

Comprehensive Note on Underfitting

1. Definition of Underfitting

Underfitting occurs when a machine learning model is too simple to capture the underlying patterns or structure of the data. It fails to learn the relationships in the training data adequately, resulting in poor performance on both the training and unseen (test) datasets. This typically happens when the model has high bias.

2. Causes of Underfitting

- **Overly Simplistic Model:** Using a model that lacks the complexity to represent the data, such as a linear model for a non-linear problem.
- **Insufficient Training Time:** Training a model for too few epochs can prevent it from learning the data patterns fully.
- **Inadequate Features:** When the input features are insufficient or do not capture the relevant information needed for prediction.
- **High Regularization:** Overuse of regularization techniques can excessively constrain the model, leading to underfitting.
- **Improper Data Preparation:** Poor preprocessing, such as not normalizing or encoding features properly, can result in underfitting.

3. Signs of Underfitting

- **High Training Error:** The model performs poorly even on the training data.
- **Similar Training and Validation Errors:** Both training and validation errors are high and close in value, indicating that the model has not learned enough.
- **Flat Learning Curve:** The training and validation loss curves do not decrease significantly with time or iterations.

4. Examples

- Using a linear regression model to predict data that follows a quadratic relationship.
- Training a shallow decision tree on a complex dataset.

5. Techniques to Address Underfitting

1. **Increase Model Complexity:**
 - Use more complex models, such as increasing the degree of a polynomial or adding layers to a neural network.
 - Select algorithms better suited to the problem (e.g., switching from linear regression to a tree-based method).
2. **Increase Training Time:**

- Train the model for more epochs, ensuring it has sufficient time to learn the data patterns.
3. **Add Relevant Features:**
 - Include additional informative features or engineer new ones to better capture the relationships in the data.
 - Use domain knowledge to identify missing features that are critical for prediction.
 4. **Reduce Regularization:**
 - Lower the regularization parameters (e.g., reduce λ in L1/L2 regularization) to allow the model more flexibility.
 5. **Optimize Hyperparameters:**
 - Fine-tune hyperparameters like learning rate, depth of trees, or number of layers to achieve better performance.
 6. **Use Appropriate Data Preparation:**
 - Ensure the data is preprocessed correctly, such as normalizing numerical features and encoding categorical variables.
 - Address issues like missing data and scaling discrepancies.
 7. **Increase Training Data Quality:**
 - Use cleaner, more representative datasets to help the model learn the true relationships.
 - Remove irrelevant or redundant data that might mislead the model.

6. Measuring Underfitting

- **Learning Curves:** Underfitting is indicated by high and flat training and validation loss curves.
- **Evaluation Metrics:** Metrics like accuracy, F1-score, or RMSE will remain poor on both the training and validation datasets.

7. Balancing Underfitting and Overfitting

Achieving optimal model performance involves striking a balance between underfitting (high bias) and overfitting (high variance). This balance is often referred to as the **bias-variance tradeoff**:

- **Bias:** Represents the error introduced by approximating a complex problem with a simple model (underfitting).
- **Variance:** Represents the sensitivity of the model to small fluctuations in the training data (overfitting).

8. Conclusion

Underfitting limits the ability of a model to make accurate predictions, as it fails to learn enough from the training data. By increasing model complexity, extending training, and optimizing features and hyperparameters, underfitting can be addressed effectively. Regular monitoring using validation datasets and

diagnostic tools like learning curves is essential to build a model that generalizes well.