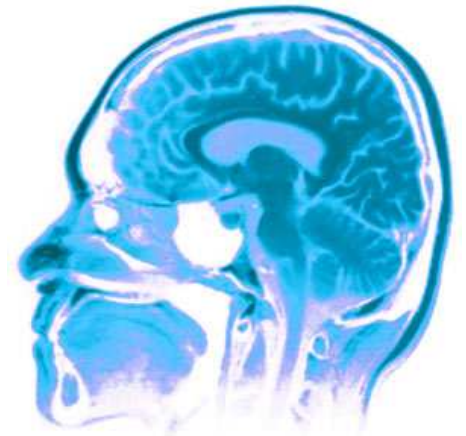# CPSC540

# Regularization, nonlinear prediction and generalization

Nando de Freitas

*Janurary, 2013*

*University of British Columbia*

# Outline of the lecture

This lecture will teach you how to fit nonlinear functions by using bases functions and how to control model complexity. The goal is for you to:

❑ Learn how to derive **ridge regression**.

❑ Understand the trade-off of fitting the data and **regularizing** it.

❑ Learn **polynomial regression**.

❑ Understand that, if basis functions are given, the problem of learning the parameters is still linear.

❑Learn **cross-validation.**

❑ Understand the effects of the number of data and the number of basis functions on **generalization**.

# Regularization

All the answers so far are of the form

$$\widehat{\boldsymbol{\theta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

*(handwritten annotations: $d{\times}1$ under $\widehat{\boldsymbol{\theta}}$; $d{\times}n$, $n{\times}d$ under $\mathbf{X}^T\mathbf{X}$; $d{\times}n$, $n{\times}1$ under $\mathbf{X}^T\mathbf{y}$)*

They require the inversion of $\mathbf{X}^T\mathbf{X}$. This can lead to problems if the system of equations is poorly conditioned. A solution is to add a small element to the diagonal:

$$\widehat{\boldsymbol{\theta}} = (\mathbf{X}^T\mathbf{X} + \delta^2 I_d)^{-1}\mathbf{X}^T\mathbf{y}$$

*(handwritten: R)*

This is the ridge regression estimate. It is the solution to the following **regularised quadratic cost function**

*(handwritten: scalar)*

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \delta^2\boldsymbol{\theta}^T\boldsymbol{\theta}$$

*(handwritten: Penalty regulariter)*

# Derivation

$$J(\theta) = (Y - X\theta)^T (Y - X\theta) + \delta^2 \theta^T \theta$$

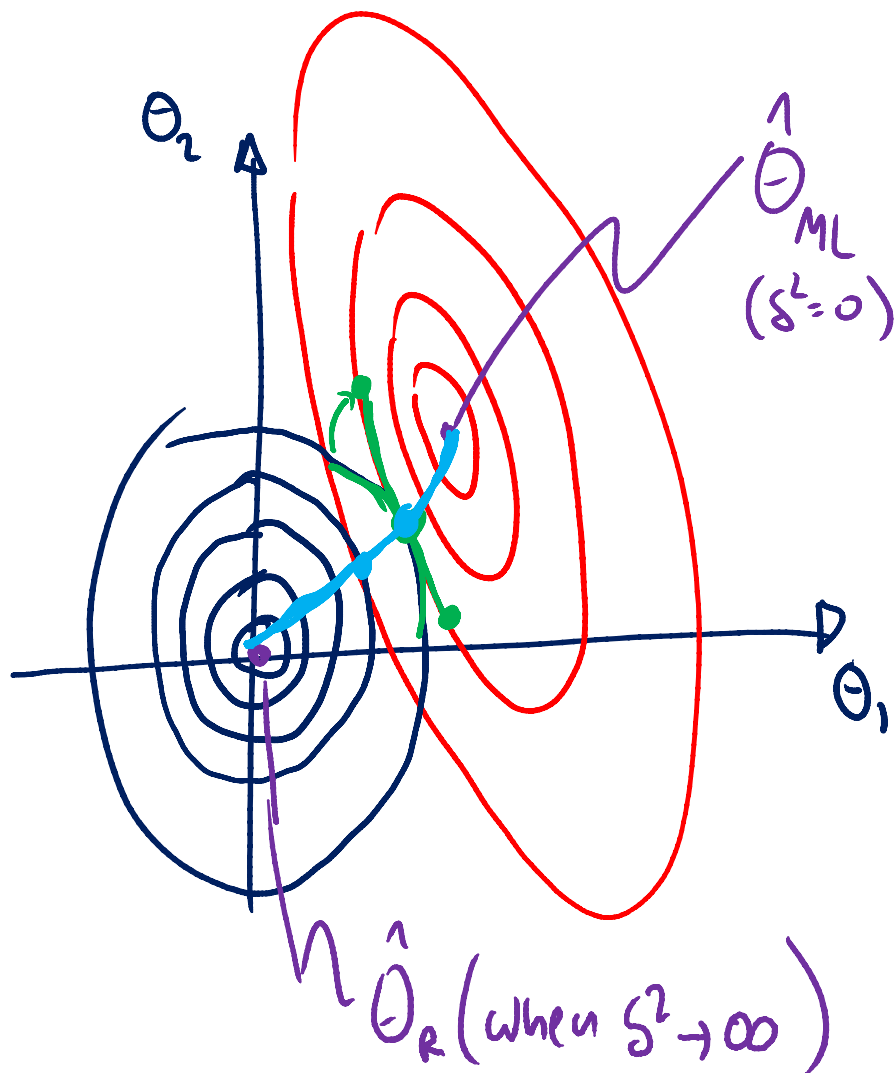$$\frac{\partial J(\theta)}{\partial \theta} = 2X^T X \theta - 2X^T y + 2\delta^2 I \theta$$

equating to zero

$$\left(X^T X + \delta^2 I\right)\theta = X^T y$$

# Ridge regression as constrained optimization

ellipses

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \delta^2\boldsymbol{\theta}^T\boldsymbol{\theta}$$

$$\min_{\boldsymbol{\theta} \,:\, \boxed{\boldsymbol{\theta}^T\boldsymbol{\theta} \le t(\delta)}} \left\{(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\right\}$$



in 2D

$$\underline{\theta} = (\theta_1, \theta_2)$$

$$\underline{\theta}^T\underline{\theta} = \theta_1^2 + \theta_2^2$$

$$\theta_1^2 + \theta_2^2 \le t^2$$

$\hat{\Theta}_{ML}$ ($\delta^2 = 0$)

$\hat{\Theta}_R$ (when $\delta^2 \to \infty$)

# Regularization paths

*As $\delta$ increases, $t(\delta)$ decreases and each $\theta_i$ goes to zero.*



[Hastie, Tibshirani & Friedman book]

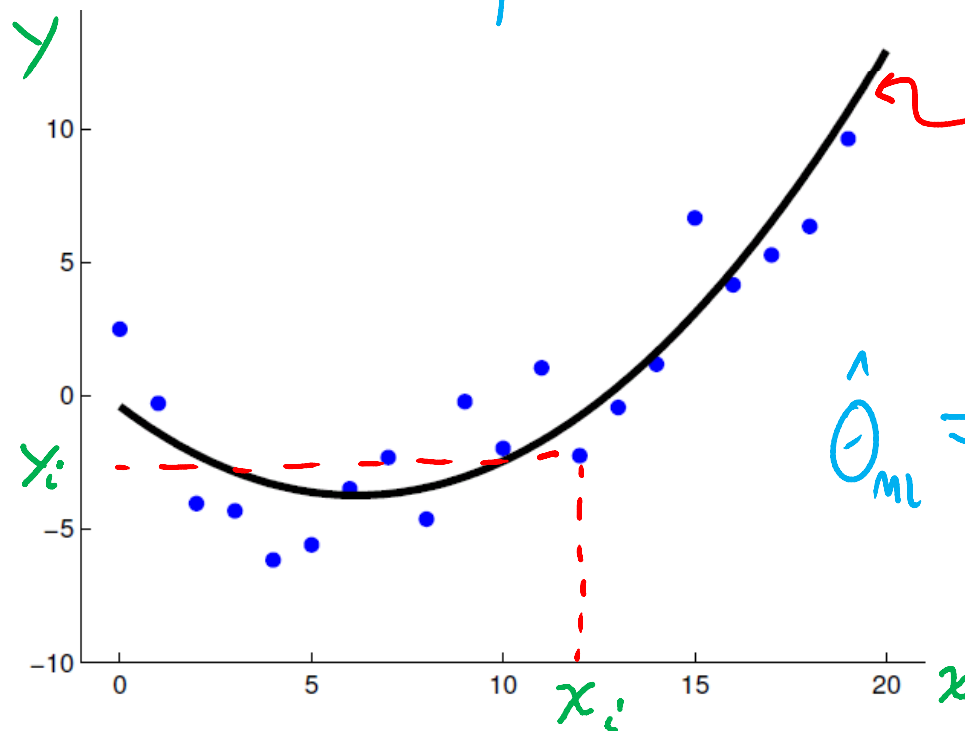# Going nonlinear via basis functions

We introduce basis functions $\phi(\cdot)$ to deal with nonlinearity:

$$y(\mathbf{x}) = \phi(\mathbf{x})\boldsymbol{\theta} + \epsilon$$

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



$\hat{Y} = \phi(x)\Theta$

$= \Theta_0 + x\Theta_1 + x^2\Theta_2$

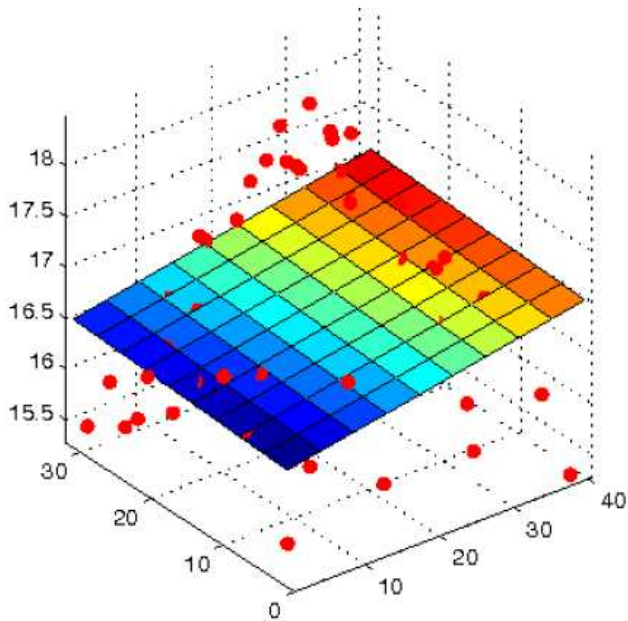$\hat{\Theta}_{ML} = \left[\phi_{(x)}^T \phi_{(x)}\right]^{-1} \phi_{(x)}^T Y$

# Going nonlinear via basis functions

$$y(\mathbf{x}) = \phi(\mathbf{x})\boldsymbol{\theta} + \epsilon$$

$$\phi(\mathbf{x}) = [1, x_1, x_2]$$

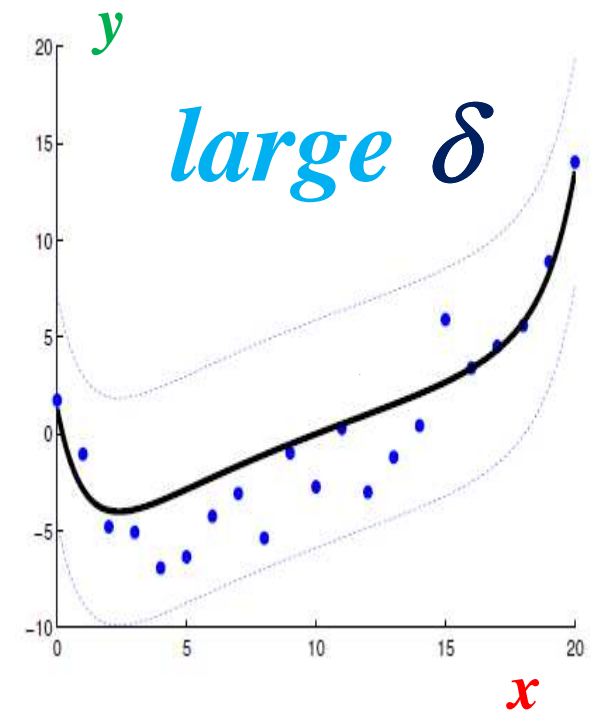$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2]$$
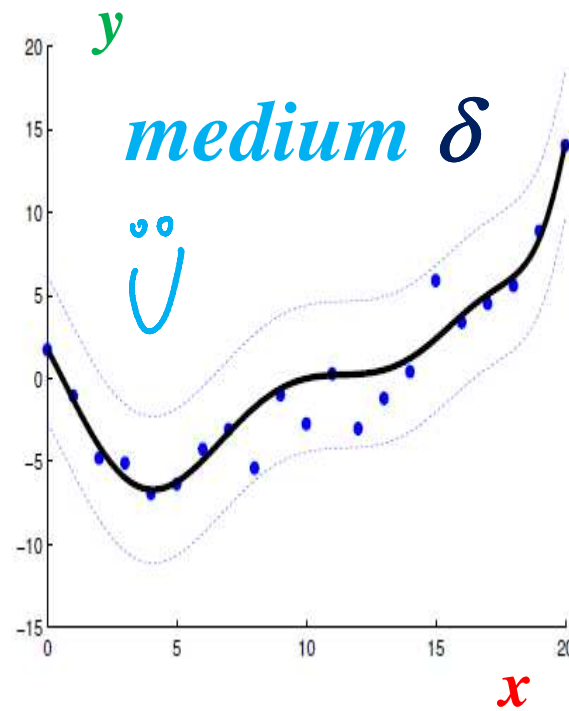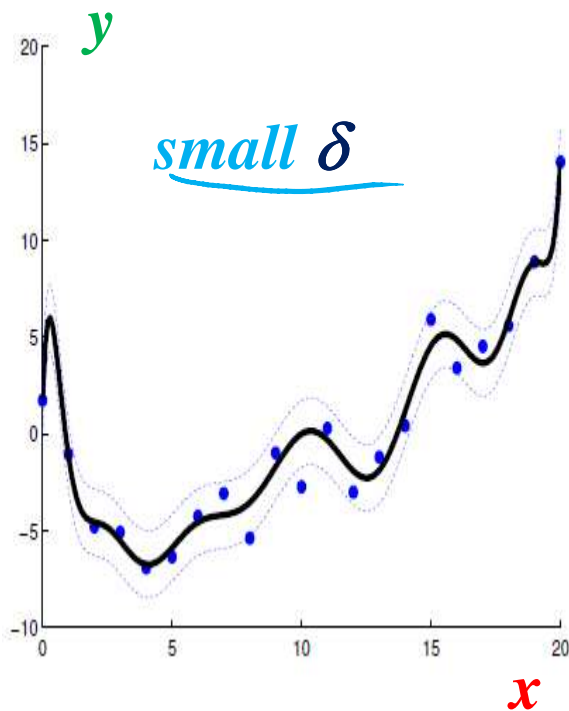
$x_1 x_2$

# Example: Ridge regression with a polynomial of degree 14

$$\hat{y}(x_i) = 1\,\theta_0 + x_i\,\theta_1 + x_i^2\,\theta_2 + \ldots + x_i^{13}\,\theta_{13} + x_i^{14}\,\theta_{14}$$

$$\Phi = [\,1 \quad x_i \quad x_i^2 \quad \ldots \quad x_i^{13} \quad x_i^{14}\,]\, x_i^{15} \ldots$$

$$J(\theta) = (\,y - \Phi\theta\,)^T(\,y - \Phi\theta\,) + \delta^2\,\theta^T\theta$$
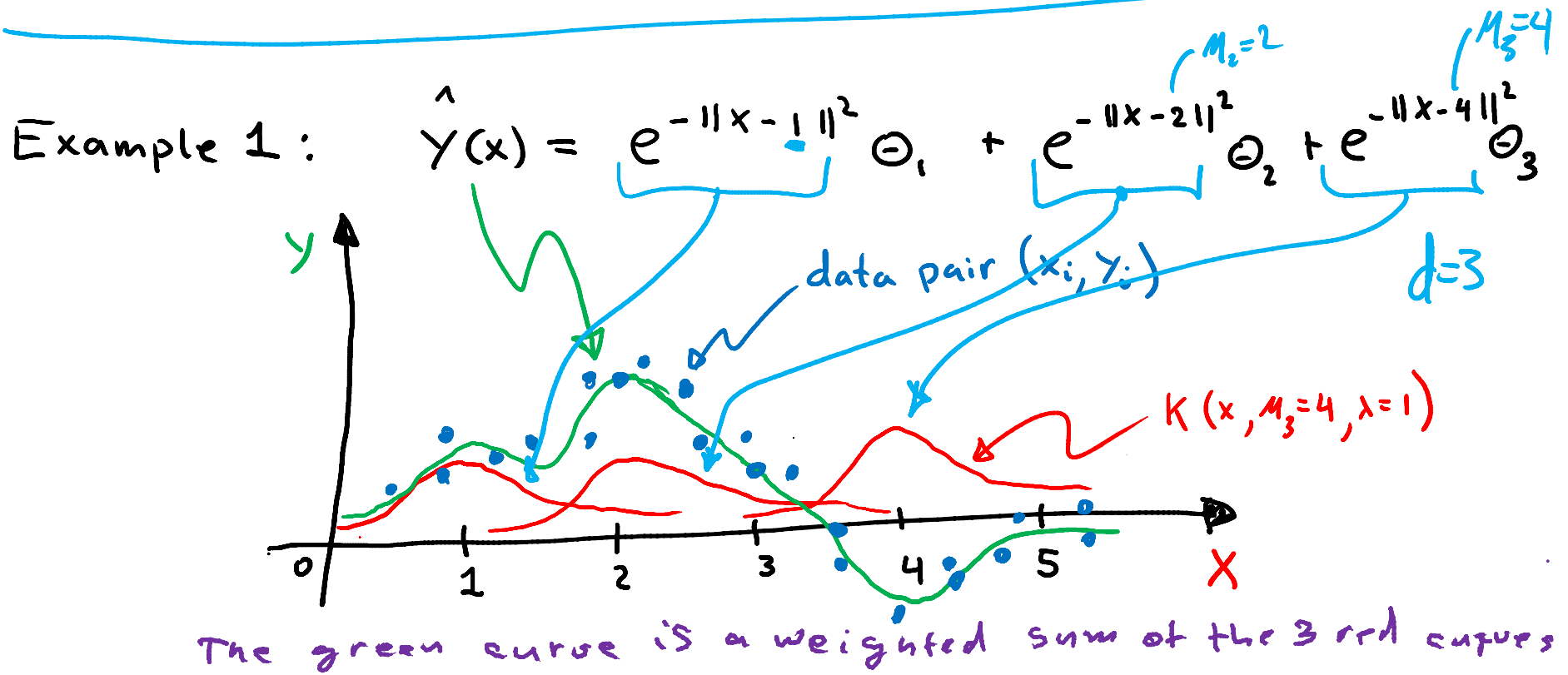


small $\delta$

medium $\delta$

large $\delta$

# Kernel regression and RBFs

We can use kernels or radial basis functions (RBFs) as features:

$$\phi(\mathbf{x}) = [\kappa(\mathbf{x}, \boldsymbol{\mu}_1, \lambda), \ldots, \kappa(\mathbf{x}, \boldsymbol{\mu}_d, \lambda)], \quad e.g. \quad \kappa(\mathbf{x}, \boldsymbol{\mu}_i, \lambda) = e^{(-\frac{1}{\lambda}\|\mathbf{x}-\boldsymbol{\mu}_i\|^2)}$$

$$\hat{y}(x_i) = \phi(x_i)\,\theta = 1\theta_0 + k(x_i, \mu_1, \lambda)\theta_1 + \ldots + k(x_i, \mu_d, \lambda)\theta_d$$

Example 1:

$$\hat{Y}(x) = e^{-\|x-1\|^2}\Theta_1 + e^{-\|x-2\|^2}\Theta_2 + e^{-\|x-4\|^2}\Theta_3$$

$\mu_2 = 2$    $\mu_3 = 4$

data pair $(x_i, y_i)$

$d = 3$

$K(x, \mu_3 = 4, \lambda = 1)$

The green curve is a weighted sum of the 3 red curves

$$\phi(x_i) = \begin{bmatrix} 1 & K(x_i, m_1, \lambda) & K(x_i, m_2, \lambda) & K(x_i, m_3, \lambda) \end{bmatrix}$$

<span style="color:red">$\phi(x_i)$ is a vector with 4 entries. There are 3 bases.</span>

<span style="color:blue">The corresponding vector of parameters is $\underline{\Theta} = \begin{bmatrix} \Theta_0 & \Theta_1 & \Theta_2 & \Theta_3 \end{bmatrix}^T$</span>

$$\hat{Y}_i = \phi(x_i)\underline{\Theta}$$

If we have $i = 1, \ldots, N$ data, let

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_N \end{bmatrix} \qquad \underline{\Phi} = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_N) \end{bmatrix}$$
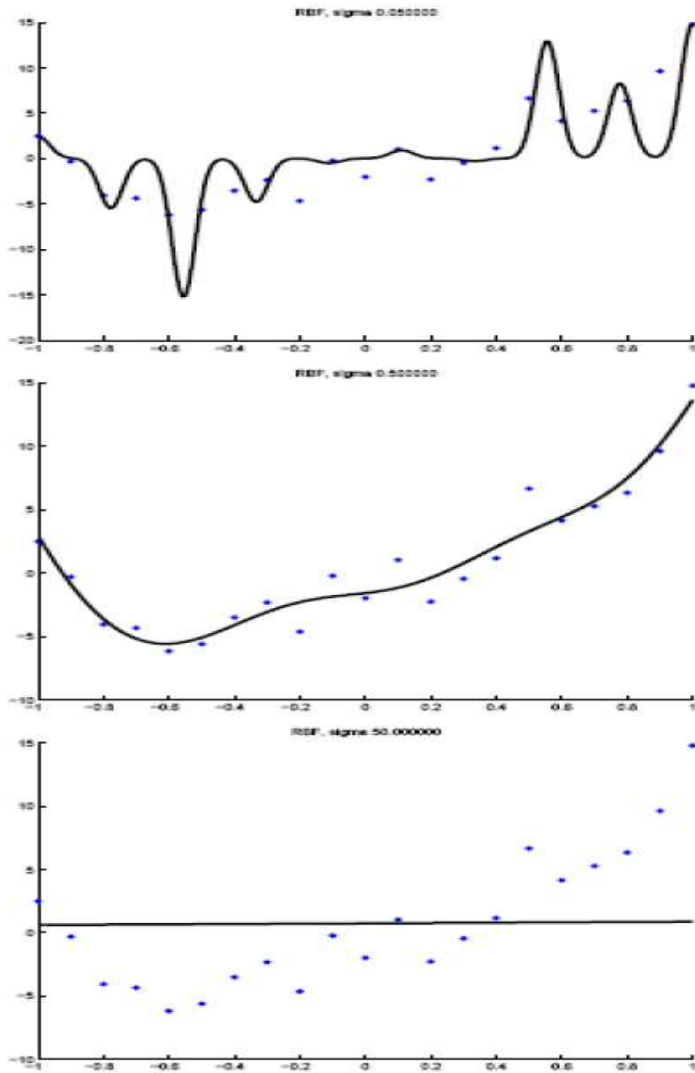
Then

$$\hat{Y} = \bar{\Phi} \Theta$$

and

$$\hat{\Theta}_{LS} = \left( \Phi^T \Phi \right)^{-1} \Phi^T Y$$

or

$$\hat{\Theta}_{ridge} = \left( \Phi^T \Phi + \delta^2 I \right)^{-1} \Phi^T Y$$

Hence, this is still linear regression, with $X$ replaced by $\Phi$.

We can choose the locations $\mu$ of the **basis functions** to be the inputs. That is, $\mu_i = x_i$. These basis functions are the known as **kernels**. The choice of width $\lambda$ is tricky, as illustrated below.

**kernels**



Too small $\lambda$

Right $\lambda$

Too large $\lambda$

The big question is how do we choose the regularization coefficient, the width of the kernels or the polynomial order?
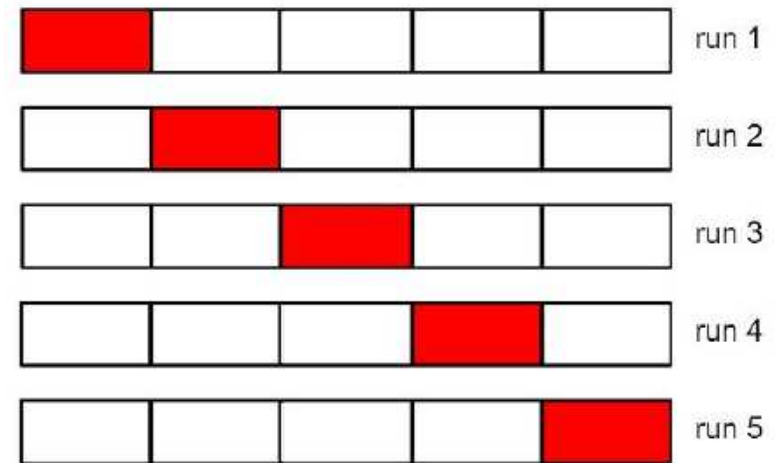
# One Solution: **cross-validation**

① Given training data $(X, Y_{train})$, and some $\delta^2$ guess, compute $\hat{\Theta}$

② $\hat{Y}_{train} = X_{train} \hat{\Theta}$ (compute training set predictions)

③ $\hat{Y}_{test} = X_{test} \hat{\Theta}$

| train | test |

| validation |

| | $\delta^2$ | Train error $\sum\limits_{i \in train}(Y_i - \hat{Y}_i)^2$ | Test error $\sum\limits_{i \in test}(Y_i - \hat{Y}_i)^2$ | max | min-max | avg |
|---|---|---|---|---|---|---|
| 0 | 0.1 | 100 | 2 | 100 | | 51 |
| 35 | 1 | 10 | 11 | 11 | X | 10.5 |
| 3 | 10 | 1 | 19 | 19 | | 10 ✗ |
| 21 | 50 | 20 | 0 | 20 | | 10 |
| 0 | 100 | 100 | 1000 | 1000 | | 550 |

# K-fold crossvalidation



The idea is simple: we split the training data into $K$ **folds**; then, for each fold $k \in \{1, \ldots, K\}$, we train on all the folds but the $k$'th, and test on the $k$'th, in a round-robin fashion.

It is common to use $K = 5$; this is called 5-fold CV.

If we set $K = N$, then we get a method called **leave-one out cross validation**, or **LOOCV**, since in fold $i$, we train on all the data cases except for $i$, and then test on $i$.

# Example: Ridge regression with polynomial of degree 14



5-fold crossvalidation error

test set error

train set error

$Log \delta^2$

$Log \delta^2$

The larger $\delta$, the larger the train set error. However the test set error improves. For good generalization we want to do well in all test sets.

# Effect of data when we have the right model

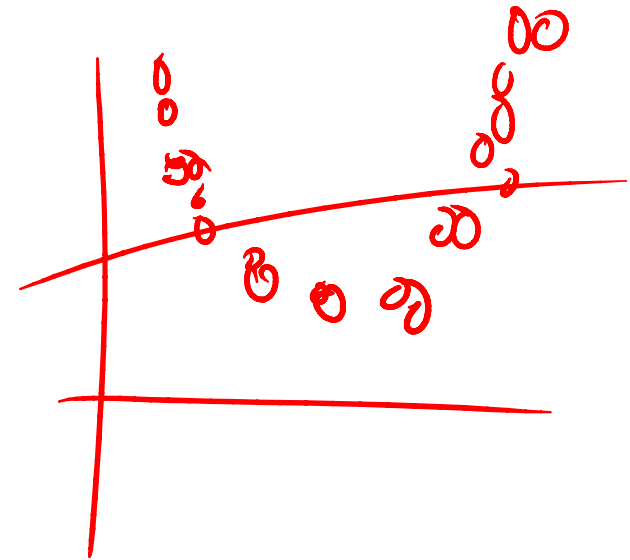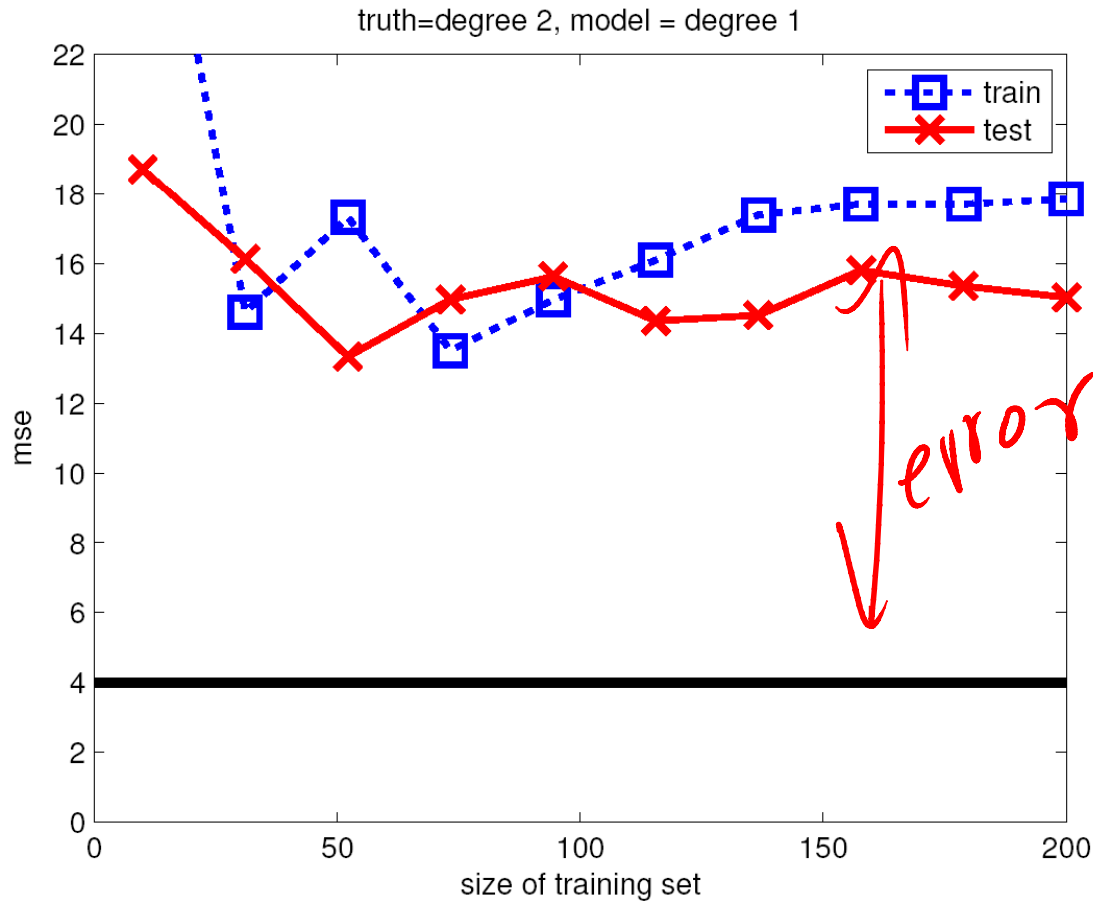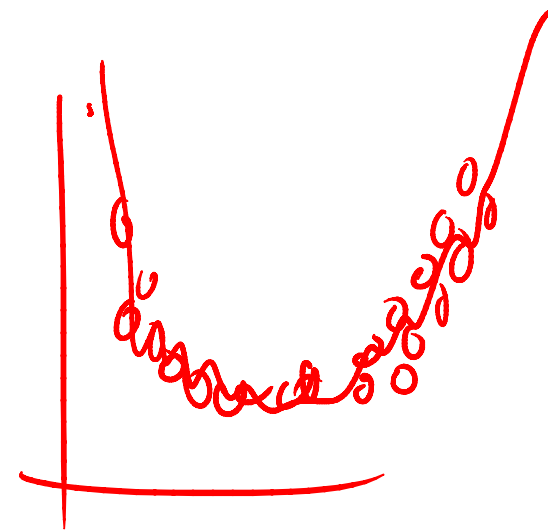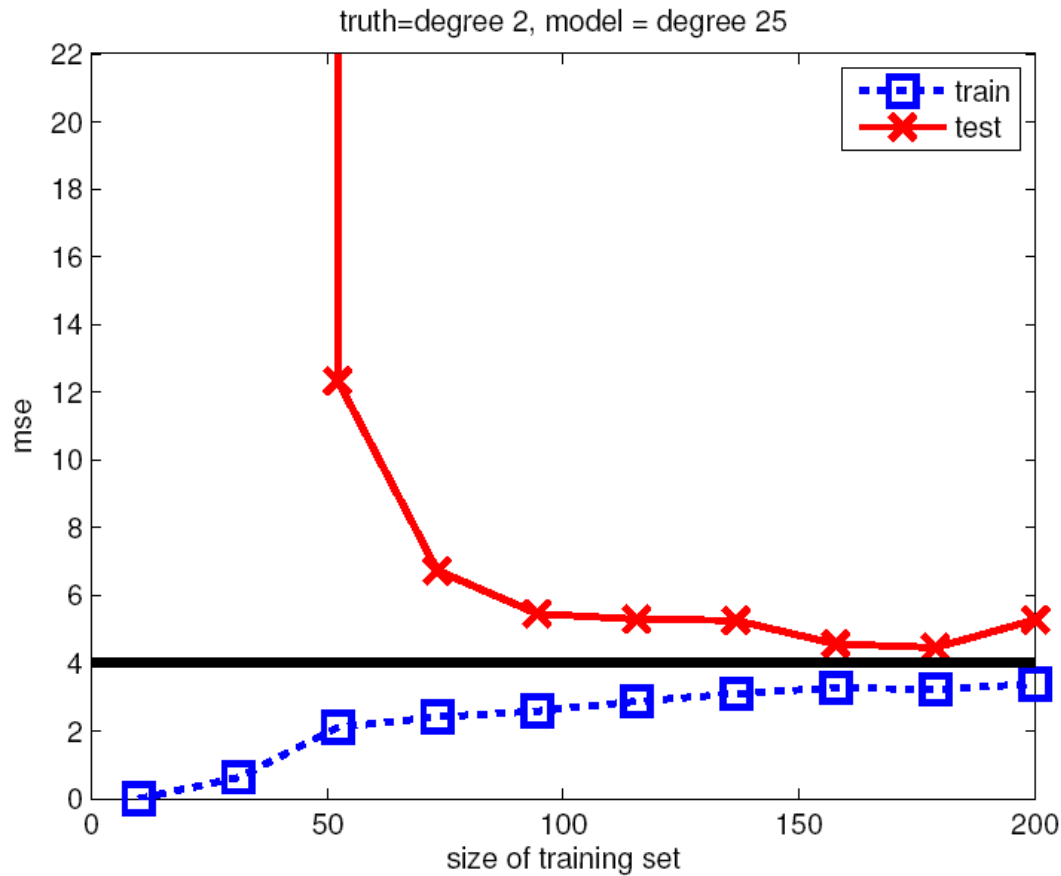$$y_i = \theta_0 + x_i\,\theta_1 + x_i^2\,\theta_2 + \mathcal{N}(\,0\,,\,\sigma^2\,)$$

$$\hat{y}_i = \hat{\Theta}_0 * + x_i\,\hat{\Theta}_1 + x_i^2\,\hat{\Theta}_2 \quad (model)$$



truth=degree 2, model = degree 2

# Effect of data when the model is too simple

$$y_i = \theta_0 + x_i\,\theta_1 + x_i^2\,\theta_2 + \mathcal{N}(\,0\,,\,\sigma^2\,)$$

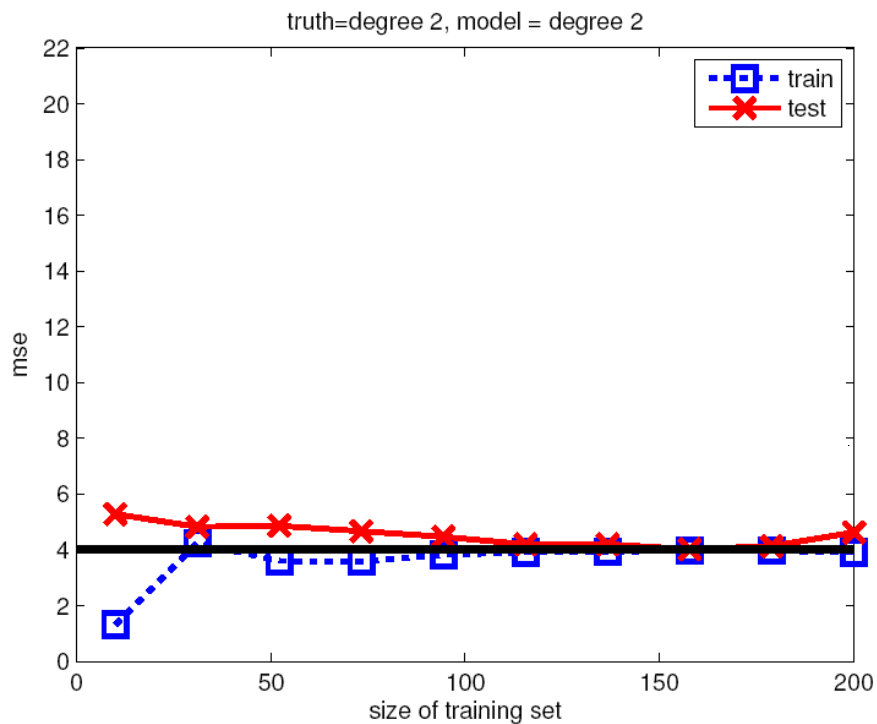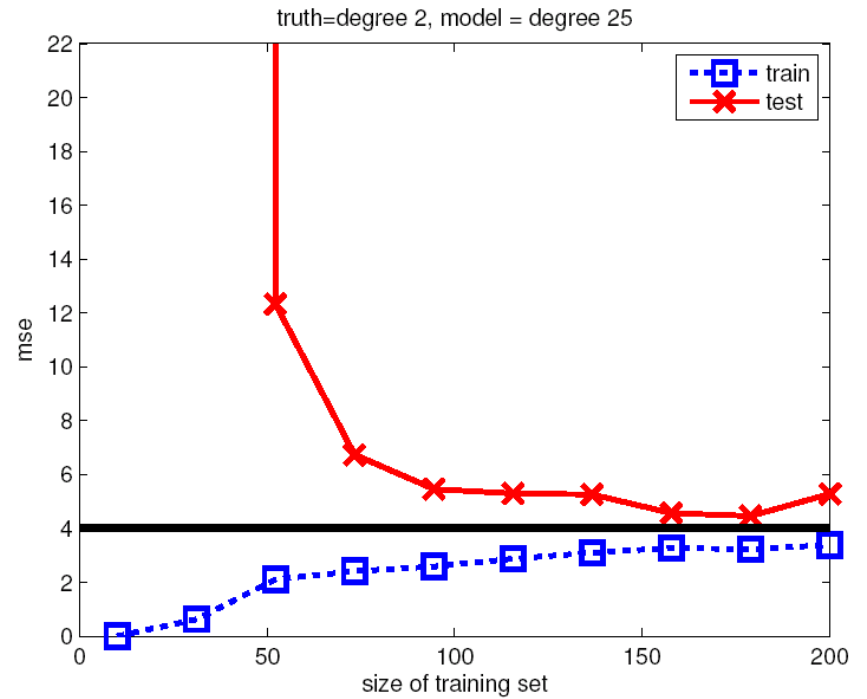$$\hat{y}_i = \hat{\theta}_0 + x_i\,\hat{\theta}_1$$



truth=degree 2, model = degree 1

error

# Effect of data when the model is very complex

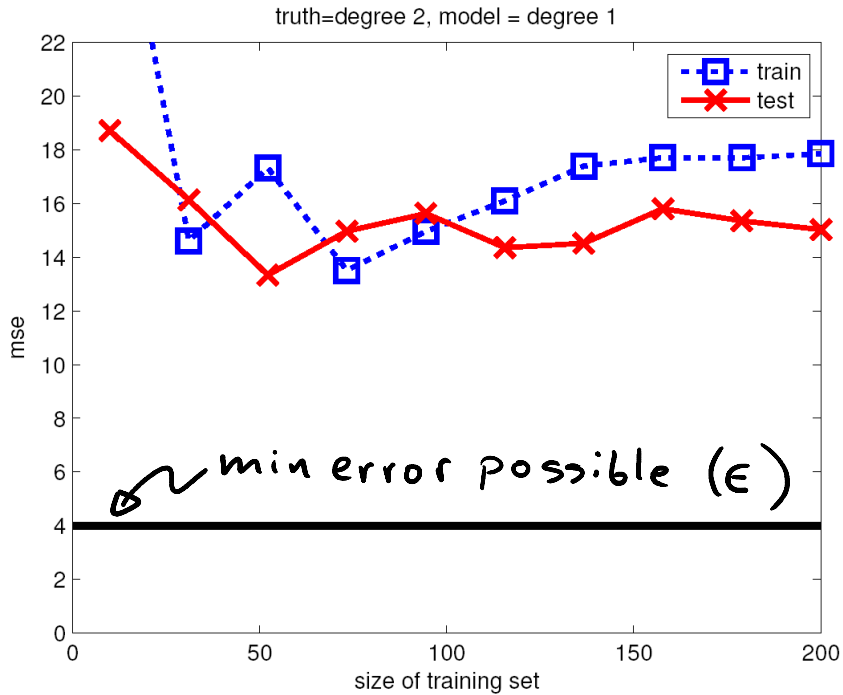$$y_i = \theta_0 + x_i\, \theta_1 + x_i^2\, \theta_2 + \mathcal{N}(\,0\,,\,\sigma^2\,)$$



truth=degree 2, model = degree 25

truth=degree 2, model = degree 1

min error possible (ε)

truth=degree 2, model = degree 25

truth=degree 2, model = degree 2

More data improves results, but only if the model has the right complexity.

# Confidence in the predictions

# Next lecture

In the next lecture, we introduce Bayesian inference, and show how it can provide us with an alternative way of learning a model from data.