

Bias-Variance Tradeoff

Cornell CS 3/5780 · Spring 2026

2. Expected Test Error (given h)

- Apply machine learning algorithm \mathcal{A} to learn a hypothesis h
- **Notation:** $h_D = \mathcal{A}(D)$
- For a specific hypothesis h_D learned on dataset D :

$$E_{(\mathbf{x}, y) \sim P} [(h_D(\mathbf{x}) - y)^2] = \int_x \int_y (h_D(\mathbf{x}) - y)^2 \Pr(\mathbf{x}, y) \partial y \partial \mathbf{x}$$

- This measures: How well does this particular hypothesis generalize?
- **Key observation:** h_D is a random variable!
- Different training sets D lead to different hypotheses
- The hypothesis depends on which data points were sampled

1. Setting

- Training data $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ drawn i.i.d. from $P(X, Y)$
- Regression: $y \in \mathbb{R}$ with squared loss
- Today's question is about generalization: what is my expected test error? (after training on D)
- **Definition:** The *expected label* given $\mathbf{x} \in \mathbb{R}^d$ (recall Bayes optimal prediction)

$$\bar{y}(\mathbf{x}) = E_{y|\mathbf{x}}[Y] = \int_y y \Pr(y|\mathbf{x}) \partial y$$

- Question: is $\bar{y}(\mathbf{x})$ a perfect prediction? When is it better/worse?

3. Expected Test Error (Given \mathcal{A})

- Taking expectation over both test data **and** training data:

$$E_{\substack{(\mathbf{x}, y) \sim P \\ D \sim P^n}} [(h_D(\mathbf{x}) - y)^2] = \int_D \int_{\mathbf{x}} \int_y (h_D(\mathbf{x}) - y)^2 P(\mathbf{x}, y) P(D) \partial \mathbf{x} \partial y \partial D$$

- Evaluates the quality of algorithm \mathcal{A} given the distribution $P(X, Y)$
- Note: D = training points and (\mathbf{x}, y) = test point
- It is also useful to compute the *average hypothesis* over all possible training sets:

$$\bar{h}(\mathbf{x}) = E_{D \sim P^n} [h_D(\mathbf{x})] = \int_D h_D(\mathbf{x}) \Pr(D) \partial D$$

- "Average predictor" across all possible training datasets (weighted average with weight = probability)

4. Decomposition part 1

Goal:
$$E_{\mathbf{x},y,D} [(h_D(\mathbf{x}) - y)^2] = \underbrace{E_{\mathbf{x},D} [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2]}_{\text{Variance}} + E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - y)^2]$$

Fill in below steps: Add and subtract, then expand, then simplify cross term

$$\begin{aligned} E_{\mathbf{x},y,D} [(h_D(\mathbf{x}) - y)^2] &= E_{\mathbf{x},y,D} [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x})) + (\bar{h}(\mathbf{x}) - y)]^2 \\ &= E_{\mathbf{x},D} [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2] + 2E_{\mathbf{x},y,D} [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)] + E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - y)^2] \end{aligned}$$

$$\begin{aligned} E_{\mathbf{x},y,D} [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)] &= E_{\mathbf{x}} [E_y [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)]] \\ &= E_{\mathbf{x}} [(E_y [h_D(\mathbf{x})] - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)] \\ &= E_{\mathbf{x}} [(\bar{h}(\mathbf{x}) - \bar{h}(\mathbf{x}))(\bar{h}(\mathbf{x}) - y)] \\ &= 0 \end{aligned}$$

5

5. Decomposition part 2

Goal:
$$E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - y)^2] = \underbrace{E_{\mathbf{x},y} [(\bar{y}(\mathbf{x}) - y)^2]}_{\text{Noise}} + \underbrace{E_{\mathbf{x}} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2]}_{\text{Bias}^2}$$

Fill in below steps: Add and subtract, then expand, then simplify cross term

$$\begin{aligned} E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - y)^2] &= E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) + (\bar{y}(\mathbf{x}) - y)]^2 \\ &= E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2] + E_{\mathbf{x}} [(\bar{y}(\mathbf{x}) - y)^2] + 2E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))(\bar{y}(\mathbf{x}) - y)] \end{aligned}$$

$$\begin{aligned} E_{\mathbf{x},y} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))(\bar{y}(\mathbf{x}) - y)] &= \\ &= \\ &= 0 \end{aligned}$$

6

7

8

6. Bias-Variance Decomposition

$$\underbrace{E_{\mathbf{x},y,D} [(h_D(\mathbf{x}) - y)^2]}_{\text{Expected Test Error}} = \underbrace{E_{\mathbf{x},D} [(h_D(\mathbf{x}) - \bar{h}(\mathbf{x}))^2]}_{\text{Variance}} + \underbrace{E_{\mathbf{x},y} [(\bar{y}(\mathbf{x}) - y)^2]}_{\text{Noise}} + \underbrace{E_{\mathbf{x}} [(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2]}_{\text{Bias}^2}$$

- Variance:** How much does h_D vary across different training sets?
 - Model "overfits" to particular training examples
 - Cause: Model is too complex relative to amount of training data
- Bias:** How far is the average hypothesis from the true expected label?
 - Model "underfits" the data, model class is not expressive enough
 - Cause: Model is too simple to capture the true pattern
- Noise:** Inherent unpredictability in the labels
 - Performance of the Bayes optimal hypothesis
 - Cause: inherent uncertainty and/or uninformative features

8. Summary

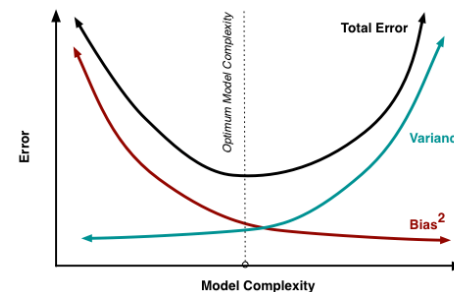
Bias-Variance Decomposition:

$$\text{Expected Test Error} = \text{Bias}^2 + \text{Variance} + \text{Noise}$$

- Bias:** Error from wrong model assumptions
- Variance:** Error from sensitivity to training data
- Noise:** Irreducible error from label randomness
- Tradeoff:** Complex models have low bias but high variance
- Goal:** Find the sweet spot that minimizes total error

7. Bias-Variance Tradeoffs

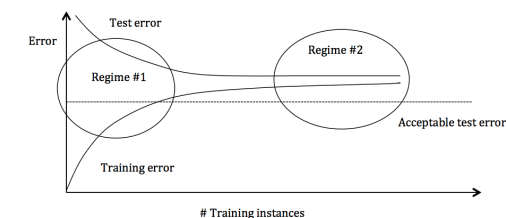
Insight: tune the model complexity to trade off Variance and Bias



9

Diagnosis: where does poor performance come from?

- High Variance (Regime 1)**
 $\text{train err} < \epsilon < \text{test err}$
 indicates overfitting
- High Bias (Regime 2)**
 $\epsilon < \text{train err} \approx \text{test err}$
 indicates underfitting



10