# Universidad Politécnica de Yucatán

# IRC9B - Machine Learning

# Solution to most common problems in ML - Portfolio evidence

## Fredy Eduardo Alonzo Mondragon

## Professor: Victor Ortiz

## Date: September 15th, 2023

- **Define the concepts of: Overfitting and Underfitting.**

**Overfitting:**

Overfitting occurs when a machine learning model is too complex relative to the complexity of the underlying data. In other words, the model learns to fit not only the underlying patterns in the data but also the noise or random fluctuations present in the training data.

Signs of overfitting include a model that performs extremely well on the training data but poorly on new, unseen data. This suggests that the model has essentially memorized the training data rather than learning the underlying relationships.
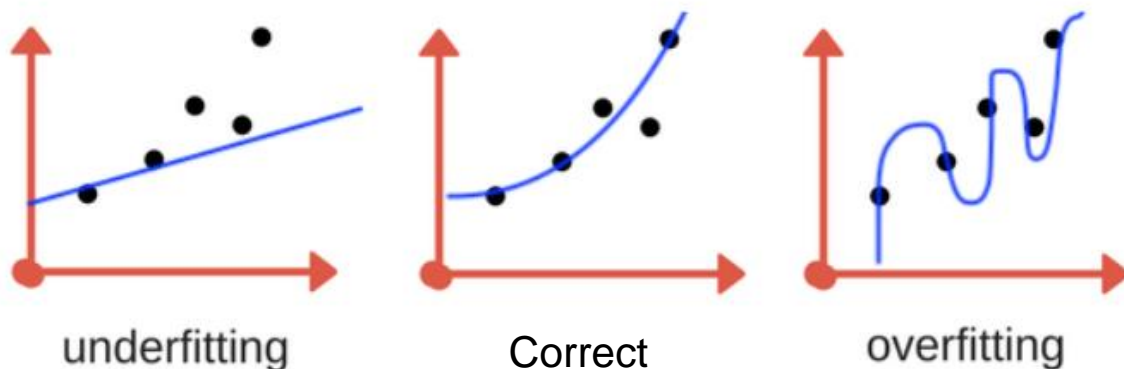
Overfitting can be thought of as "overly flexible" modeling, where the model becomes too sensitive to small variations in the training data, capturing noise rather than the actual patterns.

**Underfitting:**

Underfitting occurs when a machine learning model is too simple or lacks the capacity to capture the underlying patterns in the data. In this case, the model fails to learn the relationships between the input features and the target variable effectively.

Signs of underfitting include poor performance on both the training data and new, unseen data. The model's performance remains consistently low because it hasn't learned the essential patterns in the data.

Underfitting can be thought of as "overly rigid" modeling, where the model is too simplistic to capture the complexities present in the data.



underfitting     Correct     overfitting

- **Define and distinguish the characteristics of outliers.**

Outliers are data points that significantly deviate from most of the data in a dataset. They are observations that are either unusually high or low compared to the typical values in the dataset. Outliers can arise due to various reasons, including measurement errors, data entry errors, natural variability, or genuinely exceptional cases.

Extreme Values: Outliers are data points that have values that are significantly different from the values of most other data points in the dataset. They are typically located far from the central bulk of the data when visualized on a graph, such as a histogram or scatter plot.

Unusual Patterns: Outliers often exhibit unusual patterns or behaviors compared to the rest of the data. For example, in a scatter plot, an outlier may be an isolated point that doesn't follow the general trend or relationship observed in the data.

Influence on Statistics: Outliers can have a significant impact on summary statistics such as the mean and standard deviation. A single extreme outlier can greatly skew these statistics, potentially leading to incorrect interpretations of the data.

Context-Dependent: Whether a data point is considered an outlier can be context-dependent. In some cases, an extreme value may be a valid and important data point, while in other cases, it may be a clear error or anomaly.

Methods for Detection: Various statistical and graphical methods can be used to detect outliers, including the use of z-scores, the interquartile range (IQR), box plots, scatter plots, and machine learning-based anomaly detection techniques.

Implications: Outliers can have different implications depending on the context. They may be of interest for further investigation if they represent unique or critical cases. Alternatively, they may be removed or transformed to improve the performance of certain data analysis or modeling techniques.

Robustness: The impact of outliers on data analysis and modeling depends on the robustness of the chosen method. Some methods, like the median or robust statistical techniques, are less affected by outliers, while others, like linear regression, can be highly influenced by them.

- **Discuss the most common solutions for overfitting, underfitting and presence of outliers in datasets.**

**1. Overfitting:**

Overfitting occurs when a model learns the noise in the training data rather than the underlying patterns. To combat overfitting, you can:

Cross-validation: Use techniques like k-fold cross-validation to evaluate your model's performance on multiple subsets of your data. Cross-validation helps you understand how well your model generalizes to unseen data.

Regularization: Apply regularization techniques such as L1 (Lasso) or L2 (Ridge) regularization to penalize overly complex models. These techniques add a penalty term to the loss function, discouraging the model from fitting the noise in the data.

Simplify the model: Reduce the complexity of your model by reducing the number of features, reducing the model's capacity (e.g., decreasing the number of layers or nodes in a neural network), or using simpler algorithms. Simpler models are less prone to overfitting.

Feature selection: Choose the most relevant features for your model and discard irrelevant or noisy ones. Feature selection can help reduce the dimensionality of your data and improve model generalization.

Early stopping: Monitor the model's performance on a validation set during training and stop training when the performance starts to degrade. This prevents the model from overfitting the training data.

## 2. Underfitting:

Underfitting occurs when a model is too simple to capture the underlying patterns in the data. To address underfitting, you can:

Increase model complexity: If your model is too simple, consider increasing its complexity. This may involve adding more features, increasing the depth of a neural network, or using a more complex algorithm.

Feature engineering: Create more informative features by transforming or combining existing ones. Feature engineering can help the model better capture the relationships in the data.

Use a different algorithm: Experiment with different machine learning algorithms to find one that is better suited to your data. Some algorithms may be more capable of capturing complex patterns.

Collect more data: In some cases, collecting more data can help mitigate underfitting, especially if the current dataset is too small to capture the underlying patterns adequately.

**3. Presence of Outliers:**

Outliers can adversely affect model performance and data analysis. To handle outliers, you can:

Data preprocessing: Identify and handle outliers using techniques like the Z-score or IQR method. You can remove outliers, replace them with more typical values (e.g., using mean or median imputation), or transform the data to make it less sensitive to outliers (e.g., using log transformations).

Robust algorithms: Use algorithms that are less sensitive to outliers, such as robust regression techniques or decision tree-based models.

Feature engineering: Create new features that are less affected by outliers. For example, you can use rank-based features or percentiles instead of raw values.

Domain knowledge: Consult domain experts to understand the context of outliers. In some cases, outliers may represent important or rare events that should not be discarded.

- **Describe the dimensionality problem.**

The dimensionality problem, often referred to as the "curse of dimensionality," is a challenge that arises in various fields, including machine learning, statistics, and data analysis, when dealing with datasets that have a large number of features or dimensions. This problem encompasses several issues and difficulties that arise as the dimensionality of the data increases:

Increased Computational Complexity: As the number of dimensions in a dataset grows, the computational resources required for data processing, analysis, and modeling increase exponentially. Many algorithms become computationally infeasible or extremely slow when applied to high-dimensional data.

Data Sparsity: In high-dimensional spaces, data points become sparse, meaning that there are fewer data points available relative to the number of dimensions. Sparse data can make it challenging to estimate statistical properties and relationships accurately.

Overfitting: High-dimensional data increases the risk of overfitting in machine learning models. With many features, models can find spurious correlations or fit noise in the data, resulting in poor generalization to new, unseen data.

Reduced Model Interpretability: Models trained on high-dimensional data can become difficult to interpret because they involve many features. Understanding the

importance of individual features and their contributions to the model's predictions becomes more complex.

Increased Data Requirements: To build accurate models in high-dimensional spaces, you may need a large amount of data to adequately cover the feature space. Obtaining such large datasets can be challenging and costly.

Curse of Distance: In high-dimensional spaces, the concept of distance between data points becomes less meaningful. Due to the phenomenon known as "concentration of measure," most data points are at nearly the same distance from each other's center, making it difficult to distinguish between them.

Dimension Reduction: To mitigate the dimensionality problem, dimension reduction techniques like Principal Component Analysis (PCA) or feature selection methods can be employed. These methods aim to reduce the number of features while retaining the most important information.

Visualization Challenges: Visualizing high-dimensional data is challenging. While we can visualize data in two or three dimensions, representing data in higher-dimensional spaces requires specialized techniques such as dimensionality reduction for visualization purposes.

- **Describe the dimensionality reduction process.**

1. Data Preparation:

Begin by preparing your dataset, ensuring that it is cleaned and preprocessed. This may involve handling missing values, scaling or normalizing features, and encoding categorical variables.

2. Understanding the Data:

Before applying dimensionality reduction, it's essential to understand your data and the relationships between features. Exploratory data analysis (EDA) can help you gain insights into which features may be redundant or less informative.

## 3. Principal Component Analysis (PCA):

PCA is a linear dimensionality reduction technique that aims to find a new set of orthogonal (uncorrelated) features called principal components. These components are linear combinations of the original features and are ranked by their ability to capture the most variance in the data.

The PCA process involves the following steps:

Centering the data: Subtract the mean of each feature from the data to have a zero mean.

Computing the covariance matrix of the centered data.

Calculating the eigenvalues and eigenvectors of the covariance matrix.

Selecting a subset of the top-ranked eigenvectors (principal components) to represent the data. The number of components chosen determines the reduced dimensionality.

## 4. Feature Selection:

Feature selection involves choosing a subset of the most relevant features from the original set. Various methods can be used for feature selection, including filter methods, wrapper methods, and embedded methods. Some common techniques include mutual information, correlation analysis, and recursive feature elimination (RFE).

Feature selection methods can be based on statistical tests, machine learning algorithms, or domain knowledge.

## 5. Evaluation:

After reducing the dimensionality, it's important to evaluate the impact on your analysis or machine learning task. This may involve comparing model performance (e.g., accuracy, F1 score) before and after dimensionality reduction.

Visualizations can also help assess the separability of data points in the reduced-dimensional space.

6. Model Building and Analysis:

Once dimensionality reduction is applied, you can build your machine learning models or conduct data analysis using the reduced feature set. This often leads to improved model training times, reduced overfitting, and sometimes better generalization to new data.

7. Interpretation:

If using PCA, you can interpret the principal components to understand which original features contribute most to each component. This can provide insights into the underlying patterns in the data.

8. Fine-Tuning: Depending on the results, you may iterate and fine-tune the dimensionality reduction process, adjusting the number of retained components or selected features based on the specific requirements of your analysis or modeling task.

## • **Explain the bias-variance trade-off.**

The bias-variance trade-off is a fundamental concept in machine learning and statistical modeling that refers to the balance between two types of errors that a model can make when learning from data: bias and variance. Achieving an optimal trade-off between bias and variance is crucial for building models that generalize well to unseen data. Let's explore these concepts in detail:

Bias (Underfitting):

Bias refers to the error introduced by approximating a real-world problem that may be complex by a too-simplistic model. A model with high bias tends to oversimplify the underlying relationships in the data.

High bias models typically have low complexity and make strong assumptions about the data. These models may struggle to capture intricate patterns, leading to systematic errors across different datasets.

Models with high bias are said to underfit the data because they are not flexible enough to fit even the training data well.

Variance (Overfitting):

Variance refers to the error introduced by the model's sensitivity to small fluctuations or noise in the training data. A model with high variance is overly complex and captures random noise as genuine patterns.

High variance models are often highly flexible and can adapt to even noisy or irregular data points. However, they are prone to fitting noise and exhibit poor performance on new, unseen data.

Models with high variance are said to overfit the data because they fit the training data very closely but fail to generalize to new data.

The bias-variance trade-off can be summarized as follows:

Low Complexity Models (High Bias): These models are simple and make strong assumptions about the data. They tend to underfit, meaning they have high bias but low variance. They may not capture the underlying patterns well.

High Complexity Models (High Variance): These models are complex and have a high capacity to fit the training data. They tend to overfit, meaning they have low bias but high variance. They capture noise and may not generalize well to new data.

Achieving the right balance between bias and variance is the key to building models that generalize effectively:

Optimal Model: The goal is to find a model that achieves a reasonable trade-off between bias and variance, often referred to as the "Goldilocks zone." This model captures the essential patterns in the data while avoiding overfitting or underfitting.

Model Selection and Hyperparameter Tuning: The choice of algorithms, model complexity, and hyperparameters plays a critical role in finding the right balance. Techniques like cross-validation can help identify suitable models and hyperparameters.

Ensemble Methods: Ensemble methods, such as bagging (e.g., Random Forests) and boosting (e.g., Gradient Boosting), combine multiple models to mitigate the bias-variance trade-off. They reduce overfitting and improve generalization.

- **References:**

[1] "Qué es overfitting y underfitting y cómo solucionarlo". Aprende Machine Learning. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://www.aprendemachinelearning.com/que-es-overfitting-y-underfitting-y-como-solucionarlo/

[2] D. Lemonaki. "What is an outlier? Definition and how to find outliers in statistics". freeCodeCamp.org. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://www.freecodecamp.org/news/what-is-an-outlier-definition-and-how-to-find-outliers-in-statistics/#:~:text=What%20is%20an%20Outlier%20in,dataset%20you're%20working%20with.

[3] "A step-by-step explanation of principal component analysis (PCA)". Built In. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://builtin.com/data-science/step-step-explanation-principal-component-analysis

[4] "What is dimensionality reduction - techniques, methods, components - dataflair". DataFlair. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://data-flair.training/blogs/dimensionality-reduction-tutorial/

[5] "Overfitting and underfitting with machine learning algorithms - machinelearningmastery.com". MachineLearningMastery.com. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://machinelearningmastery.com/overfitting-and-underfitting-with-machine-learning-algorithms/

[6] "What is the curse of dimensionality?" Built In. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://builtin.com/data-science/curse-dimensionality

[7] "Lecture 12: Bias variance tradeoff". Home | Department of Computer Science. Accedido el 15 de septiembre de 2023. [En línea]. Disponible: https://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/lecturenote12.html