# A Machine Learning Handbook on Overfitting

Kartikey Vijayakumar Hebbar
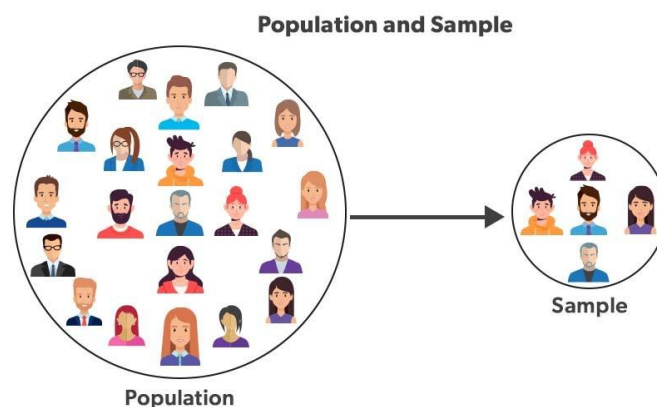
# Overfitting



Before diving into the concept of overfitting in Machine Learning, let's understand about few important terminologies like Population & Sample.

## Population

In statistics, the population refers to the entire set of possible observations that could ever be observed or collected. In machine learning, the population represents the entirety of data you would ideally like to make predictions about. This includes all possible instances that could ever be encountered in the real-world scenario the model is intended for.

In Machine Learning, the population is the ultimate target for a model's predictions. Ideally, a machine learning model should be able to generalize well across the entire population. However, since it's usually impossible to collect data on the entire population, we rely on samples.

**Sample**

A sample is a subset of the population, selected to represent the population in a manageable way. In machine learning, the sample consists of the data that we actually have available — the data on which the model is trained and validated.

A machine learning model is trained on a sample of the population data, and its performance is usually evaluated on a separate test sample. The key challenge is ensuring that these samples are representative of the overall population. If the sample is biased or too small, the model might not learn the underlying patterns effectively but instead learn the noise and specific peculiarities of the sample.

## 1.1 Definition and Causes

Consider a Machine Learning model which is being trained on a particular dataset. Imagine that it learns the noise in the training data to the extent that it negatively affects the performance of the model on new data. This is called overfitting of a Machine Learning model.

Overfitting is an undesirable machine learning behavior that occurs when the machine learning model gives accurate predictions for training data but not for new data. When data scientists use machine learning models for making predictions, they first train the model on a known data set. Then, based on this information, the model tries to predict outcomes for new data sets. An overfit model can give inaccurate predictions and cannot perform well for all types of new data.



When does overfitting occur? There are few cases where overfitting can occur:

1. **Too Complex Model:** When the model has too many parameters relative to the number of observations, it can capture excessive noise and anomalies in the training data that do not generalize to new data.
2. **Inadequate Training Data:** Having too few data points can make the variability in the data appear like meaningful patterns to the model, which are actually just random noise.
3. **Lack of Regularization:** Regularization techniques add penalties on the size of the coefficients in regression models. Without these penalties, models may become overly complex.
4. **Poor Choice of Model:** Choosing the wrong type of model for the data can also lead to overfitting. For example, using a highly flexible model for simple problems.
5. **Training for Too Long:** Especially in the context of neural networks, training a model for too many epochs can lead it to learn details that do not generalize well.
6. **Lack of Model Validation:** Not using techniques like cross-validation to check the model's performance can result in choosing a model that fits the training data very well but performs poorly on any unseen data.

In an overview, overfitting is often a result of an excessively fine-tuned model that loses its ability to generalize, hence failing on any new or unseen data.

## 1.2 Detection of Overfitting

### 1.2.1 Performance Metrics

Consider a Martial Arts professional named Jackie training for an upcoming national level competition. His coach wants him to win the competition and bring the trophy. What would he do apart from conducting training sessions for Jackie? He would analyze Jackie's performances in the training and previous competitions to understand his strengths and weaknesses. By performing this analysis, the coach will then instruct Jackie to work on his weaknesses to bring the best out of his performances.
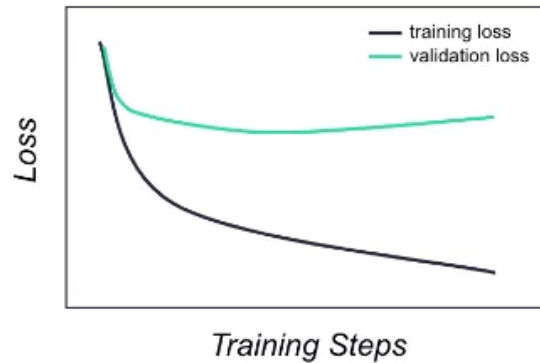


Now, let's bring this analogy to Machine Learning models to understand what Performance Metrics are and how they can be used to detect overfitting. Performance metrics are crucial in detecting overfitting by comparing how the model performs on

training data versus unseen test or validation data. Here's how they can be effectively used:
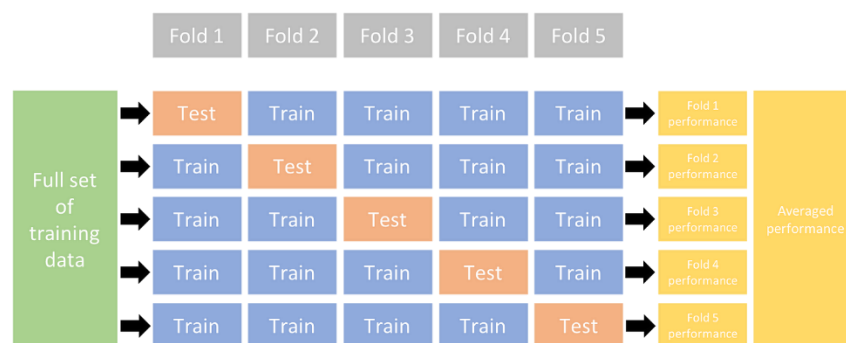
1. **Training vs. Validation Metrics:** The most direct approach to detect overfitting is to monitor performance metrics on both training and validation datasets during training. Common metrics include accuracy, precision, recall, F1 score for classification tasks, and mean squared error, mean absolute error for regression tasks.

   💡 **Signs of Overfitting:** If a model's performance on the training data improves over time but its performance on the validation data starts to worsen, this is a clear indication of overfitting.



2. **Cross-Validation:** Utilizing cross-validation techniques, especially k-fold cross-validation, helps in assessing the model's ability to generalize to an independent dataset. It involves dividing the dataset into 'k' subsets and repeatedly training the model 'k' times, each time using a different subset as the test set and the remaining as the training set.

   💡 **Consistent Performance:** A model that shows consistent performance across all folds is less likely to be overfitting. Significant variations in performance across folds suggest that the model may be capturing noise and anomalies in specific parts of the data rather than generalizing from the actual signal.



3. **Learning Curves:** Plotting learning curves, which are graphs of the model's performance on the training and validation sets over time (or over the number of training instances), can provide visual insight into whether overfitting is occurring.
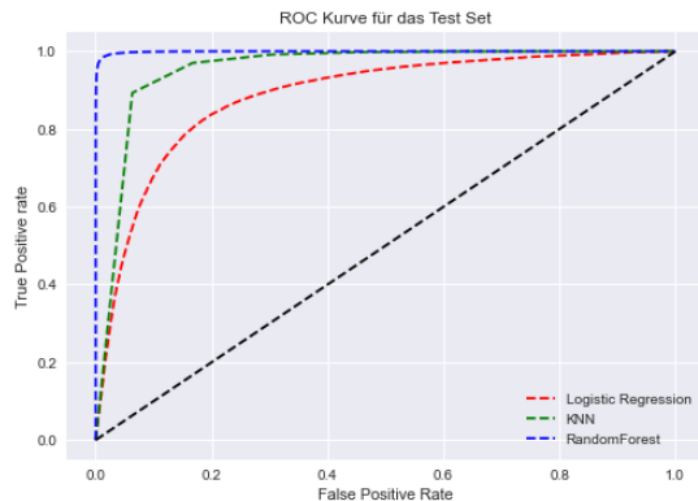
   💡 **Diverging Curves:** If the training curve continues to improve while the validation curve begins to degrade, it's likely that the model is overfitting.

4. **AUC-ROC Curve:** For classification problems, analyzing the area under the receiver operating characteristic curve (AUC-ROC) for both training and validation sets can indicate overfitting.

   💡 **Difference in AUC-ROC:** A high AUC-ROC on the training set but significantly lower on the validation set suggests overfitting.

```
ROC_AUC_SCORE of models on Test Set

Logistic Regression:  89
KNeighbors:  95
Random Forest:  100
```



5. **Adjustment with Complexity:** Performance metrics can be evaluated in response to changes in the complexity of the model. By adjusting the complexity (such as the depth of a decision tree, the number of layers in a neural network, or the degree of the polynomial in regression), and observing how performance metrics change, you can detect overfitting.

   💡 **Optimal Complexity:** The point at which increasing complexity does not significantly improve validation performance anymore, or starts to worsen it, typically indicates the onset of overfitting.

A significant performance gap between training and validation/test sets usually indicates overfitting.


### 1.2.2 Visualization

When analyzing data, it is always helpful to have a visual representation of the data. This makes the process much clearer and understandable. Visualization is a powerful tool for detecting overfitting in machine learning models, providing intuitive insights into how well a model is generalizing from its training data. Here are several visualization techniques that can help identify overfitting:

1. **Learning Curves:** These are plots of model performance on the training set and the validation set as a function of training epochs or the number of training samples.
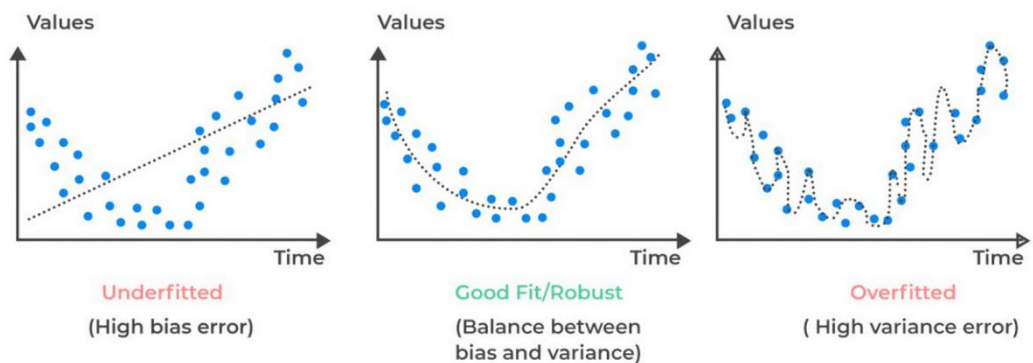
✓ **How They Help:** If a learning curve shows the training error decreasing and the validation error decreasing initially but then starting to increase, this is a classic sign of overfitting. The model is learning the training data too well, including the noise and outliers, which does not generalize to new data.

2. **Validation Curves:** These plots show the model's performance on the training and validation sets over a range of values for a model hyperparameter.

   ✓ **How They Help:** Validation curves can help you understand the relationship between model complexity and model performance. If increasing a hyperparameter (like the depth of a tree or the number of layers in a neural network) improves performance on the training data but worsens performance on validation data, overfitting might occur.
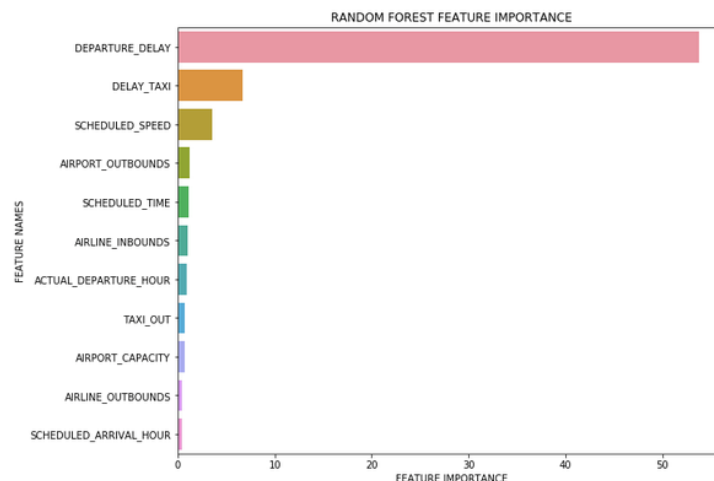
3. **Residual Plots:** These are plots that show the difference between the predicted and actual values. For a well-fitting model, the residuals should be randomly distributed around the centerline (zero residual).

   ✓ **How They Help:** If the residuals display patterns, such as systematic deviations from zero for certain types of inputs or growing dispersion with larger fitted values, it may indicate that the model is either overfitting or underfitting specific segments of the data.



4. **Feature Importance Plots:** These can be used to visualize the importance or contributions of different features to the model.

   ✓ **How They Help:** If minor features are showing disproportionately high importance, it might suggest that the model is picking up noise and overfitting. Simplifying the model by removing or combining features can sometimes help in reducing overfitting.
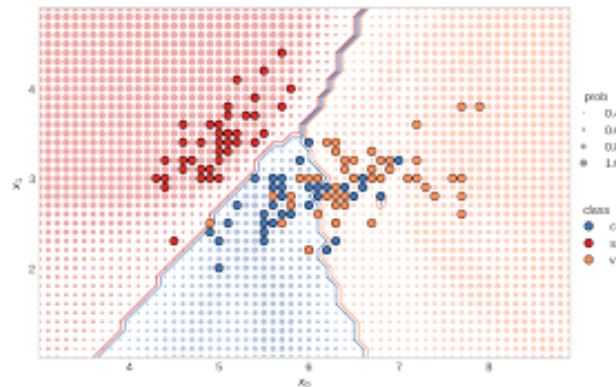
5. **ROC Curve Analysis:** For classification problems, plotting the receiver operating characteristic (ROC) curve and calculating the area under the curve (AUC) for both the training and validation datasets can provide insights.

  ✓ **How They Help:** A high AUC for the training dataset coupled with a much lower AUC for the validation dataset typically indicates overfitting.

6. **Decision Boundary Plots:** For classification problems, plotting the decision boundaries created by the model can provide visual insights into how the model is classifying different regions.

  ✓ **How They Help:** Overly complex decision boundaries that appear to be tailored to specific training examples rather than capturing the general trend suggest overfitting.



These visualization techniques not only help in detecting overfitting but also offer guidance on how parameters might be tuned to improve model performance. They make the abstract concept of model fitting more tangible and can guide non-technical stakeholders in understanding model behavior.

# 1.3 Preventing Overfitting

Since overfitting is undesired, we need to prevent that from happening. There are several ways which can be used to prevent overfitting:

A. **Simplifying the Model**
   Simplifying a machine learning model is a crucial approach to preventing overfitting, especially when a model is too complex relative to the size and diversity of the training data. Overfitting occurs when a model learns not just the underlying patterns in the data but also the noise and anomalies, which do not generalize well to new, unseen data.

   Why simplify a model?
   **Reduce Complexity:** A simpler model has fewer parameters or lower capacity and is less likely to fit the noise in the training data.
   **Improve Generalization:** By focusing on the most significant features and interactions, a simpler model tends to generalize better to new data from the same population.
   **Enhance Interpretability:** Simpler models are often easier to understand and interpret, which is valuable in many practical applications.

   Methods to simplify a model?

   - **Feature Selection:**
     **Manual Selection:** Based on domain knowledge, select only the most relevant features for use in the model.
     **Automated Techniques:** Use algorithms like backward elimination, forward selection, or recursive feature elimination to systematically remove less important features.
   - **Reducing Model Parameters:**
     **Pruning:** In tree-based models, reduce the complexity of the model by cutting back on the number of branches and depth of the tree.
     **Dimensionality Reduction:** Techniques like Principal Component Analysis (PCA) or Linear Discriminant Analysis (LDA) can reduce the number of input variables by transforming them into a new set of variables with reduced dimensions.
   - **Choosing Simpler Models:**
     **Prefer Linear Models:** For certain problems, linear models might be sufficient to capture the underlying patterns without fitting the noise.
     **Avoid High-Polynomial Features:** Higher degree polynomials can create complex decision boundaries which might lead to overfitting.
   - **Regularization:**
     **L1 Regularization (Lasso):** Adds a penalty equivalent to the absolute value of the magnitude of coefficients. This can lead to some coefficients being zero, effectively reducing the number of features.
     **L2 Regularization (Ridge):** Adds a penalty equivalent to the square of the magnitude of coefficients. This discourages large coefficients but does not set them to zero.
   - **Early Stopping:**
     **During Training:** Monitor the model's performance on a validation set and stop training when performance begins to degrade, despite improvements on the training set.
   - **Cross-validation:**
     **Validate Model Choices:** Use techniques like k-fold cross-validation to ensure that changes intended to simplify the model improve its performance on unseen data, rather than just the training set.

By strategically simplifying a machine learning model, we can enhance its ability to perform well not just on the training data but more importantly on new, unseen datasets, thereby achieving better and more reliable predictions by overcoming overfitting.
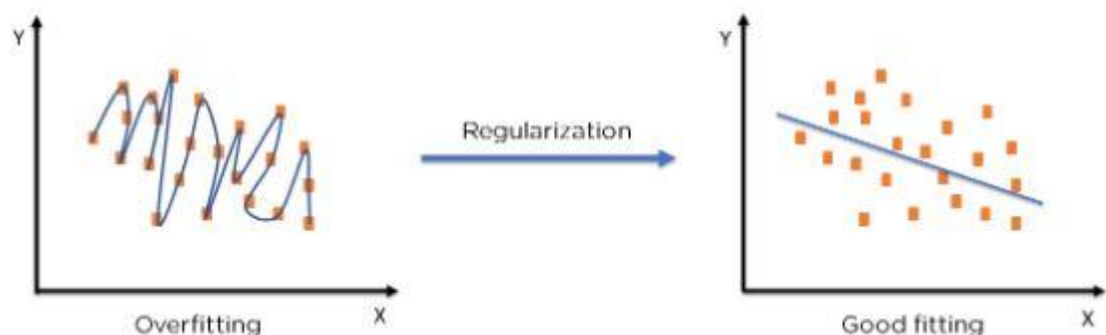
**B. Increasing the Training Data**

Increasing the amount of training data is a straightforward and often effective method to prevent overfitting in machine learning models. Here's how increasing the training data can help mitigate this problem:

- Improving Generalization: More data provides a more comprehensive representation of the underlying reality, thereby helping the model to generalize better. This means the model is less likely to be swayed by noise or outliers present in smaller datasets and can learn the true patterns more effectively.
- Enhanced Learning: With more examples, the model can learn more nuanced relationships without relying heavily on assumptions. For complex models, particularly deep learning networks that require large amounts of data, having more data points covers a broader range of possibilities and scenarios, reducing the model's need to interpolate between points or extrapolate from insufficient samples.
- Reduction of Variance: Increasing the training data reduces the variance part of the model's error. This means the model's predictions will be less sensitive to fluctuations in the training set, leading to more stable performance across different sets of data.
- Dilution of Noise: As the volume of data increases, the influence of noise (errors or irrelevant information) within the data is diluted. More data points make the real signal in the data clearer and more dominant over the noise, which a well-designed model can learn to ignore.

**C. Regularization**

Regularization methods are techniques used to reduce the model complexity and prevent overfitting, which enhances the model's generalizability. These methods work by adding a penalty term to the loss function during the training of the model. The main types of regularization techniques used for linear models are Lasso Regression and Ridge Regression.

1. **Lasso Regression (L1 Regularization)**

   Lasso Regression (Least Absolute Shrinkage and Selection Operator) is a type of linear regression that includes a regularization component. The technique adds a penalty equal to the absolute value of the magnitude of the coefficients to the loss function. This penalty is known as the L1 penalty. The key characteristic of Lasso Regression is that it can lead to sparse models with few coefficients; some coefficients can become exactly zero. This happens because the penalty effectively compresses some coefficients to zero while others remain non-zero, which is useful in feature selection in high-dimensional datasets.

$$L_{lasso}(\hat{\beta}) = \sum_{i=1}^{n}(y_i - x_i^T\hat{\beta})^2 + \lambda\sum_{j=1}^{m}|\hat{\beta}_j|$$

Where:

$y_i$ are the observed values.
$x_i$ are the feature values.
$\beta_j$ are the coefficients to be learned.
$\lambda$ is the regularization parameter that controls the strength of the penalty.

**How Lasso Helps Avoid Overfitting?**

**Regularization:** The L1 penalty encourages a model to choose simpler, sparser models over complex, dense models. By reducing the number of non-zero coefficients, Lasso helps in reducing model complexity, which inherently decreases the risk of overfitting.

**Feature Selection:** Lasso can be seen as an embedded feature selection method. It effectively reduces the number of features by setting the coefficients for less important variables to zero. This selection process directly targets the problem of overfitting by excluding irrelevant or less important features from the model.

**Enhancing Generalization:** By simplifying the model, Lasso enhances the model's ability to generalize from the training data to unseen data. This is because simpler models extrapolate better and are less likely to capture noise as signal.

**Optimization of the Regularization Parameter:** The parameter $\lambda$ provides a means to control overfitting. A larger value of $\lambda$ increases the penalty, leading to simpler models. However, setting $\lambda$ too high can lead to underfitting, where the model is too simple to capture the underlying trend. Therefore, finding the right balance for $\lambda$ is crucial, and it is usually done via cross-validation.

Lasso regression is particularly useful in scenarios where the number of predictors (features) is very high relative to the number of observations. However, Lasso might not perform well if there are highly correlated variables, as it tends to select one of the correlated terms and set others to zero. Also, if the number of predictors is greater than the number of observations, Lasso selects at most n predictors as non-zero, even if all predictors are relevant.

Overall, Lasso Regression is a powerful method for both regularization and feature selection, helping to improve the robustness and performance of regression models by mitigating overfitting.

2. **Ridge Regression (L2 Regularization)**
   Ridge Regression is a technique used for analyzing multiple regression data that suffer from multicollinearity. When multicollinearity occurs, least squares estimates are unbiased, but their variances are large, which might deviate the observed values far from the true values. By adding a degree of bias to the regression estimates, ridge regression reduces the standard errors.

$$SSE_{L_2} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{P} \beta_j^2$$

where:
$y_i$ are the observed values.
cap($y_i$) are the feature values.
$\beta_j$ are the coefficients to be learned.
$\lambda$ is the regularization parameter that controls the strength of the penalty.

**How Ridge Regression Helps Avoid Overfitting?**
**Shrinkage of Coefficients:** Unlike Lasso, which can reduce some coefficients completely to zero, Ridge Regression only shrinks the size of the coefficients. This shrinkage helps to reduce model complexity and multicollinearity, but keeps all variables in the model, thus preserving the information provided by all features.
**Bias-Variance Trade-Off:** By introducing the L2 penalty, Ridge Regression increases the bias of the model slightly (because it is a biased estimator), but significantly reduces the variance of the predictions. This trade-off typically leads to better performance on test data, as it generalizes better than an unregularized model.
**Handling Multicollinearity:** Ridge regression is particularly useful when dealing with multicollinearity (i.e., when independent variables are highly correlated). In the presence of highly correlated features, the least squares are unbiased but their variances are large which deviates the observed value far from the true value. Ridge regression stabilizes the regression estimates in such cases.
**Improvement in Prediction Accuracy:** While Ridge Regression does not reduce the number of variables as it does not lead to zero coefficients, it will still improve the prediction accuracy by a significant reduction in overfitting, especially in scenarios where the number of parameters exceeds the number of observations.
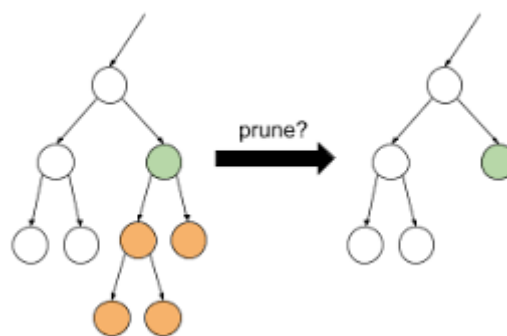**Tuning of $\lambda$:** The strength of the penalty, represented by $\lambda$, is a critical parameter in Ridge Regression. It can be tuned via cross-validation to find the model that best generalizes across unseen data.

Ridge Regression is broadly used in many fields where regression is applied, especially when the goal is to predict or infer from data with complex structures and many collinear inputs. It is less aggressive than Lasso in terms of variable reduction, maintaining all the variables but distributing the coefficient shrinkage across them. However, one of the limitations of Ridge Regression is that it includes all predictors in the final model, which may not be useful if the goal is feature elimination or selection.

Overall, Ridge Regression is a robust method for enhancing regression analysis, particularly in situations prone to overfitting or when the dataset features high multicollinearity.

## D. Pruning

Pruning is a technique used primarily in tree-based machine learning models, such as decision trees and ensemble methods (like gradient boosting trees and random forests), to reduce the size of the tree and prevent overfitting. The basic idea is to remove parts of the tree that provide little power in terms of predicting the target variable, effectively simplifying the model.



There are two types of Pruning, pre-pruning and post-pruning. Pre-pruning (Early Stopping) decisions might include stopping the tree growing when the nodes have fewer than a minimum number of points or when the depth reaches a predetermined level. The idea is to prevent the tree from becoming overly complex by limiting its growth preemptively. Post-pruning (Cost Complexity Pruning) technique involves first allowing the tree to grow to its full size and complexity and then simplifying it by removing nodes. The decision to remove a node is based on whether the removal results in a minimal increase in the overall model error on a validation set, often adjusted by a complexity parameter that penalizes a higher number of leaves.

**How Pruning Helps Prevent Overfitting?**
**Reduces Model Complexity:** Both pre-pruning and post-pruning reduce the complexity of the final model. They remove branches that have little importance, which might have been capturing noise and specific anomalies in the training data rather than underlying trends.
**Improves Generalization:** By reducing the complexity of the model, pruning helps to enhance the model's ability to generalize from the training data to unseen data. Overly

complex trees often perform well on training data but poorly on new, unseen data because they fit to quirks in the training set rather than to the underlying data structure. **Enhances Interpretability:** Simpler models are easier to understand and interpret. Pruned trees are generally more straightforward and make logical sense, as they focus on the major splits that have more statistical significance.

**Optimization of Error Trade-offs:** Post-pruning helps in optimizing the trade-off between the tree's depth and the model's predictive accuracy. It carefully evaluates whether extending or trimming parts of the tree will lead to better long-term performance.

While pruning is effective, it requires careful tuning to find the right balance between underfitting and overfitting. Over-pruning can lead to overly simplistic models that do not capture sufficient variability in the data (underfitting), whereas insufficient pruning might leave the model too complex, leading to overfitting. Choosing the right parameters for pruning, such as the minimum node size or the depth of the tree, often involves cross-validation to evaluate model performance across different potential configurations.

Pruning is a critical technique in managing the complexity of tree-based models, helping to ensure that they not only fit the training data well but also perform robustly on new, unseen data.

E. **Cross-Validation**

Cross-validation is a statistical method used to estimate the skill of machine learning models. It is primarily used to assess how the results of a statistical analysis will generalize to an independent data set. Essentially, it involves partitioning a sample of data into complementary subsets, performing the analysis on one subset (called the training set), and validating the analysis on the other subset (called the validation set or testing set).
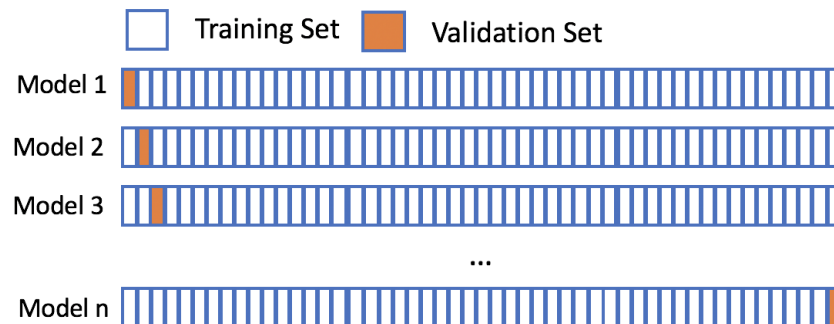
**What are the major two types of Cross-Validation?**

**k-Fold Cross-Validation:** The most common form of cross-validation. The data set is divided into k smaller sets (or folds). The model is trained on k-1 of these folds, with the remaining part used as the test set to compute a performance measure such as accuracy. This process is repeated k times, with each of the k folds used exactly once as the validation data.



4-fold validation (k=4)

**Leave-One-Out Cross-Validation (LOOCV):** A special case of k-fold cross-validation where the number of folds equals the number of instances in the dataset. This means that each fold contains exactly one data point. This method is computationally expensive but can provide detailed insights, especially for small datasets.



**How Cross-Validation Helps Prevent Overfitting?**

**Model Assessment:** Cross-validation provides a more accurate measure of a model's predictive performance. By using multiple train-test splits, it evaluates the model's ability to perform on unseen data, mitigating the risk of overfitting to the training set.

**Model Selection:** Cross-validation can be used to compare the performances of different predictive modeling procedures and choose the one that best generalizes to unseen data. For instance, it allows the comparison of different parameterizations of a model to find the best tuning parameters (like choosing the best degree for a polynomial regression).

**Reduces Bias:** Since each instance of the dataset gets to be in a training set and a validation set, it ensures that the model does not get biased towards a specific subset of data.

**Parameter Tuning:** It is often used in conjunction with grid search or random search to determine the optimal parameters for a model. This process involves running a model many times with different parameters and using cross-validation to determine which parameters perform the best.

**Improvement in Robustness:** The use of multiple training and validation sets reduces the variability of the model assessment, providing a more robust and stable estimate of model performance.

Cross-validation is a technique in machine learning for both assessing and improving the performance and generalizability of models. It is crucial for ensuring that a model not only fits the data well but also performs effectively and consistently on new, unseen data, thus helping to prevent overfitting.

## 1.4 Other Techniques
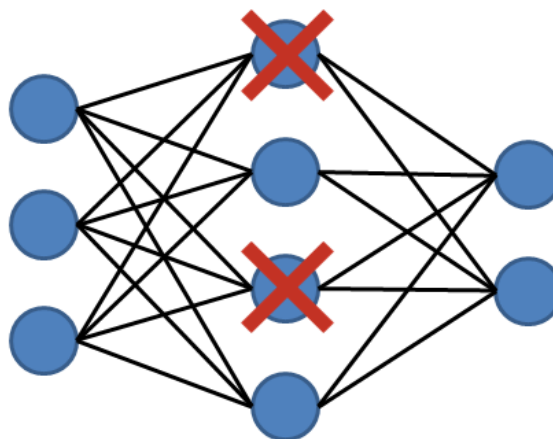
### 1.4.1 Dropout Technique

The dropout technique in neural networks is a regularization method used to prevent overfitting during training. Below are the steps involved in it:

**Random Deactivation:** During training, dropouts randomly deactivate (sets to zero) a fraction of the neurons in the network layers for each batch of data processed. The "dropout rate" specifies the percentage of neurons deactivated.

**Reduces Interdependency:** By deactivating different subsets of neurons, dropout reduces the reliance of any one neuron on the output of other specific neurons. This encourages the network to develop more robust features that are useful in conjunction with many different random subsets of the other neurons.

**Model Averaging:** Each pass through the training data uses a different "thinned" network. At test time, the whole network is used, but with the weights of the neurons scaled down proportionately to the dropout rate. This is akin to averaging the predictions of a large ensemble of different models.

By incorporating dropout technique, neural networks tend to learn more generalized representations of the data, rather than overfitting to the noise and details in the training dataset. This improves the model's performance on unseen data, making dropout particularly effective for complex networks prone to overfitting.
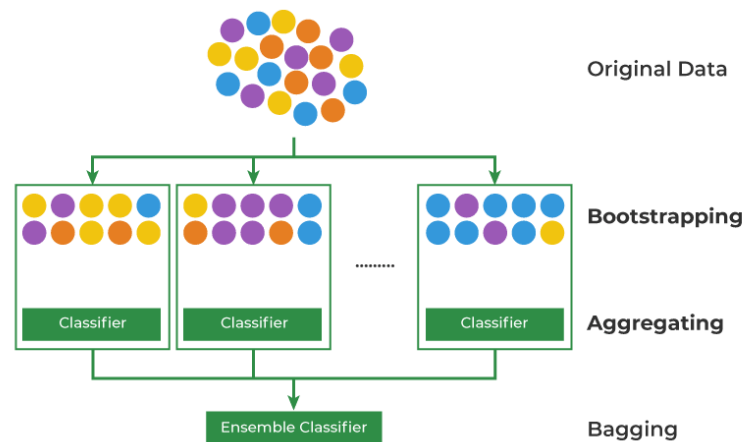


### 1.4.2 Ensemble Methods

Ensemble methods are powerful techniques in machine learning that combine multiple models to improve the robustness and accuracy of predictions. These methods can be particularly effective in reducing overfitting, enhancing stability, and improving prediction accuracy. There are two main Ensemble methods called Bagging and Boosting. Let's learn more about them.

**Bagging (Bootstrap Aggregating)**

Bagging involves training multiple models in parallel. Each model is trained on a random subset of the data, which is selected with replacement (bootstrap sample).

By using bootstrapped datasets, bagging helps in reducing the variance of the prediction model. Each individual model might overfit the data, but by averaging these models, the ensemble tends to cancel out individual errors and biases, leading to a more robust prediction.

Since it involves averaging predictions from multiple models, bagging can reduce overfitting, especially in high-variance models (like decision trees). It's less effective if the base models are biased.



## Boosting

Boosting builds models sequentially by focusing on training instances that previous models misclassified. Each new model in the sequence focuses more on the difficult cases that were handled incorrectly in previous rounds.

Boosting is effective in reducing bias (also variance to some extent), as it combines weak learners (models that do slightly better than random guessing) to create a strong overall model. By focusing on errors of preceding models, boosting seeks to produce a very accurate prediction rule by combining many weak and inaccurate rules.

While boosting can sometimes lead to overfitting if not carefully tuned, it generally performs well even with noise. The sequential nature and focus on correcting errors make it less likely to overfit compared to more simplistic models.

Both methods improve model generalization by combining the strengths of multiple models and mitigating their weaknesses. They are effective across a wide range of datasets and problems, from regression to classification. They can be used with most types of prediction models and are adaptable to any statistical learning approach that allows for iterative learning.

In summary, ensemble methods enhance prediction performance by assembling multiple models and typically result in more accurate and robust solutions than individual models alone. They can average out errors, focus on hard-to-classify instances, and reduce key problems like overfitting and variance.
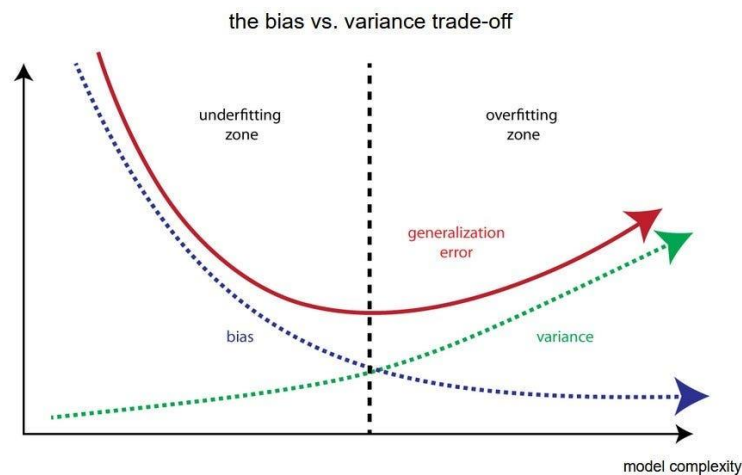


## 1.5 Bias-Variance Trade-off

The bias–variance dilemma or bias–variance problem is the conflict in trying to simultaneously minimize these two sources of error that prevent supervised learning algorithms from generalizing beyond their training set.

**Bias:** The bias error is an error from erroneous assumptions in the learning algorithm. High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting).

**Variance:** The variance is an error from sensitivity to small fluctuations in the training set. High variance may result from an algorithm modeling the random noise in the training data (overfitting).

**Tradeoff:** The tradeoff between bias and variance is challenging because typically, reducing one will increase the other. A model with very low bias must be complex enough to represent the data's underlying structure well, which can make it highly sensitive to fluctuations in the training set (high variance). Conversely, a model with very low variance might be too simple to capture the pattern in the data accurately (high bias).



the bias vs. variance trade-off

Overfitting is directly related to variance. When a model is overfitted, it is because it has a high variance. The model has learned not only the underlying pattern but also the noise in the training data. Such models perform well on the training data but do not generalize well to unseen data. The model essentially memorizes the data rather than learning to generalize from it.

To manage the bias-variance tradeoff,

- Use techniques like cross-validation to estimate how well their model is likely to perform on unseen data.
- Employ regularization techniques (like L1, L2 regularization) that can help to reduce model complexity and variance without substantial increase in bias.
- Choose model complexity based on the amount of data available; more data can support more complex models, potentially reducing bias without a corresponding increase in variance.

## 1.6 Overfitting in Supervised and Unsupervised Learning

Overfitting is a core challenge in both supervised and unsupervised learning, though its manifestations and impacts can differ between the two paradigms due to their fundamental differences in goals and methods.

In supervised learning, overfitting occurs when a model captures not only the underlying pattern in the training data but also the noise or random fluctuations. This results in a model that performs very well on training data but poorly on unseen data (test data).

Detection and Mitigation:

- Use techniques like cross-validation to detect overfitting by testing the model's performance on unseen data.
- Employ regularization techniques (like L1 and L2) to penalize overly complex models.
- Prune decision trees or use dropout in neural networks to reduce complexity.
- Early stopping during training when performance on a validation set starts to degrade.

Overfitting in unsupervised learning involves a model capturing detail that is due to noise or idiosyncrasies in the data rather than underlying structures. Since unsupervised learning doesn't use output labels, overfitting might not always be about predictive accuracy but could be about capturing spurious patterns that aren't generally useful.

Detection and Mitigation:

- Use techniques like silhouette scores to assess the appropriateness of clustering rather than just minimizing within-cluster variance which might lead to overfitting.
- For PCA, choosing the number of components to keep based on explained variance ratio, ensuring not to retain noise components.
- Applying regularization in methods like autoencoders to prevent them from merely learning to replicate the input data perfectly.

Common Challenges and Solutions:

- In both paradigms, the challenge is to discern between the signal (underlying patterns) and the noise (random fluctuations) in the data.
- Both supervised and unsupervised learning benefit from a better understanding of the data, including its source, nature of noise, and the kind of variability that is meaningful versus random.
- Selecting the right complexity of the model and tuning hyperparameters carefully based on validation metrics are universal strategies to mitigate overfitting.

## 1.7 Model Accuracy & Model Complexity in Overfitting

Model accuracy and model complexity are two critical aspects of machine learning that have a direct relationship with the phenomenon of overfitting. Understanding their interaction can help in designing better models that are both accurate and generalizable.

**Model Accuracy**

Model accuracy refers to how well a model performs on a given task; for predictive models, this is typically measured against some test data that was not used during the training process. Accuracy is determined by comparing the model's predictions against true outcomes. The goal in machine learning is often to maximize this accuracy, but achieving high accuracy on training data does not necessarily mean a model will perform well on new, unseen data.

## Model Complexity

Model complexity refers to the number of parameters in a model or the amount of change the model can accommodate with respect to its input data. A more complex model has a higher capacity to learn detailed information from the training data, including noise and outliers, which might not be relevant for generalizing to new data. Common examples of complex models include deep neural networks with many layers and parameters, and highly branched decision trees.

## Relationship with Overfitting

**Higher Complexity and Overfitting:** As model complexity increases, there's a higher risk of overfitting. Overfitting occurs when a model learns the details and noise in the training data to an extent that it negatively impacts the performance of the model on new data. Essentially, the model is too tailored to the training data, and thus unable to generalize. This is often seen with very deep neural networks or overly complex regression models where the model fits the training data almost perfectly but performs poorly on any external validation set.

**Optimal Model Complexity:** The challenge is to find a level of complexity that is just right for training data available and the noise level within it. This optimal point maximizes model accuracy on new, unseen data by balancing the bias-variance tradeoff:

**Low complexity models** may have high bias and low variance, possibly leading to underfitting, where the model is too simple to capture underlying patterns in the data.

**High complexity models** might have low bias and high variance, leading to overfitting, where the model captures the noise in the training data, mistaking it for a useful signal.

To mitigate overfitting while attempting to maintain or increase model accuracy, several strategies like Regularization, Pruning and Cross-validation are used.



A Book by Kartikey Vijayakumar Hebbar
Northeastern University, Boston – April 2024

# References

1. https://aws.amazon.com/what-is/overfitting/

2. https://chat.openai.com/

3. "Pattern Recognition and Machine Learning" by Christopher M. Bishop

4. "An Introduction to Statistical Learning" by Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani

5. https://scikit-learn.org/stable/user_guide.html

6. "The Elements of Statistical Learning" by Trevor Hastie, Robert Tibshirani, and Jerome Friedman

7. "Deep Learning" by Ian Goodfellow, Yoshua Bengio, and Aaron Courville

8. https://en.wikipedia.org/wiki/Bias%E2%80%93variance_tradeoff