# 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

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

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:

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

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

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

# Derivation

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

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta} \left( \theta^T X^T X \theta - 2 Y^T X \theta + Y^T Y + \delta^2 \theta^T \theta \right)$$

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

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

Equate to zero

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

# Ridge regression as constrained optimization

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



$$\boldsymbol{\Theta}^T = \begin{bmatrix} \Theta_1 & \Theta_2 \end{bmatrix}$$

$$\Theta^T\Theta = \begin{bmatrix} \Theta_1 & \Theta_2 \end{bmatrix} \begin{bmatrix} \Theta_1 \\ \Theta_2 \end{bmatrix}$$

$$= \Theta_1^2 + \Theta_2^2 = const$$

$\Theta_{ML}$

$\delta^2 = 0$

$\delta^2 \to \infty$

$\Theta_2$

$\Theta_1$

# Regularization paths

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



[Hastie, Tibshirani & Friedman book]

# Ridge regression and Maximum a Posteriori (MAP) learning

$$\boxed{\text{BAyes Rule}}$$

$$J(\boldsymbol{\theta}) = \underbrace{(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}_{\overline{E}(\theta)X,Y)} + \underbrace{\delta^2\boldsymbol{\theta}^T\boldsymbol{\theta}}$$

$$P(Y|X,\Theta) = \frac{1}{Z} e^{-E(\Theta,X,Y)}$$

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

$$P(Y|X) = \int P(Y|X,\Theta)P(\Theta) \, d\Theta$$

$$\max_{\Theta} \text{const} \, P(Y|X,\Theta)P(\Theta) = \max_{\Theta} \frac{P(Y|X,\Theta)P(\Theta)}{P(Y|X)}$$

# Ridge regression and Maximum a Posteriori (MAP) learning

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

$$J(\boldsymbol{\theta}) = \underbrace{(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}_{\text{Likelihood}} + \underbrace{\delta^2 \boldsymbol{\theta}^T \boldsymbol{\theta}}_{\text{prior}}$$

$$\underbrace{P(\underline{\theta}|x,y)}_{\text{Posterior}} = \frac{P(y|x,\theta)\, P(\theta)}{P(y|x)}$$

word

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

# Going nonlinear via basis functions

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

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

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

features

$$\hat{\Theta} = \left(\phi^\top \phi + s^2 I\right)^{-1} \phi^\top y$$

| x | y |
|---|---|
| 2.2 | 1 |
| 6.3 | 2.6 |
| ⋮ | ⋮ |

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

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

$= \begin{bmatrix} 1 & x & x^2 \end{bmatrix} \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \end{bmatrix}$

$x_i$

# Going nonlinear via basis functions

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

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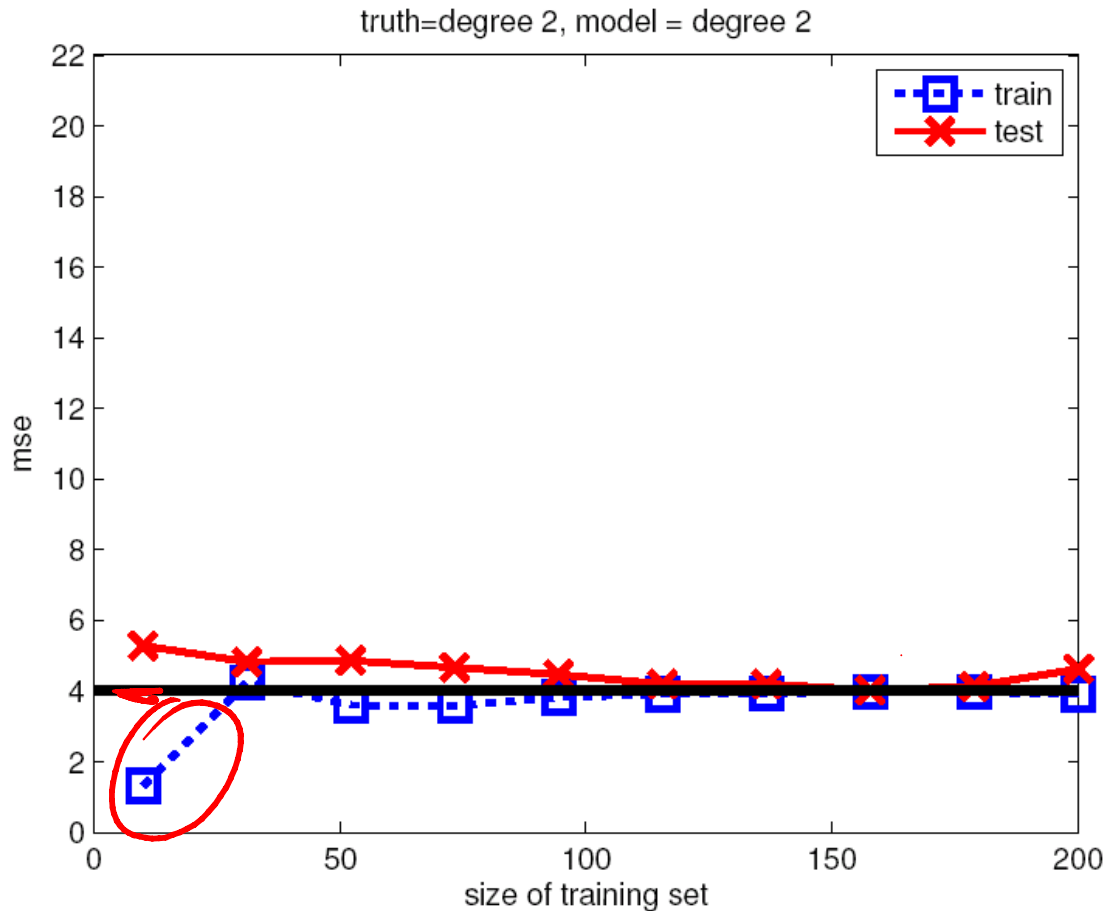
$x_1 x_2$

# Effect of data when we have the right model

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

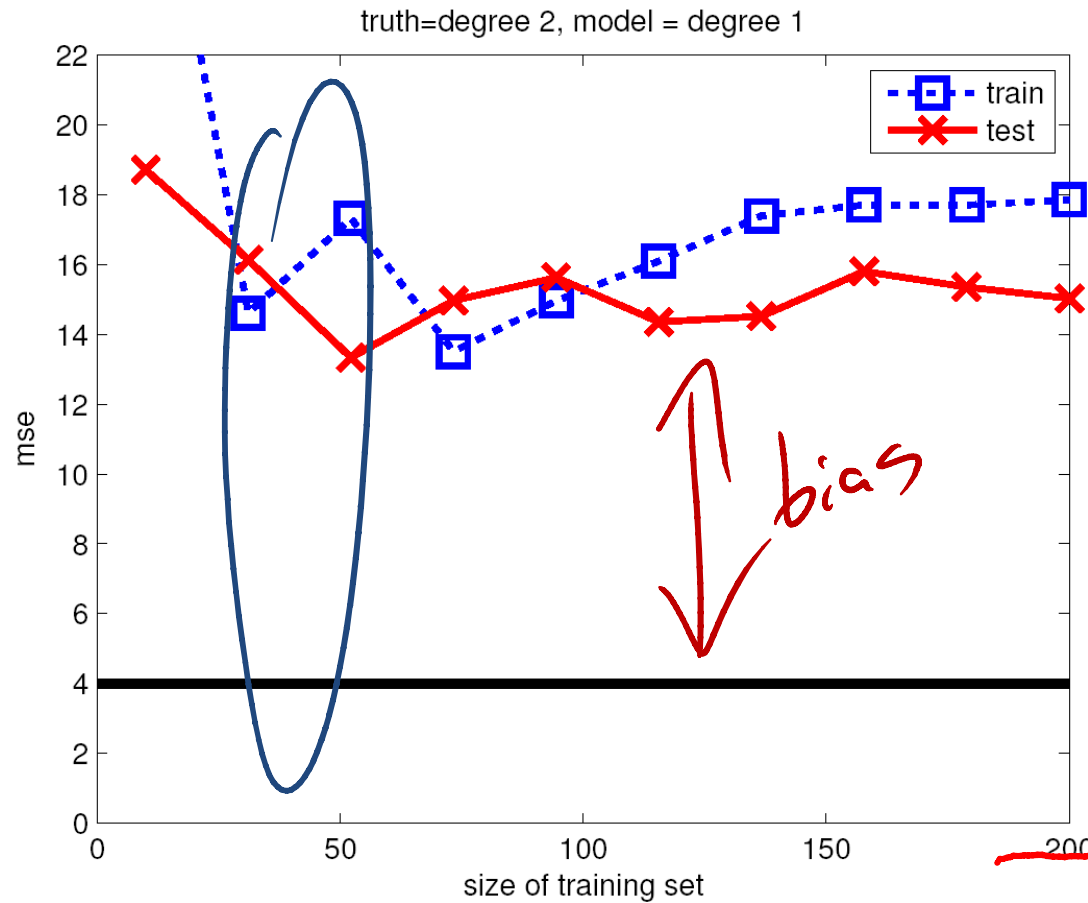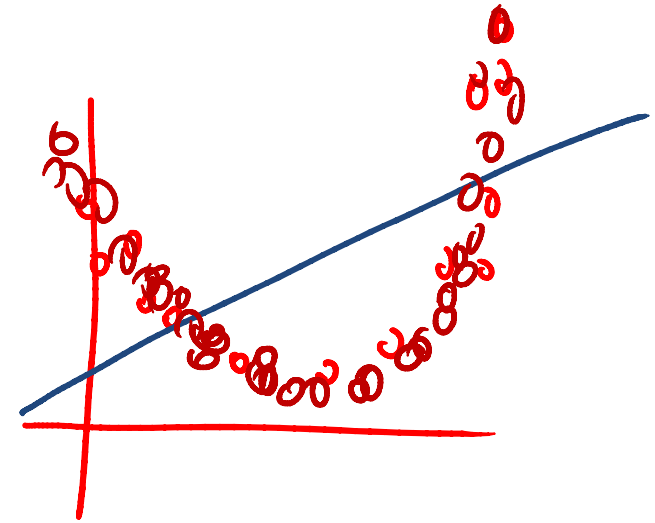$$\hat{y} = \theta_0 + x_i \theta_1 + x_i^2 \theta_2$$

truth=degree 2, model = degree 2



O Train data

● Test data

noise
Error

# Effect of data when the model is too simple

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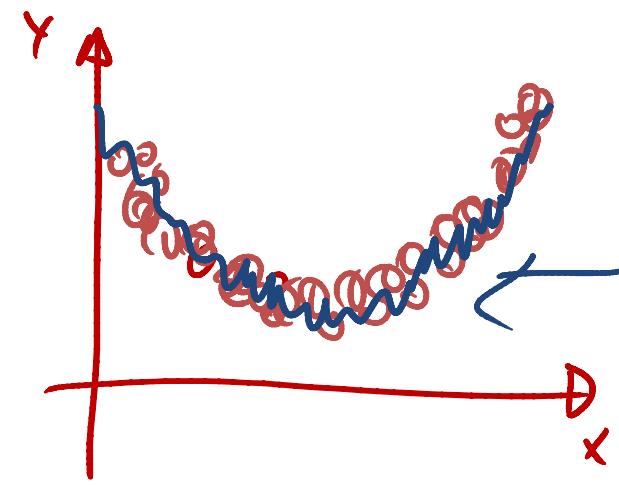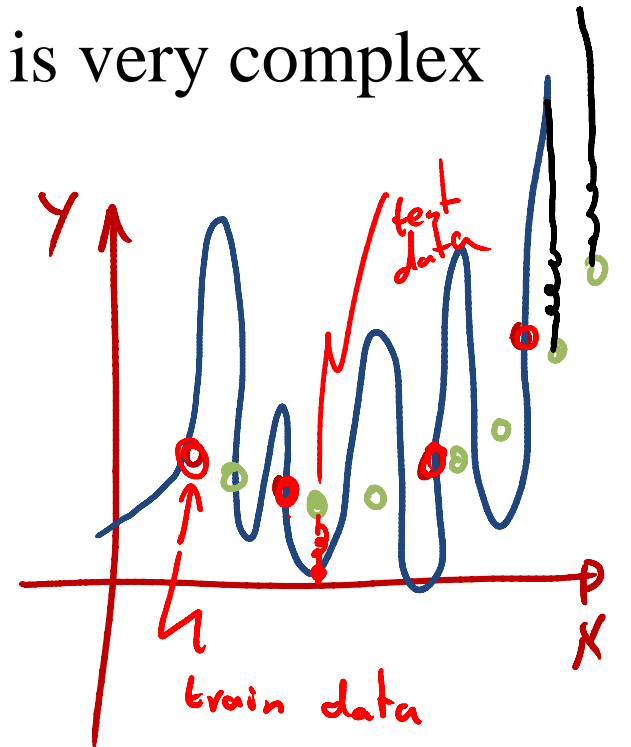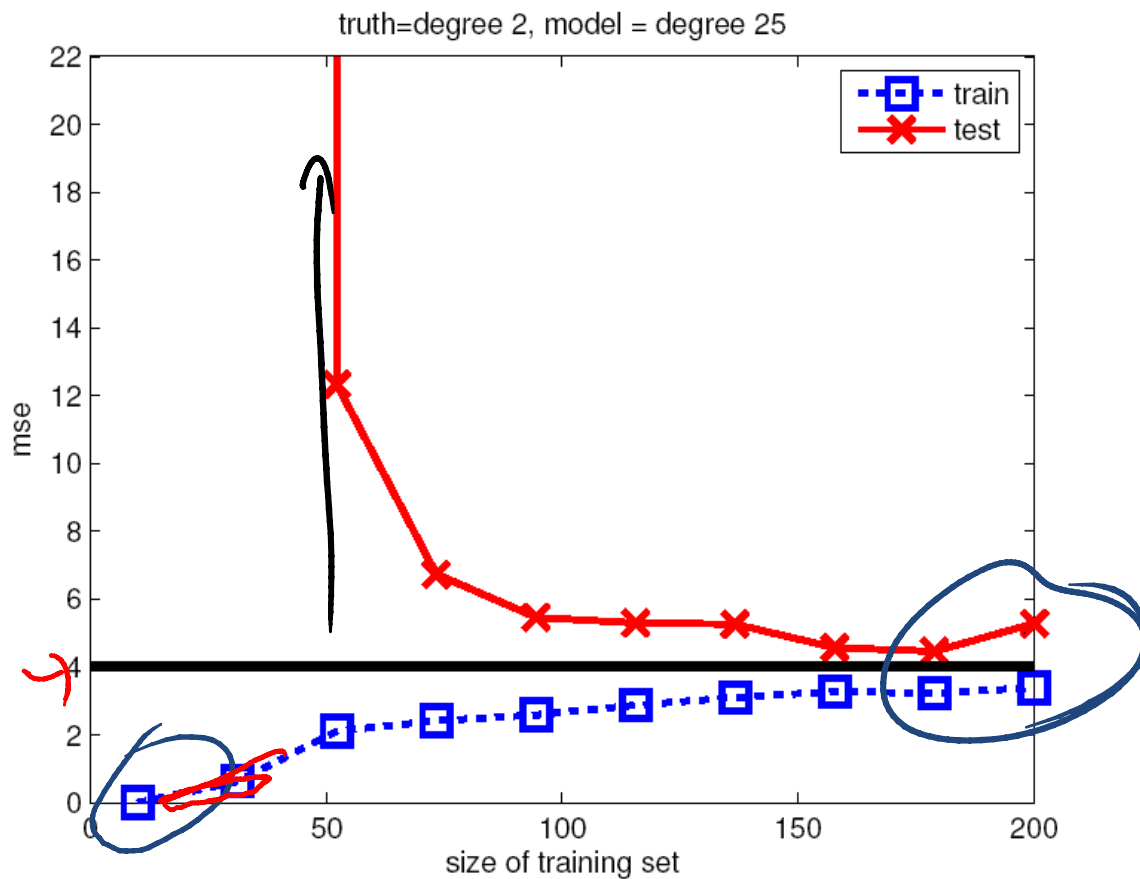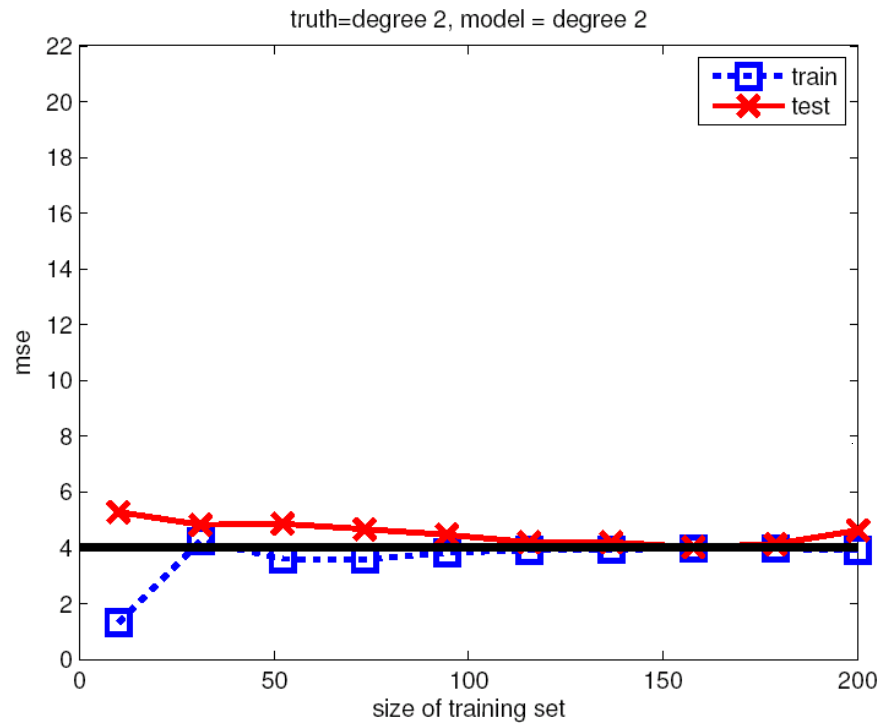
$$\hat{y} = \theta_0 + x_i\,\theta_1$$

truth=degree 2, model = degree 1



bias
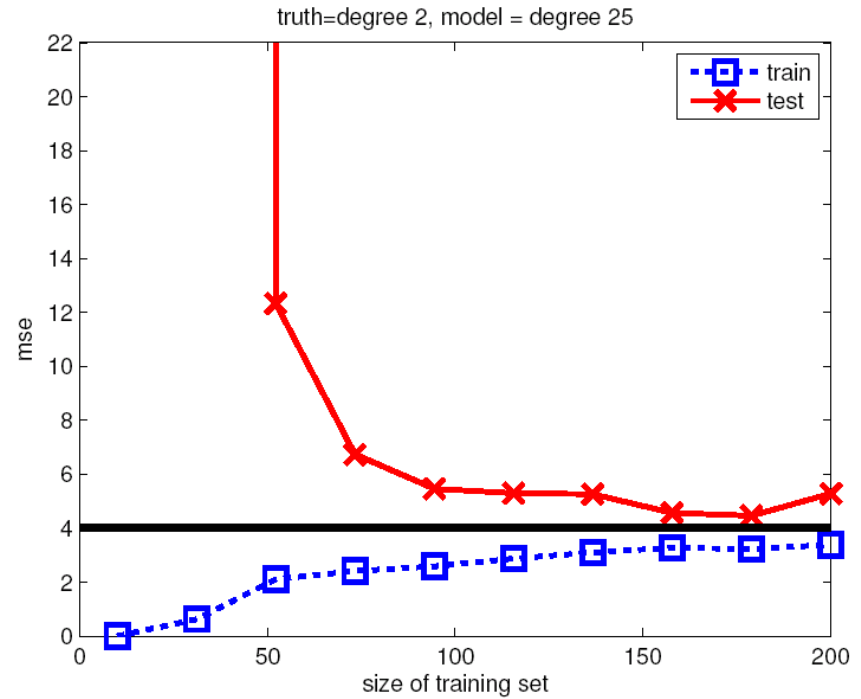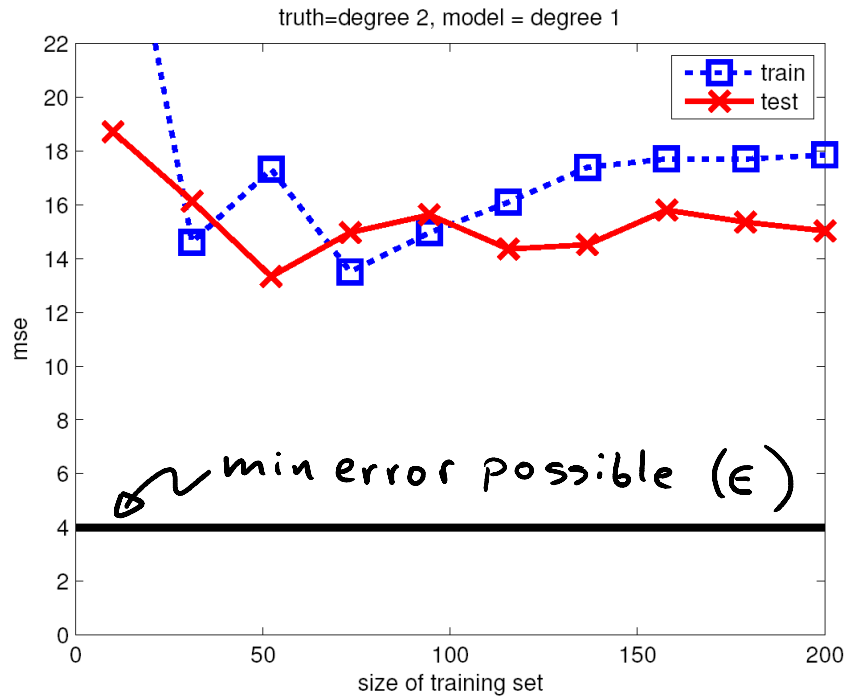
$N \to$ increases

$i = 1 : N$

# Effect of data when the model is very complex

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

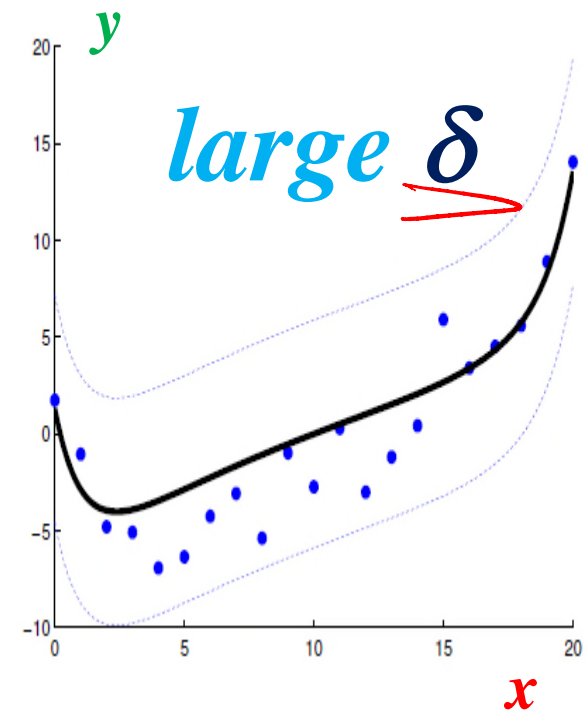$$\hat{y}_i = \theta_0 + x_i\,\theta_1 + \cdots + x_i^{25}\,\theta_{25}$$
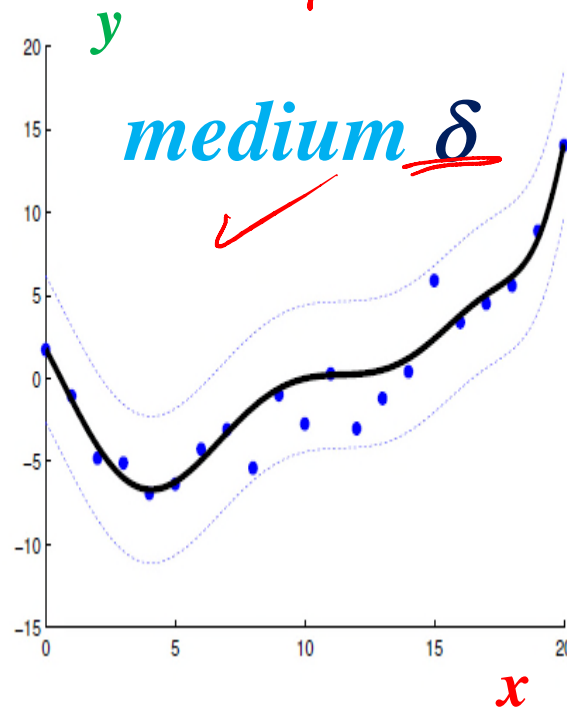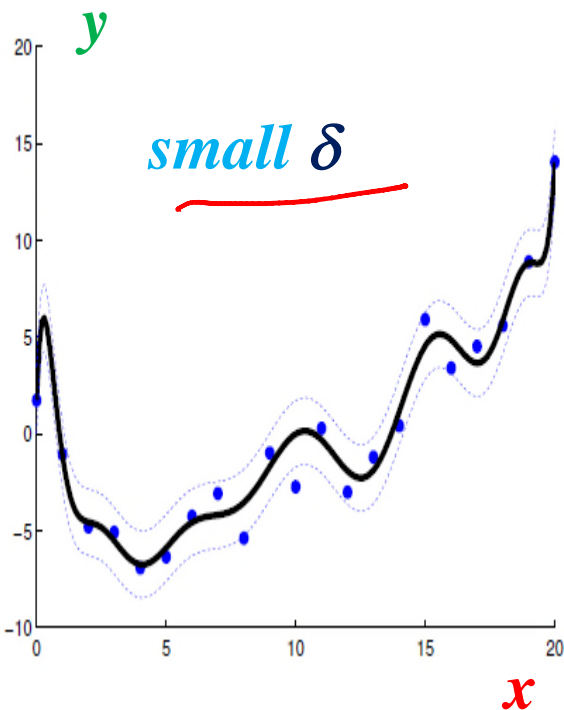
truth=degree 2, model = degree 25

truth=degree 2, model = degree 1

min error possible ($\epsilon$)

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.

**_Example_**: *Ridge regression* **with a** *polynomial of degree 14*

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

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

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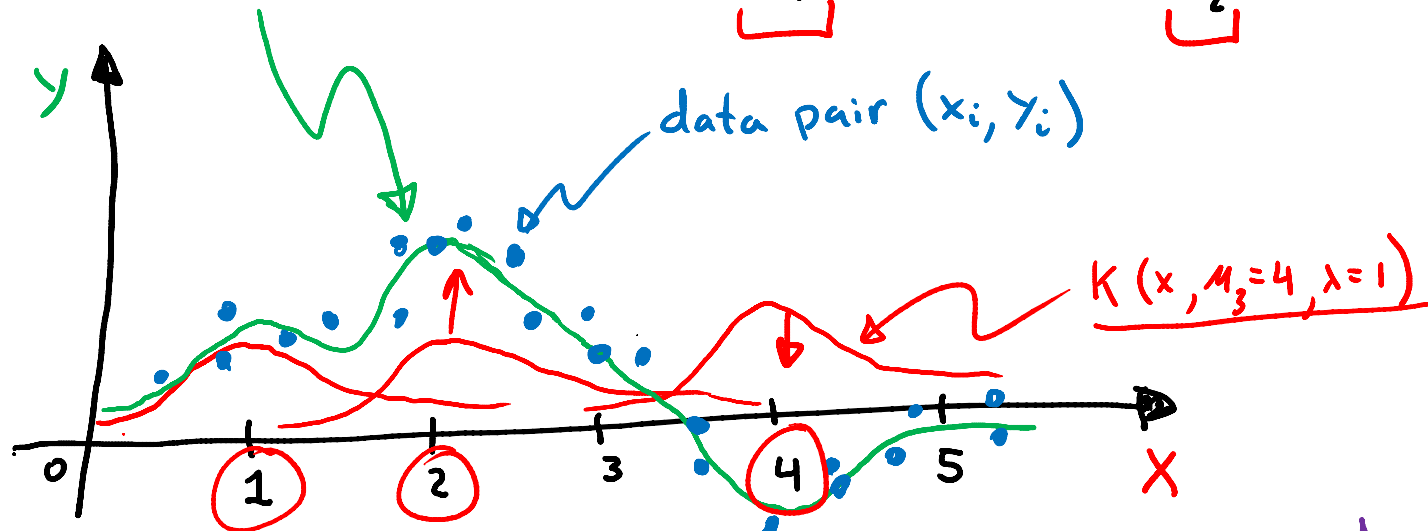


small $\delta$

medium $\delta$

large $\delta$

# Kernel regression and RBFs

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

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

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

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



data pair $(x_i, Y_i)$

$K(x, M_3=4, \lambda=1)$

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

$$\phi(x_i) = [\boxed{1} \; \kappa(x_i, m_1, \lambda) \; \kappa(x_i, m_2, \lambda) \; \kappa(x_i, m_3, \lambda)]$$

$\phi(x_i)$ 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 $i = 1, \ldots, N$ data, let

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

$N \times 1$ $\qquad$ $N \times 4$

Then

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

and

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

or

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

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

# 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

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

$$\bar{\Phi} = \begin{bmatrix} \phi_{11} \cdots \phi_{1n} \\ \\ \phi_{nn} \end{bmatrix}$$

```
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
```

$\sim \delta^2$

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

# 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,'-'})
```
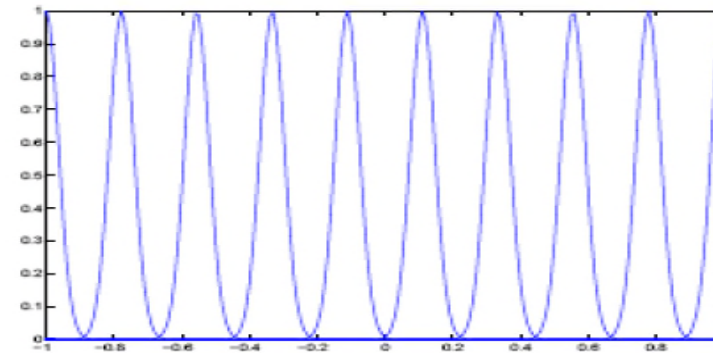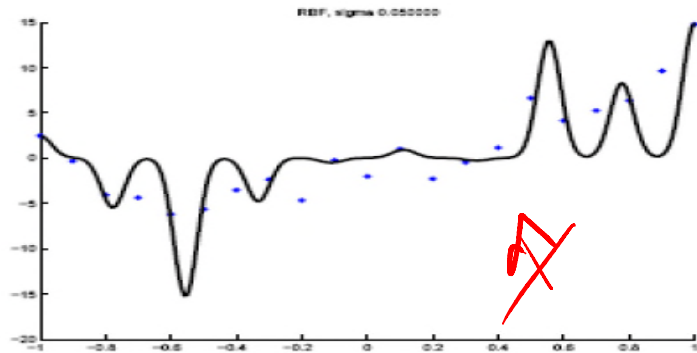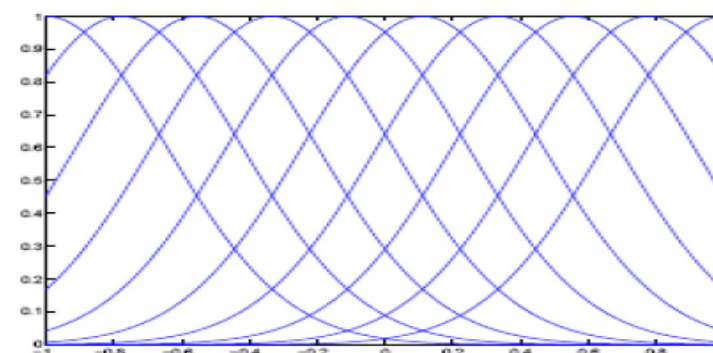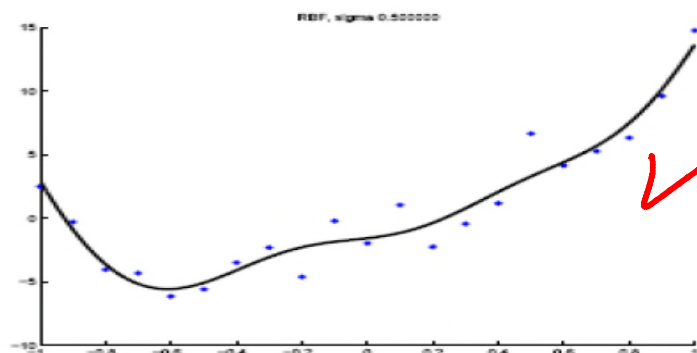
*We can choose the locations **μ** of the **basis functions** to be the inputs. That is, **μ$_i$** = **x$_i$**. These basis functions are the known as **kernels**. The choice of width **λ** is tricky, as illustrated below.*

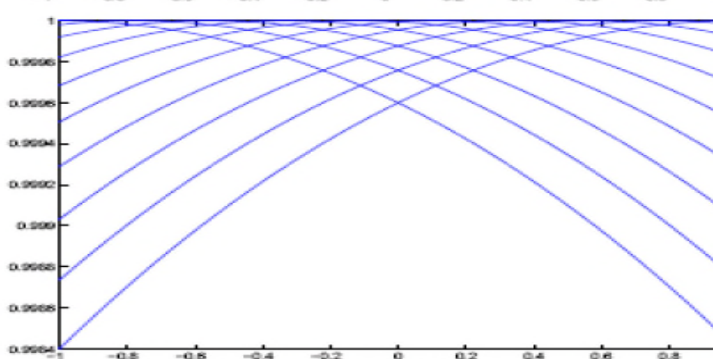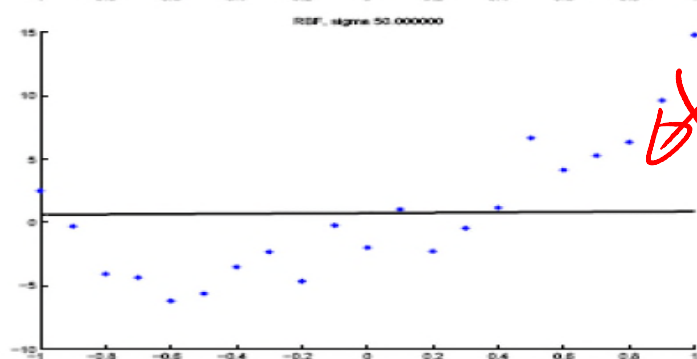**kernels**



Too small λ

Right λ
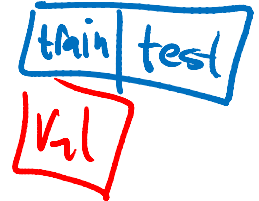
Too large λ

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 $(X, Y_b)_{train}$ 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
val

validation

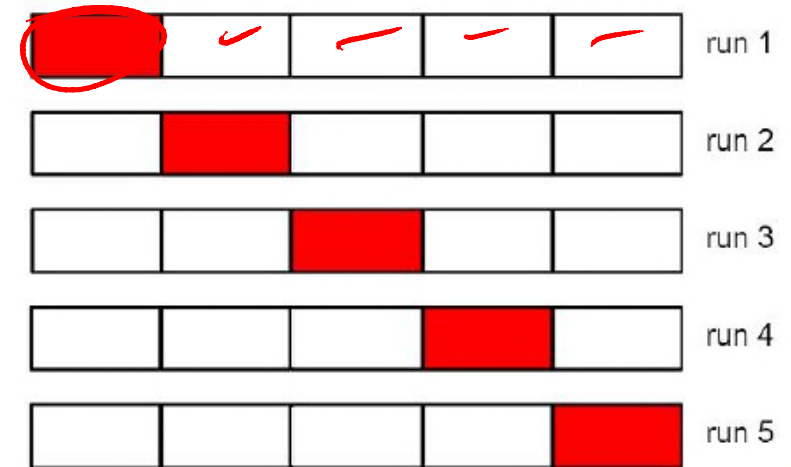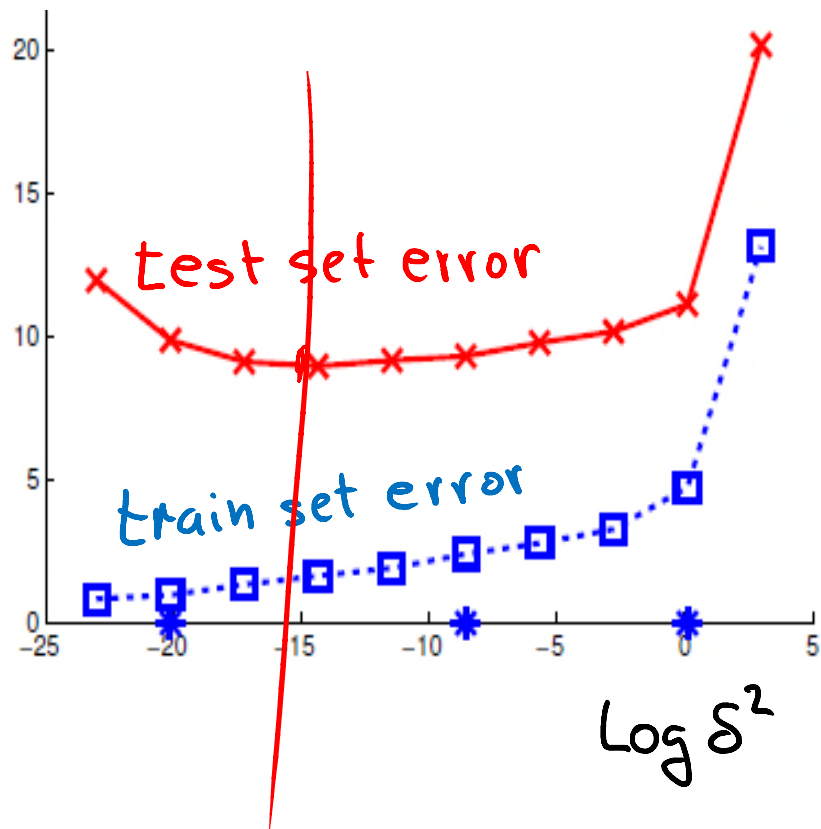| $q_1 q_2 q_3 \lambda, \delta^2$ | Train error $\sum_{i \in train}(Y_i - \hat{Y}_i)^2$ | Test error $\sum_{i \in test}(Y_i - \hat{Y}_i)^2$ | max | min-max | avg |
|---|---|---|---|---|---|
| ①→ 0.1 | 100 | 2 | 100 | | |
| ③ → 1 | 10 | 11 | 11 | 11 | |
| →\|→ 10 | 1 | 19 | 19 | | 10.5 ✗ ✗ 10 |
| \|\| → 50 | 20 | 0 | 20 | | ✗ 10 |
| \| → 100 | 100 | 1000 | 1000 | | |

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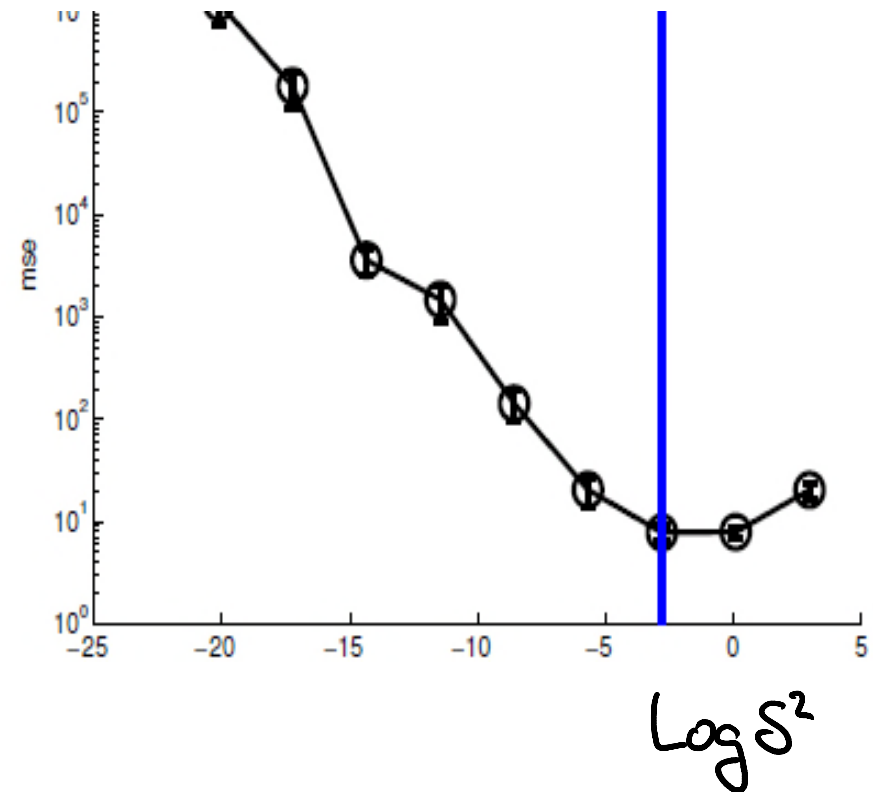
It is common to use $K = 5$; this is called 5-fold CV.

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 $i$.
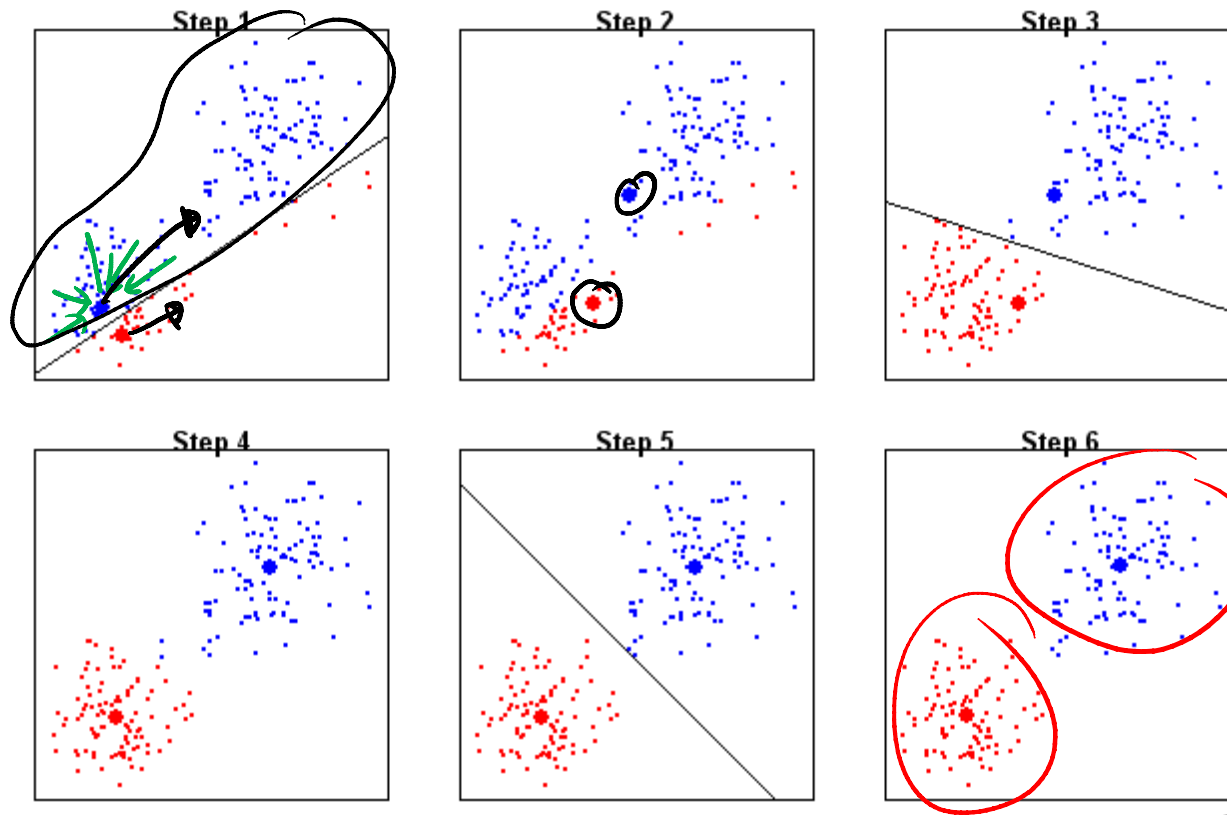
# Example: Ridge regression with polynomial of degree 14



5-fold crossvalidation error

test set error

train set error

$Log \delta^2$

mse

$Log \delta^2$

# Where cross-validation fails) (K-means)



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

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

$$z_i = \operatorname{argmin}_k ||\mathbf{x}_i - \boldsymbol{\mu}_k||_2^2$$

# 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**