# Introduction to K-Fold Cross-Validation

## A Model Evaluation Technique

Jean-Michel Marin

September 19, 2024

# What is Cross-Validation?

- Cross-validation is a statistical method used to estimate the performance of a statistical model
- It helps in assessing how well a model will generalize to an independent dataset

# Why Do We Need Cross-Validation?

- **Avoid Overfitting** Helps in detecting when a model performs well on training data but poorly on unseen data
- Provides a better estimate of model performance compared to train-test split, which might not capture variability

# What is K-Fold Cross-Validation?

- A type of cross-validation where data is divided into **K subsets** (or "folds")
- The model is trained on **K-1 folds** and tested on the remaining fold
- The process is repeated **K times**, with each fold used once as a test set
- The final performance is the **average** of all K trials

# How Does K-Fold Work?

1. Split the dataset randomly into K equal parts (folds)
2. Train the model on K-1 parts and test it on the remaining part
3. Repeat this process K times, each time using a different part as the test set
4. Calculate the average performance score (accuracy, precision, etc.) across all K trials

# Choosing K in K-Fold

- Commonly used values for K are 5 or 10
- A small K value (e.g., 5) reduces computation time but might lead to higher bias
- A large K value (e.g., 10) reduces bias but may increase variance and computational cost

# Advantages and Disadvantages

**Advantages**

- ▶ Provides a better estimate of model performance
- ▶ Helps in selecting the best model by comparing different models
- ▶ More efficient use of data compared to a single train-test split

**Disadvantages**

- ▶ Computationally expensive for large datasets
- ▶ Can lead to higher variance in performance metrics

# Special Case: Stratified K-Fold

▶ Stratified K-Fold is a variation where each fold maintains the same proportion of classes as the original dataset

▶ Useful for imbalanced datasets to ensure that each fold represents the data distribution

# Conclusion

- ▶ K-Fold Cross-Validation is a powerful technique to assess model generalization
- ▶ It reduces overfitting and provides a more reliable performance estimate
- ▶ Choosing the right K value and considering computational cost is important