

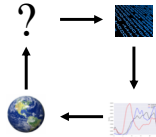
The Bias Variance Tradeoff and Regularization

Slides by:

Joseph E. Gonzalez jegonzal@cs.berkeley.edu

Spring '18 updates:

Fernando Perez fernando.perez@berkeley.edu



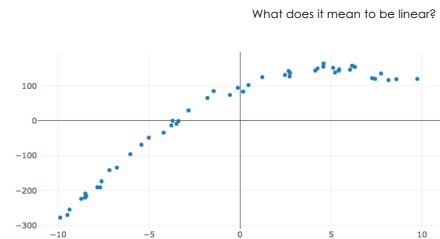
Quick announcements

- Please **be respectful** on Piazza
 - Both of your fellow students and of your teaching staff.
 - The teaching team monitors Piazza, but you can report any incidents directly to Profs. Gonzalez and/or Perez.
- Our infrastructure isn't perfect
 - We're working hard on improving it.
 - We're building the plane while we fly it, full of passengers.
- We have a textbook: textbook.ds100.org
 - It's a **work in progress!**

Linear models for non-linear relationships

Advice for people who are dealing with non-linear relationship issues but would really prefer the simplicity of a linear relationship.

Is this data Linear?



What does it mean to be a linear model?

$$f_{\theta}(\phi(x)) = \phi(x)^T \theta = \sum_{j=1}^k \phi(x)_j \theta_j$$

In what sense is the above **model linear**?

Are linear models linear in the

1. the features?
2. the parameters?

Introducing Non-linear Feature Functions

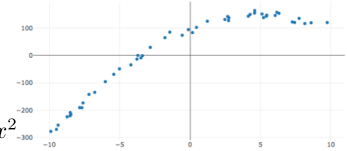
- One reasonable feature function might be:

$$\phi(x) = [1, x, x^2]$$

- That is:

$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

- This is **still a linear model**, in the parameters θ



What are the fundamental challenges in learning?

Fundamental Challenges in Learning?

- **Fit the Data**
 - Provide an explanation for what we observe
- **Generalize to the World**
 - Predict the future
 - Explain the unobserved



Is this cat grumpy or are we overfitting to human faces?

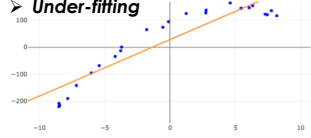
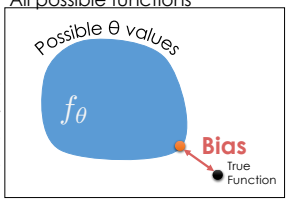
Fundamental Challenges in Learning?

- **Bias:** the expected deviation between the predicted value and the true value
- **Variance:** two sources
 - **Observation Variance:** the variability of the random noise in the process we are trying to model.
 - **Estimated Model Variance:** the variability in the predicted value across different training datasets.

Bias

The expected deviation between the predicted value and the true value

- Depends on both the:
 - choice of f
 - learning procedure
- **Under-fitting**

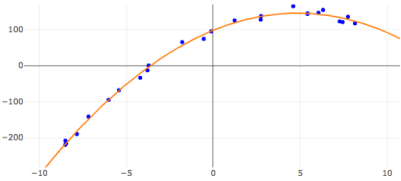



Observation Variance

the variability of the random noise in the process we are trying to model

- measurement variability
- stochasticity
- missing informati

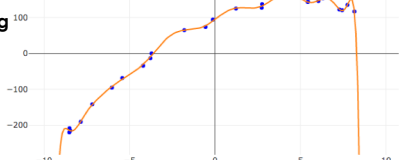
Beyond our control (usually)



Estimated Model Variance

variability in the predicted value across different training datasets

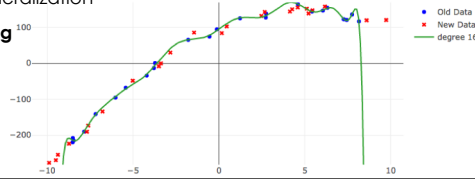
- Sensitivity to variation in the training data
- Poor generalization
- **Overfitting**



Estimated Model Variance

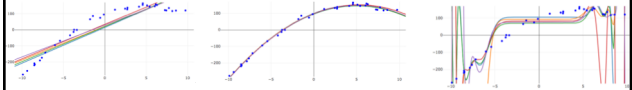
variability in the predicted value across different training datasets

- > Sensitivity to variation in the training data
- > Poor generalization
- > **Overfitting**



The Bias-Variance Tradeoff

Estimated Model Variance →



← Bias

Demo

Analysis of the Bias-Variance Trade-off

Analysis of Squared Error

> For the test point x the expected error:

> Random variables are **red**
Assume noisy observations
→ y is a random variable

True Function
 $y = h(x) + \epsilon$

Noise term:
 $\mathbf{E}[\epsilon] = 0$
 $\mathbf{Var}[\epsilon] = \sigma^2$

$$\mathbf{E} \left[\left(y - f_{\hat{\theta}}(x) \right)^2 \right]$$

Assume **training data** is random
→ θ is a random variable

Analysis of Squared Error

Goal:

$$\mathbf{E} \left[\left(y - f_{\hat{\theta}}(x) \right)^2 \right] =$$

Obs. Var. + **(Bias)²** + **Mod. Var.**

Other terminology:

“Noise” + **(Bias)²** + **Variance**

$$\mathbf{E} \left[(y - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[(y - \underbrace{h(x) + h(x)} - f_{\hat{\theta}}(x))^2 \right]$$

Subtracting and adding $h(x)$

Useful Eqns:
 $y = h(x) + \epsilon$
 $\mathbf{E}[\epsilon] = 0$
 $\text{Var}[\epsilon] = \sigma^2$

$$\mathbf{E} \left[(y - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[(\underbrace{y - h(x)}_a + \underbrace{h(x) - f_{\hat{\theta}}(x)}_b)^2 \right]$$

Expanding in terms of a and b : $(a + b)^2 = a^2 + b^2 + 2ab$

$$= \mathbf{E} \left[(y - h(x))^2 \right] + \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] + 2\mathbf{E} \left[(y - h(x)) (h(x) - f_{\hat{\theta}}(x)) \right]$$

$\underbrace{\hspace{10em}}_{\epsilon}$

$$+ 2\mathbf{E} \left[\epsilon (h(x) - f_{\hat{\theta}}(x)) \right]$$

Useful Eqns:
 $y = h(x) + \epsilon$
 $\mathbf{E}[\epsilon] = 0$
 $\text{Var}[\epsilon] = \sigma^2$

$$\mathbf{E} \left[(y - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[(y - h(x) + h(x) - f_{\hat{\theta}}(x))^2 \right]$$

Expanding in terms of a and b :

$$= \mathbf{E} \left[(y - h(x))^2 \right] + \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] + 2\mathbf{E} \left[\epsilon (h(x) - f_{\hat{\theta}}(x)) \right]$$

Independence of ϵ and θ

$$+ 2\mathbf{E}[\epsilon] \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x)) \right]$$

0

Useful Eqns:
 $y = h(x) + \epsilon$
 $\mathbf{E}[\epsilon] = 0$
 $\text{Var}[\epsilon] = \sigma^2$

$$\mathbf{E} \left[(y - f_{\hat{\theta}}(x))^2 \right] = \sigma^2$$

$\mathbf{E} \left[(y - h(x))^2 \right]$
Obs. Value

$\mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right]$
True Value
Pred. Value

Obs. Variance
"Noise" Term

Model Estimation Error

Useful Eqns:
 $y = h(x) + \epsilon$
 $\mathbf{E}[\epsilon] = 0$
 $\text{Var}[\epsilon] = \sigma^2$

$$\mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] = \text{Next we will show....}$$

$$(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 + \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right]$$

(Bias)² Model Variance

➤How?
 ➤Adding and Subtracting what?

$$\mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] = \mathbf{E} \left[(\underbrace{h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]}_a + \underbrace{\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x)}_b)^2 \right]$$

Expanding in terms of a and b : $(a + b)^2 = a^2 + b^2 + 2ab$

$$\mathbf{E} \left[(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 \right] + \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right] + 2\mathbf{E} \left[(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]) (\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x)) \right]$$

$\underbrace{\hspace{10em}}_{2ab}$

$$\begin{aligned} \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] &= \\ \mathbf{E} \left[(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 \right] &+ \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right] \\ + 2\mathbf{E} \left[(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]) (\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x)) \right] & \\ \underbrace{\hspace{10em}}_{\text{Constant}} & \\ + 2(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]) \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x)) \right] & \end{aligned}$$

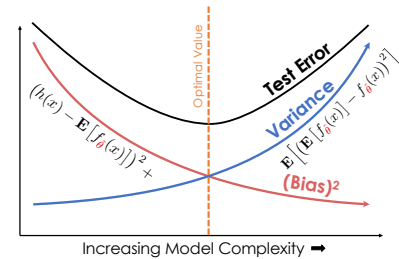
$$\begin{aligned} \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] &= \\ \mathbf{E} \left[(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 \right] &+ \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right] \\ + 2(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]) \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x)) \right] & \\ \underbrace{\hspace{10em}}_{\text{Constant}} & \\ + 2(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)]) \underbrace{\left(\mathbf{E} [f_{\hat{\theta}}(x)] - \mathbf{E} [f_{\hat{\theta}}(x)] \right)}_0 & \end{aligned}$$

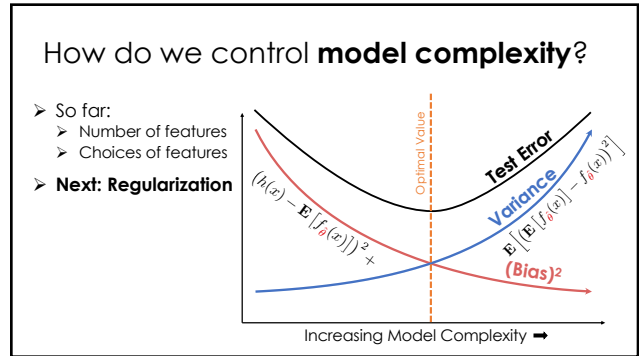
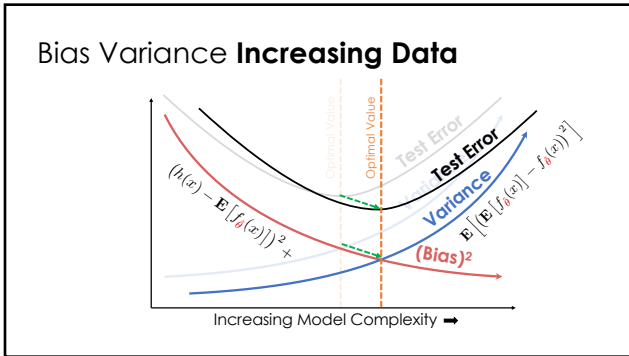
$$\begin{aligned} \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] &= \\ \mathbf{E} \left[(h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 \right] &+ \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right] \\ \underbrace{\hspace{10em}}_{\text{Constant}} & \\ (h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 &+ \end{aligned}$$

$$\begin{aligned} \mathbf{E} \left[(h(x) - f_{\hat{\theta}}(x))^2 \right] &= \\ (h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 &+ \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right] \\ \mathbf{(Bias)^2} &\quad \mathbf{Model Variance} \end{aligned}$$

$$\begin{aligned} \mathbf{E} \left[(y - f_{\theta}(x))^2 \right] &= \\ \mathbf{E} \left[(y - h(x))^2 \right] &\stackrel{= \sigma^2}{=} \mathbf{Obs. Variance} \\ &\text{"Noise"} \\ (h(x) - \mathbf{E} [f_{\hat{\theta}}(x)])^2 &+ \mathbf{(Bias)^2} \\ \mathbf{E} \left[(\mathbf{E} [f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2 \right] &\mathbf{Model Variance} \end{aligned}$$

Bias Variance Plot





Bias Variance Derivation Quiz

> Match each of the following: <http://bit.ly/ds100-sp18-bvt>

(1) $\mathbf{E}[y]$	A. 0
(2) $\mathbf{E}[\epsilon^2]$	B. Bias ²
(3) $\mathbf{E}[(h(x) - \mathbf{E}[f_{\hat{\theta}}(x)])^2]$	C. Model Variance
(4) $\mathbf{E}[\epsilon(h(x) - f_{\hat{\theta}}(x))]$	D. Obs. Variance
	E. $h(x)$
	F. $h(x) + \epsilon$

Bias Variance Derivation Quiz

> Match each of the following: <http://bit.ly/ds100-sp18-bvt>

(1) $\mathbf{E}[y]$	A. 0
(2) $\mathbf{E}[\epsilon^2]$	B. Bias ²
(3) $\mathbf{E}[(h(x) - \mathbf{E}[f_{\hat{\theta}}(x)])^2]$	C. Model Variance
(4) $\mathbf{E}[\epsilon(h(x) - f_{\hat{\theta}}(x))]$	D. Obs. Variance
	E. $h(x)$
	F. $h(x) + \epsilon$

Note: In the original image, blue lines connect (1) to E, (2) to B, (3) to C, and (4) to D.

Regularization

Parametrically Controlling the Model Complexity

> Tradeoff:
 > Increase bias
 > Decrease variance

The dial shows a red line indicating a specific value of λ between the **MIN** and **MAX** positions.

Basic Idea of Regularization

Fit the Data vs. Penalize Complex Models

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i)) + \lambda \mathbf{R}(\theta)$$

> How should we define $\mathbf{R}(\theta)$?
 > How do we determine λ ?

The **Regularization Parameter** λ is highlighted in a red box.

The Regularization Function $R(\theta)$

Goal: Penalize model complexity

Recall earlier: $\phi(x) = [x, x^2, x^3, \dots, x^p]$

- More features → overfitting ...
- How can we control overfitting through θ
- **Proposal:** set weights = 0 to remove features

Common Regularization Functions

Ridge Regression (L2-Reg) $R_{\text{ridge}}(\theta) = \sum_{i=1}^d \theta_i^2$

- Distributes weight across related features (robust)
- Analytic solution (easy to compute)
- Does not encourage sparsity → small but non-zero weights.

LASSO (L1-Reg) $R_{\text{lasso}}(\theta) = \sum_{i=1}^d |\theta_i|$

- **Encourages sparsity** by setting weights = 0
- Used to select informative features
- Does not have an analytic solution → numerical methods

Regularization and Norm Balls

Regularization and Norm Balls

Python Demo!

The shapes of the norm balls.
Maybe show reg. effects on actual models.

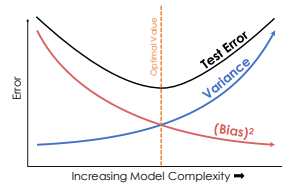
Determining the Optimal λ

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i)) + \lambda R(\theta)$$

- Value of λ determines bias-variance tradeoff
- Larger values → more regularization → more bias → less variance

Summary

$$\begin{aligned} \mathbb{E}[(y - f_{\theta}(x))^2] &= \\ &\mathbb{E}[(y - h(x))^2] + \\ &(h(x) - \mathbb{E}[f_{\theta}(x)])^2 + \\ &\mathbb{E}[(\mathbb{E}[f_{\theta}(x)] - f_{\theta}(x))^2] \end{aligned}$$



Regularization

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i)) + \lambda \mathbf{R}(\theta)$$

