

Training/Testing and Regularization

Il-Chul Moon
Dept. of Industrial and Systems Engineering
KAIST

icmoon@kaist.ac.kr

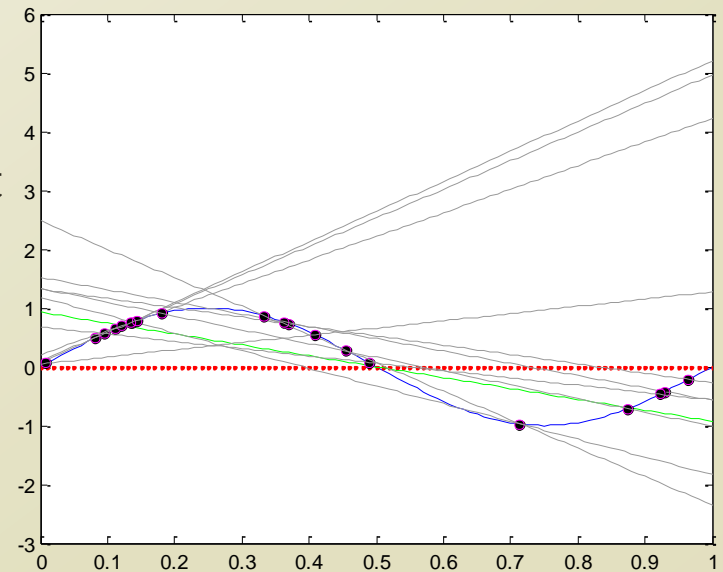
Weekly Objectives

- Understand the concept of bias and variance
 - Know the concept of over-fitting and under-fitting
 - Able to segment two sources, bias and variance, of error
- Understand the bias and variance trade-off
 - Understand the concept of Occam's razor
 - Able to perform cross-validation
 - Know various performance metrics for supervised machine learning
- Understand the concept of regularization
 - Know how to apply regularization to
 - Linear regression
 - Logistic regression
 - Support vector machine

MODEL REGULARIZATION

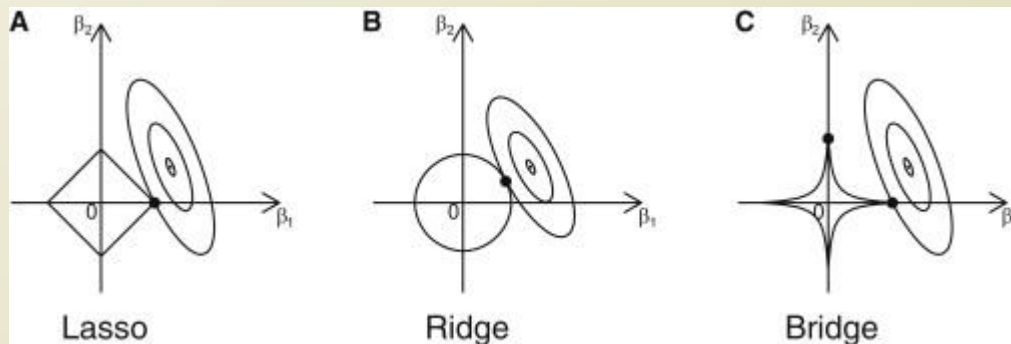
Concept of Regularization

- Disaster in terms of variance
- With regularization
 - We sacrifice the perfect fit
 - Reducing the training accuracy
 - We increase the potential fit in the test
 - Because of the increased model complexity, the bias tends to decrease a little bit
 - Eventually, regularization is another constraint for models
 - Existing constraint?
 - Minimizing error from training set
- We add a new term to the MSE



Formal Definition of Regularization

- Regularization is another constraint for the regression
 - The below $J(B)$ is the regularization function to minimize
 - B is the weight of the regression model except the constant term
- There are diverse regularization
 - L1 Regularization == Lasso regularization
 - The first order
 - L2 Regularization == Ridge regularization
 - The second order
 - Depends on the order of the regularization term
 - The order determines the shape of the loss function



$$E(w) = \frac{1}{2} \sum_{n=0}^N (\text{train}_n - g(x_n, w))^2 + \frac{\lambda}{2} \|w\|^2$$

$$E(w) = \frac{1}{2} \sum_{n=0}^N (\text{train}_n - g(x_n, w))^2 + \lambda |w|$$

Regularization of Linear Regression

- Let's apply the regularization idea to the linear regression

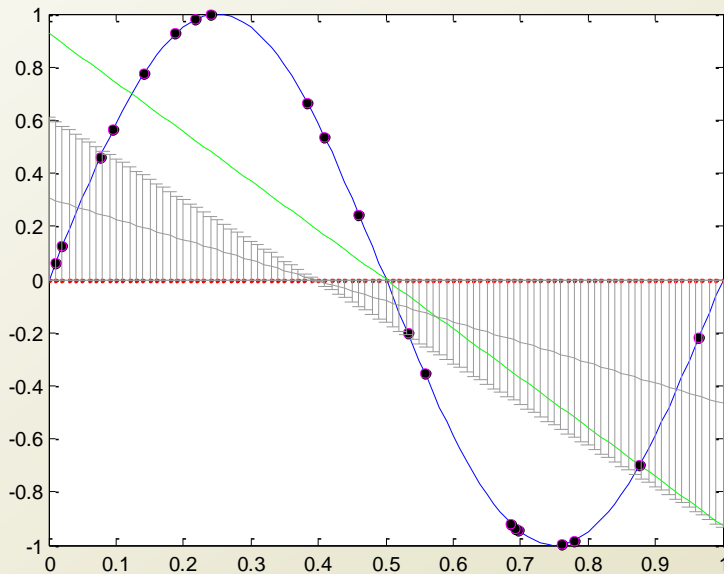
$$E(w) = \frac{1}{2} \sum_{n=0}^N (\text{train}_n - g(x_n, w))^2 + \frac{\lambda}{2} \|w\|^2$$

- We can calculate w in the closed form.

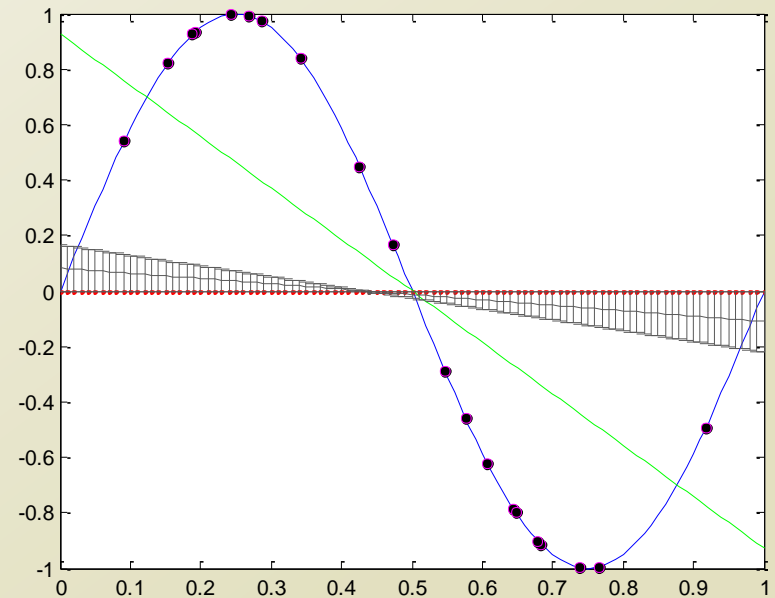
$$\begin{aligned} \frac{d}{dw} E(w) &= 0 \\ \frac{d}{dw} E(w) &= \frac{d}{dw} \left(\frac{1}{2} \|\text{train} - Xw\|^2 + \frac{\lambda}{2} \|w\|^2 \right) \\ &= \frac{d}{dw} \left(\frac{1}{2} \|\text{train} - Xw\|^T \|\text{train} - Xw\| + \frac{\lambda}{2} w^T w \right) \\ &= \frac{d}{dw} \left(\frac{1}{2} (\text{train}^T \text{train} - 2X^T w \cdot \text{train} + X^T X w^T w) + \frac{\lambda}{2} w^T w \right) \\ &= \frac{d}{dw} \left(\frac{1}{2} \text{train}^T \text{train} - X^T w \cdot \text{train} + \frac{1}{2} X^T X w^T w + \frac{\lambda}{2} w^T w \right) \\ &= -X^T \cdot \text{train} + X^T X w + \lambda w = 0 \end{aligned}$$

$$\begin{aligned} -X^T \cdot \text{train} + X^T X w + \lambda I w &= 0 \\ -X^T \cdot \text{train} + (X^T X + \lambda I) w &= 0 \\ (X^T X + \lambda I) w &= X^T \cdot \text{train} \\ w &= (X^T X + \lambda I)^{-1} X^T \cdot \text{train} \end{aligned}$$

Effect of Regularization



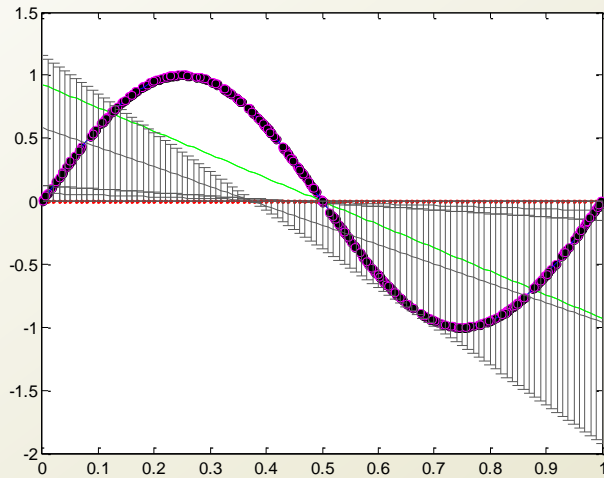
Bias = 0.3092
Var. = 2.0708



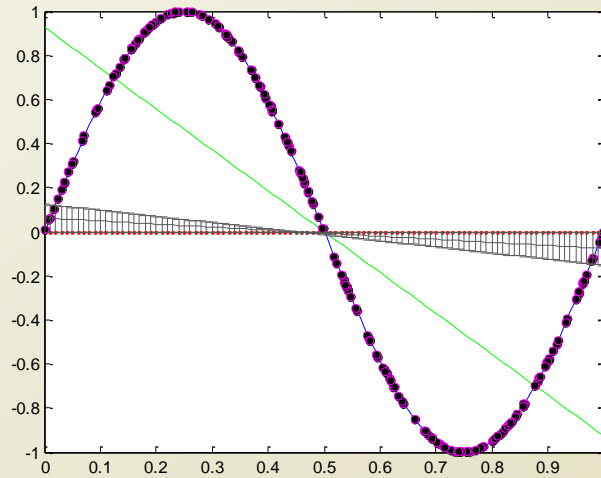
Bias = 0.4372
Var. = 0.1167

- When $\lambda = 1$
 - The bias increases a little bit
 - The variance reduces significantly

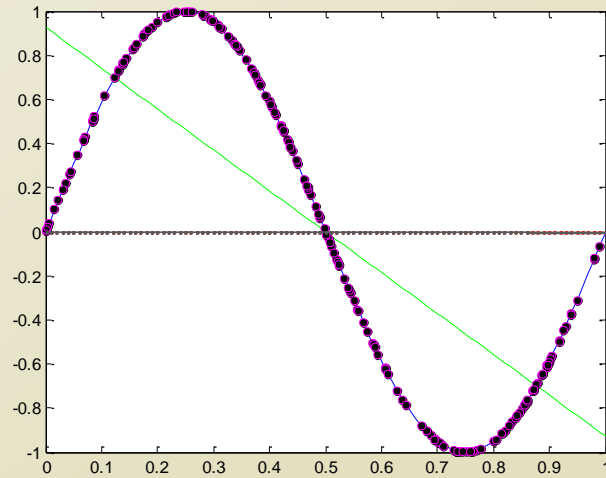
Optimizing the Regularization



$\lambda = 0$



$\lambda = 1$



$\lambda = 100$

- We need to optimize λ
 - Too low λ : Too high variance
 - Works like an unregularized model
 - Too high λ : Too low variance
 - Works like a less complex model
 - Converting the first-order model into the constant model
- How to optimize λ ?

Regularization of Logistic Regression

- Regularization is applicable to other models
 - Such as logistic regression
- You can search for the closed form and the approximate form of finding θ

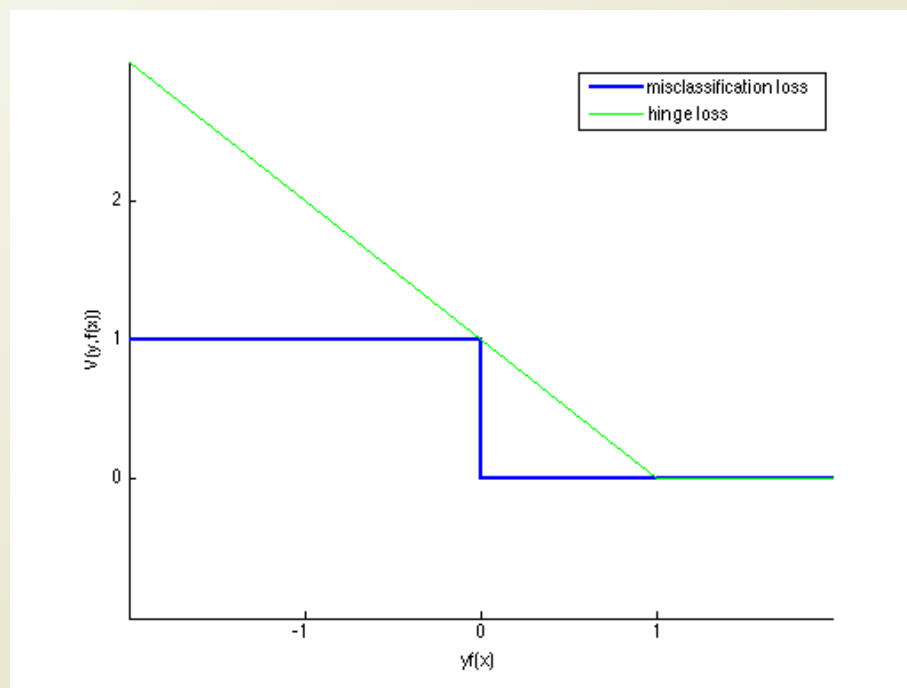
$$\arg \max_{\theta} \sum_{i=1}^m \log p(y_i | x_i, \theta) - \alpha R(\theta)$$

$$\text{L1: } R(\theta) = \|\theta\|_1 = \sum_{i=1}^n |\theta_i|$$

$$\text{L2: } R(\theta) = \|\theta\|_2^2 = \sum_{i=1}^n \theta_i^2$$

Regularization and SVM

$$f = \arg \min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2 \right\}$$



$$V(y_i, f(x_i)) = (1 - yf(x))_+ \\ (s)_+ = \max(s, 0)$$

$$f = \arg \min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n (1 - yf(x))_+ + \lambda \|f\|_{\mathcal{H}}^2 \right\}$$

$$f = \arg \min_{f \in \mathcal{H}} \left\{ C \sum_{i=1}^n (1 - yf(x))_+ + \frac{1}{2} \|f\|_{\mathcal{H}}^2 \right\}$$

$$C = \frac{1}{2\lambda n}$$

- Support vector is a special case of regularization with the hinge loss