

Bias-Variance Trade-off in ML

Sargur Srihari
srihari@cedar.buffalo.edu

Bias-Variance Decomposition

1. Model Complexity in Linear Regression

2. Point estimate

Bias-Variance in Statistics

3. Bias-Variance in Regression

- Choosing λ in maximum likelihood/least squares estimation
- Formulation for regression
- Example
- Choice of optimal λ

Model Complexity in Linear Regression

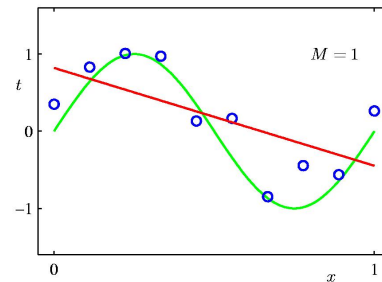
- We looked at linear regression where form of basis functions ϕ and their no. M are fixed
- Using maximum likelihood (equivalently least squares) leads to severe overfitting if complex models are trained with limited data
 - However limiting M to avoid overfitting has side effect of not capturing important trends in the data
- Regularization can control overfitting for models with many parameters
 - But seeking to minimize wrt both \mathbf{w} and λ leads to unregularized solution with $\lambda=0$

Overfitting is a property of Max Likelihood

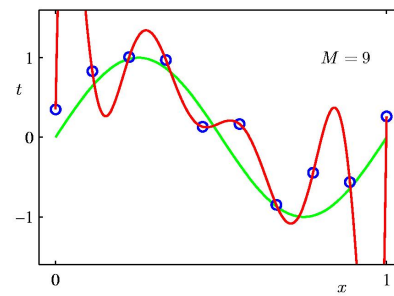
- Does not happen when we marginalize over parameters in a Bayesian setting
- Before considering Bayesian view, instructive to consider frequentist viewpoint of model complexity
- It is called *Bias-Variance* trade-off

Bias-Variance in Regression

- Low degree polynomial has high bias (fits poorly) but has low variance with different data sets



- High degree polynomial has low bias (fits well) but has high variance with different data sets



Bias-Variance in Point Estimate

True height of Chinese emperor: 200cm (6.5ft)

Poll question: “How tall is the emperor?”

Determine how wrong people are, on average



If all answer 200, ave. squared error is 0. Consider three datasets with mean=180 (or bias error = 20), but increasing variance (0, 10 and 20)

<p>200 — Bias 180 — No variance</p> <ul style="list-style-type: none">• Dataset 1• Everyone believes it is 180 (variance=0)• Answer is always 180• The error is always -20• <u>Ave squared error is 400</u>• Average bias error is 20• <u>400=400+0</u>	<p>Bias — 200 — Some variance — 180 —</p> <ul style="list-style-type: none">• Dataset 2• Normally distributed beliefs with mean 180 and std dev 10 (variance 100)• Poll two: One says 190, other 170• Bias Errors are -10 and -30<ul style="list-style-type: none">– Average bias error is -20• Squared errors: 100 and 900<ul style="list-style-type: none">– <u>Ave squared error: 500</u>• <u>500 = 400 + 100</u>	<p>Bias — 200 — More variance — 180 —</p> <ul style="list-style-type: none">• Dataset 3• Normally distributed beliefs with mean 180 and std dev 20 (variance=400)• Poll two: One says 200 and other 160• Errors: 0 and -40<ul style="list-style-type: none">– Ave error is -20• Sq. errors: 0 and 1600<ul style="list-style-type: none">– <u>Ave squared error: 800</u>• <u>800 = 400 + 400</u>
--	--	---

Average Squared Error = (Bias error)² + Variance

As variance increases, error increases

Prediction in Linear Regression

- Need to apply *Decision Theory*: choose a specific estimate $y(\mathbf{x})$ of the value of t for a given \mathbf{x}
- In doing so if we incur incur loss $L(t, y(\mathbf{x}))$
- Then the expected loss is

$$E[L] = \iint L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt$$

- Using squared loss function $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$

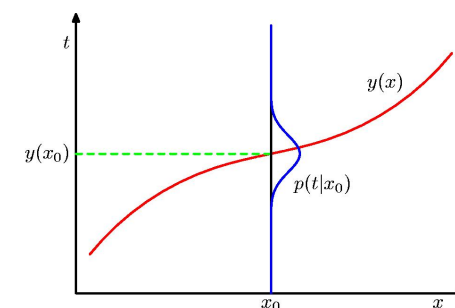
$$E[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

- Taking derivative of E wrt $y(\mathbf{x})$, using calculus of variations,

$$\frac{\delta E[L]}{\delta y(\mathbf{x})} = 2 \int (y(\mathbf{x}) - t) p(\mathbf{x}, t) dt$$

- Setting equal to zero, solving for $y(\mathbf{x})$ and using sum and product rules

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t | \mathbf{x}) dt = E_t[t | \mathbf{x}]$$



Regression function $y(\mathbf{x})$, which minimizes the expected squared loss, is given by the mean of the conditional distribution $p(t|x)$

Alternative Derivation

- We can show that the optimal prediction is equal to the conditional mean in another way
- First we have

$$\{y(\mathbf{x}) - t\}^2 = \{y(\mathbf{x}) - E[t | \mathbf{x}] + E[t | \mathbf{x}] - t\}^2$$

- Substituting into the loss function, we obtain expression for the loss function as

$$E[L] = \int \{y(\mathbf{x}) - E[t | \mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}(t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

- The function $y(\mathbf{x})$ we seek to determine enters only in the first term, which will be minimum when $y(\mathbf{x}) = E[t | \mathbf{x}]$

Bias -Variance in Regression

- $y(\mathbf{x})$: regression function using some method
- $h(\mathbf{x})$: optimal prediction (using squared loss)

$$h(\mathbf{x}) = E[t | \mathbf{x}] = \int tp(t | \mathbf{x})dt$$
- If we assume loss function $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- $E[L]$ for a particular data set D can be written as
expected loss = (bias)² + variance + noise
- where

$$(\text{bias})^2 = \int \{E_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 p(\mathbf{x}) dx$$

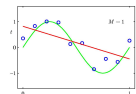
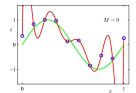
Difference between expected value and optimal

$$\text{variance} = \int E_D[\{y(\mathbf{x}; D) - E_D[y(\mathbf{x}; D)]\}^2] p(\mathbf{x}) dx$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) dx dt$$

Goal: Minimize Expected Loss

- We have decomposed expected loss into sum of (squared) bias, a variance and a constant noise term
- There is a trade-off between bias and variance
 - Very flexible models have low bias and high variance
 - Rigid models have high bias and low variance
 - Optimal model has has the best balance



Dependence of Bias-Variance on Model Complexity

- $h(x) = \sin(2\pi x)$
- Regularization parameter λ
- $L=100$ data sets
- Each with $N=25$
- 24 Gaussian Basis functions
 - No of parameters $M=25$
- Total Error function:

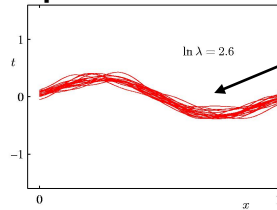
$$\frac{1}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \phi(x_n) \}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Where ϕ is a vector of basis functions

Result of averaging multiple solutions with complex model gives good fit

Weighted averaging of multiple solutions is at heart of Bayesian approach: not wrt multiple data sets but wrt posterior distribution of parameters

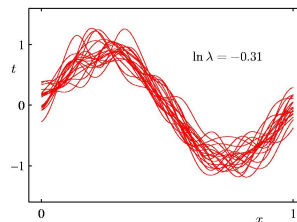
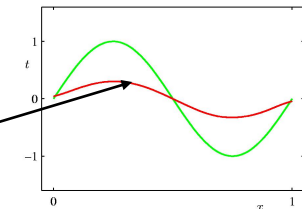
20 Fits for
25 data
points each



High λ

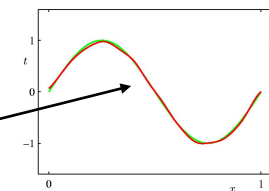
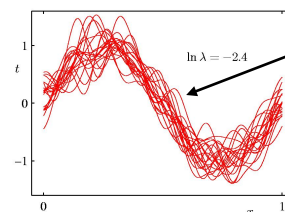
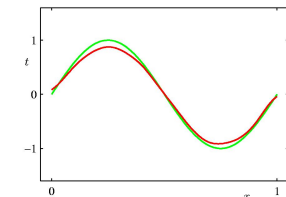
Low
Variance
High bias

Red: Average of Fits
Green: Sinusoid from which
data was generated



Low λ

High
Variance
Low bias



Determining optimal λ

- Average Prediction

$$\bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x)$$

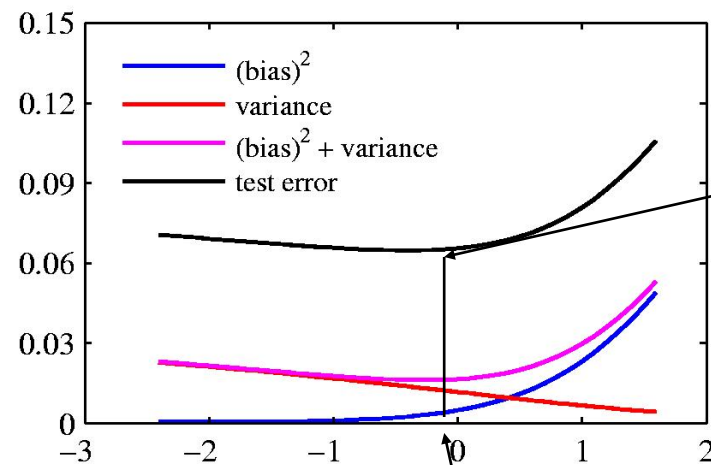
- Squared Bias

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \left\{ \bar{y}(x_n) - h(x_n) \right\}^2$$

- Variance

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \left\{ y^{(l)}(x_n) - \bar{y}(x_n) \right\}^2$$

Squared Bias and Variance vs λ



Test error minimum occurs close to minimum of $(\text{bias}^2 + \text{variance})$

$$\ln \lambda = -0.31$$

Small values of λ allow model to become finely tuned to noise leading to large variance

Large values of λ pull weight parameters to zero leading to large bias

Bias-Variance vs Bayesian

- Bias-Variance decomposition provides insight into model complexity issue
- Limited practical value since it is based on ensembles of data sets
 - In practice there is only a single observed data set
 - If there are many training samples then combine them
 - which would reduce over-fitting for a given model complexity
- Bayesian approach gives useful insights into over-fitting and is also practical