

Underfitting and Overfitting in Machine Learning Models

Pournami P N

October 22, 2024

Underfitting and Overfitting

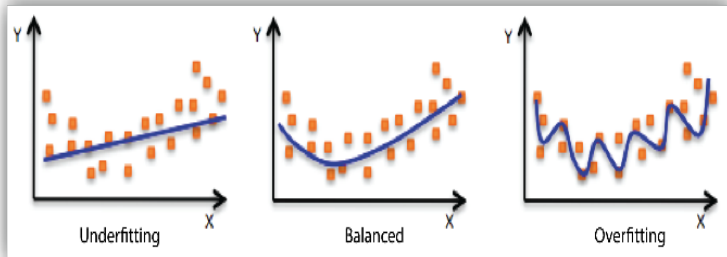
- ▶ **Underfitting** occurs when a model is too simple to capture the underlying structure of the data.
- ▶ The model has high bias and performs poorly on both training and test data.
- ▶ The input features which are used to train the model are not the adequate representations of underlying factors influencing the target variable.
- ▶ The size of the training dataset used is not enough.
- ▶ Example: Linear regression for a highly non-linear problem.
- ▶ **Overfitting** occurs when a model is too complex and fits the noise in the training data.
- ▶ The model has high variance, performs well on the training data but poorly on unseen test data.
- ▶ Example: A high-degree polynomial that fits all points exactly.

Bias and Variance

Bias is the difference between the average prediction of our model and the correct value which we are trying to predict. Model with high bias pays very little attention to the training data and oversimplifies the model. It always leads to high error on training and test data.

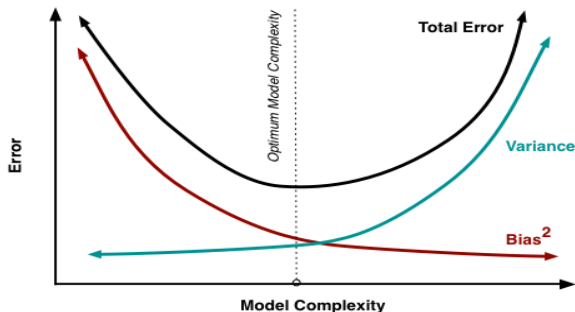
Variance is the variability of model prediction for a given data point or a value which tells us spread of our data. Model with high variance pays a lot of attention to training data and does not generalize on the data which it hasn't seen before. As a result, such models perform very well on training data but has high error rates on test data.

Underfitting and Overfitting



Bias Variance Trade-off

- ▶ If our model is too simple and has very few parameters then it may have high bias and low variance.
- ▶ On the other hand if our model has large number of parameters then it's going to have high variance and low bias.
- ▶ To build a good model, we need to find a good balance between bias and variance such that it minimizes the total error.



Techniques to Avoid Underfitting

- ▶ Increase model complexity.
- ▶ Increase the number of features, performing feature engineering.
- ▶ Remove noise from the data.
- ▶ Increase the number of epochs or increase the duration of training to get better results.

Techniques to Avoid Overfitting

- ▶ Regularization
- ▶ Early Stopping
- ▶ Dataset Augmentation
- ▶ Dropout
- ▶ Batch Normalization

Regularization

- ▶ Regularization introduces a penalty for larger coefficients in the model.
- ▶ Common techniques:
 - ▶ **L1 Regularization (Lasso)**: Encourages sparsity by penalizing the sum of the absolute values of the weights.
 - ▶ **L2 Regularization (Ridge)**: Penalizes the sum of the squared weights, keeping the weights small.
- ▶ Helps the model generalize better by avoiding large weights.

$$\text{Loss} = \text{MSE} + \lambda \sum_i w_i^2 \quad (\text{for L2})$$

Early Stopping

- ▶ **Early stopping** is used to halt training when the model's performance on validation data starts to degrade.
- ▶ It prevents the model from overfitting by stopping at the point of best generalization.
- ▶ Useful in training neural networks, where monitoring validation loss can prevent unnecessary extra epochs.

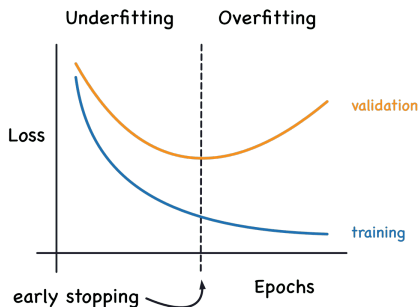


Figure: Early Stopping: Stop when validation error increases

Dataset Augmentation

- ▶ **Data Augmentation** artificially increases the size of the training data by applying transformations to the dataset.
- ▶ Common techniques include - Flipping, rotating, cropping, and zooming (for images), Adding noise (for audio, images).
- ▶ Helps prevent overfitting by providing more varied training examples.

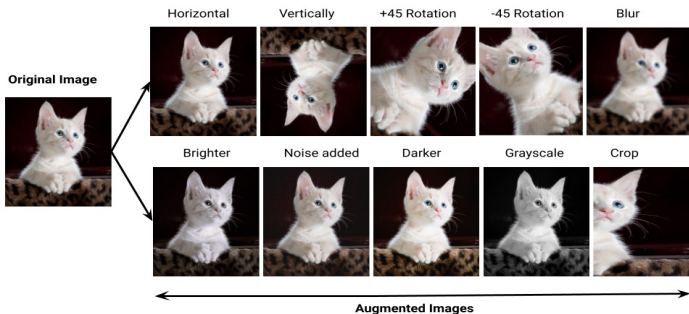


Figure: Data Augmentation: Creating variety in the dataset

Dropout

- ▶ **Dropout** is a regularization technique used in neural networks.
- ▶ During training, randomly drop a fraction of neurons (typically 20% to 50%) from each layer.
- ▶ Prevents co-adaptation of neurons and forces the network to generalize better.

Batch Normalization

- ▶ **Batch Normalization** normalizes the activations in each layer during training, stabilizing the learning process.
- ▶ Helps reduce internal covariate shift, allowing the use of higher learning rates and improving model generalization.
- ▶ It also acts as a regularizer, similar to dropout.

$$\hat{x} = \frac{x - \mu}{\sigma}$$

Conclusion

- ▶ Overfitting and underfitting are common problems in machine learning models.
- ▶ Techniques like regularization, early stopping, dataset augmentation, dropout, and batch normalization help in improving the generalization of models.
- ▶ Proper tuning and careful model selection are essential to achieving the right balance.