

Bias variance trade-off

Monday, 9 February 2026 12:29 PM

Bias-Variance Tradeoff

In supervised learning, the goal of a model is to learn an underlying true function $f(x)$ from data and make accurate predictions on unseen samples. However, prediction error arises due to multiple sources, primarily bias and variance.

Error Decomposition

For a regression problem with squared loss, the expected prediction error at a point x can be decomposed as:

$$\mathbb{E}[(y - f(x))^2] = \frac{\text{Bias}^2}{\text{systematic error}} + \frac{\text{Variance}}{\text{sensitivity to data}} + \frac{\sigma^2}{\text{irreducible noise}}$$

where σ^2 represents noise inherent in the data that no model can eliminate.

Bias

Bias measures the error introduced by approximating a real-world problem with a simplified model.

$$\text{Bias}(x) = \mathbb{E}[f(x)] - f(x)$$

Characteristics

- High bias models make strong assumptions
- Often too simple to capture complex patterns
- Lends to underfitting
- Example: linear regression on highly nonlinear data

Variance

Variance measures how much the model's predictions change when trained on different datasets.

$$\text{Variance}(x) = \mathbb{V}[f(x) - \mathbb{E}[f(x)]]^2$$

Characteristics

- High variance models are overly complex
- Highly sensitive to training data
- Lends to overfitting
- Example: high-degree polynomial fitting noisy data

Bias-Variance Tradeoff

Bias and variance have an inverse relationship:

- Increasing model complexity \downarrow bias \uparrow variance
- Decreasing model complexity \uparrow bias \downarrow variance

The objective is not to minimize bias or variance individually, but to find an optimal balance that minimizes total generalization error.

Regularization

Regularization is a technique used to control model complexity by adding a penalty term to the loss function. It discourages overly large model parameters, helping reduce variance and improve generalization.

Ordinary Least Squares (OLS)

$$\mathcal{L}_{OLS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

OLS minimizes training error only, which can lead to overfitting.

L2 Regularization (Ridge Regression)

$$\mathcal{L}_{Ridge} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p w_j^2$$

Effect

- Penalizes squared magnitude of weights
- Shrinks coefficients smoothly toward zero
- Reduces variance without eliminating features
- Increases bias slightly to lower total error

$$w_1 + w_2 + w_3$$

tuning knobs

L1 Regularization (Lasso Regression)

$$\mathcal{L}_{Lasso} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |w_j|$$

Effect

- Encourages sparse solutions
- Forces some weights exactly to zero
- Performs feature selection
- Can increase bias more than Ridge

$$\text{cost func}$$

Elastic Net Regularization

$$\mathcal{L}_{ElasticNet} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (\alpha \sum_{j=1}^p |w_j| + (1 - \alpha) \sum_{j=1}^p w_j^2)$$

Elastic Net Regularization

$$\mathcal{L}_{ElasticNet} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (\alpha \sum_{j=1}^p |w_j| + (1 - \alpha) \sum_{j=1}^p w_j^2)$$

Effect

- Combines strengths of L1 and L2
- Handles correlated features well
- Balances sparsity and stability

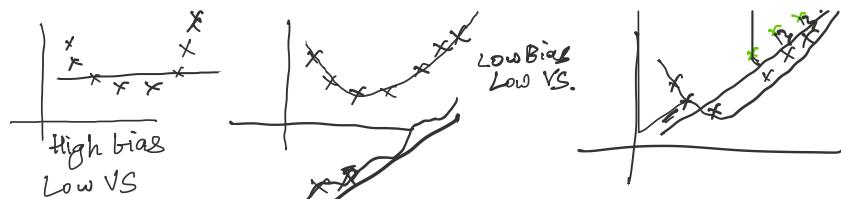
Impact of Regularization on Bias and Variance

Model	Bias	Variance
No regularization	Low	High
Ridge (L2)	Slightly higher	Lower
Lasso (L1)	Higher	Much lower
Elastic Net	Balanced	Balanced

Regularization intentionally increases bias to achieve a larger reduction in variance, resulting in better generalization performance.

Key Takeaways

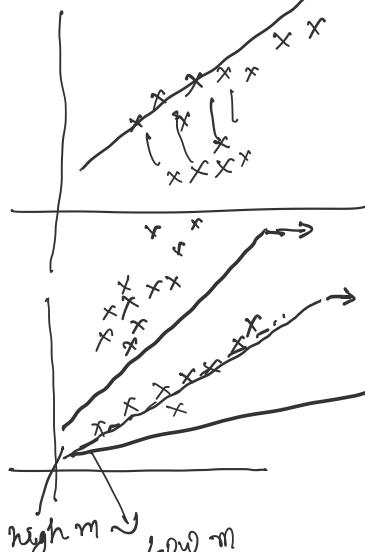
- Bias-variance tradeoff is fundamental to machine learning
- Overfitting and underfitting are direct consequences of imbalance
- Regularization is a principled way to control complexity
- The best model minimizes total expected error $\mathbb{E}[\text{error}]$



Inability to capture the relationship in the training data

overfitting \downarrow bias \uparrow variance

\rightarrow bias \uparrow variance \downarrow



$\beta_0 + \beta_1 x_1 + \beta_2 x_2$
High weights w regularization

$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

If $\lambda \uparrow$ $\beta/m \downarrow$
 $\lambda \rightarrow \infty$ shrinks to zero
 $m \rightarrow 0$

$\beta_1 = 0$
 $\beta_2 \neq 0$
 No role in predicting y

2) Higher coefficients are impacted more

3) Bias variance trade off

Elastic Net

Ridge

Lasso
 $n, \cdot, ^2, 1, \dots, 1$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda |w|^2$$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i) + \lambda |w|$$

when we don't know if all the features are important or not.

EN

$$\Rightarrow L = \sum (y_i - \hat{y}_i)^2 + \alpha |w|^2 + \beta |w|$$

Also used in multicollinearity.

$\left\{ \begin{array}{l} \lambda = a + b \\ L\text{-ratio} = \frac{\alpha}{a+b} \\ d_1 = \frac{a}{\lambda} \\ a = L \times \lambda \\ b = \lambda - a \end{array} \right.$