

Feature Engineering

- **What is feature engineering?**
 - Creating and transforming features to maximize predictive power.
(it covers topics like generating features, feature extraction, feature normalization, feature selection, etc.)
- **Feature generation:** Generating features for the given machine learning task.
Example: for time-series, features are generated to highlight key aspects of the time-series data.
- **Dimensionality reduction:** Dimensionality reduction reduces the number of features for a machine learning task to simplify analysis, improve model performance, and reduce computation time. E.g., feature selection and feature extraction.
- **What is feature selection?** Process of reducing features to improve efficiency and model performance
(not all of the generated features are relevant or necessary for the given task).
- **Feature extraction:** is the process of transforming raw data into a set of features (variables) that better represent the underlying patterns in the data, making it more suitable for machine learning models.
Example: **Principal Component Analysis (PCA)**
- **Selection vs Extraction**
Feature selection keeps the most relevant features, discarding redundant ones.
Feature extraction creates new features that capture essential information (e.g., PCA, t-SNE).
Both help address the curse of dimensionality by retaining only the most informative aspects of data.
- **Importance of feature engineering and selection:**
 - Boosts model interpretability, computational efficiency, and performance.
- **Question** What are the drawbacks of feature extraction compared to feature selection?
 - Tracability?
 - Computational Cost?
 - What else?

▼ Feature Generation for Time-series

In the example below, we can see the concept of feature generation and selection.

First hand, we have a noisy signal!

Then some features are extracted.

Then we select a few of them and show that the initial time-series is actually reconstructable with those few selected features!

A few toolboxes and libraries for time-series feature generation will be introduced in the future.

```

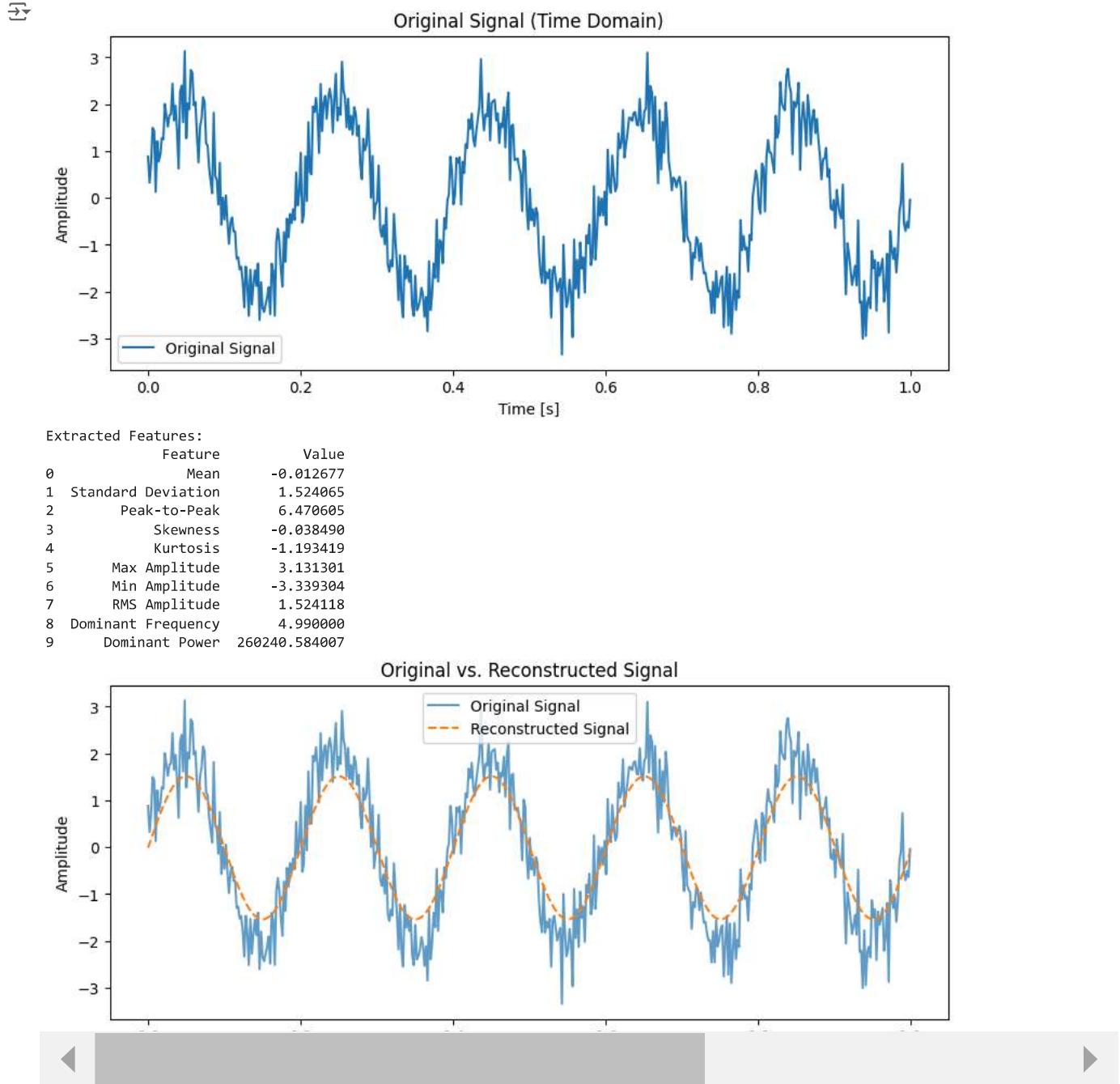
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from scipy.fft import fft, ifft, fftfreq
5 from scipy.stats import skew, kurtosis
6
7 # Generate an artificial time-series signal
8 np.random.seed(0)
9 time = np.linspace(0, 1, 500) # 500 time points over 1 second
10 frequency = 5 # 5 Hz signal frequency
11 amplitude = 2
12 noise = np.random.normal(0, 0.5, time.shape) # Gaussian noise
13 signal = amplitude * np.sin(2 * np.pi * frequency * time) + noise
14
15 # Plot the original signal
16 plt.figure(figsize=(10, 4))
17 plt.plot(time, signal, label='Original Signal')
18 plt.xlabel("Time [s]")
19 plt.ylabel("Amplitude")
20 plt.title("Original Signal (Time Domain)")
21 plt.legend()
22 plt.show()
23
24 # Feature Generation
25 # 1. Time Domain Features
26 peak_to_peak = np.ptp(signal) # Difference between max and min
27
28 # 2. Statistical Features
29 mean_val = np.mean(signal)
30 std_dev = np.std(signal)
31 signal_skewness = skew(signal)
32 signal_kurtosis = kurtosis(signal)
33
34 # 3. Amplitude Domain Features
35 max_amplitude = np.max(signal)
36 min_amplitude = np.min(signal)
37 rms_amplitude = np.sqrt(np.mean(signal**2)) # Root mean square amplitude
38
39 # 4. Spectral Domain Features
40 # Perform FFT to get frequency components
41 signal_fft = fft(signal)
42 signal_freq = fftfreq(len(signal), time[1] - time[0])
43 power_spectrum = np.abs(signal_fft)**2
44
45 # Dominant frequency (excluding zero-frequency component)
46 dominant_freq = signal_freq[np.argmax(power_spectrum[1:]) + 1]
47 dominant_power = np.max(power_spectrum[1:])

```

```

48
49 # Display extracted features
50 features = {
51     "Mean": mean_val,
52     "Standard Deviation": std_dev,
53     "Peak-to-Peak": peak_to_peak,
54     "Skewness": signal_skewness,
55     "Kurtosis": signal_kurtosis,
56     "Max Amplitude": max_amplitude,
57     "Min Amplitude": min_amplitude,
58     "RMS Amplitude": rms_amplitude,
59     "Dominant Frequency": dominant_freq,
60     "Dominant Power": dominant_power,
61 }
62
63 # Converting to DataFrame
64 features_df = pd.DataFrame(list(features.items()), columns=["Feature", "Value"])
65 print("Extracted Features:")
66 print(features_df)
67
68 # Reconstructing the signal using key features
69 reconstructed_signal = rms_amplitude * np.sin(2 * np.pi * dominant_freq * time) + mean_val
70
71 # Plot the reconstructed signal
72 plt.figure(figsize=(10, 4))
73 plt.plot(time, signal, label="Original Signal", alpha=0.7)
74 plt.plot(time, reconstructed_signal, label="Reconstructed Signal", linestyle="--")
75 plt.xlabel("Time [s]")
76 plt.ylabel("Amplitude")
77 plt.title("Original vs. Reconstructed Signal")
78 plt.legend()
79 plt.show()
80

```



Understanding the Curse of Dimensionality

- What is the Curse of Dimensionality?

- The curse of dimensionality refers to the various challenges that arise when working with high-dimensional data. As the number of features (dimensions) in a dataset increases, **each data point becomes increasingly sparse in the feature space**, making it harder to identify meaningful patterns.
- **Example:** In image recognition, a single image can contain thousands of pixels as features. While some pixels contribute to identifying objects, many do not, adding complexity and potentially introducing noise rather than useful information.

- Impact on Models:

- Distance-Based Models (e.g., k-NN, k-Means):

- Distance-based models are highly affected by high dimensionality because as dimensions increase, the distances between points become less distinguishable. In k-NN, for example, as the number of irrelevant features increases, it becomes harder to find truly "nearest" neighbors, as all points end up roughly equidistant in high-dimensional space. This results in poor model accuracy, particularly on sparse datasets.

- Linear and Logistic Regression:

- High-dimensional data often **leads to overfitting** in models like linear or logistic regression, where each additional feature introduces a new parameter to learn. If irrelevant or noisy features are included, the model may "learn" from noise, resulting in poor generalization on new data.

- Decision Trees and Ensemble Models (e.g., Random Forest, XGBoost):

- Decision trees and ensemble models tend to perform well on high-dimensional data, but their performance decreases when many irrelevant features are present. Each split in a decision tree depends on feature values; if features are mostly irrelevant, splits become noisy and add unnecessary complexity, leading to overfitting.

- Neural Networks:

- Neural networks, especially deep networks, can theoretically handle high-dimensional data by learning hierarchical representations. However, if the data includes a large number of irrelevant features, **training becomes slower**, and the network may converge to **suboptimal solutions** or require more epochs and **larger training datasets** to distinguish signal from noise effectively.

- Computational Cost and Training Time:

- High-dimensional data significantly **increases computational costs** and **training time** because **each added feature increases the complexity of the model**. Training a model with many features involves more parameters, more calculations per iteration, and a greater memory requirement. This is especially problematic for large datasets, where high dimensionality can make training prohibitively slow, particularly for complex models like neural networks and ensemble models.
- Storage and Memory:** High-dimensional datasets also **require more storage and memory** resources for data loading and preprocessing, as each feature needs to be loaded and processed in each batch, which can slow down training and require specialized hardware in extreme cases.

- Theoretical Solutions - More Data vs. Feature Selection:

- In theory, if a model has access to a much larger number of training examples, it can learn to ignore noise in high-dimensional data and focus on relevant patterns. For instance, in neural networks, with enough data, irrelevant features can be ignored over time. However, this solution has **limitations**:
 - Training Cost:** Collecting and using more data significantly increases the training time and computational cost, requiring more powerful hardware and extended training sessions.
 - Data Quality:** More data may not always be feasible or effective if the additional samples do not represent meaningful patterns or are difficult to acquire in certain domains.
- Feature Selection as a Practical Solution:** Instead of expanding the dataset, feature selection techniques provide a more practical approach by identifying and keeping only the most informative features, reducing the dimensionality and allowing the model to focus on the most relevant aspects of the data. This leads to improved model interpretability, lower computational costs, and faster training times.

How Feature Dimensionality Impacts Model Performance

- Balancing the Right Features

- Underfitting:** Caused by too few features, leading to missing essential patterns.
- Overfitting:** Caused by too many features, causing noise to be learned as patterns.

```

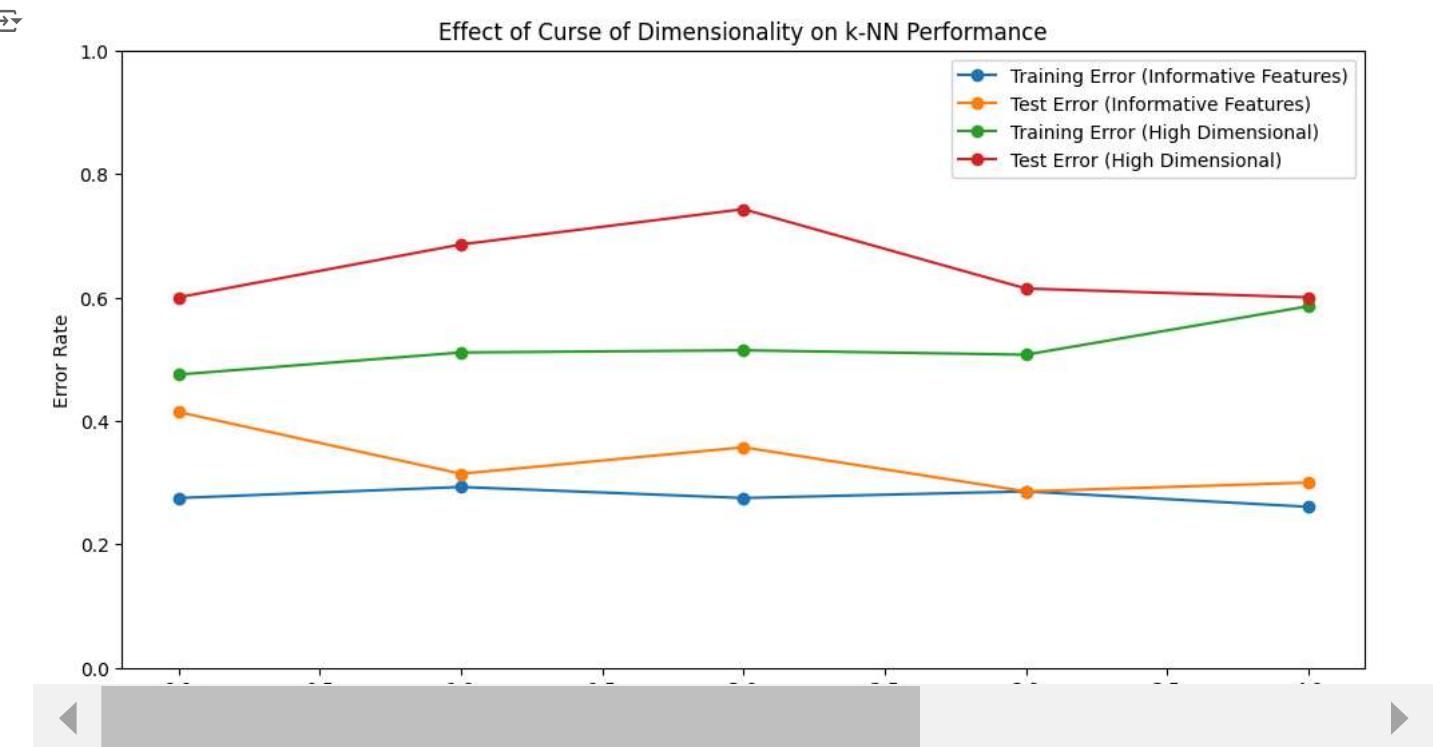
1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import make_classification
4 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
5 from sklearn.model_selection import StratifiedKFold
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.metrics import accuracy_score
8 import matplotlib.pyplot as plt
9
10 # Generate a synthetic dataset with 5 informative features and no redundant or irrelevant features initially
11 X_informative, y = make_classification(
12     n_samples=500, n_features=5, n_informative=4, n_redundant=0, n_classes=8, n_clusters_per_class=1, random_state=0
13 )
14
15 # Split the informative dataset into training and test sets
16 X_train_inf, X_test_inf, y_train, y_test = train_test_split(X_informative, y, test_size=0.3, random_state=42)
17
18 # Now add irrelevant features to create a high-dimensional dataset
19 X_high_dim = np.hstack([X_informative, np.random.normal(0, 1, (X_informative.shape[0], 100))])
20 X_train_high, X_test_high, _, _ = train_test_split(X_high_dim, y, test_size=0.3, random_state=42)
21
22 # Define k-NN and hyperparameter tuning for both datasets
23 knn = KNeighborsClassifier()
24 param_grid = {'n_neighbors': np.arange(1, 31)}
25
26 # Nested cross-validation for the informative dataset
27 cv_outer = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)
28 train_errors_inf, test_errors_inf = [], []
29
30 for train_idx, test_idx in cv_outer.split(X_train_inf, y_train):
31     X_train_outer, X_test_outer = X_train_inf[train_idx], X_train_inf[test_idx]
32     y_train_outer, y_test_outer = y_train[train_idx], y_train[test_idx]
33

```

```

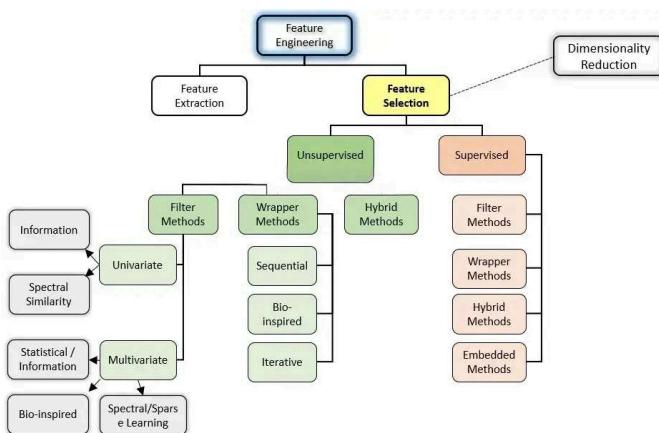
34     # Inner cross-validation for hyperparameter tuning
35     cv_inner = StratifiedKFold(n_splits=3, shuffle=True, random_state=0)
36     grid_search = GridSearchCV(knn, param_grid, scoring='f1_weighted', cv=cv_inner)
37     grid_search.fit(X_train_outer, y_train_outer)
38     best_k = grid_search.best_params_['n_neighbors']
39
40     # Evaluate with best hyperparameters
41     knn_best = KNeighborsClassifier(n_neighbors=best_k)
42     knn_best.fit(X_train_outer, y_train_outer)
43     train_errors_inf.append(1 - knn_best.score(X_train_outer, y_train_outer))
44     test_errors_inf.append(1 - knn_best.score(X_test_outer, y_test_outer))
45
46 # Nested cross-validation for the high-dimensional dataset
47 train_errors_high, test_errors_high = [], []
48
49 for train_idx, test_idx in cv_outer.split(X_train_high, y_train):
50     X_train_outer, X_test_outer = X_train_high[train_idx], X_train_high[test_idx]
51     y_train_outer, y_test_outer = y_train[train_idx], y_train[test_idx]
52
53     # Inner cross-validation for hyperparameter tuning
54     grid_search = GridSearchCV(knn, param_grid, scoring='f1_weighted', cv=cv_inner)
55     grid_search.fit(X_train_outer, y_train_outer)
56     best_k = grid_search.best_params_['n_neighbors']
57
58     # Evaluate with best hyperparameters
59     knn_best = KNeighborsClassifier(n_neighbors=best_k)
60     knn_best.fit(X_train_outer, y_train_outer)
61     train_errors_high.append(1 - knn_best.score(X_train_outer, y_train_outer))
62     test_errors_high.append(1 - knn_best.score(X_test_outer, y_test_outer))
63
64 # Plotting the training and test errors for both datasets
65 plt.figure(figsize=(12, 6))
66 plt.plot(train_errors_inf, label="Training Error (Informative Features)", marker='o')
67 plt.plot(test_errors_inf, label="Test Error (Informative Features)", marker='o')
68 plt.plot(train_errors_high, label="Training Error (High Dimensional)", marker='o')
69 plt.plot(test_errors_high, label="Test Error (High Dimensional)", marker='o')
70 plt.xlabel("Cross-validation Fold")
71 plt.ylabel("Error Rate")
72 plt.ylim(0, 1) # Set y-axis range from 0 to 1
73 plt.title("Effect of Curse of Dimensionality on k-NN Performance")
74 plt.legend()
75 plt.show()
76

```



▼ Types of Feature Selection Algorithms

- **Overview:**
 - **Filter Methods:** Select features based on statistical/mathematical metrics, independent of models.
 - **Wrapper Methods:** Select based on model performance.
 - **Embedded Methods:** Feature selection happens during model training.
- **Choosing the Right Method:** Each method has advantages based on dataset size, computational resources, and model needs.



source: <https://medium.com/analytics-vidhya/feature-selection-extended-overview-b58f1d524c1c>

Filter Methods

- Overview:**
 - Independent statistical/mathematical evaluation of each feature for relevance.
- Pros and Cons:**
 - Pros:** Computationally efficient, model-agnostic.
 - Cons:** Usually does not account for feature interactions.
- When to Use:** Good for initial reduction on large datasets.
- Common Algorithms:** Here's a more detailed list of filter-based feature selection methods, including mathematical formulations and Python libraries commonly used to implement them.

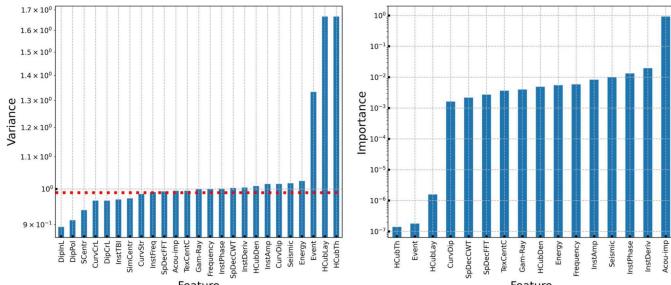
1. Variance Thresholding

- Formulation:** Variance thresholding removes features with variance below a given threshold:

$$\text{Var}(X_j) < \text{threshold}$$

where X_j is the feature and $\text{Var}(X_j)$ is its variance.

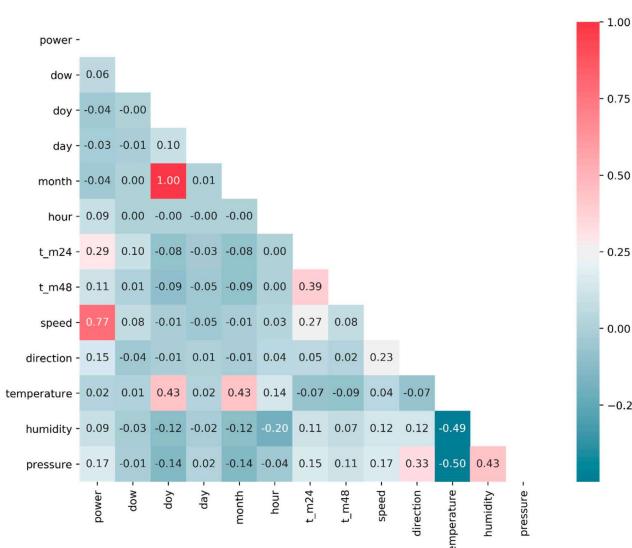
- Python Library:** `sklearn.feature_selection.VarianceThreshold`



source: <https://doi.org/10.1029/2023EA003101>

2. Correlation Analysis

- Formulation:** Features are removed if their correlation coefficient exceeds a threshold. For linear relationships, the Pearson correlation is:
$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$
where $\text{Cov}(X, Y)$ is the covariance between X and Y , and σ represents standard deviation.
- Python Library:** `pandas.DataFrame.corr` or `scipy.stats.pearsonr`



Source:

https://www.researchgate.net/publication/334711620_A_XGBoost_Model_with_Weather_Similarity_Analysis_and_Feature_Engineering_for_Short-Term_Wind_Power_Forecasting/figures?lo=1

3. Chi-Square Test

- **Formulation:** Measures the association between categorical features and target using the chi-square statistic:
- **Python Library:** `sklearn.feature_selection.chi2`

4. ANOVA (Analysis of Variance)

- **Formulation:** Analyzes variance across feature values by calculating the F-statistic:

$$F = \frac{\text{Between-group variance}}{\text{Within-group variance}}$$

where a high F-value indicates significant variance between classes.

- **Python Library:** `sklearn.feature_selection.f_classif`

5. Mutual Information (MI)

- **Formulation:** Measures dependency between feature X and target Y (Mutual Information is a concept from information theory that measures how much knowing one variable tells you about another. For feature selection, mutual information helps us understand how much a feature tells us about the target variable (what we're trying to predict).):

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

where $p(x, y)$ is the joint probability, and $p(x)$ and $p(y)$ are the marginal probabilities.

- **Python Library:** `sklearn.feature_selection.mutual_info_classif`

6. ReliefF Algorithm

Nearest Neighbors: For each data point, ReliefF finds its nearest neighbors:

One neighbor from the same class (called the "near hit"). One neighbor from a different class (called the "near miss"). Feature Scoring:

ReliefF compares the values of each feature for the data point and its neighbors. If a feature has similar values between the data point and its "near hit" but different values between the data point and its "near miss," it's likely a useful feature. The feature's importance score increases if it consistently helps distinguish between classes. Repeat: This process is repeated for many points in the dataset, and the final score for each feature reflects how well it consistently helps separate classes across the data.

- **Formulation:** Assigns weights to features based on nearest neighbor distances. The weight for feature X_j is adjusted as:

$$W_j = W_j - \frac{\text{dist}(X_j^{\text{near hit}}, X_j)}{m} + \frac{\text{dist}(X_j^{\text{near miss}}, X_j)}{m}$$

where dist represents distance, m is the number of samples, and "near hit" and "near miss" refer to closest samples of the same and different classes, respectively.

- **Python Library:** `skrebate` (Relief-based feature selection)

7. Information Gain

- **Formulation:** Measures the decrease in entropy from knowing feature X for predicting Y :

$$IG(Y, X) = H(Y) - H(Y|X)$$

where $H(Y)$ is the entropy of Y and $H(Y|X)$ is the conditional entropy of Y given X .

- **Python Library:** `sklearn.feature_selection.mutual_info_classif` (can calculate information gain by computing entropy reduction)

8. Minimum Redundancy Maximum Relevance (mRMR)

- **Formulation:** Balances feature relevance (correlation with target) and redundancy (correlation among features):

$$\text{mRMR} = \max \left(\text{Relevance}(X_j, Y) - \frac{1}{|S|} \sum_{X_k \in S} \text{Redundancy}(X_j, X_k) \right)$$

where S is the set of selected features.

- **Python Library:** `pymrmr` (Python package for mRMR)

Wrapper Methods

• Overview:

- Wrapper methods evaluate subsets of features based on how they impact model performance, often by using a model's prediction performance as a selection criterion. This iterative process identifies the subset of features that maximizes model performance, accounting for interactions between features.
- These methods can lead to highly optimized models but are generally computationally intensive because they require repeated training and evaluation of the model.

• Pros and Cons:

- **Pros:** Accounts for feature interactions, often providing highly accurate feature subsets for model training.
- **Cons:** Computationally expensive, especially with large datasets or complex models, due to repeated training across multiple feature subsets. Also they are optimized for the model under study and not other models.

- **When to Use:** Wrapper methods are especially suitable for smaller datasets and cases where high model performance is critical. They are also useful when interactions between features might play a significant role in model accuracy.

Common Algorithms

1. Sequential Feature Selection (SFS)

SFS is a greedy search algorithm that builds an optimal feature subset sequentially, either by adding features (forward selection) or removing features (backward elimination), and stopping based on model performance.

- **Forward Selection:**

- **Description:** Starts with an empty feature set and adds one feature at a time that most improves model performance, continuing until adding additional features does not improve the model.

- **Formulation:**

$$\text{Forward Step: } S_{k+1} = S_k \cup \{\arg \max_{j \notin S_k} \text{Perf}(S_k \cup \{j\})\}$$

where S_k is the current feature subset, and `Perf` is a performance metric (like accuracy or F1-score).

- **Pros:** Less computationally demanding than exhaustive methods; good for finding feature subsets in moderate-sized datasets.

- **Cons:** Can miss optimal feature combinations due to the greedy approach.

- **Python Library:** `mlxtend.feature_selection.SequentialFeatureSelector`

- **Backward Elimination:**

- **Description:** Starts with the full feature set, removes the least significant feature one at a time, and re-evaluates until removing additional features reduces performance.

- **Formulation:**

$$\text{Optimal Subset} = \arg \max_{S \subseteq \{1, \dots, n\}} \text{Perf}(S)$$

where S_k is the current feature subset, and `Perf` represents the model performance metric.

- **Pros:** Useful when starting with a larger feature set and aiming to reduce dimensionality.

- **Cons:** More computationally intensive than forward selection, especially with high-dimensional datasets.

- **Python Library:** `mlxtend.feature_selection.SequentialFeatureSelector`

2. Exhaustive Feature Selection

- **Description:** Evaluates all possible feature combinations to find the subset that provides the best model performance. Since it evaluates every subset, it guarantees finding the optimal feature set for a given model.

- **Formulation:** For n features, the exhaustive approach evaluates each subset S of the power set 2^n :

$$\text{Optimal Subset} = \arg \max_{S \subseteq \{1, \dots, n\}} \text{Perf}(S)$$

where `Perf` is a performance metric like accuracy, F1-score, or AUC.

- **Pros:** Provides the best possible feature subset.

- **Cons:** Extremely computationally expensive and impractical for large datasets; limited to small feature sets due to the exponential growth in combinations.

- **Python Library:** `mlxtend.feature_selection.ExhaustiveFeatureSelector`

- **Real-World Example:** In healthcare prediction models, where the importance of specific features (such as test results) is critical, Recursive Feature Elimination (RFE) is often used to select the most predictive medical test results for diagnostics. For example, in cancer diagnostics, RFE might identify a small number of biomarkers from a large set, improving model interpretability and focusing on the most relevant indicators.

Embedded Methods

- **Overview:**

- Embedded methods select features as part of the model training process. These methods incorporate regularization or feature importance scoring within the model itself, thereby eliminating the need for an external feature selection step.

- **Pros and Cons:**

- **Pros:** Embedded methods balance performance and computational efficiency since feature selection is integrated into the training process, and they can account for feature interactions.
- **Cons:** They are often model-specific and rely on selection mechanisms inherent to the chosen algorithm, meaning they may not be transferrable across different models.

- **When to Use:** Embedded methods are ideal when using models with built-in feature selection capabilities, such as linear models with regularization, decision trees, and ensemble methods.

Common Algorithms

1. LASSO (L1 Regularization)

- **Description:** LASSO (Least Absolute Shrinkage and Selection Operator) adds an L_1 -norm penalty to the model's cost function. This forces some feature coefficients to become exactly zero, effectively removing them from the model and resulting in a sparse model with only the most important features.

- **Formulation:** For a linear regression model with parameters β , the LASSO objective function is:

$$\text{Minimize} \quad \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p X_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

where λ is the regularization parameter controlling the strength of feature selection.

- **Pros:** Effective in high-dimensional data, performs feature selection and regularization simultaneously.
- **Cons:** May select only one feature from a set of highly correlated features.
- **Python Library:** `sklearn.linear_model.Lasso`

2. Ridge Regression (L2 Regularization)

- **Description:** Ridge regression adds an $L2$ -norm penalty to the model, shrinking coefficients but not setting them to zero, making it useful for identifying smaller yet still influential features.
- **Formulation:** The Ridge regression objective function is:

$$\text{Minimize} \quad \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p X_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

where λ controls the strength of the penalty.

- **Pros:** Helps prevent overfitting; retains all features, which can be advantageous when removing features entirely is undesirable.
- **Cons:** Does not create a sparse model (i.e., does not fully eliminate features).
- **Python Library:** `sklearn.linear_model.Ridge`

3. Elastic Net Regularization

- **Description:** Elastic Net combines both $L1$ and $L2$ regularization to retain the benefits of LASSO (sparsity) and Ridge regression (stability with correlated features). It is suitable for cases with highly correlated features.
- **Formulation:** For parameters β , the Elastic Net objective function is:

$$\text{Minimize} \quad \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p X_{ij}\beta_j \right)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2$$

where λ_1 and λ_2 control the $L1$ and $L2$ penalties, respectively.

- **Pros:** Can handle correlated features and offers flexible regularization.
- **Cons:** Requires tuning two regularization parameters.
- **Python Library:** `sklearn.linear_model.ElasticNet`

4. Tree-Based Feature Importance

Impurity in decision trees is a measure of how mixed the classes are within a node (or group) of the tree. For **feature selection**, impurity helps to decide which features are most helpful in separating the data into distinct classes. Decision trees use this measure to choose the best features step-by-step as the tree grows.

1. Impurity Measures:

- **Gini Impurity** and **Entropy** are common impurity measures. They assess how "pure" or "impure" a group is:
 - **Pure:** If a node has instances from only one class, it's "pure," meaning it has zero impurity.
 - **Impure:** If a node has instances from multiple classes, it's "impure" (higher impurity score).

2. Feature Selection by Reducing Impurity:

- The tree tests each feature to see how much it can **reduce impurity** by splitting the data.
- A feature that reduces impurity the most is chosen as the next split in the tree, making it an important feature.

3. Feature Importance Scores:

- As the tree grows, it calculates how much each feature reduces impurity across all splits.
- Features that consistently reduce impurity by a lot get high importance scores, making them strong predictors.

Why It's Useful:

- **Higher scores** mean a feature is more useful for classification, while **lower scores** indicate that the feature doesn't help much and might be dropped.

For example, in a tree predicting loan approval, "income" might have a high importance score if it consistently helps separate approved from non-approved cases, while "zip code" might have a lower score if it doesn't reduce impurity much.

- **Pros:** Intuitive and useful for both classification and regression tasks; accounts for interactions between features.
- **Cons:** Biased toward features with more levels (in categorical data).
- **Python Library:** `sklearn.ensemble.RandomForestClassifier` or `DecisionTreeClassifier`

5. Gradient Boosting Feature Importance

- **Description:** Gradient boosting models rank features by the frequency and effectiveness of their splits in each decision tree, thereby providing a feature importance score based on how often and effectively a feature is used across trees.
- **Pros:** Effective for complex datasets and provides highly accurate models.
- **Cons:** Computationally intensive due to iterative training of multiple trees.
- **Python Library:** `sklearn.ensemble.GradientBoostingClassifier`

6. Penalized Logistic Regression

- **Description:** Logistic regression with an $L1$ or $L2$ penalty can be used to select features. The penalty term shrinks less important features, with $L1$ setting some coefficients exactly to zero for sparsity.
- **Formulation:** The objective function with $L1$ regularization is:

$$\text{Minimize} \quad - \sum_{i=1}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i)) + \lambda \sum_{j=1}^p |\beta_j|$$

where p_i is the predicted probability and λ controls the penalty.

- **Pros:** Effective for binary classification tasks, especially when many features are uninformative.
- **Cons:** Sensitive to multicollinearity; requires careful parameter tuning.
- **Python Library:** `sklearn.linear_model.LogisticRegression`

✓ Comparing Feature Selection Methods

- **Performance and Efficiency:**
 - **Filter Methods:** Fastest; use as a first-pass technique for dimensionality reduction.
 - **Wrapper Methods:** Best for capturing feature interactions but computationally intensive.
 - **Embedded Methods:** Balanced; efficient for models with built-in selection.
- **Choosing the Right Method:**
 - **Filter Methods:** When dealing with large datasets and quick feature selection is needed.
 - **Wrapper Methods:** For moderate datasets where accuracy is critical.
 - **Embedded Methods:** When using models that integrate feature selection.

Conclusion

- **Key Takeaways:**
 - Effective feature engineering and selection are essential for optimizing model performance.
 - Consider the curse of dimensionality, balancing too few or too many features.
 - Select methods based on dataset size, computational resources, and model type.
 - Remember other considerations like feature stability, transformations, and encoding choices.

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import make_classification
4 from sklearn.model_selection import train_test_split, GridSearchCV, StratifiedKFold
5 from sklearn.metrics import accuracy_score
6 import matplotlib.pyplot as plt
7 from xgboost import XGBClassifier
8
9 # Generate a synthetic dataset with 5 informative features and no redundant or irrelevant features initially
10 X_informative, y = make_classification(
11     n_samples=500, n_features=5, n_informative=4, n_redundant=0, n_classes=8, n_clusters_per_class=1, random_state=0
12 )
13
14 # Split the informative dataset into training and test sets
15 X_train_inf, X_test_inf, y_train, y_test = train_test_split(X_informative, y, test_size=0.3, random_state=42)
16
17 # Now add irrelevant features to create a high-dimensional dataset
18 X_high_dim = np.hstack([X_informative, np.random.normal(0, 1, (X_informative.shape[0], 100))])
19 X_train_high, X_test_high, _, _ = train_test_split(X_high_dim, y, test_size=0.3, random_state=42)
20
21 # Define XGBoost classifier and hyperparameter grid
22 xgb = XGBClassifier(eval_metric='mlogloss')
23 param_grid = {
24     'max_depth': [5, 7],
25     'learning_rate': [0.1, 0.2],
26     'n_estimators': [50, 100]
27 }
28
29 # Nested cross-validation for the informative dataset
30 cv_outer = StratifiedKFold(n_splits=5, shuffle=True, random_state=0)
31 train_errors_inf, test_errors_inf = [], []
32
33 for train_idx, test_idx in cv_outer.split(X_train_inf, y_train):
34     X_train_outer, X_test_outer = X_train_inf[train_idx], X_train_inf[test_idx]
35     y_train_outer, y_test_outer = y_train[train_idx], y_train[test_idx]
36
37     # Inner cross-validation for hyperparameter tuning
38     cv_inner = StratifiedKFold(n_splits=3, shuffle=True, random_state=0)
39     grid_search = GridSearchCV(xgb, param_grid, scoring='f1_weighted', cv=cv_inner, n_jobs=-1)
40     grid_search.fit(X_train_outer, y_train_outer)
41     best_params = grid_search.best_params_
42
43     # Evaluate with best hyperparameters
44     xgb_best = XGBClassifier(**best_params, eval_metric='mlogloss')
45     xgb_best.fit(X_train_outer, y_train_outer)
46     train_errors_inf.append(1 - xgb_best.score(X_train_outer, y_train_outer))
47     test_errors_inf.append(1 - xgb_best.score(X_test_outer, y_test_outer))
48
49 # Nested cross-validation for the high-dimensional dataset
50 train_errors_high, test_errors_high = [], []
51
52 for train_idx, test_idx in cv_outer.split(X_train_high, y_train):
53     X_train_outer, X_test_outer = X_train_high[train_idx], X_train_high[test_idx]
54     y_train_outer, y_test_outer = y_train[train_idx], y_train[test_idx]
55
56     # Inner cross-validation for hyperparameter tuning
57     grid_search = GridSearchCV(xgb, param_grid, scoring='f1_weighted', cv=cv_inner, n_jobs=-1)
58     grid_search.fit(X_train_outer, y_train_outer)
59     best_params = grid_search.best_params_
60
61     # Evaluate with best hyperparameters

```

```
62     xgb_best = XGBClassifier(**best_params, eval_metric='mlogloss')
63     xgb_best.fit(X_train_outer, y_train_outer)
64     train_errors_high.append(1 - xgb_best.score(X_train_outer, y_train_outer))
65     test_errors_high.append(1 - xgb_best.score(X_test_outer, y_test_outer))
66
67 # Plotting the training and test errors for both datasets
68 plt.figure(figsize=(12, 6))
69 plt.plot(train_errors_inf, label="Training Error (Informative Features)", marker='o')
70 plt.plot(test_errors_inf, label="Test Error (Informative Features)", marker='o')
71 plt.plot(train_errors_high, label="Training Error (High Dimensional)", marker='o')
72 plt.plot(test_errors_high, label="Test Error (High Dimensional)", marker='o')
73 plt.xlabel("Cross-validation Fold")
74 plt.ylabel("Error Rate")
75 plt.ylim(0, 1) # Set y-axis range from 0 to 1
76 plt.title("Effect of Curse of Dimensionality on XGBoost Performance")
77 plt.legend()
78 plt.show()
79
```

