

A photograph of a classic brick clock tower. The tower features a large white clock face with black Roman numerals and hands. It is topped with a small white dome and a weather vane. The tower is made of red brick and has white decorative elements like cornices and columns. The background is a clear blue sky.

Introduction to Machine Learning

Loss Functions, Bias and Variance

Loss functions

- 0/1 loss
 - Perceptron

Any hope?

- Consider modifying function to add a regularizer
 - $R(w,b)$
-
- Optimize trade-off between low training error and simplicity

Relationship between margin and loss

- Small changes to w, b have large impact on objective function

Relationship between margin and loss

- Sigmoid function

Surrogate loss functions

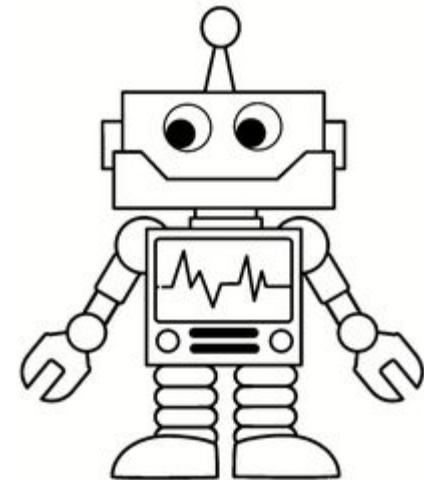
- 0/1 $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$
- Hinge $\ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$
- Logistic $\ell^{(\log)}(y, \hat{y}) = \frac{1}{\log 2} \log (1 + \exp[-y\hat{y}])$
- Exponential $\ell^{(\exp)}(y, \hat{y}) = \exp[-y\hat{y}]$
- Squared $\ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$

Surrogate loss functions

- 0/1 $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$
- Hinge $\ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$
- Logistic $\ell^{(\log)}(y, \hat{y}) = \frac{1}{\log 2} \log (1 + \exp[-y\hat{y}])$
- Exponential $\ell^{(\exp)}(y, \hat{y}) = \exp[-y\hat{y}]$
- Squared $\ell^{(\text{sqr})}(y, \hat{y}) = (y - \hat{y})^2$

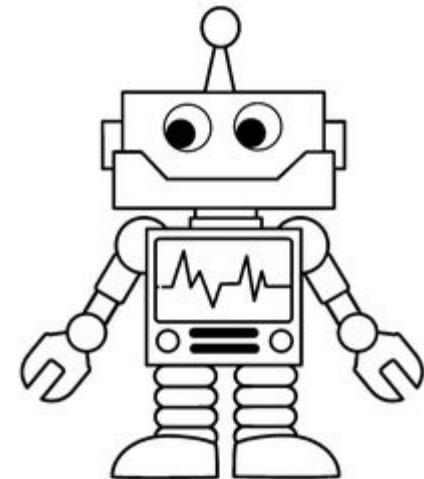
Regularizer

- Helps to prevent overfit
- Perceptron $\min_{w,b} \sum_n \mathbf{1}[y_n(w \cdot x_n + b) \leq 0] + \lambda R(w, b)$
- What are good choices of $R(w,b)$ for hyperplanes?



Regularizer

- Helps to prevent overfit
- Perceptron $\min_{w,b} \sum_n \mathbf{1}[y_n(w \cdot x_n + b) \leq 0] + \lambda R(w, b)$
- What are good choices of $R(w,b)$ for hyperplanes?



Norm

- Popular choice
- Keeps function (weights) from changing too quickly
- $R^{(\text{norm})}(w, b) = ||w|| = \sqrt{\sum_d w_d^2}$
 - Convex
 - Smooth
- l-norm (p-norm)

$$||w||_p = \left(\sum_d |w_d|^p \right)^{\frac{1}{p}}$$

P-norm

$$\|w\|_p = \left(\sum_d |w_d|^p \right)^{\frac{1}{p}}$$

Bias and fairness



Bias and Fairness



Business

Apple Card algorithm sparks gender bias allegations against Goldman Sachs

COMPUTING

Racial Bias Found in a Major Health Care Risk Algorithm

When a Computer Program Keeps You in Jail

What types of bias can exist?

- Adaptation

Fairness and data bias

- Labeling
- Sample selection
- Task
- Feedback loops

Error

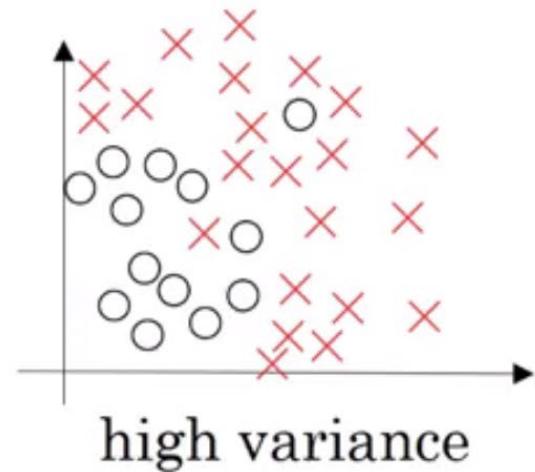
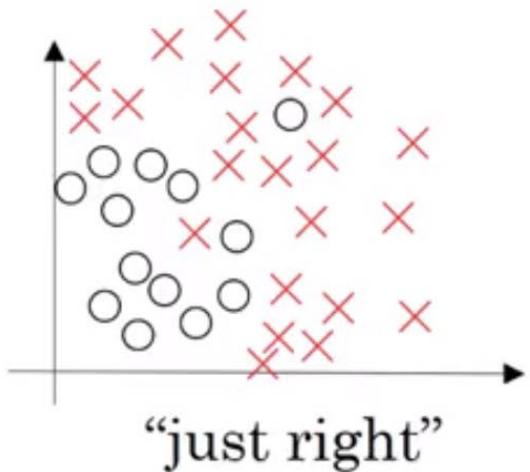
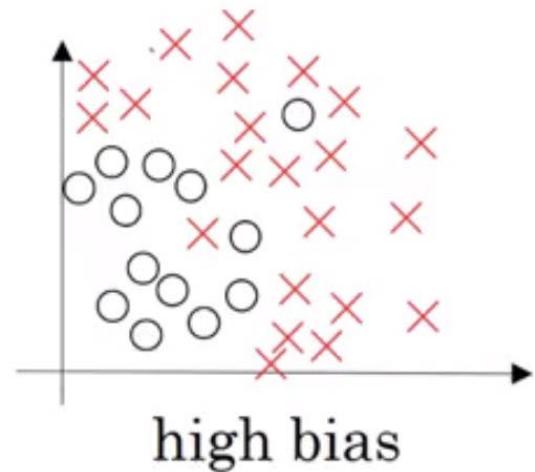
$$\text{error}(f) = \underbrace{\left[\text{error}(f) - \min_{f^* \in \mathcal{F}} \text{error}(f^*) \right]}_{\text{estimation error}} + \underbrace{\left[\min_{f^* \in \mathcal{F}} \text{error}(f) \right]}_{\text{approximation error}}$$

Error

- Bias error
 - Simplifying assumptions that make target function easier to learn
- Variance error
 - Amount that estimated target function will change if change training data

Bias-Variance tradeoff

- Goal: low bias and low variance
- As we increase complexity
 - Bias decreases (better fit to data)
 - Variance increases (fit varies more with data)



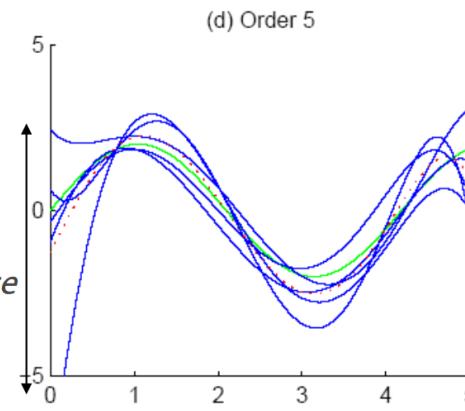
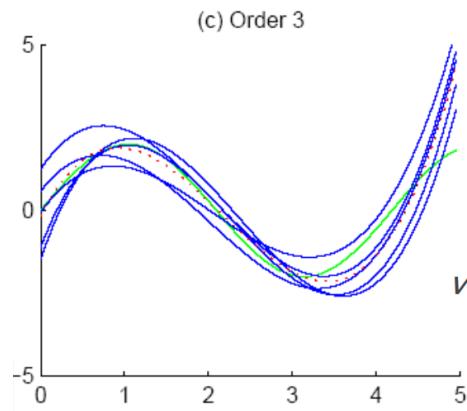
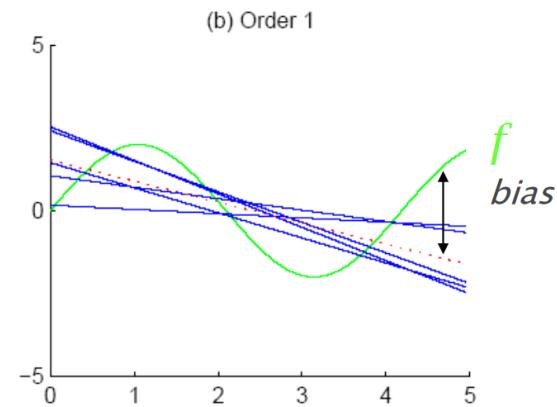
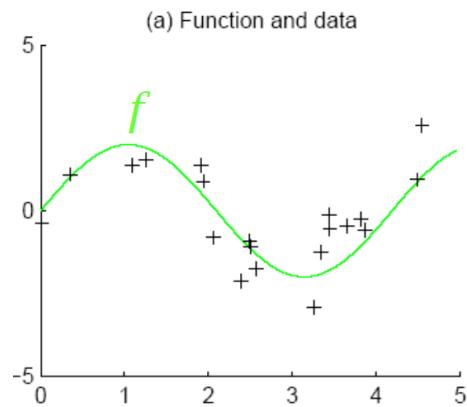
Bias and variance

Cat classification



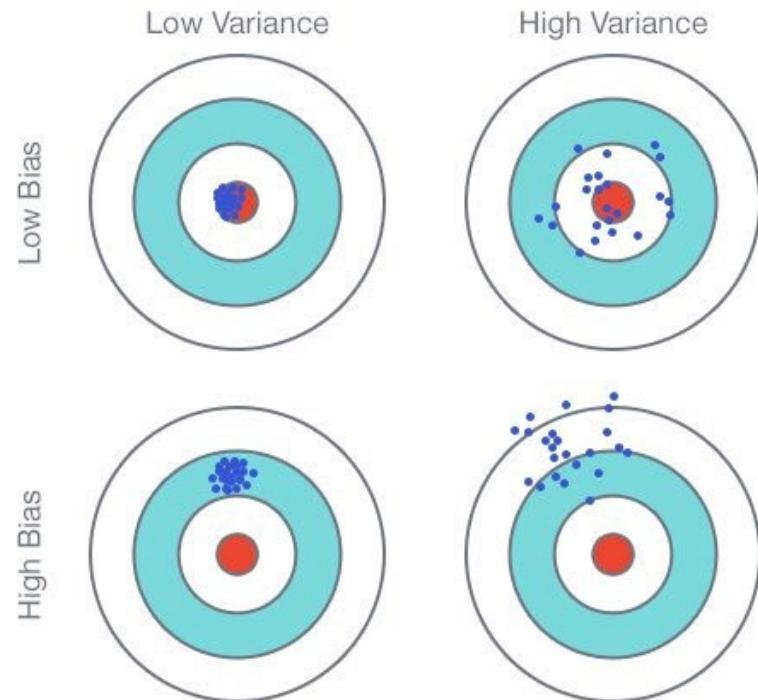
Example

- $f(x) = 2 \sin(1.5x)$



Regularization

- Path to bias-variance trade-off



DL Bias and Variance



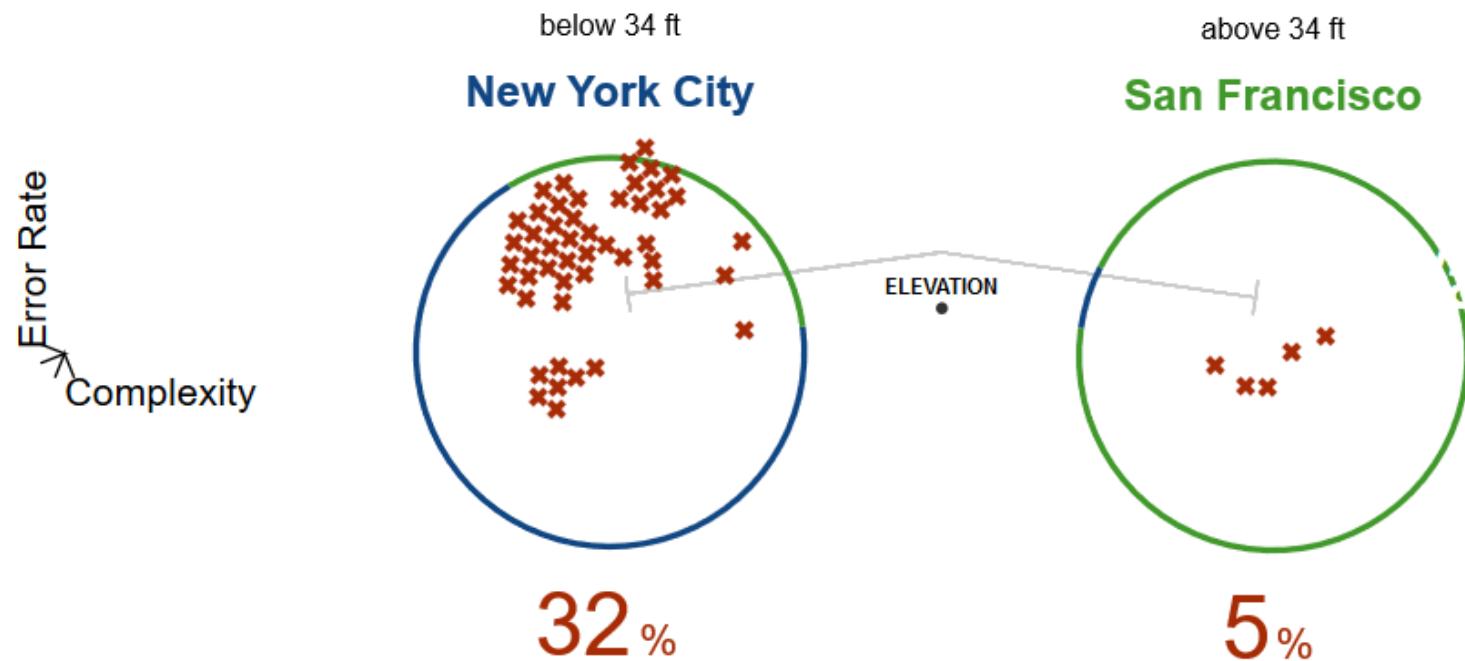
Decision stump



Decision stump



Errors due to bias



Decrease error due to bias



Example of var



Example of var



Example of varia



Example of variance

155/250

92/155

34/92

30/34

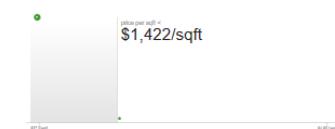
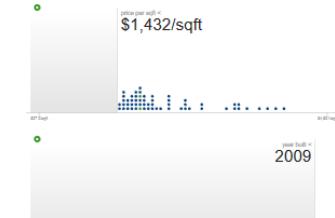
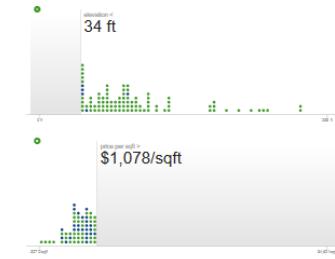
28/30

27/28

2/27

1/2

Error Rate
Complexity



Flukes are nor



Bias-variance tradeoff



Bias-variance tradeoff



Bias-var tradeoff

