

Unit 3-Resampling Methods, Linear Model Selection and Regularization(ML)

Resampling Methods:

Resampling methods are techniques used in statistics and machine learning to estimate the performance or assess the uncertainty of a model. They involve repeatedly drawing samples from a given dataset and analyzing the samples to make inferences about the population or to evaluate the model's performance. **The two commonly used resampling methods are:**

1. **Cross-Validation:** Cross-validation is a technique used to assess the performance of a model on an independent dataset. It involves dividing the dataset into several subsets or folds. The model is trained on a subset of the data (training set) and evaluated on the remaining subset (validation set). This process is repeated multiple times, with different subsets serving as the validation set each time. The performance measures obtained from each iteration are then averaged to obtain a more robust estimate of the model's performance.
 - **K-Fold Cross-Validation:** In k-fold cross-validation, the dataset is divided into k equal-sized folds. The model is trained on k-1 folds and evaluated on the remaining fold. This process is repeated k times, with each fold serving as the validation set once. The performance results are then averaged across all iterations.
 - **Leave-One-Out Cross-Validation (LOOCV):** LOOCV is a special case of k-fold cross-validation where k is equal to the number of samples in the dataset. In LOOCV, the model is trained on all but one sample and evaluated on the single left-out sample. This process is repeated for each sample in the dataset, and the performance results are averaged.
2. **Bootstrap:** Bootstrap is a resampling method that involves sampling with replacement from the original dataset to generate multiple bootstrap samples. Each bootstrap sample has the same size as the original dataset, but it may contain duplicate instances and exclude some original instances. The model is trained on each bootstrap sample, and the aggregated results are used to estimate the performance or uncertainty associated with the model. Bootstrap is particularly useful when the dataset is limited or when the underlying distribution is unknown.

Linear Model Selection and Regularization:

Linear model selection and regularization techniques are used to improve the performance of linear regression models by preventing overfitting and selecting the most relevant features. These techniques are especially valuable when dealing with high-dimensional data.

1. **Subset Selection:** Subset selection methods aim to find the best subset of features from a larger pool of candidates. There are two main approaches:
 - **Best Subset Selection:** Best subset selection involves fitting the linear regression model with all possible combinations of predictors and selecting the one that yields the best fit based on a criterion such as the lowest residual sum of squares (RSS) or highest adjusted R-squared.
 - **Stepwise Selection:** Stepwise selection methods, including forward selection, backward elimination, and stepwise regression, start with an empty model or a full model and iteratively add or remove predictors based on a criterion like Akaike information criterion (AIC) or Bayesian information criterion (BIC).
2. **Regularization:** Regularization methods introduce a penalty term to the linear regression objective function to shrink the coefficients of the predictors towards zero, effectively reducing the complexity of the model. This helps to avoid overfitting and improve generalization. Two commonly used regularization techniques are:
 - **Ridge Regression (L2 Regularization):** Ridge regression adds the L2 norm of the coefficient vector to the objective function, penalizing large coefficient values. This leads to a more stable and robust model. The regularization strength is controlled by a hyperparameter (lambda) that determines the amount of shrinkage applied to the coefficients.
 - **Lasso Regression (L1 Regularization):** Lasso regression adds the L1 norm of the coefficient vector to the objective function, promoting sparsity and encouraging some coefficients to become exactly zero. Lasso can perform both feature selection and coefficient shrinkage. The regularization strength is controlled by the lambda parameter.

ResamplingMethods: Cross-Validation

Cross-validation is a popular resampling method used to evaluate the performance of a machine learning model on a dataset. It helps estimate how well the model will generalize to unseen data and provides insights into its robustness and potential overfitting.

The most commonly used form of cross-validation is k-fold cross-validation. Here's how it works:

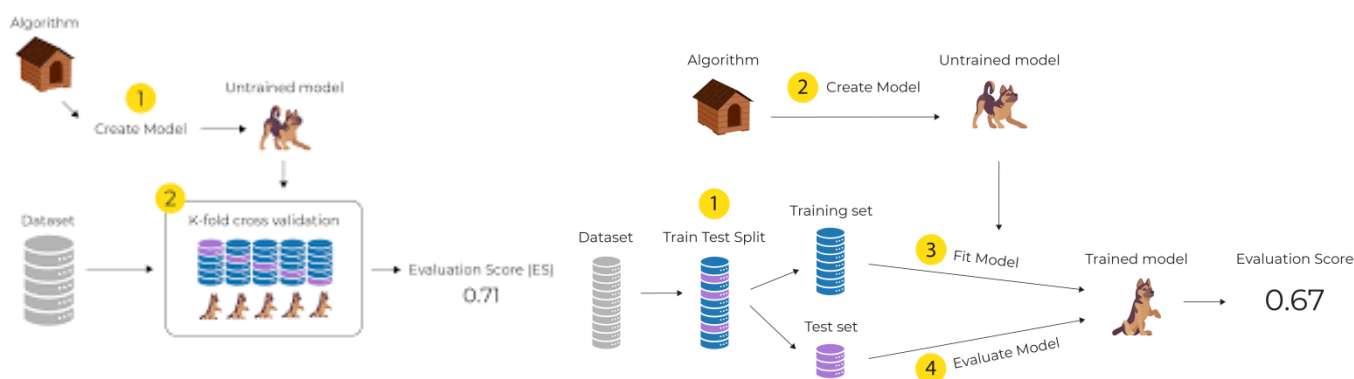
1. **Splitting the data:** The original dataset is divided into k equal-sized subsets or folds.
2. **Iterations:** The model is trained and evaluated k times. In each iteration, one of the folds is designated as the validation set, and the remaining k-1 folds are used as the training set.
3. **Training and evaluation:** The model is trained on the training set, and its performance is evaluated on the validation set using a predefined evaluation metric (e.g., accuracy, mean squared error, etc.).
4. **Performance aggregation:** The performance metric obtained from each iteration is recorded. Typically, the metrics are averaged across all iterations to obtain a single performance estimate.

By repeating the training and evaluation process k times with different fold configurations, cross-validation provides a more reliable estimate of the model's performance compared to a single train-test split.

The value of k can vary depending on the size of the dataset and the available computational resources. Common choices for k include 5, 10, and in some cases, even higher values like 20 or 50.

The benefits of cross-validation are:

- 1) **Reliable performance estimation:** Cross-validation provides a more accurate estimate of a model's performance compared to a single train-test split, as it utilizes multiple validation sets.
- 2) **Robustness assessment:** Cross-validation helps assess the stability and consistency of a model's performance across different subsets of data.
- 3) **Hyperparameter tuning:** Cross-validation is often used in conjunction with hyperparameter tuning techniques like grid search or randomized search to find the optimal hyperparameter values that yield the best model performance.
- 4) **Data utilization:** Cross-validation ensures that each data point is used for both training and validation at some point, maximizing the use of available data.



the Bootstrap

The bootstrap is a resampling method that is commonly used in statistics and machine learning to estimate uncertainty, evaluate model performance, and make inferences about population parameters. Unlike cross-validation, which involves partitioning the data into distinct training and validation sets, the bootstrap resamples the dataset itself to create multiple bootstrap samples.

Here's how the bootstrap method works:

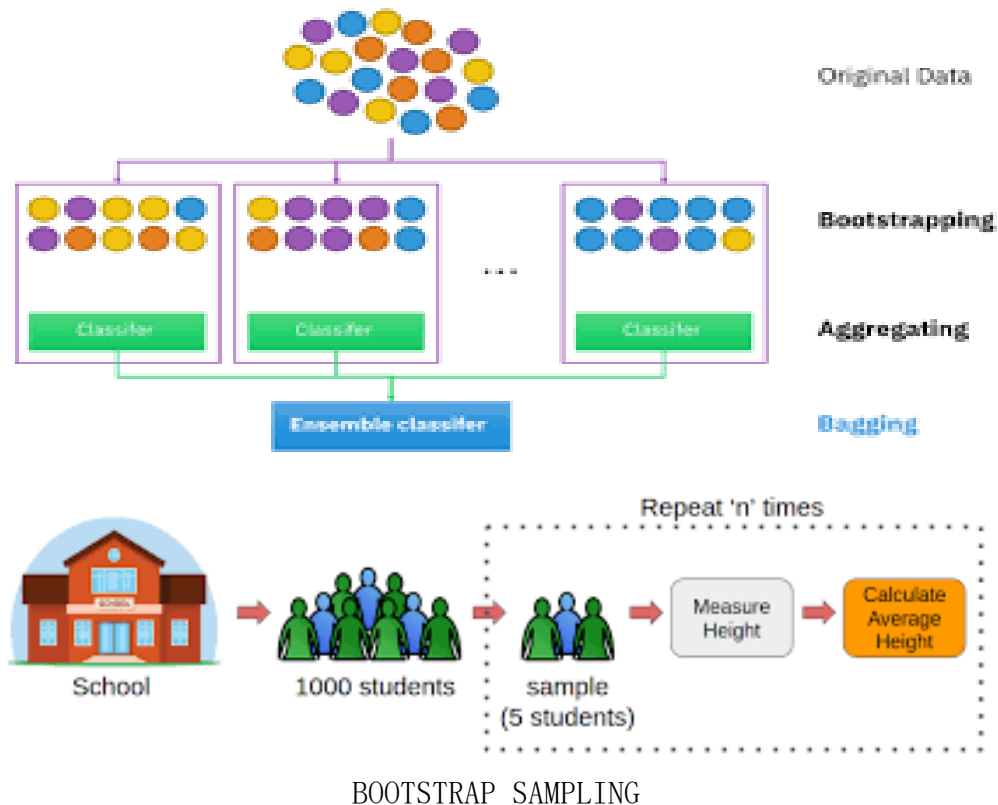
1. **Sample creation:** Given a dataset of size N , the bootstrap method generates multiple bootstrap samples by randomly selecting N instances from the original dataset with replacement. This means that each bootstrap sample can contain duplicate instances from the original dataset, while some original instances may not be included.
2. **Training and evaluation:** For each bootstrap sample, a model is trained using the resampled data. The model's performance is evaluated on the remaining instances that were not selected in the bootstrap sample. This process is repeated for each bootstrap sample.
3. **Aggregation of results:** The performance metrics obtained from each iteration (each bootstrap sample) are then aggregated to estimate the model's performance. This aggregation can be done by averaging the results, calculating percentiles, or constructing confidence intervals.

The bootstrap method allows us to estimate uncertainty and assess the variability of model performance. It provides measures such as bootstrap mean, standard error, and confidence intervals that quantify the uncertainty associated with the performance metric.

Advantages of the bootstrap method include:

- 1) **No assumptions about the underlying distribution:** The bootstrap method is distribution-free, meaning it does not require assumptions about the distribution of the data. This makes it a flexible and widely applicable resampling technique.
- 2) **Utilizes the entire dataset:** The bootstrap method effectively utilizes the available data by creating multiple resamples from the original dataset. This is especially useful when the dataset is limited or when the underlying distribution is unknown.

- 3) Estimation of uncertainty: The bootstrap provides estimates of uncertainty by generating multiple bootstrap samples and calculating statistics from these samples. This information can be used to quantify the variability of model performance or parameter estimates.
- 4) Robustness against outliers: The bootstrap method is robust against outliers because it allows for the inclusion of duplicated instances. Outliers that are present in the original dataset are likely to be included in some of the bootstrap samples, mitigating their impact on the overall analysis.



Linear Model Selection and Regularization :

Linear model selection and regularization are techniques used to improve the performance and interpretability of linear regression models, particularly when dealing with high-dimensional data or when multicollinearity is present. These methods aim to select the most relevant features and prevent overfitting, thereby improving the model's generalization to unseen data. Two commonly used techniques are:

1. Feature Selection:

Feature selection involves choosing a subset of the most relevant features from a larger set of predictors. This is essential when dealing with high-dimensional data, as it can help in reducing model complexity and improving interpretability. There are several approaches to feature selection:

a. Forward Selection: Start with an empty model and iteratively add the most significant predictor at each step until a stopping criterion is met (e.g., reaching a certain model size or significance threshold).

b. Backward Elimination: Start with a model containing all predictors and iteratively remove the least significant predictor at each step until a stopping criterion is met.

c. Stepwise Selection: A combination of forward selection and backward elimination, where predictors are added or removed based on their significance levels.

d. Lasso Regression: Lasso stands for "Least Absolute Shrinkage and Selection Operator." It adds a penalty term to the linear regression cost function, which forces some regression coefficients to become exactly zero. This leads to automatic feature selection, as some predictors will be excluded from the model.

2. Regularization:

Regularization techniques add a penalty term to the linear regression cost function, discouraging large coefficients for predictors and thus reducing model complexity. Regularization helps prevent overfitting and improves the model's ability to generalize to new data. Two common forms of regularization are:

a. L1 Regularization (Lasso): In Lasso regression, the penalty term is proportional to the absolute values of the regression coefficients. This leads to sparse coefficients and effectively performs feature selection, as some coefficients can become exactly zero.

b. b. L2 Regularization (Ridge): In Ridge regression, the penalty term is proportional to the square of the regression coefficients. While it does not result in exact feature selection like Lasso, it effectively shrinks the coefficients, making them small but non-zero.

Both feature selection and regularization help to address the issue of multicollinearity, where some predictors are highly correlated with each other, leading to unstable and unreliable coefficient estimates.

Subset selection

Subset selection is a technique used in linear model selection to identify the best subset of predictors from a larger pool of potential predictors. The goal is to find a subset of predictors that provides the best balance between model fit and complexity. Subset selection methods can be broadly categorized into two main approaches: best subset selection and stepwise selection.

1. Best Subset Selection:

Best subset selection involves fitting the linear regression model with all possible combinations of predictors and selecting the model that yields the best fit based on a specified criterion. The criterion can be chosen based on statistical metrics such as the lowest residual sum of squares (RSS), highest adjusted R-squared, or other appropriate metrics. The steps involved in best subset selection are as follows:

- a. Generate all possible subsets: Generate all possible combinations of predictors from the original pool of predictors.
- b. Fit the model for each subset: Fit a linear regression model for each subset of predictors.
- c. Evaluate model fit: Evaluate the model fit for each subset based on the chosen criterion.
- d. Select the best model: Select the model that provides the best fit according to the chosen criterion.

Best subset selection can provide an exhaustive search over all possible subsets and identify the subset of predictors that yields the best model fit. However, it becomes computationally expensive as the number of predictors increases because the number of possible combinations grows exponentially.

2. Stepwise Selection:

Stepwise selection methods iteratively build the model by adding or removing predictors based on specific criteria. Stepwise selection techniques are computationally efficient compared to best subset selection and are commonly used in practice. There are three commonly used stepwise selection techniques:

a. Forward Selection:

- Start with an empty model.
- Iteratively add predictors that result in the largest improvement in the chosen criterion until no further improvement is observed.

b. Backward Elimination:

- Start with a model that includes all predictors.
- Iteratively remove predictors that result in the least deterioration in the chosen criterion until further removals are not justified.

c. Stepwise Regression:

- Combines forward selection and backward elimination.
- Start with an empty model.
- Allow both adding and removing predictors at each step based on the chosen criterion.

Stepwise selection methods provide a trade-off between computational efficiency and model quality. While they may not always find the globally optimal subset of predictors, they often provide good subsets with reasonable computational resources.

Subset selection methods can help in improving model interpretability by selecting a subset of relevant predictors and removing irrelevant or redundant predictors. However, it's important to note that subset selection methods do not account for the collinearity among predictors, and they can lead to overfitting if the selection criterion is not appropriately chosen or if the sample size is small.

Shrinkage Methods

Shrinkage methods are a class of techniques used in linear regression to mitigate the impact of multicollinearity and overfitting by shrinking or regularizing the regression coefficients. These methods add a penalty term to the linear regression cost function that discourages large coefficient values, effectively "shrinking" the coefficients towards zero. As a result, some coefficients may become exactly zero, leading to feature selection, while others are reduced but

remain non-zero. Shrinkage methods are particularly useful when dealing with high-dimensional data where the number of predictors is larger than the number of observations.

Two common shrinkage methods are:

1. Ridge Regression (L2 Regularization):

Ridge regression adds a penalty term to the linear regression cost function, which is proportional to the sum of squared values of the regression coefficients multiplied by a tuning parameter (lambda or alpha). The cost function in Ridge regression is given by:

➤ $\text{Cost} = \text{RSS (Residual Sum of Squares)} + \lambda * \sum(\text{coefficient}^2)$

The tuning parameter λ controls the amount of regularization applied to the coefficients. A higher λ value will result in stronger regularization, shrinking the coefficients more towards zero. Ridge regression is particularly effective when dealing with multicollinearity, as it can handle cases where predictors are highly correlated with each other.

2. Lasso Regression (L1 Regularization):

Lasso regression also adds a penalty term to the linear regression cost function, but this time it is proportional to the sum of the absolute values of the regression coefficients multiplied by the tuning parameter (lambda or alpha). The cost function in Lasso regression is given by:

➤ $\text{Cost} = \text{RSS (Residual Sum of Squares)} + \lambda * \sum|\text{coefficient}|$

Lasso has the attractive property of performing feature selection, as it tends to drive some regression coefficients to exactly zero. This makes Lasso a valuable tool when there is a need to identify the most relevant predictors, effectively reducing the model's complexity and improving its interpretability.

Both Ridge and Lasso regression provide a balance between bias and variance in the model. By penalizing large coefficients, they prevent overfitting, leading to better generalization to new data. The choice between Ridge and Lasso, or the combination of both (Elastic Net), depends on the specific characteristics of the dataset and the underlying assumptions about the importance of the predictors.

Dimension Reduction Methods

Dimension reduction methods are techniques used to reduce the number of predictors (features) in a dataset while preserving as much relevant information as possible. These methods are particularly useful when dealing with high-dimensional data, as they help to address issues like multicollinearity, overfitting, and the curse of dimensionality. By reducing the number of features, dimension reduction methods can simplify models, improve interpretability, and speed up computation.

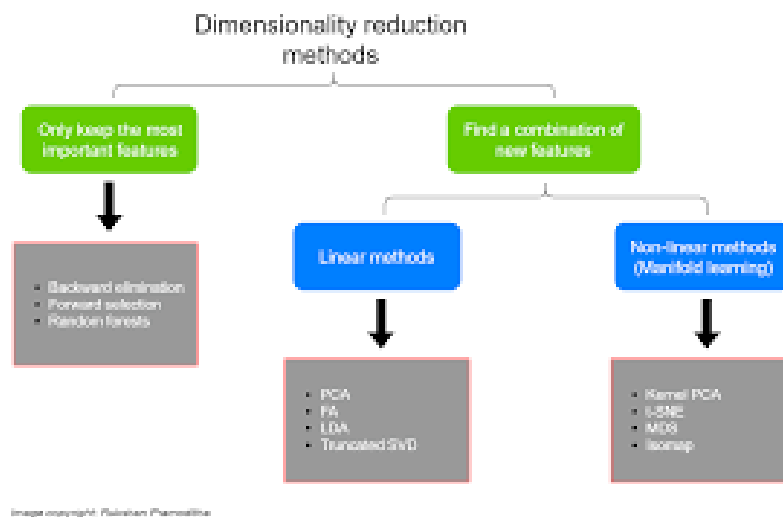


image copyright: Reinchen Perceval

There are two primary types of dimension reduction methods:

1. Principal Component Analysis (PCA):

PCA is an unsupervised dimension reduction technique that transforms the original features into a new set of orthogonal variables called principal components. These components are ordered in such a way that the first principal component explains the most variance in the data, the second component explains the second most, and so on. The number of principal components is equal to the number of original features.

By retaining only the top k principal components, where k is a user-defined parameter, we effectively reduce the dimensionality of the data. PCA ensures that the new components capture the maximum variance possible, preserving the most important information from the original features. It is important to note that PCA is a linear transformation and may not be suitable for all types of data distributions.

2. t-distributed Stochastic Neighbor Embedding (t-SNE):

t-SNE is a non-linear dimension reduction technique that is primarily used for visualization purposes. It is particularly effective at visualizing high-dimensional data in lower-dimensional space (usually 2D or 3D) while preserving the local structure and relationships between data points.

t-SNE works by modeling pairwise similarities between data points in the original high-dimensional space and the lower-dimensional space. It tries to minimize the divergence between the two probability distributions, effectively preserving the neighborhood relationships among data points. While t-SNE is excellent for visualization, it does not preserve global structures as well as PCA.

Both PCA and t-SNE are widely used in various fields, such as data science, machine learning, and data visualization. The choice between these methods depends on the specific goals of the analysis. If the goal is to reduce dimensionality while preserving as much variance as possible, PCA is often preferred. On the other hand, if the main objective is to visualize high-dimensional data in a lower-dimensional space, t-SNE is a powerful tool for revealing patterns and clusters in the data.

Considerations in High Dimensions

In high-dimensional data, where the number of predictors (features) is much larger than the number of observations, several challenges and considerations arise. High-dimensional datasets are common in various fields, including genomics, image processing, natural language processing, and machine learning. Here are some key considerations when dealing with high-dimensional data:

- 1) **Curse of Dimensionality:** High-dimensional data can suffer from the curse of dimensionality. As the number of dimensions increases, the volume of the data space expands exponentially. Consequently, the available data points become sparse, making it difficult to accurately estimate relationships and patterns in the data.
- 2) **Overfitting:** In high-dimensional settings, there is a higher risk of overfitting. When the number of features is large compared to the sample size, a model can learn noise or spurious patterns in the data, leading to poor generalization to new, unseen data.
- 3) **Multicollinearity:** High-dimensional datasets often exhibit multicollinearity, where predictors are highly correlated with each other. This can lead to unstable coefficient estimates and difficulty in interpreting the individual effects of predictors on the response variable.
- 4) **Computational Complexity:** Analyzing high-dimensional data can be computationally intensive. Many algorithms have time and memory requirements that increase dramatically with the number of features, making traditional methods impractical for large-scale datasets.
- 5) **Feature Selection:** Identifying the most relevant features becomes crucial in high-dimensional data. Including irrelevant or redundant features can lead to reduced model performance and increased computation time. Proper feature selection methods are required to obtain a parsimonious and interpretable model.
- 6) **Dimension Reduction:** Dimension reduction techniques like Principal Component Analysis (PCA) and t-distributed Stochastic Neighbor Embedding (t-SNE) can be useful for reducing the number of dimensions while preserving the most important information. This can help in data visualization and improving model performance.
- 7) **Regularization:** Regularization techniques like Ridge and Lasso regression can be employed to address overfitting and multicollinearity in high-dimensional models. These methods add penalty terms to the regression cost function to control the magnitude of regression coefficients.
- 8) **Cross-Validation:** Cross-validation is essential for assessing the generalization performance of models in high-dimensional data. It helps in estimating how well the model will perform on unseen data and can assist in tuning hyperparameters effectively.
- 9) **Feature Engineering:** In high-dimensional data, feature engineering becomes critical to create meaningful and informative features that can improve model performance. Domain knowledge and expertise play a crucial role in this process.
- 10) **Sample Size Considerations:** The sample size should be large enough relative to the number of features to achieve reliable and robust results. Insufficient data points may lead to overfitting and unreliable estimates.