



Nonlinear ridge regression Risk, regularization, and cross-validation

Nando de Freitas



UNIVERSITY OF
OXFORD

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 model complexity and **generalization**.

Regularization

All the answers so far are of the form

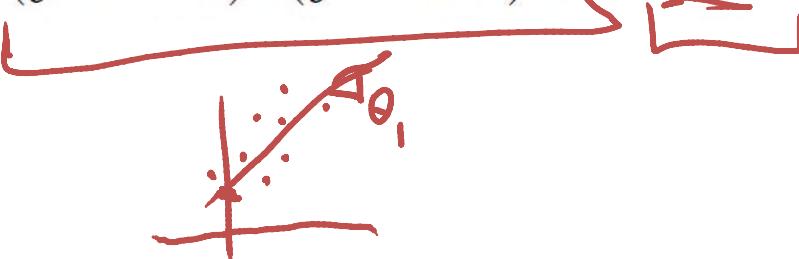
$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

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:

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

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

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



Derivation

$$J(\theta) = \frac{1}{2} \underbrace{(y - x\theta)^T (y - x\theta)}_{\text{red}} + \frac{1}{2} \theta^T \theta$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \left(\theta^T x^T x \theta - 2 y^T x \theta + y^T y + \frac{1}{2} \theta^T \theta \right)$$

$$= 2 x^T x \theta - 2 x^T y + 2 y^T y$$

$$= 2 (x^T x + \frac{1}{2} I) \theta - 2 x^T y$$

Equate to zero

$$\hat{\theta}_{\text{ridge}} = (x^T x + \frac{1}{2} I)^{-1} x^T y$$

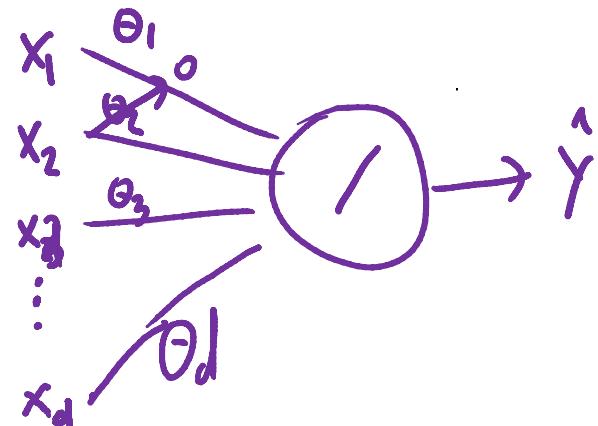
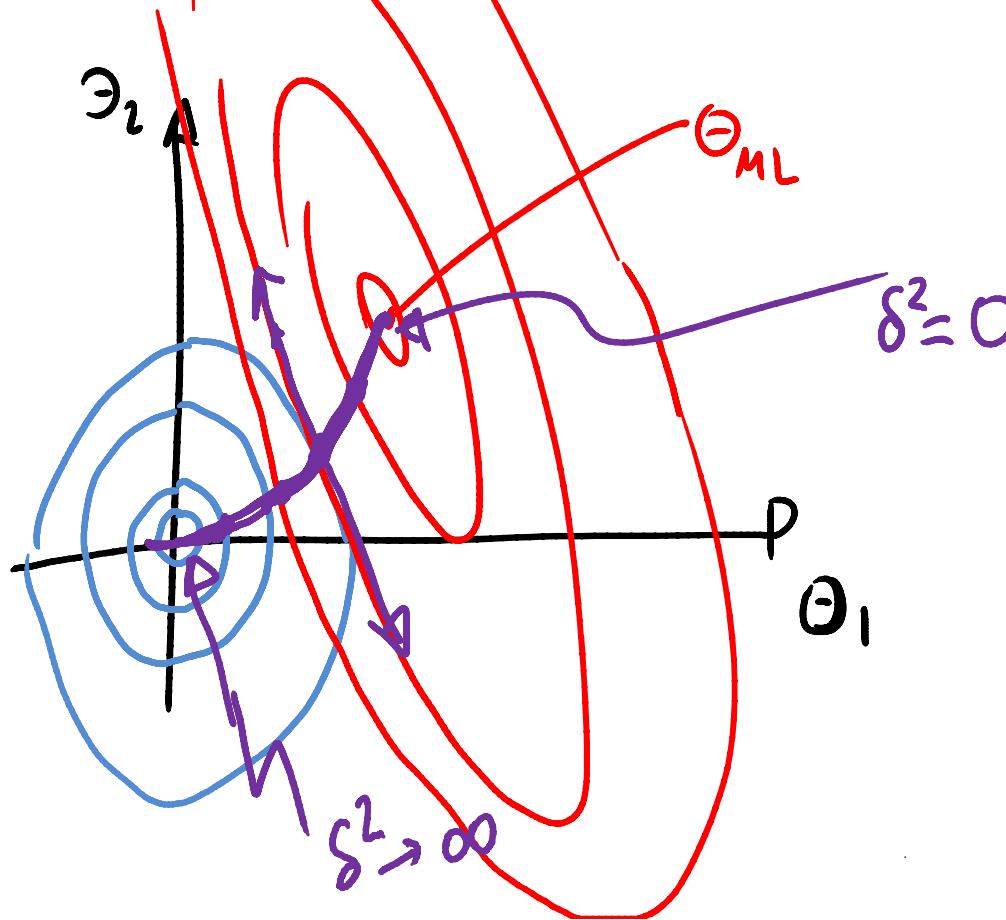
Ridge regression as constrained optimization

$$J(\theta) = (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta) + \delta^2 \theta^T \theta \quad \equiv \quad \min_{\theta : \theta^T \theta \leq t(\delta)} \{ (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta) \}$$

$$\Theta^T = [\theta_1 \quad \theta_2]$$

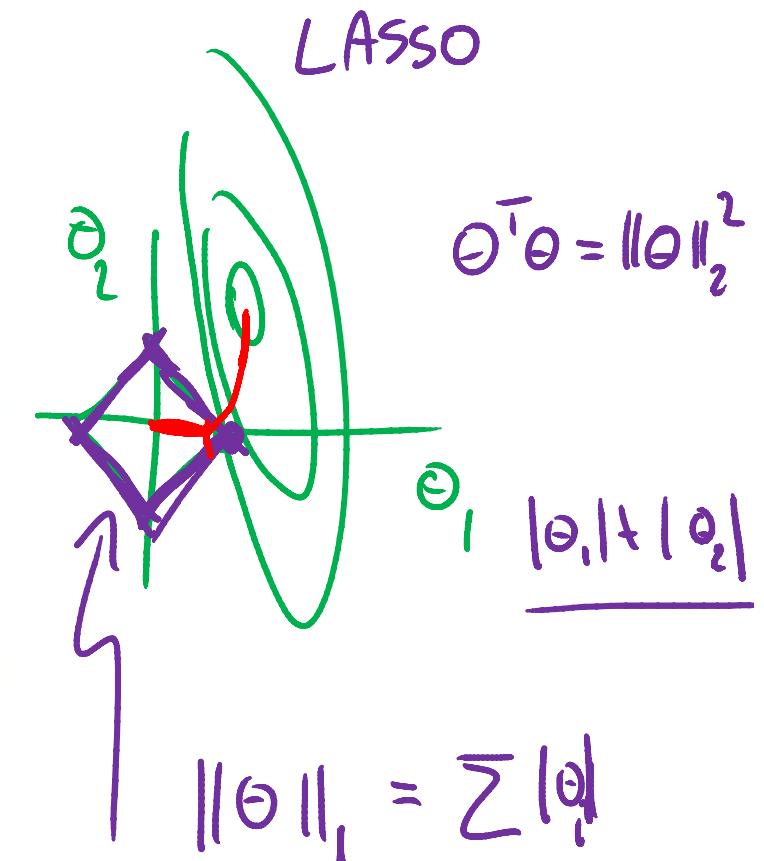
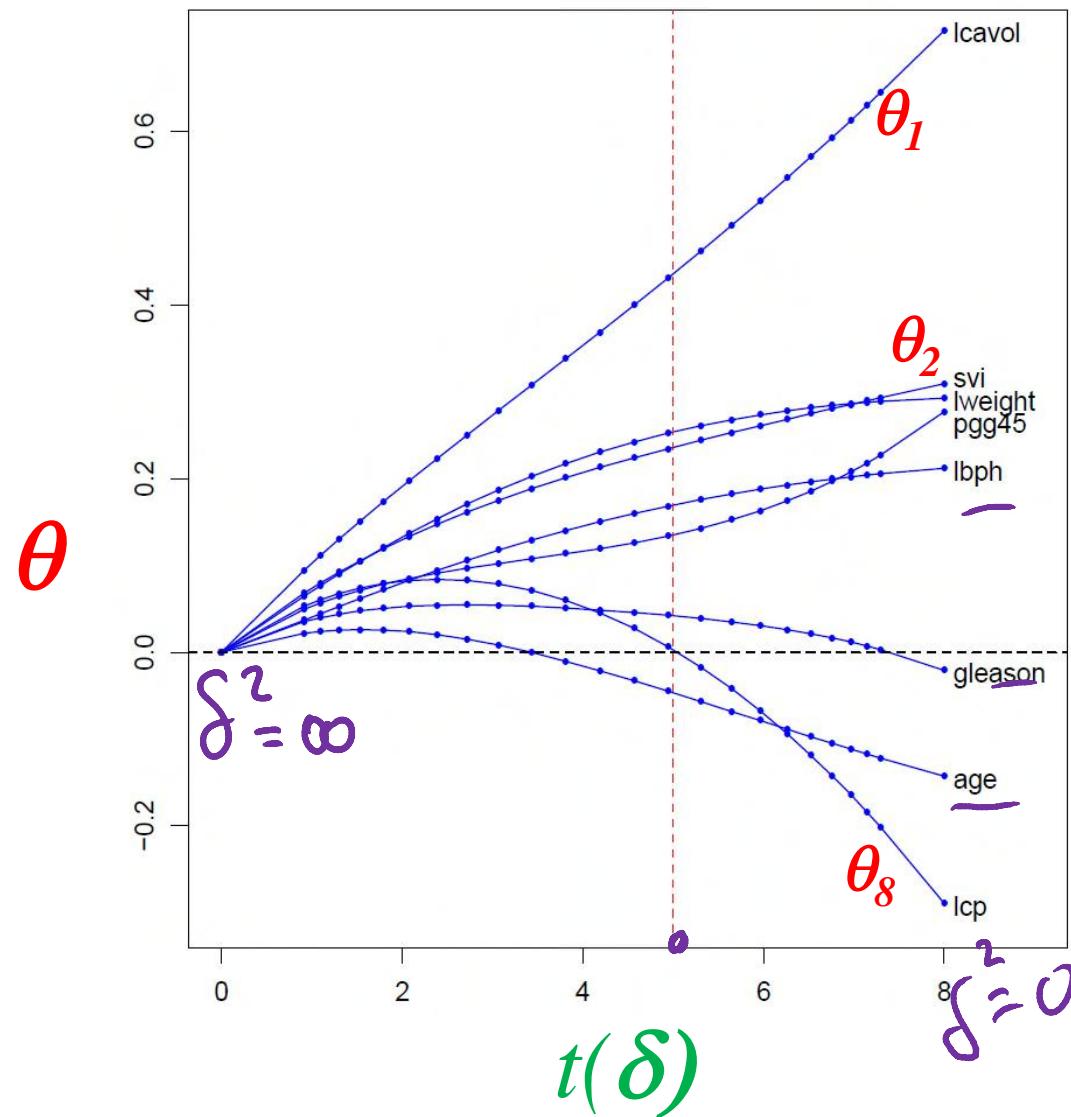
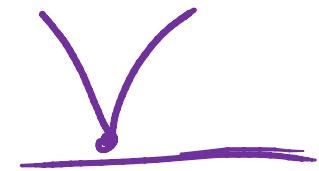
$$\Theta^T \Theta = [\theta_1 \quad \theta_2] \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$= \theta_1^2 + \theta_2^2 = \text{const}$$



Regularization paths

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



[Hastie, Tibshirani & Friedman book]

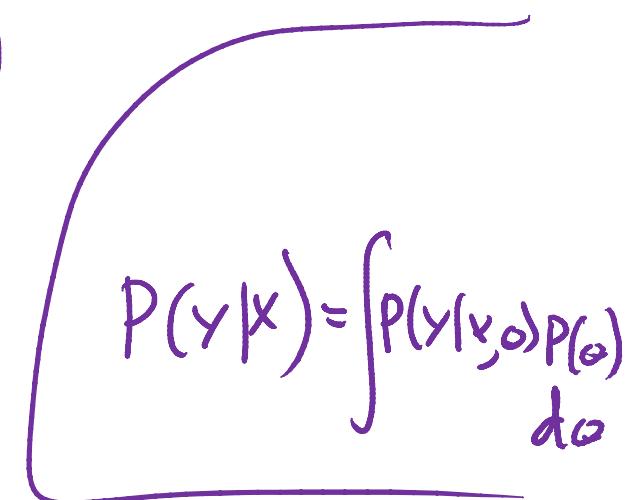
Ridge regression and Maximum a Posteriori (MAP) learning

Bayes Rule

$$J(\theta) = \underbrace{(\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta)}_{E(\theta|x,y)} + \delta^2 \underbrace{\theta^T\theta}_{\text{Prior}}$$

$$P(y|x,\theta) = \frac{1}{Z} e^{-E(\theta|x,y)}$$

Prior $\rightarrow P(\theta) = \frac{1}{Z_2} e^{-\delta^2 \theta^T \theta}$



$$P(y|x) = \int P(y|x,\theta) P(\theta) d\theta$$

$$\max_{\theta} \text{Post} P(y|x,\theta) P(\theta) \Leftrightarrow \max_{\theta} \frac{P(y|x,\theta) P(\theta)}{P(y|x)}$$

Ridge regression and Maximum a Posteriori

(MAP) learning

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$J(\theta) = \underbrace{(y - X\theta)^T (y - X\theta)}_{\text{likelihood}} + \delta^2 \theta^T \theta \quad \begin{matrix} \leftarrow \\ \text{prior} \end{matrix}$$

Posterior

$$P(\underline{\theta}|x, y) = \frac{P(y|x, \theta) P(\theta)}{P(y|x)}$$

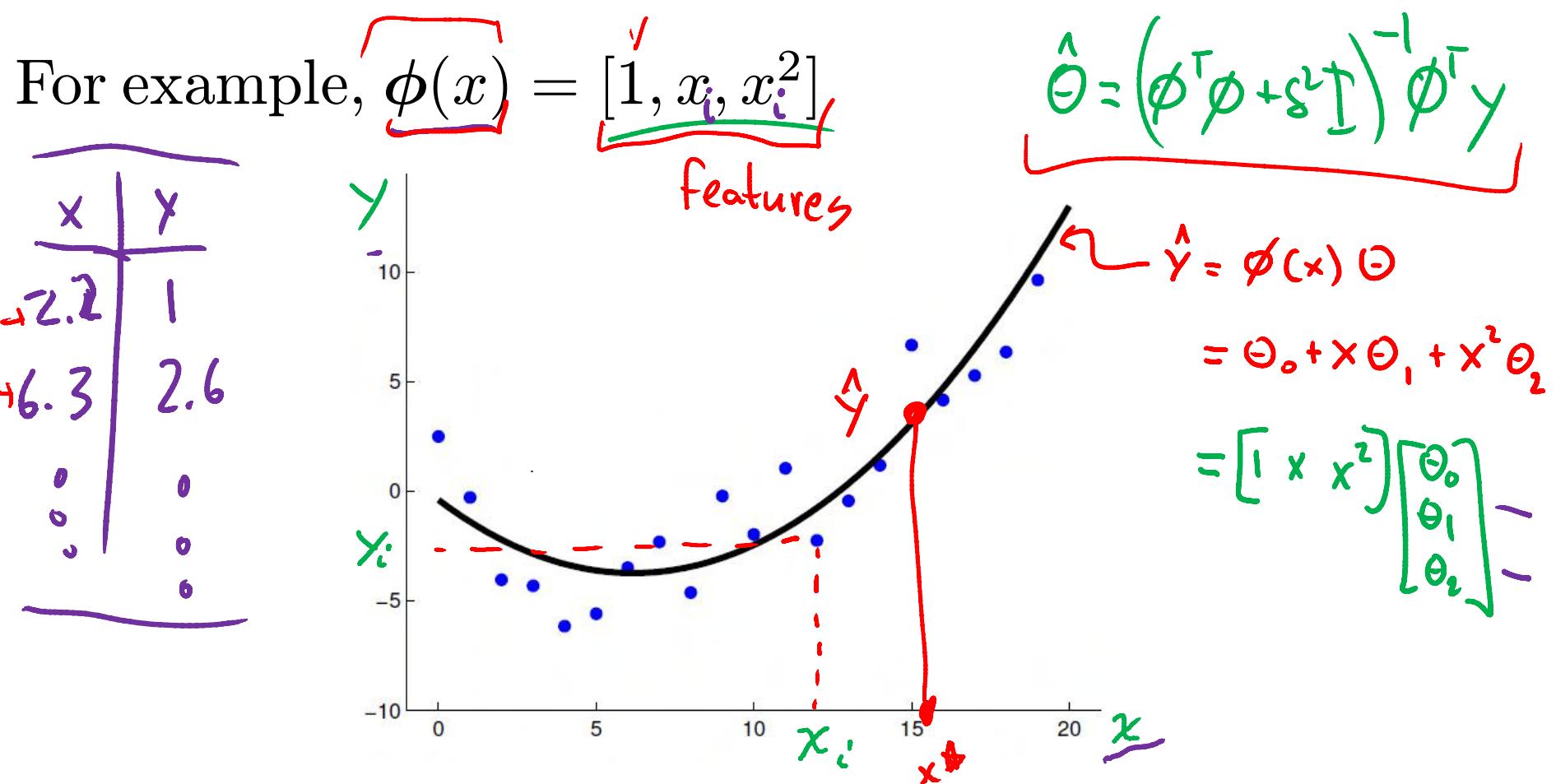
Word $\int P(y|\theta, x) P(\theta) d\theta$

$$\boxed{P(y|x, \theta) P(\theta)} = \frac{\int P(y|\theta, x) P(\theta) d\theta}{\int P(y|\theta, x) P(\theta) d\theta}$$

Going nonlinear via basis functions

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

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

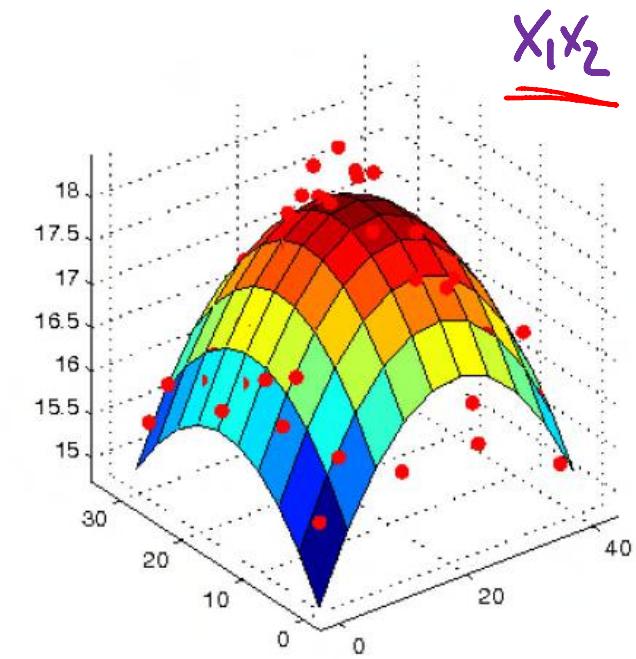
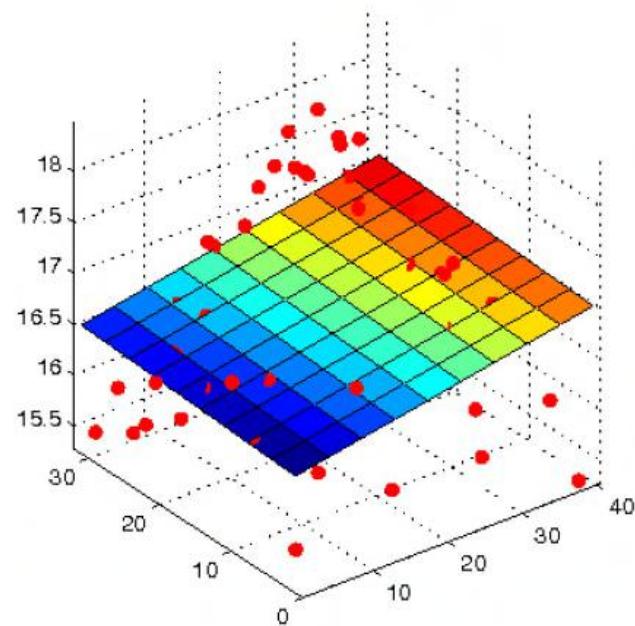


Going nonlinear via basis functions

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

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

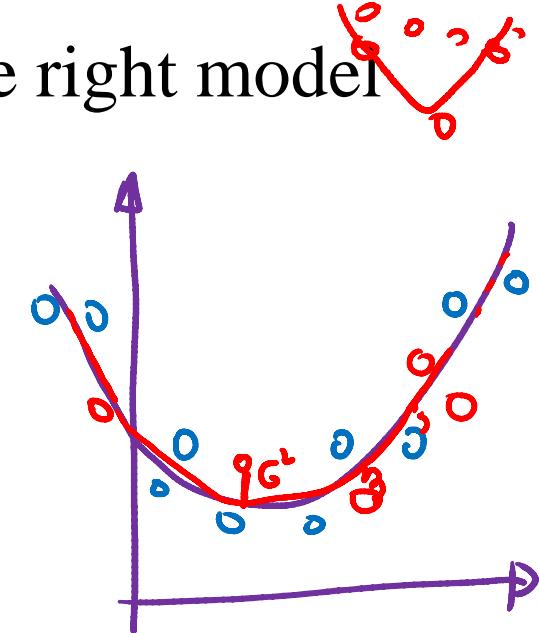
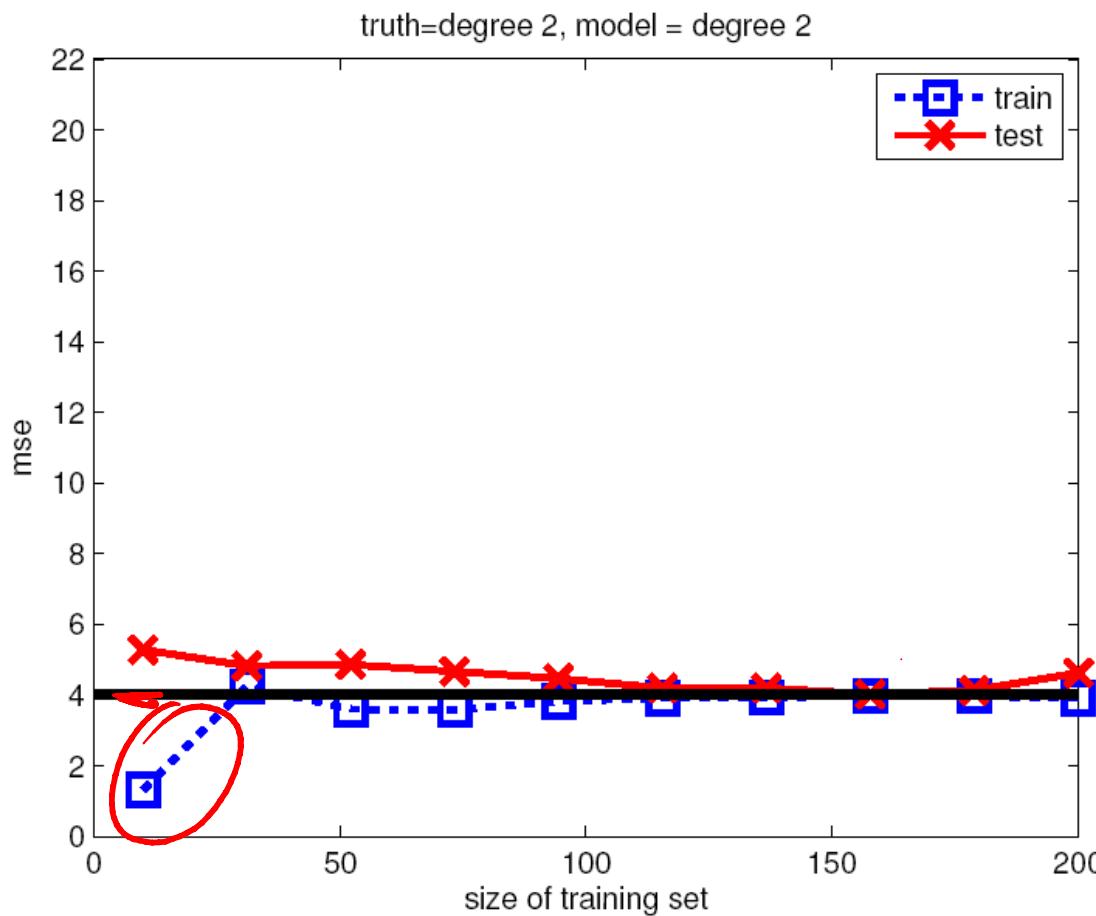
$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2]$$



Effect of data when we have the right model

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

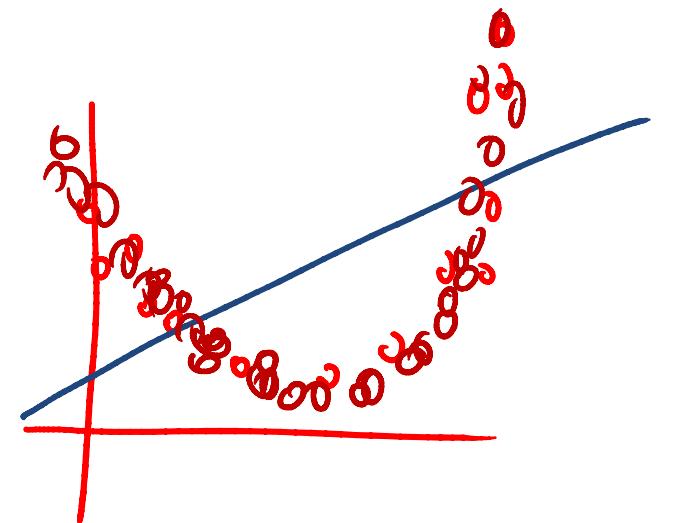
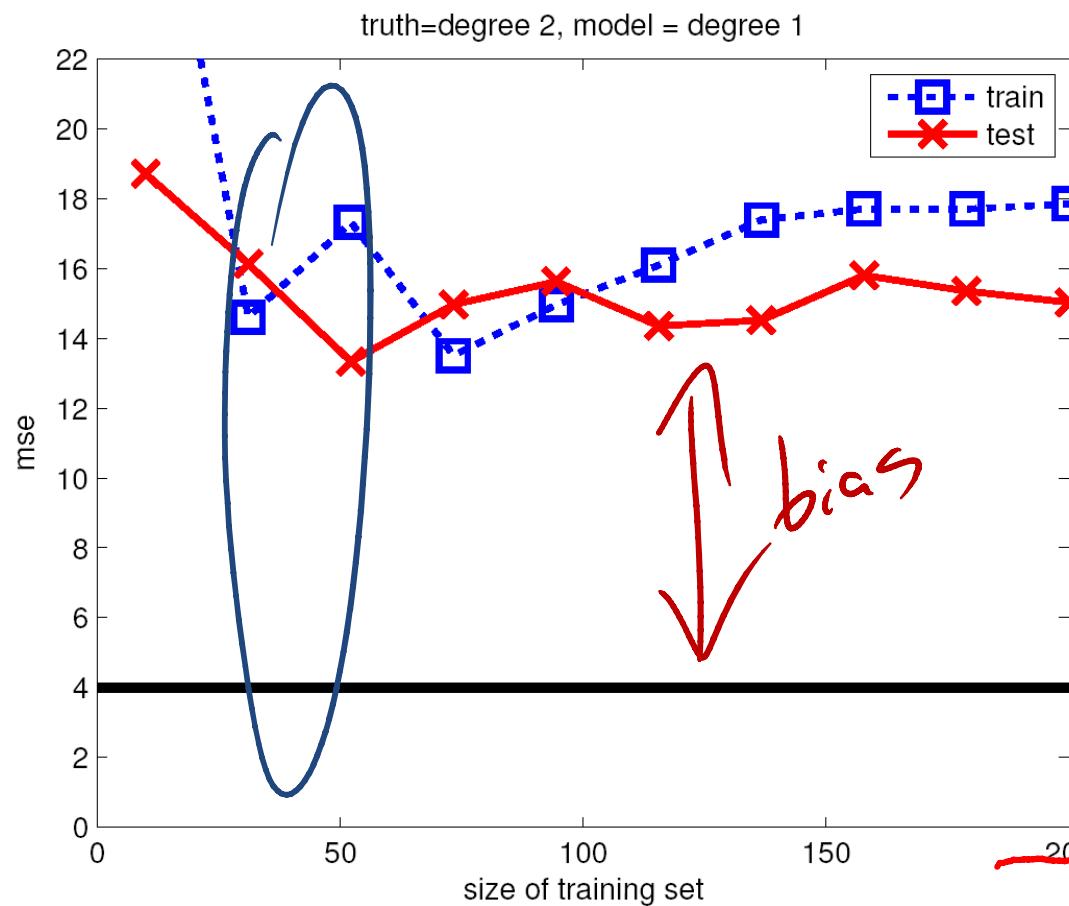
$$\hat{y} = \theta_0 + x_i \theta_1 + x_i^2 \theta_2$$



Effect of data when the model is too simple

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

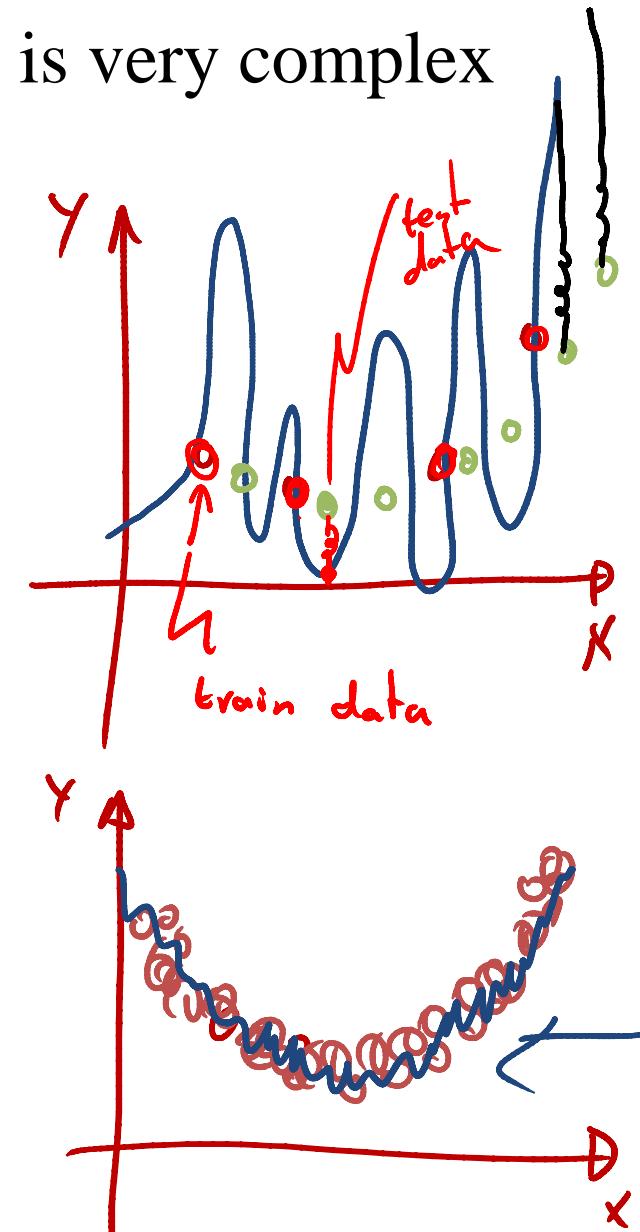
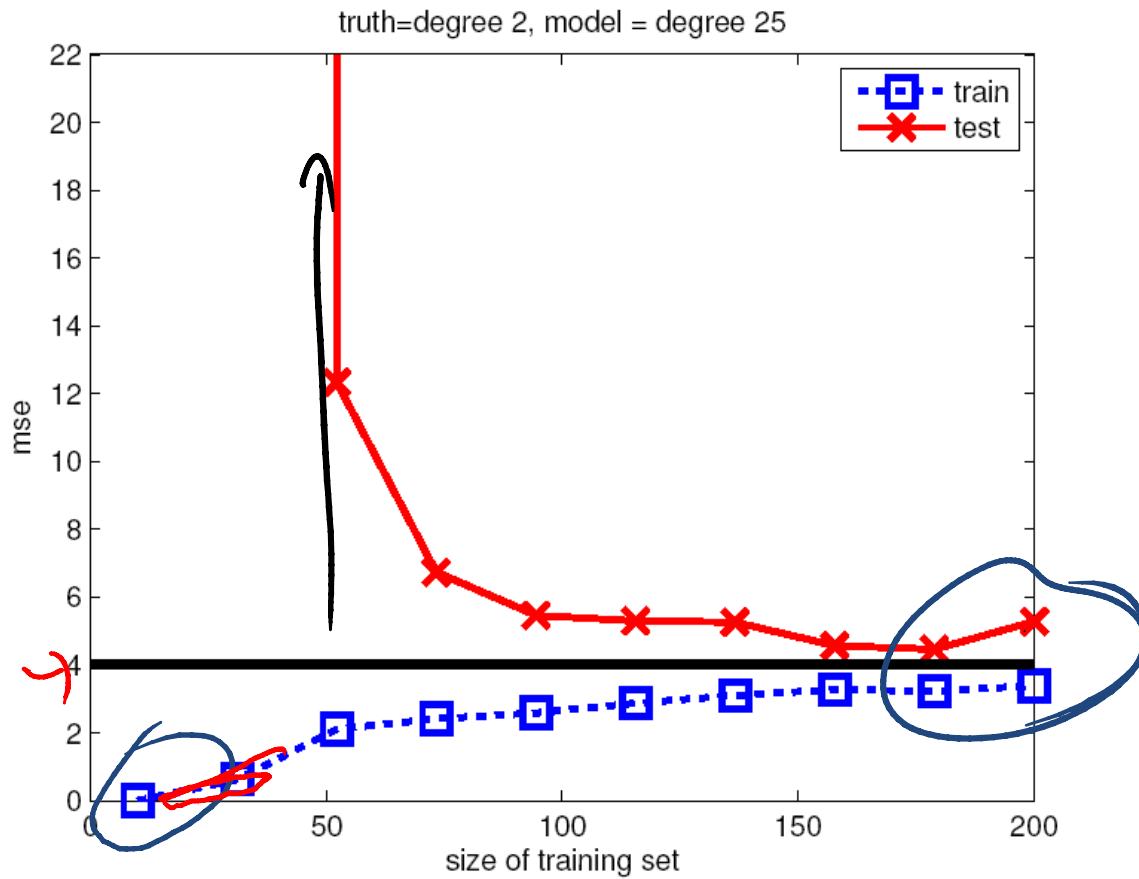
$$\hat{Y} = \theta_0 + x_i \theta_1$$

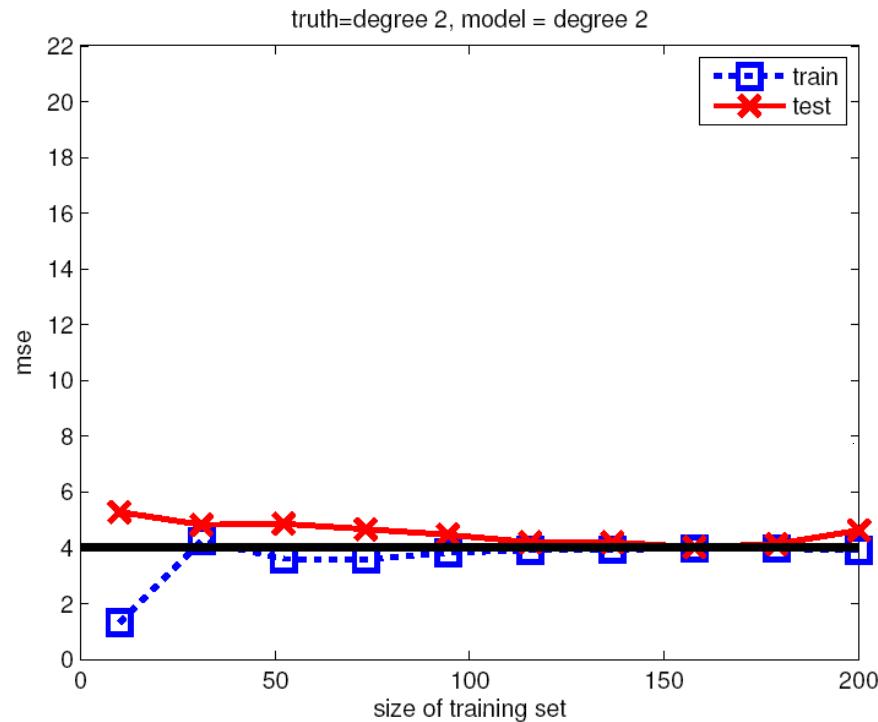
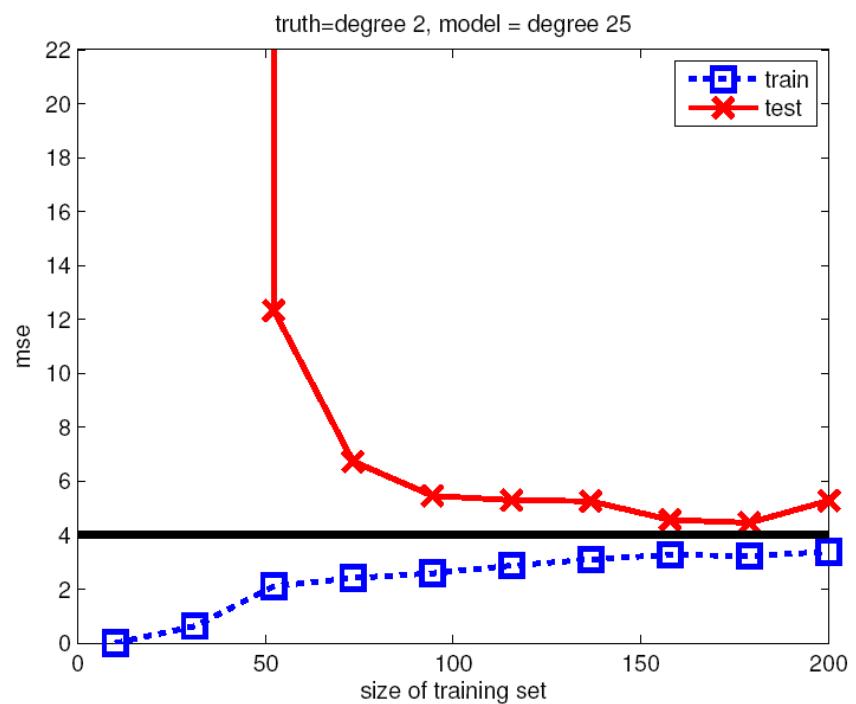
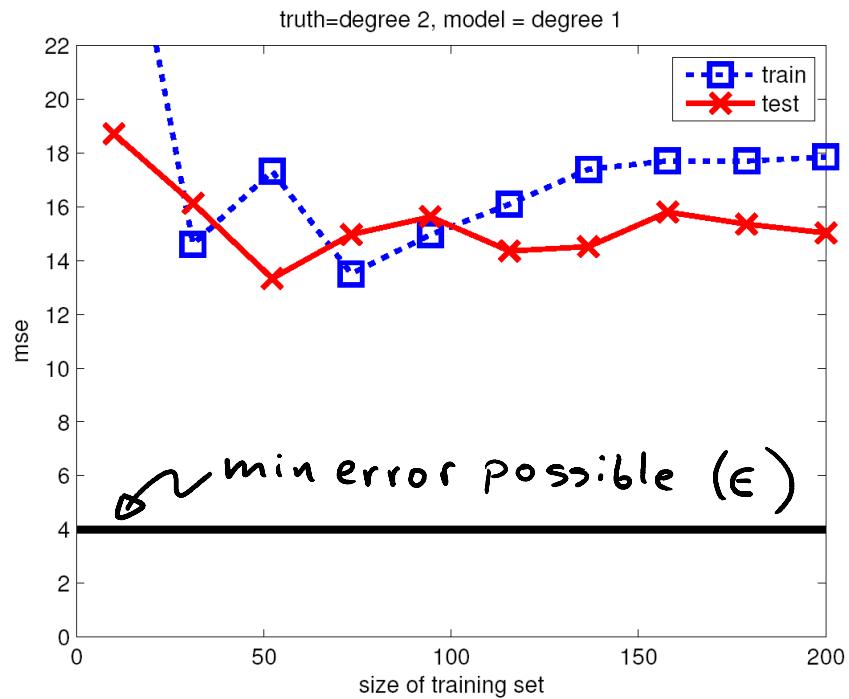


Effect of data when the model is very complex

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

$$\hat{Y}_i = \theta_0 + x_i \theta_1 + \dots + x_i^{25} \theta_{25}$$





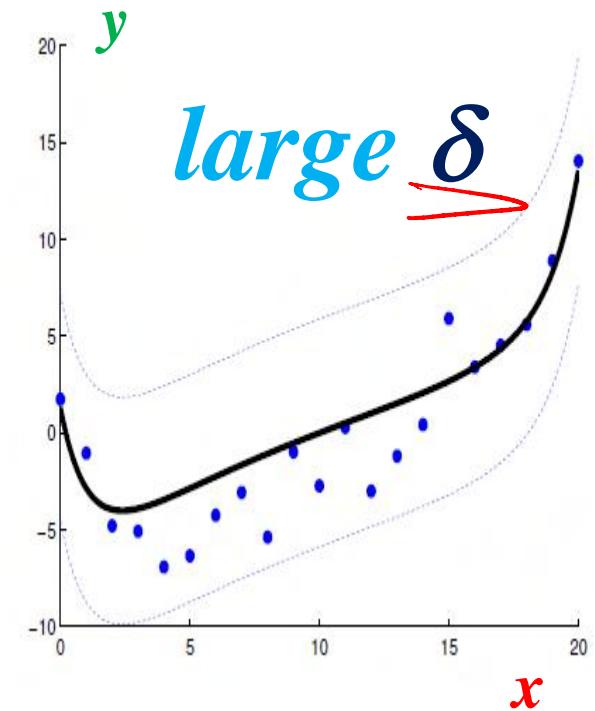
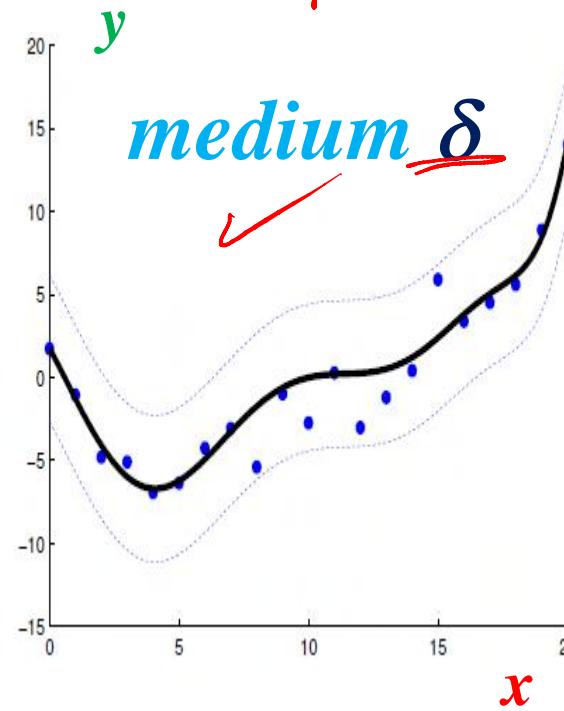
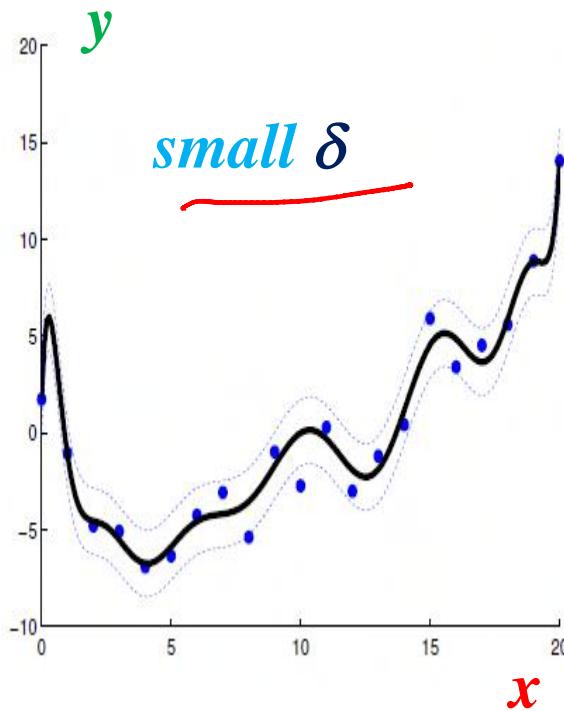
More data improves results,
but only if the model
has the right complexity.

Example: Ridge regression with a polynomial of degree 14

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

$$\Phi_i = [1 \ x_i \ x_i^2 \ \dots \ x_i^{13} \ x_i^{14}]$$

$$J(\theta) = (y - \Phi \theta)^T (y - \Phi \theta) + \underbrace{\delta^2 \theta^T \theta}_{\gamma}$$



Kernel regression and RBFs

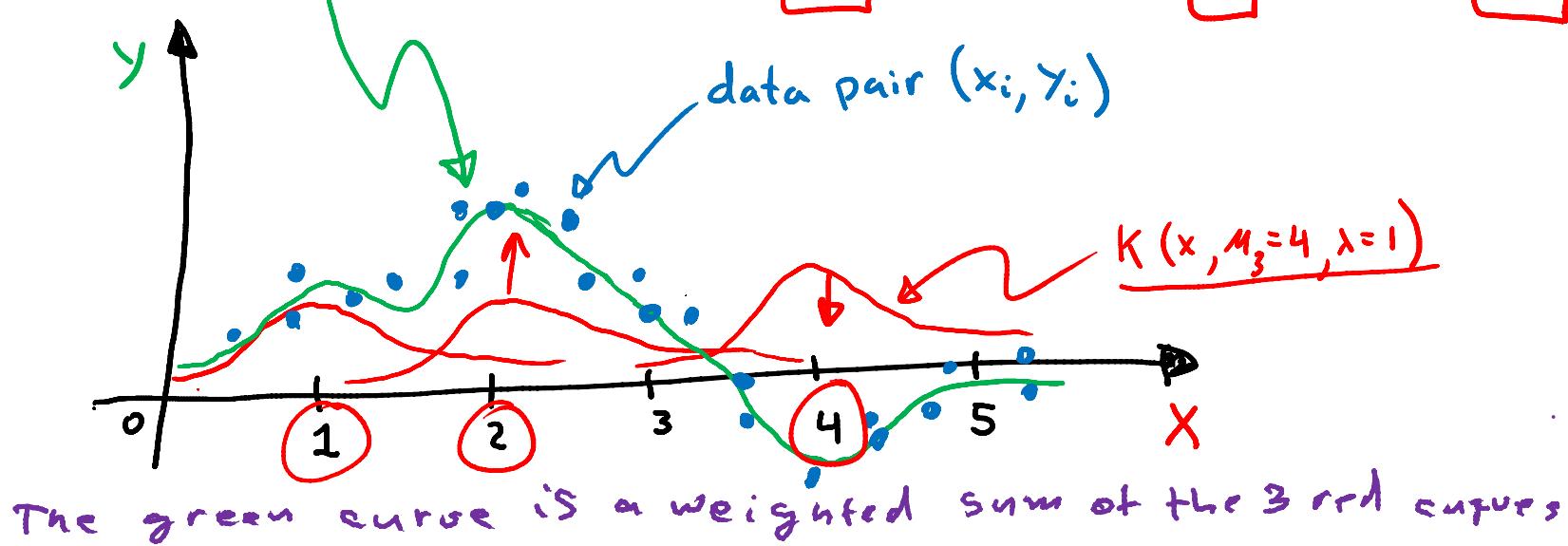
We can use kernels or radial basis functions (RBFs) as features:

$$\phi(\mathbf{x}_i) = [\kappa(\mathbf{x}_i, \mu_1, \lambda), \dots, \kappa(\mathbf{x}_i, \mu_d, \lambda)],$$

e.g. $\kappa(\mathbf{x}, \mu_i, \lambda) = e^{-\frac{1}{\lambda} \|\mathbf{x} - \mu_i\|^2}$

$$\hat{y}(\mathbf{x}_i) = \phi(\mathbf{x}_i) \theta = \theta_0 + k(\mathbf{x}_i, \mu_1, \lambda) \theta_1 + \dots + k(\mathbf{x}_i, \mu_d, \lambda) \theta_d$$

Example 1 : $\hat{y}(x) = e^{-\|x - 1\|^2} \theta_1 + e^{-\|x - 2\|^2} \theta_2 + e^{-\|x - 4\|^2} \theta_3$



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

$\phi(x_i)$ is a vector with 4 entries. There are 3 bases.

The corresponding vector of parameters is $\underline{\Theta} = [\Theta_0 \ \Theta_1 \ \Theta_2 \ \Theta_3]^T$

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

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

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad N \times 1$$

$$\underline{\Phi} = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_N) \end{bmatrix} \quad N \times 4$$

Then

$$\hat{Y} = \Phi \theta$$

and

$$\hat{\theta}_{ls} = (\Phi^T \Phi)^{-1} \Phi^T y$$

or

$$\hat{\theta}_{ridge} = (\Phi^T \Phi + \lambda^2 I)^{-1} \Phi^T y$$

Hence, this is still linear regression,
with X replaced by Φ .

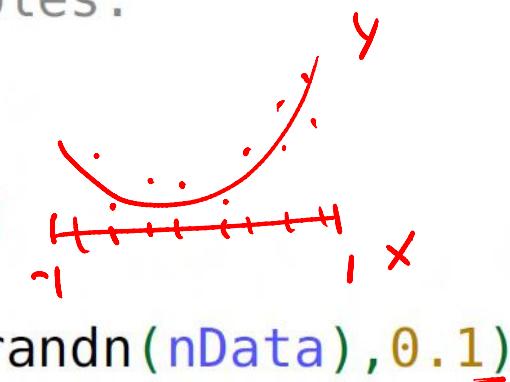
Kernel regression in Torch

```
require 'torch'  
require 'gnuplot'
```

```
local nData = 10 -- Number of data samples.  
local kWidth = 1 -- Kernel width.
```

```
local xTrain = torch.linspace(-1,1,nData)  
local yTrain = torch.pow(xTrain,2)  
local yTrain = yTrain + torch.mul(torch.randn(nData), 0.1)
```

```
local function phi(x, y)  
    return torch.exp(-(1/kWidth)*torch.sum(torch.pow(x-y,2)))  
end
```



Kernel regression in Torch

$$\phi(x_i, x_j) = e^{-\frac{1}{\lambda} \|x_i - x_j\|^2}$$

$$\underline{\Phi} = \begin{bmatrix} \underline{\phi}_{11} & \dots & \underline{\phi}_{1n} \\ & \ddots & \\ & & \underline{\phi}_{nn} \end{bmatrix}$$

```
local Phi = torch.Tensor(nData, nData)
for i=1,nData do
    for j=1,nData do
        Phi[i][j]=phi(xTrain[{{i}}], xTrain[{{j}}])
    end
end
```

```
local regularizer = torch.mul(torch.eye(nData), 0.001)
local theta = torch.inverse((Phi:t()*Phi) + regularizer) * Phi:t() * yTrain
```

$$\Theta = [\underline{\Phi}^T \underline{\Phi} + \delta^2 I]^{-1} \underline{\Phi}^T y$$

Kernel regression in Torch

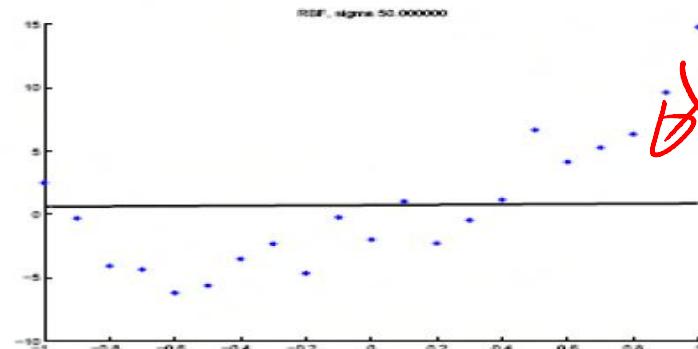
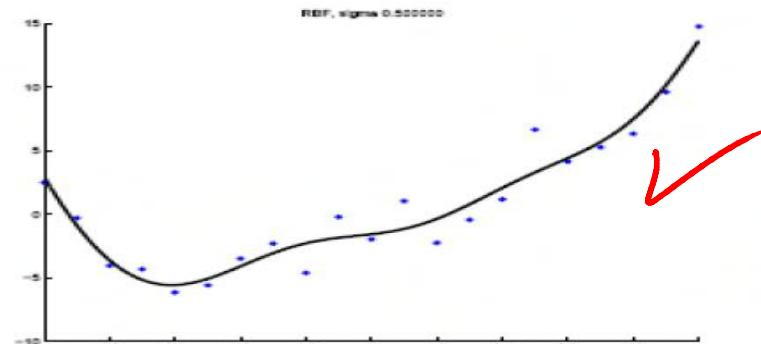
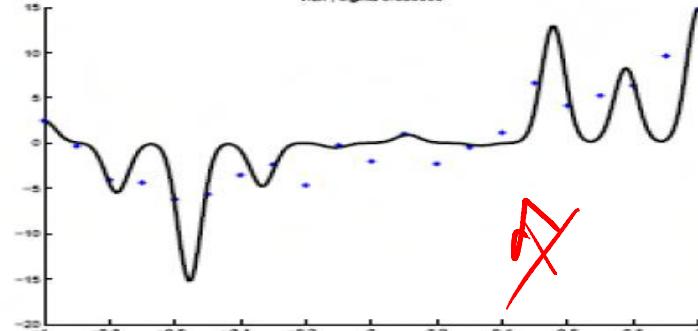
```
local nTestData = 100 -- Number of test data samples
local xTest = torch.linspace(-1,1,nTestData)

local PhiTest = torch.Tensor(nData,nTestData)
for i=1,nData do
    for j=1,nTestData do
        PhiTest[i][j]=phi(xTrain[{{i}}],xTest[{{j}}])
    end
end

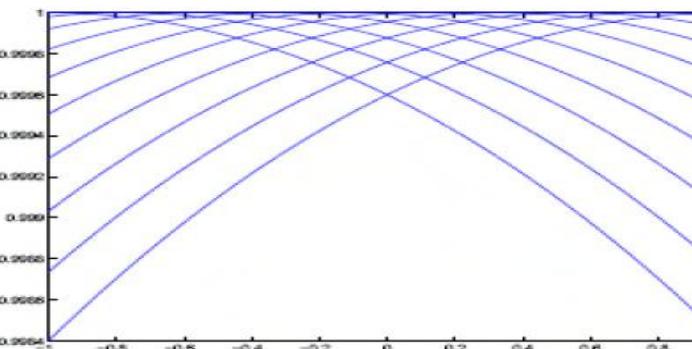
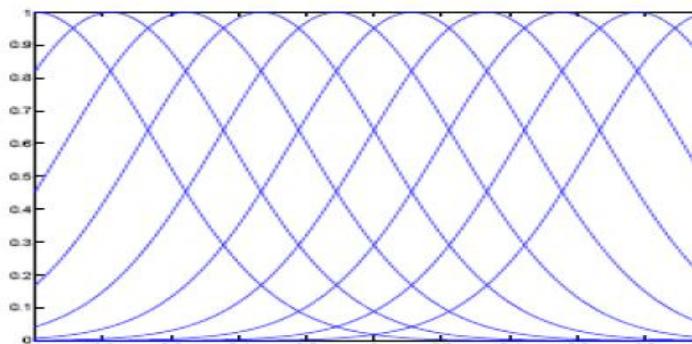
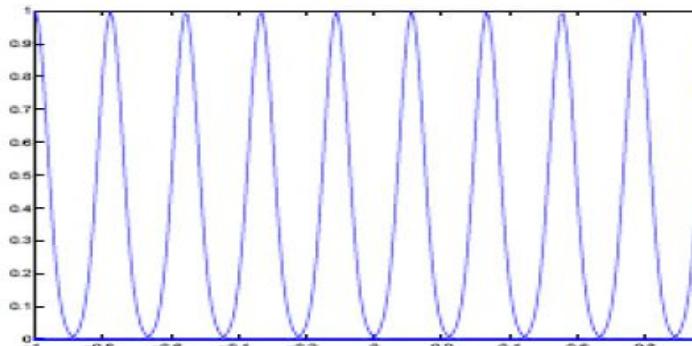
local yPred = PhiTest:t() * theta

gnuplot.plot({'Data',xTrain,yTrain,'+'},{'Prediction',xTest,yPred,'-'})
```

We can choose the locations μ of the **basis functions** to be the inputs. That is, $\mu_i = \mathbf{x}_i$. These basis functions are known as **kernels**. The choice of width λ is tricky, as illustrated below.



kernels



Too small λ

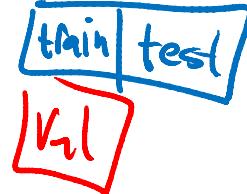
Right λ

Too large λ

The big question is how do we choose the regularization coefficient, the width of the kernels or the polynomial order?

Simple solution: cross-validation

- ① Given training data (X, Y) , and some \hat{g}^2 guess, Compute $\hat{\theta}$
- ② $\hat{Y}_{\text{train}} = X_{\text{train}} \hat{\theta}$ (compute training set predictions)
- ③ $\hat{Y}_{\text{test}} = X_{\text{test}} \hat{\theta}$



$\alpha_1, \alpha_2, \alpha_3, \lambda, \sigma^2$	Train error $\sum_{i \in \text{train}} (y_i - \hat{y}_i)^2$	Test error $\sum_{i \in \text{test}} (y_i - \hat{y}_i)^2$	Max	min-max	avg
1 → 0.1	100	2	100		
3 → 1	10	11	11	11	
71 → 10	1	19	19		10.5
11 → 50	20	0	20		x 16 *
1 → 100	100	1000	1000		x 10

K-fold crossvalidation

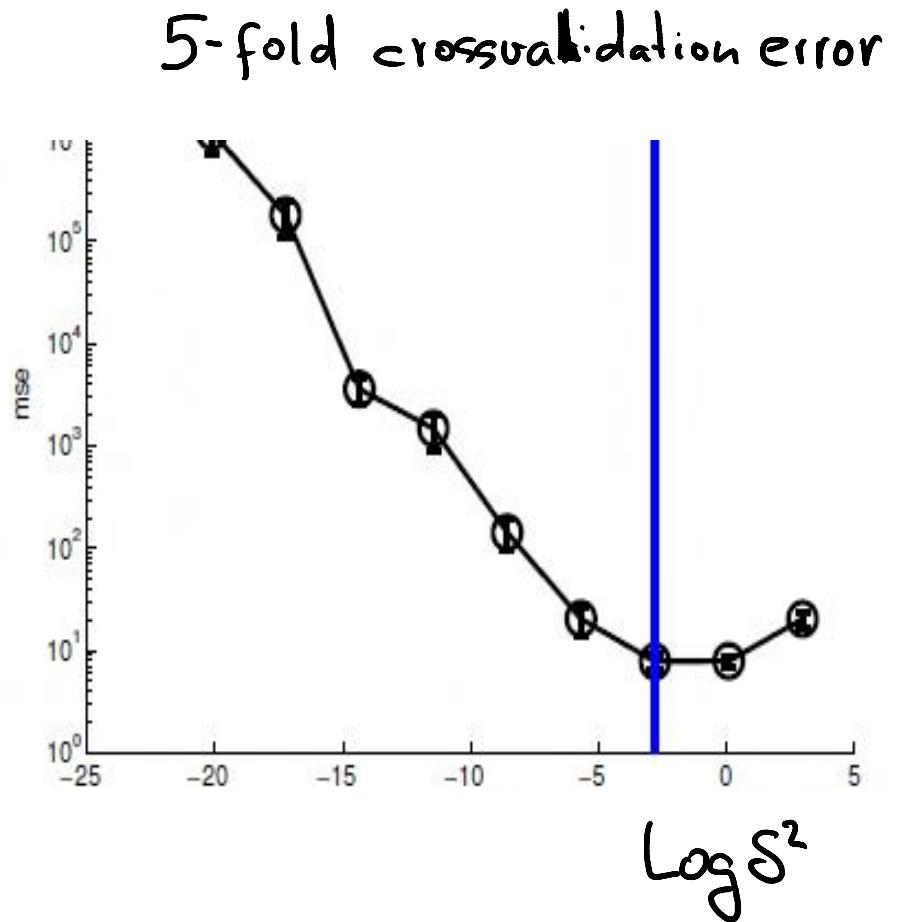
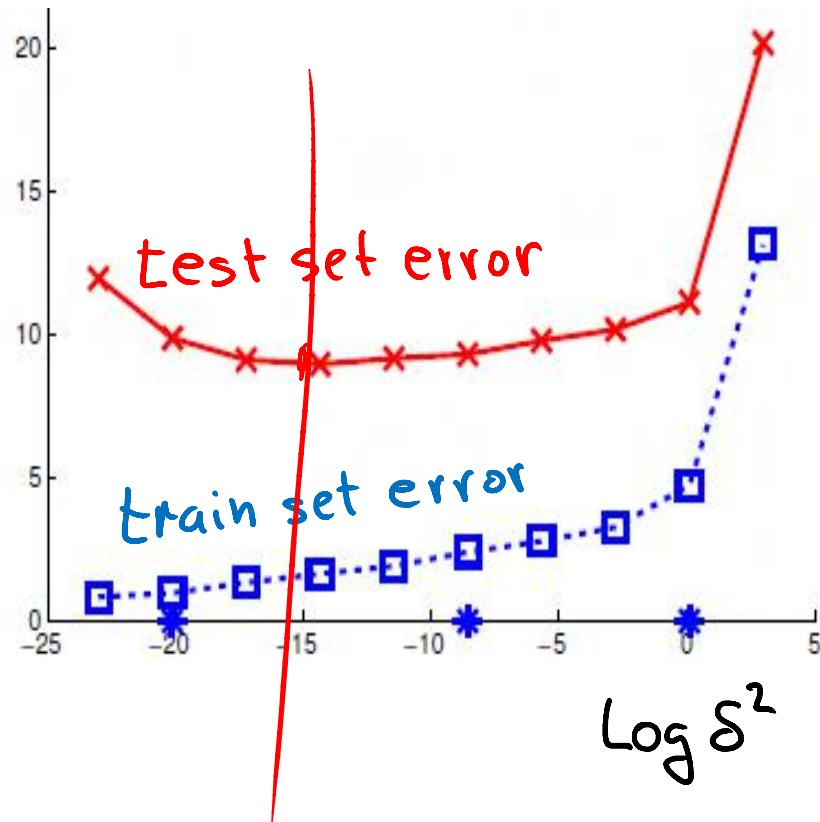


The idea is simple: we split the training data into K **folds**; then, for each fold $k \in \{1, \dots, K\}$, we train on all the folds but the k 'th, and test on the k 'th, in a round-robin fashion.

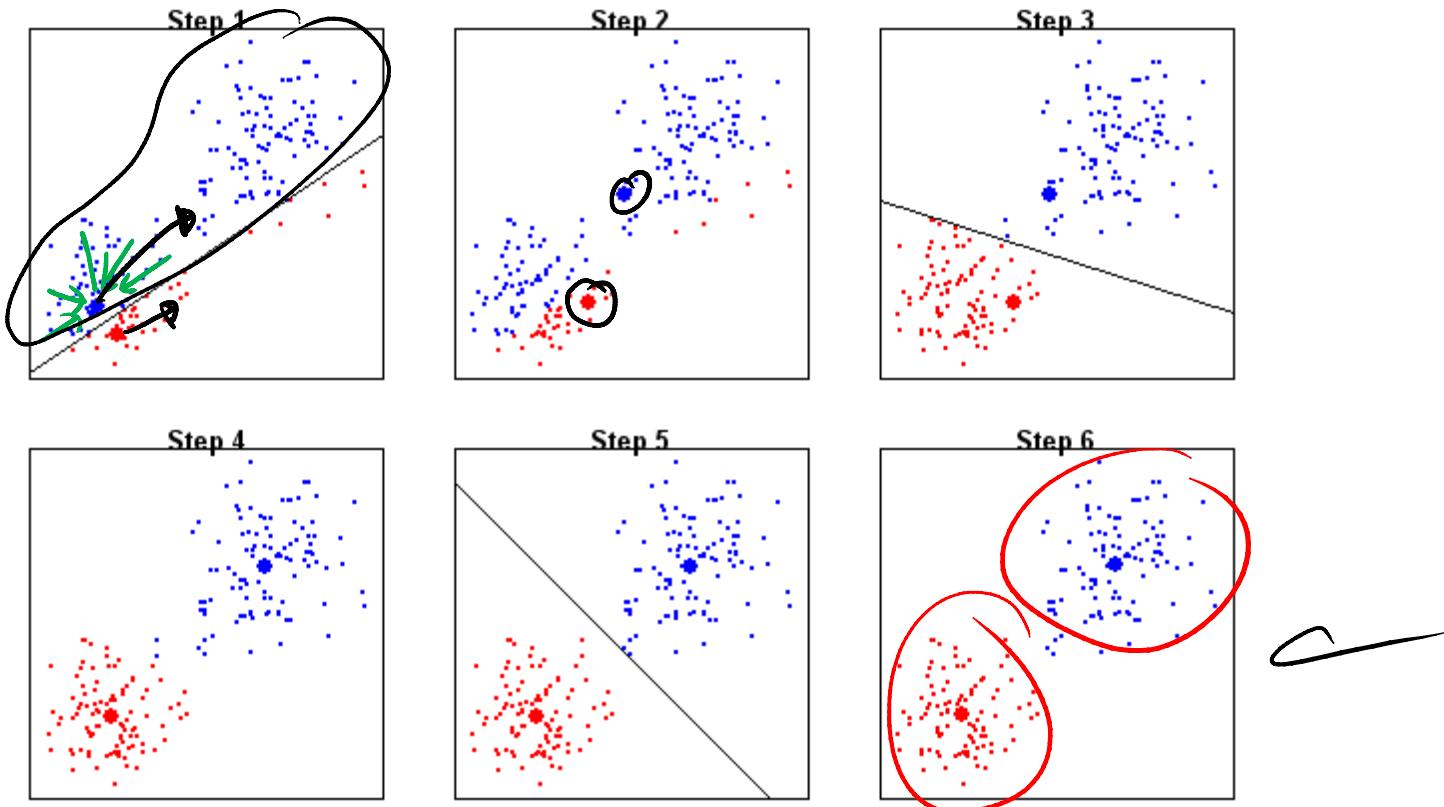
It is common to use $K = 5$; this is called 5-fold CV.

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



Where cross-validation fails) (K-means)



$$E(\mathcal{D}, L) = \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} \|\underline{\mathbf{x}}_i - \hat{\underline{\mathbf{x}}}_i\|^2$$

$$\hat{\underline{\mathbf{x}}}_i = \arg\min_k \|\underline{\mathbf{x}}_i - \underline{\boldsymbol{\mu}}_k\|_2^2$$

$$\hat{\underline{\mathbf{x}}}_i = \underline{\boldsymbol{\mu}}_{z_i}$$

Next lecture

In the next lecture, we delve into the world of optimization.

Please revise your multivariable calculus and in particular the definition of **gradient**