

# Overfitting and Underfitting

David Robinson

## Underfitting

**Underfitting** occurs when a machine learning model is too simple to capture the underlying patterns in the data. This leads to poor performance both on the training and testing data.

### Causes of Underfitting

- **Model is too simple:** The model lacks complexity to capture data relationships
- **Inadequate features:** Input features do not adequately represent the factors influencing the target variable
- **Small training dataset:** Insufficient data may prevent the model from learning key patterns
- **Excessive regularization:** Too much regularization restricts the model's ability to learn effectively.

### Solutions to Underfitting

- **Increase model complexity:** Use more advanced algorithms to capture data complexities
- **Add more features:** Perform feature engineering to improve representation
- **Remove noise:** Clean the dataset to enhance model accuracy
- **Train longer:** Increase the number of epochs or training time to allow the model to better learn from the data

## Bias in Machine Learning

**Bias** is the error introduced by overly simplistic assumptions made by the learning algorithm. A model with high bias is too simplified, failing to represent the relationship between input and output accurately and usually leads to underfitting.

## Overfitting

**Overfitting** occurs when a model learns not only the underlying patterns in the training data but also the noise and random fluctuations. This leads to poor generalization to unseen data even if performance is good on the training data.

### Causes of Overfitting

- **Model complexity:** The model is too complex, capturing noise and irrelevant patterns in the training data
- **Insufficient data:** Small training datasets can lead to models that overfit due to learning irrelevant details
- **Excessive training:** Training for too many epochs can cause the model to fit the data too closely

## Solutions to Overfitting

- **Increase training data:** Provide more examples to help the model learn better generalizations
- **Reduce model complexity:** Use simpler models or limit the parameters
- **Regularization techniques:** Lasso (L1) or Ridge (L2) regularization
- **Early stopping:** Stop training when the performance on validation data begins to degrade
- **Dropout:** Randomly drop neurons during training to avoid over-reliance on specific features

## Variance in Machine Learning

**Variance** is the model's sensitivity to small changes or fluctuations in the training data. A model with high variance are typically very complex, capturing random noise in the training data and usually leads to overfitting.

## Regularization

**Regularization** is used to reduce overfitting by penalizing overly complex models and encouraging simpler, more generalizable patterns, striking a balance between bias and variance.

### Common Techniques

- **Lasso (L1):** Adds a penalty proportional to the absolute value of coefficients, encouraging sparsity

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m |\theta_i|$$

- **Ridge (L2):** Adds a penalty proportional to the square of the coefficients, shrinking them

$$\text{Cost} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m \theta_i^2$$

- **Elastic Net (L1 + L2):** Combines both Lasso and Ridge regularization, balancing between shrinking and sparsity