# Group Discussion: Model Selection

## Objective:

This activity aims to deepen your understanding of key concepts in model selection, such as performance trade-offs, validation techniques, and model complexity. You will work in groups to analyse real-world scenarios and explore how various factors influence the choice of an optimal model.

## Discussion Points for Each Group:

- **Outcome Discussion**: What does the outcome imply about the model's performance and suitability?
- **Influencing Factors**: What factors influenced the model selection? Why was one model preferred over another?
- **Concept Impact**: How do overfitting, generalization, computational cost, and operational requirements affect the decision?
- **Trade-off Evaluation**: What compromises are made in terms of accuracy, speed, and usability when choosing one model over another?

## Expected Outcomes:

- You will better understand how to make informed decisions about model selection based on performance, complexity, and application-specific needs.
- You will learn to articulate the reasoning behind choosing specific machine learning models and strategies in varied real-world contexts.

# Examples to illustrate key concepts in Model Selection

Each of these examples illustrates different aspects of model selection, emphasizing how different scenarios might require different priorities and considerations.

## 1. Performance Trade-offs

### Scenario: Predicting Customer Churn

A telecommunications company wants to predict which customers are likely to churn based on their usage patterns, demographics, and customer service interactions. The data science team builds two models:

- A complex deep neural network (DNN) with multiple layers and parameters.

- A simpler logistic regression model.

### Outcome:

- The DNN performs exceptionally well on the training dataset, achieving near-perfect accuracy.

- However, when deployed on new, unseen data, the DNN's performance drops significantly, indicating overfitting.

- The logistic regression model, while slightly underperforming on the training data compared to the DNN, shows much better generalization on new data.

## Possible solutions

**Discuss the Outcome of the Scenario:**

- The deep neural network (DNN) excelled on the training data but performed poorly on new, unseen data, suggesting it overfitted the training data.

- The logistic regression model, while less accurate on the training data, demonstrated better generalization on new data.

- This scenario highlights the importance of considering how well a model can generalize beyond the training data. The simpler logistic regression model is ultimately chosen for its better performance on unseen data, illustrating the trade-off between complexity and practical usability.

**What Does the Outcome Imply About the Model's Performance and Suitability?**

- The DNN, despite its high accuracy on training data, is unsuitable for deployment due to poor generalization.

- The logistic regression model is more suitable as it provides a reliable performance on unseen data, which is critical for real-world applications.

**Identify Key Factors Influencing the Model Selection:**

- **Generalization:** The ability of the logistic regression model to perform well on unseen data made it preferable.

- **Overfitting:** The DNN's overfitting made it less desirable despite its high initial accuracy.

**How Do Concepts Affect the Decision?**

- **Overfitting and Generalization:** Directly influenced the decision against the DNN due to its inability to generalize beyond the training dataset.

- **Computational Cost and Operational Requirements:** Logistic regression, being simpler and faster, also likely required less computational resources, aligning better with operational efficiency.

**Evaluate the Trade-offs Involved:**

- **Compromises Made:** Accuracy was slightly compromised with the logistic regression model in favour of better generalization and lower computational cost. The trade-off was between top-notch training accuracy (DNN) and practical usability and reliability (logistic regression).

# 2. Validation Techniques

**Scenario: Developing a Diagnostic Tool**

A healthcare company is developing a diagnostic tool to detect a specific type of cancer from patient test results. The team decides to use k-fold cross-validation to evaluate the performance of three different models: a support vector machine (SVM), a random forest classifier, and a gradient boosting machine (GBM).

**Procedure:**

- The dataset is divided into 'k' subsets. Each model is trained on 'k-1' subsets and tested on the remaining subset, repeated 'k' times with each subset used as the test set once.

- Each model's accuracy, sensitivity, and specificity are recorded across all 'k' folds.

**Outcome:**

- The SVM shows high variance in performance across different folds.

- The random forest provides the best balance of sensitivity and specificity across all folds, with lower variance compared to the SVM.

- The GBM shows good performance but at a higher computational cost.

## Possible solutions

**Discuss the Outcome of the Scenario:**

- SVM displayed high variance across different data folds, indicating possible instability in its predictions.

- The random forest achieved the best balance of sensitivity and specificity with consistent performance, while GBM also performed well but was computationally more demanding.

- Using k-fold cross-validation helps in assessing not just the overall performance but also the stability of each model across different subsets of data. The random forest is selected for its consistent and reliable performance across folds.

**What Does the Outcome Imply About the Model's Performance and Suitability?**

- The random forest's consistent performance and lower variance make it most suitable for medical diagnostic tools where reliability is paramount.

**Identify Key Factors Influencing the Model Selection:**

- **Stability and Consistency:** Random forest's lower variance in performance across folds influenced its selection.

- **Sensitivity and Specificity:** Important metrics in medical diagnostics, where both the ability to detect positives (sensitivity) and true negatives (specificity) are critical.

**How Do Concepts Affect the Decision?**

- **Validation Technique (k-fold cross-validation):** Helped in identifying which model performs consistently under different subsets of data, favouring the random forest.

- **Computational Cost:** Affected the decision against GBM, which, although powerful, demanded more resources.

**Evaluate the Trade-offs Involved:**

- **Compromises Made:** In choosing the random forest over the GBM, computational efficiency and model stability were prioritized over potentially higher performance that GBM might offer on a larger computational budget.

# 3. Model Complexity

**Scenario: Real-Time Fraud Detection**

A financial institution implements a real-time system to detect fraudulent transactions. The team explores several models ranging from simple logistic regression to complex ensemble methods and neural networks.

**Considerations:**

- The simpler logistic regression model can be implemented and run with lower latency, which is critical in real-time applications.

- More complex models like ensemble methods and neural networks offer higher accuracy but require significantly more computational resources and result in slower response times.

**Outcome:**

- The logistic regression model, while slightly less accurate, provides sufficiently good performance and can evaluate transactions in real-time without causing delays.

- The complex models, despite their higher accuracy, are deemed unsuitable for real-time processing due to their computational demands.

# Possible solutions

**Discuss the Outcome of the Scenario:**

- The logistic regression model, despite being less accurate, could process transactions in real-time without delays.

- Complex models like ensemble methods and neural networks, while more accurate, were unsuitable for real-time processing due to their computational demands.

- This scenario underscores the necessity of balancing model complexity with operational requirements, especially in contexts where performance and speed are critical. The choice of logistic regression aligns with the need for fast, efficient processing crucial for real-time fraud detection.

**What Does the Outcome Imply About the Model's Performance and Suitability?**

- Logistic regression's suitability for real-time applications makes it the preferred choice despite slightly lower accuracy.

**Identify Key Factors Influencing the Model Selection:**

- **Operational Requirements:** The need for real-time processing heavily influenced the selection of a simpler model.

- **Model Complexity:** Directly impacted the decision, favouring a model that efficiently balances speed and accuracy.

**How Do Concepts Affect the Decision?**

- **Model Complexity and Operational Requirements:** Were the key drivers, prioritizing a model that meets the real-time operational demands over one that, while more accurate, fails in speed and efficiency.

**Evaluate the Trade-offs Involved:**

- **Compromises Made:** Accuracy was compromised for speed and the ability to implement the model in a real-time environment. The decision balanced the need for efficiency and sufficient accuracy, which is critical in fraud detection systems to prevent lag in transaction processing.