
21
22 class KMeansClusterer:
23     def __init__(self, data_dir='dataset', n_clusters=10):
24         self.data_dir = Path(data_dir)
25         self.n_clusters = n_clusters
26         self.model = None
27         self.pca = None
28         self.scaler = None
29         self.class_names = None
30
31     def extract_features(self, img_path):
32         """
33         Improved feature extraction with focus on clothing-specific
34         characteristics
35         """
36         try:
37             # Load and resize image
38             img = Image.open(img_path).convert('RGB')
39             img = img.resize((128, 128)) # Increase image size to preserve
40             more details
41             img_array = np.array(img)
42
43             # 1. Enhanced Color Features
44             # Use HSV color space
45             img_hsv = cv2.cvtColor(img_array, cv2.COLOR_RGB2HSV)
46
47             # Calculate HSV histograms
48             hist_h = np.histogram(img_hsv[:, :, 0], bins=16)[0]
49             hist_s = np.histogram(img_hsv[:, :, 1], bins=16)[0]
50             hist_v = np.histogram(img_hsv[:, :, 2], bins=16)[0]
51
52             # Color statistics
53             color_stats = np.concatenate([
54                 np.mean(img_hsv, axis=(0,1)),
55                 np.std(img_hsv, axis=(0,1))
56             ])
57
58             # 2. Improved texture features
59             gray = cv2.cvtColor(img_array, cv2.COLOR_RGB2GRAY)
60
61             # Use different scales of LBP
62             lbp_features = []
63             for radius in [1, 2, 3]:
64                 lbp = local_binary_pattern(gray, P=8*radius, R=radius, method=

```

```

'uniform')
    lbp_hist = np.histogram(lbp, bins=10)[0]
    lbp_features.extend(lbp_hist)

    # 3. Improved shape features
    # Use denser HOG features
    hog_features = hog(gray,
                       orientations=9,
                       pixels_per_cell=(16, 16),
                       cells_per_block=(2, 2),
                       visualize=False)

    # 4. Edge features
    edges = cv2.Canny(gray, 100, 200)
    edge_hist = np.histogram(edges, bins=16)[0]

    # Combine all features
    features = np.concatenate([
        hist_h/hist_h.sum(), hist_s/hist_s.sum(), hist_v/hist_v.sum(),
        # Normalized color histograms
        color_stats,
        # HSV statistics
        np.array(lbp_features)/sum(lbp_features),
        # Normalized LBP features
        hog_features/np.linalg.norm(hog_features),
        # Normalized HOG features
        edge_hist/edge_hist.sum(),
        # Normalized edge features
    ])

    return features

except Exception as e:
    print(f"Error processing {img_path}: {e}")
    return None

def load_data(self):
    """Load and prepare data for K-means clustering"""
    X = [] # Features
    y = [] # True labels (for evaluation only)
    img_paths = [] # Store paths for visualization

    # Get class names from train directory
    self.class_names = [d for d in os.listdir(self.data_dir / 'train')
                        if os.path.isdir(self.data_dir / 'train' / d)]

    # Create class to index mapping
    class_to_idx = {cls_name: i for i, cls_name in enumerate(self.
class_names)}

    # Process all data (train, test combined for unsupervised learning)
    for split in ['train', 'test']:
        for class_name in self.class_names:
            class_dir = self.data_dir / split / class_name

```

```

110         for img_file in os.listdir(class_dir):
111             if img_file.lower().endswith(('.png', '.jpg', '.jpeg')):
112                 img_path = class_dir / img_file
113                 features = self.extract_features(img_path)
114                 if features is not None:
115                     X.append(features)
116                     y.append(class_to_idx[class_name])
117                     img_paths.append(str(img_path))
118
119         return np.array(X), np.array(y), img_paths
120
121     def train(self):
122         """Improved training process"""
123         # Load data
124         X, y_true, img_paths = self.load_data()
125
126         # Feature selection
127         selector = SelectKBest(score_func=mutual_info_classif, k=100)
128         X_selected = selector.fit_transform(X, y_true)
129
130         # Standardization
131         self.scaler = StandardScaler()
132         X_scaled = self.scaler.fit_transform(X_selected)
133
134         # Use UMAP for dimensionality reduction
135         self.reducer = umap.UMAP(n_components=30)
136         X_reduced = self.reducer.fit_transform(X_scaled)
137
138         # Try different numbers of clusters
139         silhouette_scores = []
140         ari_scores = []
141         ch_scores = []
142         db_scores = []
143         k_range = list(range(2, 21)) # Integers from 2 to 20
144
145         # Create plot
146         plt.figure(figsize=(15, 10))
147
148         # Cluster for each K value and calculate metrics
149         for k in k_range:
150             kmeans = KMeans(n_clusters=k, random_state=42, n_init=10)
151             labels = kmeans.fit_predict(X_reduced)
152
153             # Calculate metrics
154             silhouette = silhouette_score(X_reduced, labels)
155             ari = adjusted_rand_score(y_true, labels)
156             ch = calinski_harabasz_score(X_reduced, labels)
157             db = davies_bouldin_score(X_reduced, labels)
158
159             silhouette_scores.append(silhouette)
160             ari_scores.append(ari)
161             ch_scores.append(ch)
162             db_scores.append(db)
163

```

```

164         print(f"K={k}, Silhouette Score: {silhouette:.4f}, ARI: {ari:.4f},
165         CH: {ch:.4f}, DB: {db:.4f}")
166
167
168     # Silhouette Score
169     plt.subplot(2, 2, 1)
170     plt.plot(k_range, silhouette_scores, 'bo-')
171     plt.xlabel('Number of Clusters (K)')
172     plt.ylabel('Silhouette Score')
173     plt.title('Silhouette Score vs K\n(Higher is better)')
174     plt.grid(True)
175     plt.xticks(k_range, k_range) # Set x-axis ticks to integers
176
177     # ARI Score
178     plt.subplot(2, 2, 2)
179     plt.plot(k_range, ari_scores, 'ro-')
180     plt.xlabel('Number of Clusters (K)')
181     plt.ylabel('Adjusted Rand Index')
182     plt.title('Adjusted Rand Index vs K\n(Higher is better)')
183     plt.grid(True)
184     plt.xticks(k_range, k_range)
185
186     # Calinski-Harabasz Score
187     plt.subplot(2, 2, 3)
188     plt.plot(k_range, ch_scores, 'go-')
189     plt.xlabel('Number of Clusters (K)')
190     plt.ylabel('Calinski-Harabasz Score')
191     plt.title('Calinski-Harabasz Score vs K\n(Higher is better)')
192     plt.grid(True)
193     plt.xticks(k_range, k_range)
194
195     # Davies-Bouldin Score
196     plt.subplot(2, 2, 4)
197     plt.plot(k_range, db_scores, 'mo-')
198     plt.xlabel('Number of Clusters (K)')
199     plt.ylabel('Davies-Bouldin Score')
200     plt.title('Davies-Bouldin Score vs K\n(Lower is better)')
201     plt.grid(True)
202     plt.xticks(k_range, k_range)
203
204     plt.tight_layout()
205     plt.savefig('plots/kmeans_metrics_comparison.png')
206
207     # Find best K value for each metric
208     best_k_silhouette = k_range[np.argmax(silhouette_scores)]
209     best_k_ari = k_range[np.argmax(ari_scores)]
210     best_k_ch = k_range[np.argmax(ch_scores)]
211     best_k_db = k_range[np.argmin(db_scores)] # Note: Lower DB score is
better
212
213     print("\nBest K values by different metrics:")
214     print(f"Silhouette Score: K={best_k_silhouette}")
215     print(f"Adjusted Rand Index: K={best_k_ari}")

```

```

216     print(f"Calinski-Harabasz Score: K={best_k_ch}")
217     print(f"Davies-Bouldin Score: K={best_k_db}")
218
219     # Use best K value (using best K from ARI score here)
220     self.n_clusters = best_k_ari
221     self.model = KMeans(n_clusters=self.n_clusters, random_state=42,
n_init=10)
222     cluster_labels = self.model.fit_predict(X_reduced)
223
224     # Evaluate and visualize
225     self._visualize_clusters(X_reduced, cluster_labels, y_true)
226     self._analyze_clusters(cluster_labels, y_true)
227     metrics = self.evaluate_clustering(X_reduced, cluster_labels, y_true)
228
229     return self.model
230
231 def _visualize_clusters(self, X_pca, cluster_labels, y_true):
232     """Visualize the clusters in 2D using PCA"""
233     # Further reduce to 2D for visualization
234     pca_viz = PCA(n_components=2)
235     X_pca_2d = pca_viz.fit_transform(X_pca)
236
237     # Plot clusters
238     plt.figure(figsize=(12, 10))
239
240     # Plot by cluster assignment
241     plt.subplot(1, 2, 1)
242     scatter = plt.scatter(X_pca_2d[:, 0], X_pca_2d[:, 1], c=cluster_labels
, cmap='viridis', alpha=0.6)
243     plt.colorbar(scatter)
244     plt.title('K-means Clustering Results')
245     plt.xlabel('PCA Component 1')
246     plt.ylabel('PCA Component 2')
247
248     # Plot by true labels
249     plt.subplot(1, 2, 2)
250     scatter = plt.scatter(X_pca_2d[:, 0], X_pca_2d[:, 1], c=y_true, cmap='
tab10', alpha=0.6)
251     plt.colorbar(scatter)
252     plt.title('True Class Labels')
253     plt.xlabel('PCA Component 1')
254     plt.ylabel('PCA Component 2')
255
256     plt.tight_layout()
257     plt.savefig('plots/kmeans_clusters_visualization.png')
258
259 def _analyze_clusters(self, cluster_labels, y_true):
260     """Analyze the composition of each cluster"""
261     # Create a mapping from cluster labels to true labels
262     cluster_to_label = {}
263
264     for cluster_id in range(max(cluster_labels) + 1): # Use actual
cluster labels instead of self.n_clusters
265         # Get indices of samples in this cluster

```

```

266         indices = np.where(cluster_labels == cluster_id)[0]
267
268         # Skip empty clusters
269         if len(indices) == 0:
270             print(f"\nCluster {cluster_id} is empty")
271             continue
272
273         # Get true labels of these samples
274         true_labels = y_true[indices]
275
276         # Count occurrences of each true label
277         label_counts = Counter(true_labels)
278
279         # Find the most common true label in this cluster
280         most_common_label = label_counts.most_common(1)[0][0]
281         cluster_to_label[cluster_id] = most_common_label
282
283         # Print cluster composition
284         print(f"\nCluster {cluster_id} composition:")
285         for label, count in label_counts.most_common():
286             percentage = count / len(indices) * 100
287             print(f"    {self.class_names[label]}: {count} samples ({
percentage:.2f}%)")
288
289         # Create confusion matrix-like visualization
290         n_clusters = max(cluster_labels) + 1
291         confusion = np.zeros((n_clusters, len(self.class_names)))
292
293         for i in range(len(cluster_labels)):
294             cluster = cluster_labels[i]
295             true_label = y_true[i]
296             confusion[cluster, true_label] += 1
297
298         # Normalize by cluster size
299         for i in range(n_clusters):
300             if np.sum(confusion[i, :]) > 0:
301                 confusion[i, :] = confusion[i, :] / np.sum(confusion[i, :])
302
303         # Plot heatmap
304         plt.figure(figsize=(12, 10))
305         sns.heatmap(confusion, annot=True, fmt='.2f', cmap='Blues',
306                     xticklabels=self.class_names,
307                     yticklabels=[f'Cluster {i}' for i in range(n_clusters)])
308         plt.xlabel('True Class')
309         plt.ylabel('Cluster')
310         plt.title('Cluster Composition')
311         plt.tight_layout()
312         plt.savefig('plots/kmeans_cluster_composition.png')
313
314     def evaluate_clustering(self, X_pca, cluster_labels, y_true):
315         """Evaluate the clustering performance"""
316         # 1. Calculate ARI
317         ari = adjusted_rand_score(y_true, cluster_labels)
318         print(f"Adjusted Rand Index: {ari:.4f}")

```



```

319
320     # 2. Calculate Silhouette Score
321     silhouette_avg = silhouette_score(X_pca, cluster_labels)
322     print(f"Silhouette Score: {silhouette_avg:.4f}")
323
324     # 3. Analyze the purity of each cluster
325     purities = []
326     for cluster_id in range(self.n_clusters):
327         indices = np.where(cluster_labels == cluster_id)[0]
328         if len(indices) > 0:
329             true_labels = y_true[indices]
330             most_common = Counter(true_labels).most_common(1)[0]
331             purity = most_common[1] / len(indices)
332             purities.append(purity)
333
334     avg_purity = np.mean(purities)
335     print(f"Average Cluster Purity: {avg_purity:.4f}")
336
337     return {
338         'ari': ari,
339         'silhouette': silhouette_avg,
340         'purity': avg_purity
341     }
342
343 if __name__ == "__main__":
344     clusterer = KMeansClusterer(data_dir='dataset', n_clusters=10)
345     clusterer.train()

```