# feature-engineering-2-3

August 13, 2023

```
[ ]:  # Q1. What is Min-Max scaling, and how is it used in data preprocessing?

      # Provide an example to illustrate its application.
```

Min-Max scaling, also known as feature scaling, is a data preprocessing technique used in machine learning and data analysis to transform numerical features so that they are within a specific range, typically between 0 and 1.

This method is particularly useful when working with algorithms that are sensitive to the scale of the input features, such as distance-based algorithms (e.g., k-means clustering) and gradient-based optimization algorithms (e.g., neural networks).

The formula to perform Min-Max scaling on a feature x is as follows:

x_scaled = (x - min(x)) / (max(x) - min(x))

Where:

x is the original value of the feature. min(x) is the minimum value of the feature in the dataset. max(x) is the maximum value of the feature in the dataset. x_scaled is the scaled value of the feature within the range [0, 1].

Here's an example to illustrate Min-Max scaling:

Let's say you have a dataset of exam scores with scores ranging from 60 to 95. You want to scale these scores using Min-Max scaling.

Original scores:

Student 1: 60 Student 2: 75 Student 3: 85 Student 4: 95 First, calculate the minimum and maximum values:

min(score) = 60 max(score) = 95

Now, apply the Min-Max scaling formula to each student's score:

Student 1 scaled score = (60 - 60) / (95 - 60) = 0.0 Student 2 scaled score = (75 - 60) / (95 - 60) = 0.375 Student 3 scaled score = (85 - 60) / (95 - 60) = 0.625 Student 4 scaled score = (95 - 60) / (95 - 60) = 1.0

After Min-Max scaling, the scores are transformed to the range [0, 1] while preserving their relative differences.

1

```
[ ]: # Q2. What is the Unit Vector technique in feature scaling, and how does it␣
     ↪differ from Min-Max scaling?

     # Provide an example to illustrate its application.
```

The Unit Vector technique, also known as Normalization, is another data preprocessing technique used to scale numerical features.

Mathematically, the normalization of a feature vector x is performed as follows:

x_normalized = x / ||x||

Where:

x is the original feature vector. ||x|| represents the Euclidean norm of the vector, calculated as the square root of the sum of squared elements.

Unlike Min-Max scaling that scales features within a specific range (e.g., [0, 1]), the Unit Vector technique scales features to have a unit norm, meaning that the vector representing the feature values will have a length of 1 while maintaining the direction of the original vector.

Here's an example to illustrate the Unit Vector technique:

Consider a dataset with two features representing height in centimeters and weight in kilograms:

Original data:

Person 1: Height = 160 cm, Weight = 55 kg Person 2: Height = 175 cm, Weight = 70 kg Person 3: Height = 180 cm, Weight = 65 kg Let's calculate the Euclidean norms for each person's feature vector:

Person 1: $||x1|| = sqrt(160^2 + 55^2)$   167.91 Person 2: $||x2|| = sqrt(175^2 + 70^2)$   188.46 Person 3: $||x3|| = sqrt(180^2 + 65^2)$   193.29 Now, normalize each feature vector:

Person 1 normalized: (160 / 167.91, 55 / 167.91)   (0.953, 0.302) Person 2 normalized: (175 / 188.46, 70 / 188.46)   (0.928, 0.372) Person 3 normalized: (180 / 193.29, 65 / 193.29)   (0.932, 0.336) In this case, the normalization process ensures that the length of each feature vector is 1, while the relative relationships between height and weight for each person are maintained.

In summary, the key difference between Min-Max scaling and the Unit Vector technique (Normalization) lies in their objectives. Min-Max scaling scales features to a specific range, while the Unit Vector technique normalizes features to have a unit norm, maintaining their direction and ensuring equal contribution to distance calculations. The choice between these techniques depends on the requirements of the specific machine learning algorithm you're using and the nature of your data.

```
[ ]: # Q3. What is PCA (Principle Component Analysis), and how is it used in␣
     ↪dimensionality reduction?

     # Provide an  example to illustrate its application
```

Principal Component Analysis (PCA) is a dimensionality reduction technique used in the field of machine learning and statistics to transform high-dimensional data into a lower-dimensional representation while retaining as much of the original data's variability as possible. PCA achieves this

by identifying the principal components, which are orthogonal (uncorrelated) linear combinations of the original features that capture the most significant variations in the data.

he primary steps of PCA are as follows:

Standardize the Data: Before performing PCA, it's common to standardize the data by subtracting the mean and dividing by the standard deviation of each feature. This ensures that each feature contributes equally to the analysis.

Compute Covariance Matrix: The covariance matrix is computed based on the standardized data. It shows how features vary together.

Compute Eigenvectors and Eigenvalues: Eigenvectors and eigenvalues are calculated from the covariance matrix. Eigenvectors represent the directions of maximum variance, and eigenvalues indicate the magnitude of variance in those directions.

Sort Eigenvectors: Eigenvectors are sorted based on their corresponding eigenvalues in decreasing order. The eigenvector with the highest eigenvalue represents the first principal component, the second highest represents the second principal component, and so on.

Select Principal Components: You can choose a subset of the top eigenvectors (principal components) that capture a desired amount of variance. This determines the dimensionality of the reduced space.

Transform Data: The original data is transformed into the new reduced-dimensional space by projecting it onto the selected principal components.

Here's an example to illustrate PCA's application:

Suppose you have a dataset of 2D points in space:

Original data:

Point 1: (2, 3) Point 2: (5, 7) Point 3: (8, 10) Point 4: (11, 14) Standardize the Data: If necessary, you'd standardize the data by subtracting the mean and dividing by the standard deviation of each dimension.

Compute Covariance Matrix: Compute the covariance matrix based on the standardized data.

Compute Eigenvectors and Eigenvalues: Calculate the eigenvectors and eigenvalues of the covariance matrix.

Sort Eigenvectors: Assume the eigenvectors are sorted as [eig1, eig2], where eig1 has the higher eigenvalue.

Select Principal Components: Let's say you decide to keep only the first principal component (eig1).

Transform Data: Project the original data onto the first principal component.

After the transformation, your data would now be represented along a single dimension, capturing the most significant variance. This lower-dimensional representation is the result of PCA's dimensionality reduction process.

In practice, PCA can be particularly useful for visualization, noise reduction, and improving the efficiency of machine learning algorithms by reducing the number of features while retaining the most important information.

```
# Q4. What is the relationship between PCA and Feature Extraction, and how can␣
↪PCA be used for Feature Extraction?


# Provide an example to illustrate this concept
```

PCA (Principal Component Analysis) and Feature Extraction are closely related concepts, with PCA being a technique commonly used for feature extraction. Feature extraction involves transforming the original features of a dataset into a new set of features that capture the most important information, while reducing dimensionality. PCA achieves this by identifying the principal components, which are linear combinations of the original features that capture the most significant variability in the data.

Here's how PCA can be used for feature extraction:

1. Original Data: Consider a dataset with multiple features, which might represent various measurements or attributes of the data instances.

2. Standardize the Data: If necessary, standardize the data to ensure all features have comparable scales.

3. Compute Covariance Matrix: Calculate the covariance matrix based on the standardized data.

4. Compute Eigenvectors and Eigenvalues: Compute the eigenvectors and eigenvalues of the covariance matrix.

5. Sort Eigenvectors: Sort the eigenvectors in decreasing order based on their corresponding eigenvalues.

6. Select Principal Components: Choose a subset of the top eigenvectors (principal components) that capture a desired amount of variance. These selected components will be used as the new features.

7. Transform Data: Transform the original data into the reduced-dimensional space by projecting it onto the selected principal components.

Here's how PCA can be used for feature extraction:

1. Original Data: Consider a dataset with multiple features, which might represent various measurements or attributes of the data instances.

2. Standardize the Data: If necessary, standardize the data to ensure all features have comparable scales.

3. Compute Covariance Matrix: Calculate the covariance matrix based on the standardized data.

4. Compute Eigenvectors and Eigenvalues: Compute the eigenvectors and eigenvalues of the covariance matrix.

5. Sort Eigenvectors: Sort the eigenvectors in decreasing order based on their corresponding eigenvalues.

6. Select Principal Components: Choose a subset of the top eigenvectors (principal components) that capture a desired amount of variance. These selected components will be used as the new features.

7. Transform Data: Transform the original data into the reduced-dimensional space by projecting it onto the selected principal components.

Here's an example to illustrate how PCA is used for feature extraction:

Original Data: Suppose you have a dataset with three features: height, weight, and age, of individuals.

Standardize the Data: Standardize the data to ensure all features have a mean of 0 and a standard deviation of 1.

Compute Covariance Matrix: Calculate the covariance matrix based on the standardized data.

Compute Eigenvectors and Eigenvalues: Calculate the eigenvectors and eigenvalues of the covariance matrix.

Sort Eigenvectors: Assume the sorted eigenvectors are [eig1, eig2, eig3], with eig1 having the highest eigenvalue.

Select Principal Components: Decide to keep only the first two eigenvectors (eig1 and eig2), capturing a significant portion of the variance.

Transform Data: Project the original data onto the first two principal components.

In this case, the original three features have been transformed into a new two-dimensional space represented by the first two principal components. These new features, derived from the eigenvectors, capture the most important information from the original features while reducing dimensionality. The transformed data can now be used for analysis or machine learning tasks.

Feature extraction through PCA can be especially beneficial when working with high-dimensional datasets, as it allows you to represent the data in a more compact form while preserving the key characteristics.

```
# Q5. You are working on a project to build a recommendation system for a food␣
 ↪delivery service.

# The dataset contains features such as price, rating, and delivery time.

#  Explain how you would use Min-Max scaling to preprocess the data.
```

In the context of building a recommendation system for a food delivery service, you can use Min-Max scaling to preprocess the data and ensure that the features are within a common range. This is important because recommendation systems often rely on similarity or distance measures between items, and having features on a consistent scale can improve the accuracy of the recommendations.

Here's how you would use Min-Max scaling to preprocess the features in your dataset:

Understand the Features: Review the features available in your dataset, which in this case are price, rating, and delivery time. Understand the range and distribution of each feature.

Calculate Min and Max: For each feature, calculate the minimum and maximum values present in the dataset. For example, let's say:

Price: Minimum = $5, Maximum = $30 Rating: Minimum = 2.0, Maximum = 5.0 Delivery Time: Minimum = 15 minutes, Maximum = 60 minutes

Apply Min-Max Scaling: For each feature, apply the Min-Max scaling formula to transform the values to the [0, 1] range:

x_scaled = (x - min(x)) / (max(x) - min(x))

Scaled Price = (Price - $5) / ($30 - $5) Scaled Rating = (Rating - 2.0) / (5.0 - 2.0) Scaled Delivery Time = (Delivery Time - 15) / (60 - 15)

Preprocessed Data: After applying Min-Max scaling, your dataset's features will now be in the [0, 1] range:

Scaled Price: 0.0 to 1.0 Scaled Rating: 0.0 to 1.0 Scaled Delivery Time: 0.0 to 1.0

Use Preprocessed Data: The preprocessed data with scaled features can now be used to build your recommendation system. When calculating similarities or distances between items (food options in this case), the consistent scaling ensures that no feature dominates others due to its original scale.

By using Min-Max scaling, you have transformed the features in a way that maintains their relative relationships while bringing them to a common scale. This preprocessing step can contribute to the accuracy and effectiveness of your recommendation system by enabling more meaningful comparisons between different food options based on their features.

```
# Q6. You are working on a project to build a model to predict stock prices.

# The dataset contains many features, such as company financial data and market
    ↪trends.

# Explain how you would use PCA to reduce the dimensionality of the dataset.
```

When working with a dataset that contains numerous features, such as company financial data and market trends, using PCA (Principal Component Analysis) can be a valuable technique to reduce the dimensionality of the data while retaining the most important information. Reducing dimensionality can help in simplifying the model, improving its efficiency, and potentially mitigating the curse of dimensionality. Here's how you could use PCA for this stock price prediction project:

Understand the Dataset: First, gain a thorough understanding of the dataset's features, which might include financial indicators (e.g., revenue, profit, debt), market indicators (e.g., market sentiment, interest rates), and potentially other external factors affecting stock prices.

Data Preprocessing: Standardize the data to ensure that all features have comparable scales. This step is important for PCA as it's sensitive to the scale of the features.

Compute Covariance Matrix: Calculate the covariance matrix based on the standardized data. The covariance matrix provides insight into how features are correlated and how they vary together.

Compute Eigenvectors and Eigenvalues: Calculate the eigenvectors and eigenvalues of the covariance matrix. These eigenvectors represent the directions of maximum variance in the data.

Sort Eigenvectors: Sort the eigenvectors in decreasing order of their corresponding eigenvalues. The eigenvectors with higher eigenvalues capture more variance and are considered the principal components.

Select Principal Components: Decide on the number of principal components you want to retain. This decision can be based on the cumulative explained variance, where you aim to retain a certain percentage of the total variance in the data. This determines how many dimensions your reduced dataset will have.

Transform Data: Project the original standardized data onto the selected principal components. This results in a new dataset with reduced dimensions, as the principal components serve as the new features.

Build Prediction Model: Use the transformed dataset, which has reduced dimensionality, to build your stock price prediction model. You can use various machine learning algorithms for regression, taking advantage of the simplified dataset.

It's important to note that while PCA reduces dimensionality, it comes with the trade-off of losing some interpretability of the individual features. However, the retained principal components are combinations of the original features that capture the most significant variability in the data.

```
[ ]:  # Q7. For a dataset containing the following values: [1, 5, 10, 15, 20],␣
      ↪perform Min-Max scaling to transform the values to a range of -1 to 1
```

To perform Min-Max scaling on a dataset and transform the values to a range of -1 to 1, you can use the following formula:

x_scaled = 2 * (x - min(x)) / (max(x) - min(x)) - 1

Where:

x is the original value. min(x) is the minimum value in the dataset. max(x) is the maximum value in the dataset. x_scaled is the scaled value within the range of -1 to 1.

Let's apply this formula to your dataset: [1, 5, 10, 15, 20]

min(x) = 1 max(x) = 20

Now, let's calculate the scaled values for each element:

For x = 1: x_scaled = 2 * (1 - 1) / (20 - 1) - 1 = -1

For x = 5: x_scaled = 2 * (5 - 1) / (20 - 1) - 1 = -0.5

For x = 10: x_scaled = 2 * (10 - 1) / (20 - 1) - 1 = 0

After applying Min-Max scaling with the specified formula, the dataset [1, 5, 10, 15, 20] is transformed to the range of -1 to 1 as follows:

Scaled dataset: [-1, -0.5, 0, 0.5, 1]

```
[1]:  import numpy as np

      # Original dataset
      dataset = np.array([1, 5, 10, 15, 20])
```

```python
# Min-Max scaling
min_val = dataset.min()
max_val = dataset.max()

scaled_dataset = 2 * (dataset - min_val) / (max_val - min_val) - 1

print("Original dataset:", dataset)
print("Scaled dataset:", scaled_dataset)
```

```
Original dataset: [ 1  5 10 15 20]
Scaled dataset: [-1.         -0.57894737 -0.05263158  0.47368421  1.         ]
```

```
[ ]: # Q8. For a dataset containing the following features: [height, weight, age,␣
     ↪gender, blood pressure], perform Feature Extraction using PCA.

     # How many principal components would you choose to retain, and why?
```

Choosing the number of principal components to retain in PCA depends on the variance explained by those components and the desired level of dimensionality reduction.

Let's break down the steps to determine how many principal components to retain for the given dataset with features: [height, weight, age, gender, blood pressure].

Standardize the Data: Start by standardizing the data so that all features have comparable scales. This step is crucial for PCA.

Compute Covariance Matrix: Calculate the covariance matrix based on the standardized data.

Compute Eigenvectors and Eigenvalues: Calculate the eigenvectors and eigenvalues of the covariance matrix.

Sort Eigenvectors: Sort the eigenvectors in decreasing order based on their corresponding eigenvalues.

Explained Variance: Calculate the explained variance ratio for each principal component. This ratio tells you the proportion of the total variance in the data that is captured by each component.

Cumulative Explained Variance: Calculate the cumulative explained variance by summing up the explained variance ratios. This will help you understand how much variance is retained when considering a certain number of principal components.

Choose Number of Components: Based on the cumulative explained variance, you can choose the number of principal components that retain a desired amount of variance. A common approach is to retain components that collectively capture a high percentage (e.g., 95% or more) of the total variance.

The decision of how many principal components to retain depends on your specific goals and the trade-off between dimensionality reduction and preserving variance. Retaining a higher number of components retains more information but might lead to a higher-dimensional transformed dataset. On the other hand, retaining fewer components simplifies the dataset but might lose some information.

It's worth mentioning that in the context of the given features [height, weight, age, gender, blood pressure], some of these features might be less relevant for dimensionality reduction. For instance, "gender" might not contribute significantly to variance compared to continuous numerical features like "height," "weight," and "blood pressure."

In summary, the choice of how many principal components to retain involves a balance between retaining variance and reducing dimensionality while considering the specific characteristics of your dataset and the goals of your analysis.

```python
import numpy as np
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

# Example dataset with features [height, weight, age, gender, blood pressure]
dataset = np.array([
    [170, 65, 30, 1, 120],
    [165, 55, 25, 0, 110],
    [180, 75, 40, 1, 130],
    [155, 45, 22, 0, 100],
])

# Standardize the data
scaler = StandardScaler()
scaled_data = scaler.fit_transform(dataset)

# Perform PCA with 2 components (you can choose the number of components)
num_components = 4
pca = PCA(n_components=num_components)
pca_result = pca.fit_transform(scaled_data)

# Print the transformed data
print("Transformed data:")
print(pca_result)
```

```
Transformed data:
[[ 1.00406417e+00  6.70269784e-01  1.77664948e-01  8.86256479e-17]
 [-1.23315205e+00 -5.26264780e-01  2.57265124e-01  8.86256479e-17]
 [ 2.97898198e+00 -2.99876773e-01 -1.84530372e-01  8.86256479e-17]
 [-2.74989410e+00  1.55871769e-01 -2.50399700e-01  8.86256479e-17]]
```