# Linear Regression

Alice Oh
[alice.oh@kaist.edu](mailto:alice.oh@kaist.edu)

# Linear Regression

degree 1



$$y = ax + b$$
$$y = w_1 x + w_0$$

(1, 6)

(9, -0.5)

# Linear Regression

- Response (real number) is a linear function of the inputs

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \epsilon$$

- Assume that \epsilon (the residual error) has a Gaussian distribution

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(y|\mu(\mathbf{x}), \sigma^2(\mathbf{x}))$$

where

$$\mu(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = w_0 + w_1 x$$

# Modeling non-linear relationships

- Simply take

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(y|\mathbf{w}^T\mathbf{x}, \sigma^2)$$

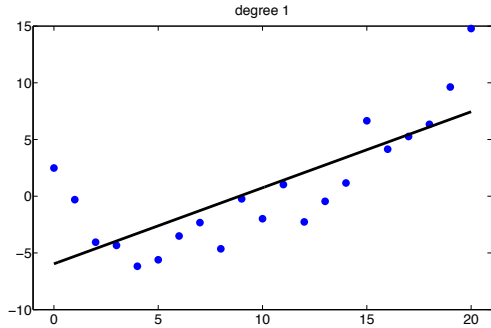and replace x with some non-linear function of the inputs

$$\phi(x) = X^2$$

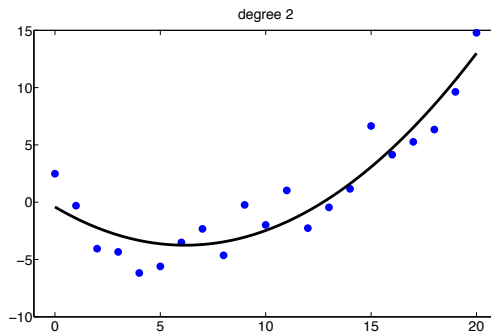$$p(y|\mathbf{x}, \theta) = \mathcal{N}(y|\mathbf{w}^T\phi(\mathbf{x}), \sigma^2)$$

$$2X + 3X^2$$

This is called the basis function expansion.

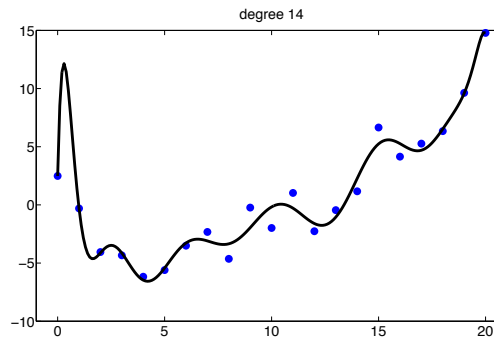(But the model is still called linear regression because it is linear in the parameters w)
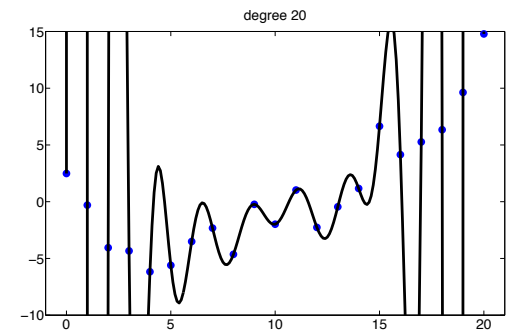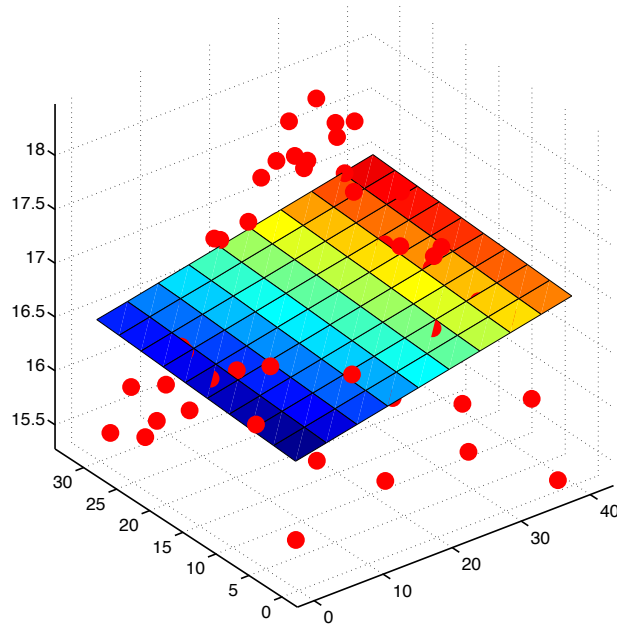
# Polynomial Regression

# Multivariate linear regression



$$w_0 + w_1 x_1 + w_2 x_2$$

$$w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2$$

Figures from Kevin Murphy's book -- Machine Learning: A Probabilistic Perspective

# Maximum likelihood estimation

- Using MLE, arguments \theta can be computed by

$$\arg \max \log p(D|\theta)$$

$$= \sum_{i=1}^{N} log p(y_i | \mathbf{x}_i, \theta)$$

- If we plug in the Gaussian formulation

$$p(y|\mathbf{x}, \theta) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \sigma^2)$$

and put it into the log likelihood above, we get

$$= \sum_{i=1}^{N} log[(\frac{1}{2\pi\sigma^2})^{\frac{1}{2}} exp(-\frac{1}{2\sigma^2}(y_i - \mathbf{w}^T \mathbf{x}_i)^2)]$$

$$= -\frac{N}{2} log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{N}(y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

$$y = \mathbf{w} \mathbf{x} + \mathbf{w_0}$$

# Residual Sum of Squares

- Log likelihood

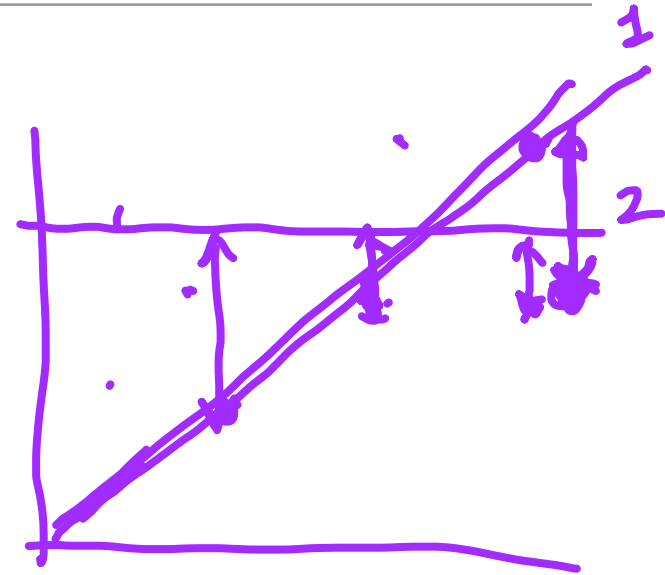$$-\frac{N}{2}log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2$$

- Negative log likelihood (NLL)

$$\frac{N}{2}log(2\pi\sigma^2) + \frac{1}{2\sigma^2}\sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2$$
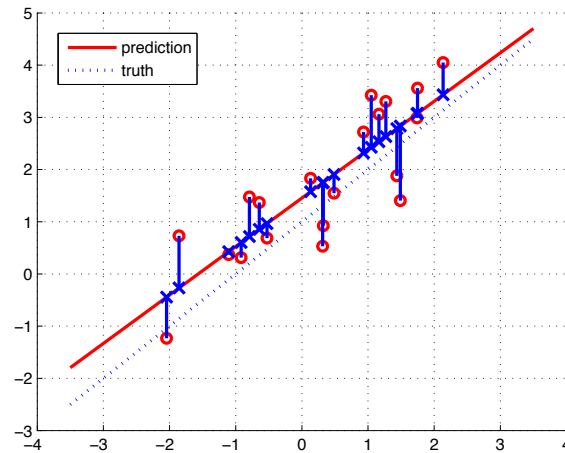
- To minimize NLL, we minimize this term

$$\sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2$$

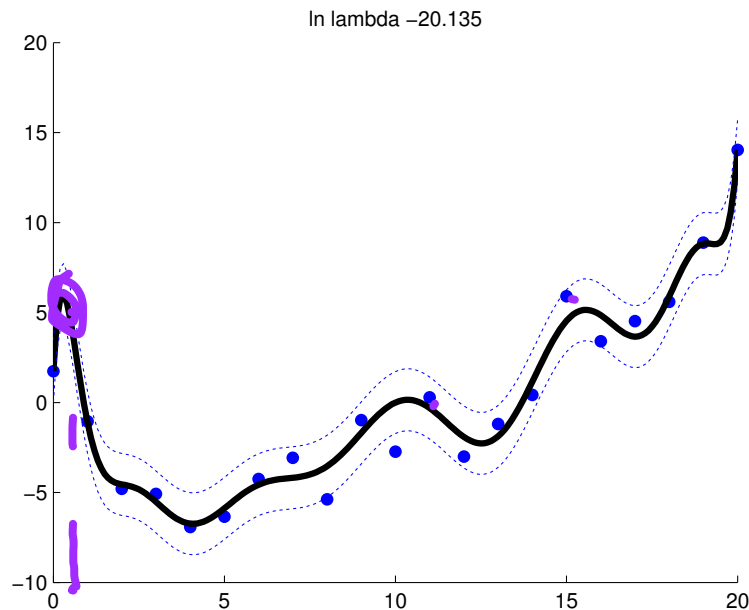called residual sum of squares (RSS)

# Least Squares

- MLE for w is the one that minimizes the RSS

# Ridge Regression

- MLE can overfit

  - For linear regression, this means the weights can become large



ln lambda −20.135

**Regularization**

$W_5 x^6$

$\frac{Na}{2} 0$

$W = W_2 x^2, W_3 x^3$

W =

6.560, -36.934, -109.255,
543.452, 1022.561,
-3046.224, -3768.031,
8524.540 …

- We can encourage the weights to be small by putting a zero-mean Gaussian prior on the weights – $l_2$ Regularization!

Figure from Murphy's MLPP book

# Ridge Regression

- Zero-mean Gaussian prior on the weights

$$p(\mathbf{w}) = \prod_j \mathcal{N}(w_j | 0, \tau^2)$$

- MAP estimation problem

$$\mathrm{argmax} \sum_{i=1}^{N} log\mathcal{N}(y_i | w_0 + \mathbf{w}^T \mathbf{x}_i, \sigma^2) + \sum_{j=1}^{D} log\mathcal{N}(w_j | 0, \tau^2)$$

- Compare with the MLE problem

$$\arg\max \log p(D|\theta) = \sum_{i=1}^{N} logp(y_i | \mathbf{x}_i, \theta)$$

- To solve the MAP estimation, minimize

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^{N} (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda ||\mathbf{w}||_2^2$$

$$\text{where } \lambda = \frac{\sigma^2}{\tau^2} \text{ and } ||\mathbf{w}||_2^2 = \sum_j w_j^2 = \mathbf{w}^T \mathbf{w}$$

# Ridge Regression

$$J(\mathbf{w}) = \frac{1}{N}\sum_{i=1}^{N}(y_i - (w_0 + \mathbf{w}^T\mathbf{x}_i))^2 + \lambda||\mathbf{w}||_2^2$$

- Compare with NLL before
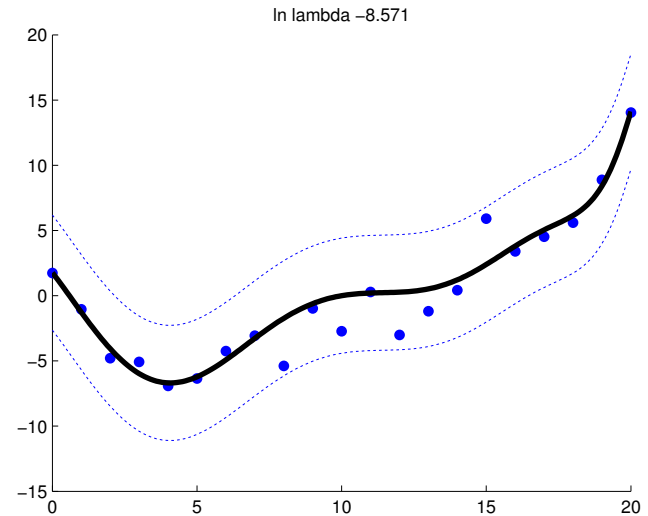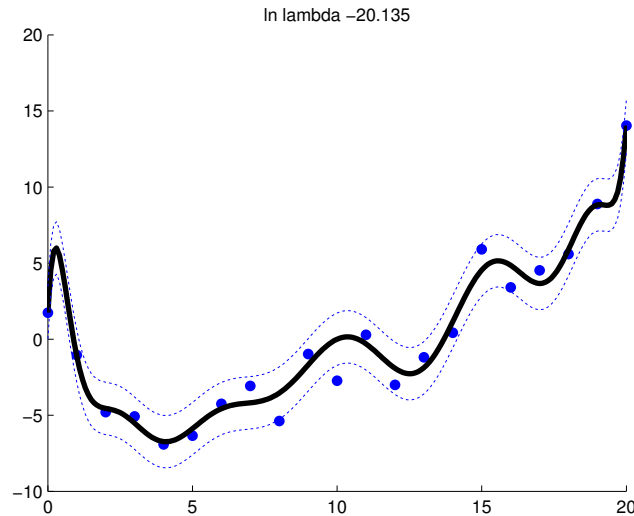
$$\sum_{i=1}^{N}(y_i - \mathbf{w}^T\mathbf{x}_i)^2$$

- So the first term of ridge regression is same as NLL, and the second term is the complexity penalty (when \lamba > 0)

- Corresponding solution is

$$\hat{\mathbf{w}}_{ridge} = (\lambda\mathbf{I}_D + \mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \quad \text{(from Murphy's MLPP book Section 7.5.1)}$$

$$\hat{\mathbf{w}}_{OLS} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \quad \text{(from Murphy's MLPP book Section 7.3.1)}$$

# Ridge Regression



$$\hat{\mathbf{w}}_{ridge} = (\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$
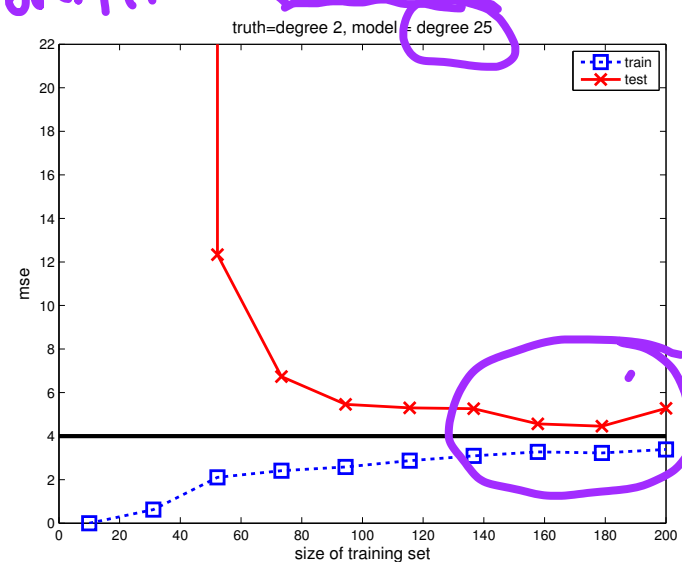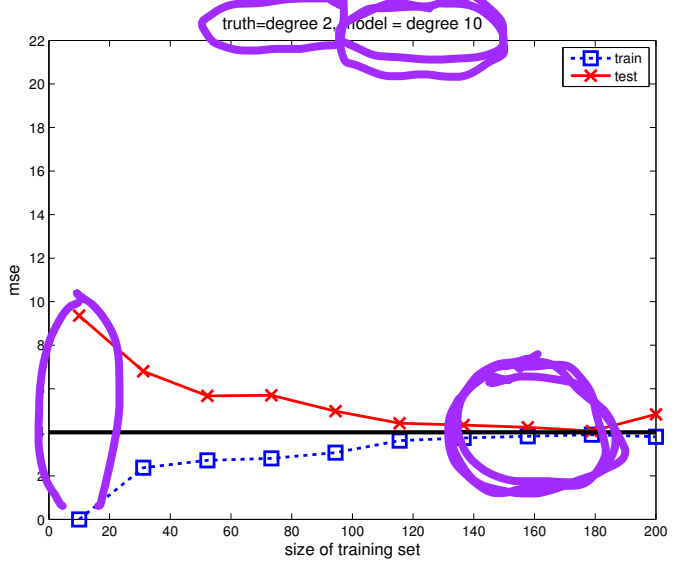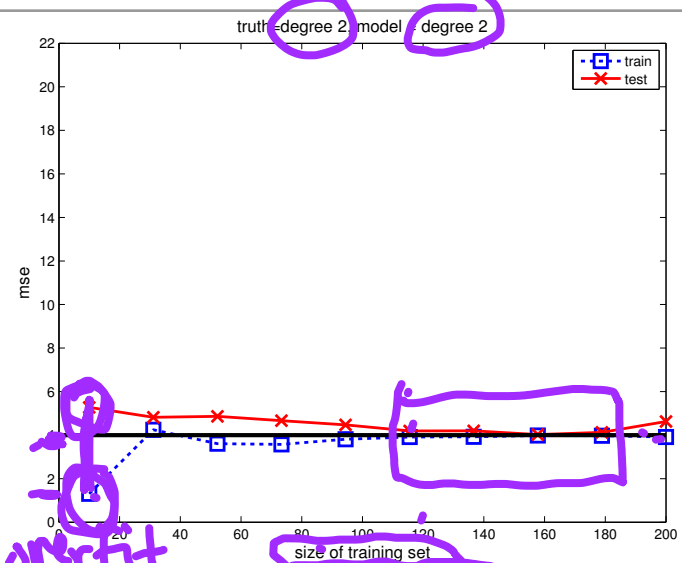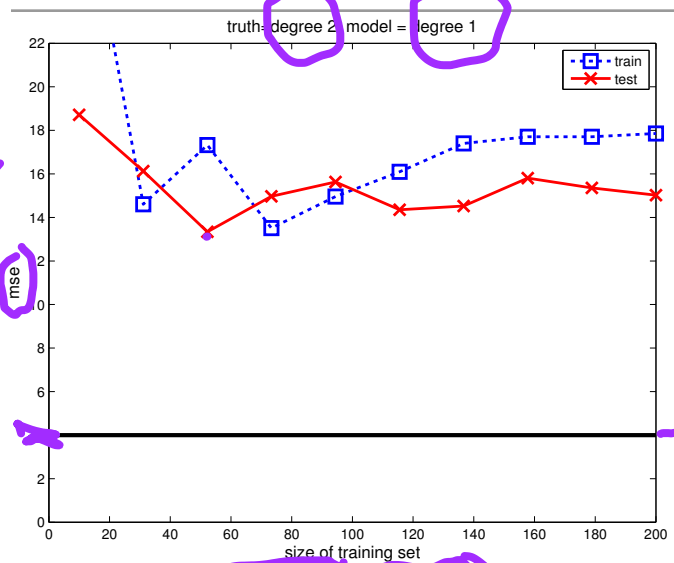
$$\hat{\mathbf{w}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

What happens when \lambda is 0?
What happens as \lambda increases?
What happens when \lambda is infinity?

Figures from Murphy's MLPP book

# Regularization effects of big data (Murphy MLPP Section 7.5.4)



$\sum |w|^2$

$\ell_2$ norm

$\ell_2$ reg.

overfit