

# EECS 545: Machine Learning

## Lecture 2. Linear Regression

Honglak Lee

1/10/2011

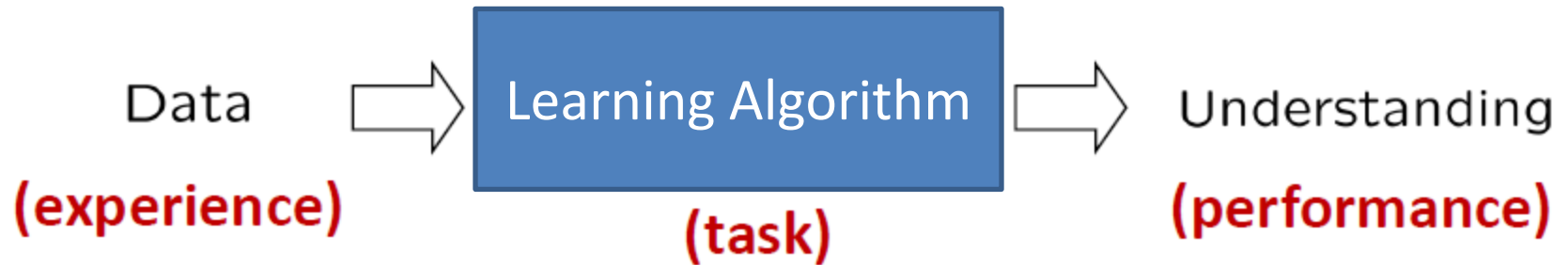


# Outline

- Recap: ML, Supervised Learning
- Linear Regression

# Informal definition of ML

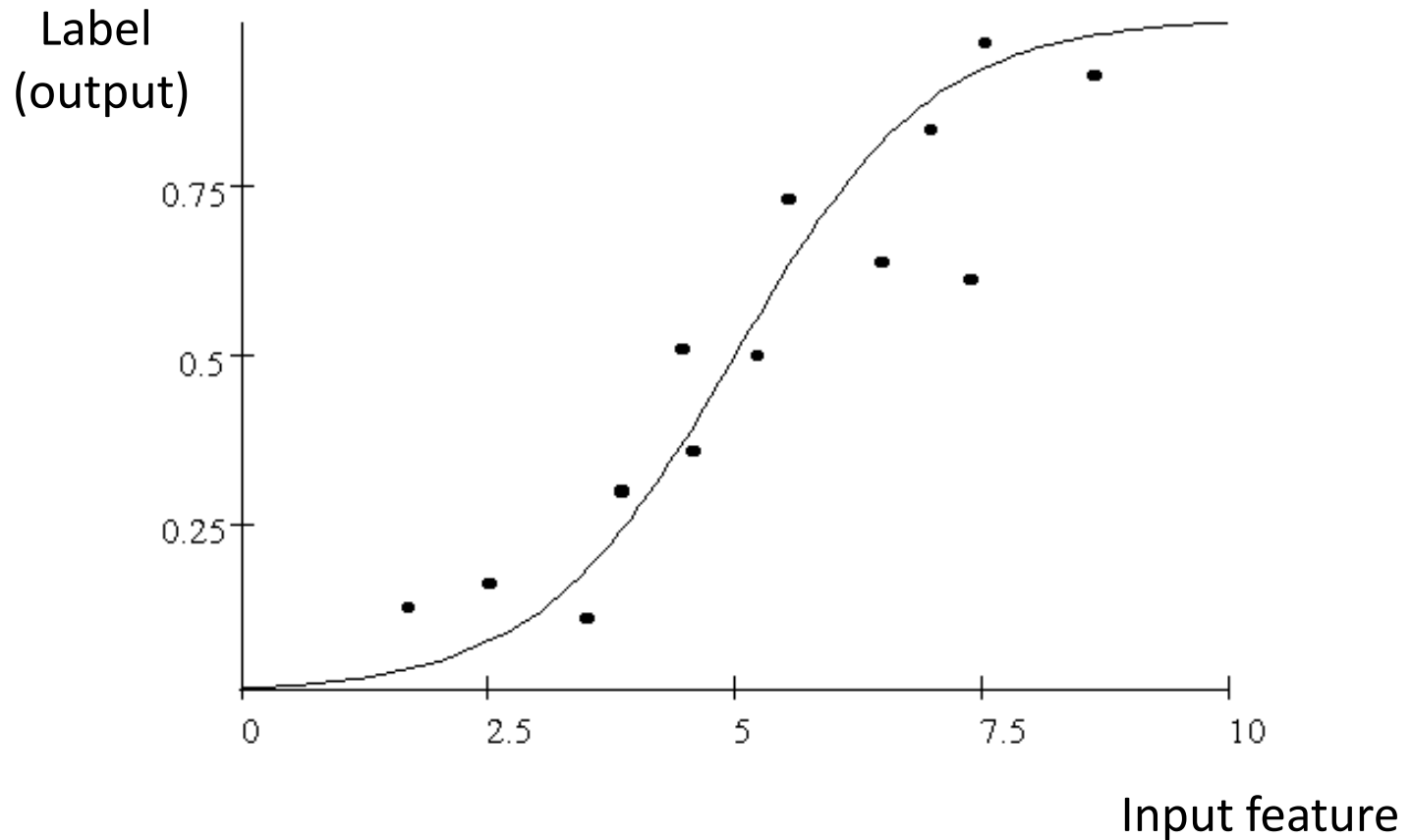
- Algorithms that improve their performance at some task with experience.



# Supervised Learning

- Goal:
  - Given data  $X$  in feature space and the labels  $Y$
  - Learn to predict  $Y$  from  $X$
- Labels could be discrete or continuous
  - Discrete labels: classification
  - Continuous labels: regression

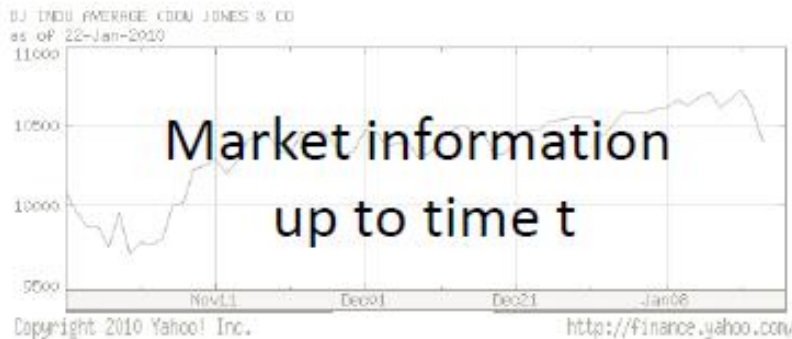
# Supervised Learning - Regression



“Learning regression function  $f(X)$ ”

# Supervised Learning - Regression

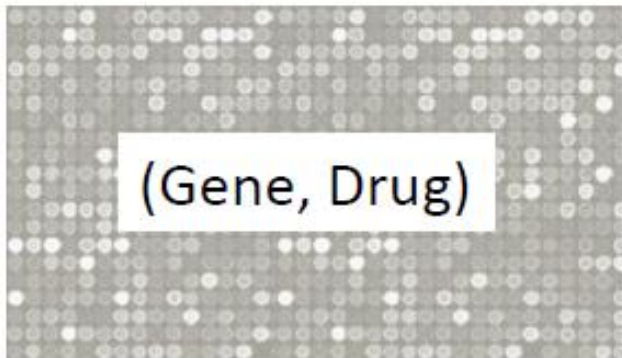
**Feature Space  $\mathcal{X}$**



**Label Space  $\mathcal{Y}$**



Share Price  
"\$ 24.50"



Expression level  
"0.01"

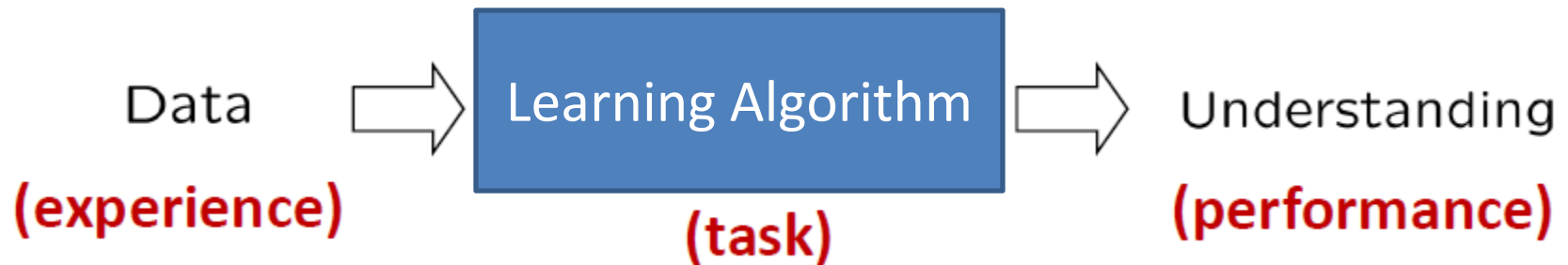
**Continuous Labels**

# Example: Housing price prediction

- **Given statistics about houses in a local area, predict median value of homes.**
- **Features:**
  - 1. Average number of rooms per dwelling
  - 2. Average area (in square foot)
  - 3. Per capita crime rate by town
  - 4. Proportion of residential land zoned for lots
  - 6. Proportion of non-retail business acres per town
  - 7. Nitric oxides concentration (parts per 10 million)
  - .....
- **Label: median value of the houses**

# Overview of linear regression

- In this lecture, we will assume
  - Data (or features): vector or scalar representation
  - Learning algorithm: linear regression to predict  $y$  from  $X$  (parameterized by  $w$ )
  - Performance: measured by sum of squared errors between prediction and labels (objective function)





# Outline

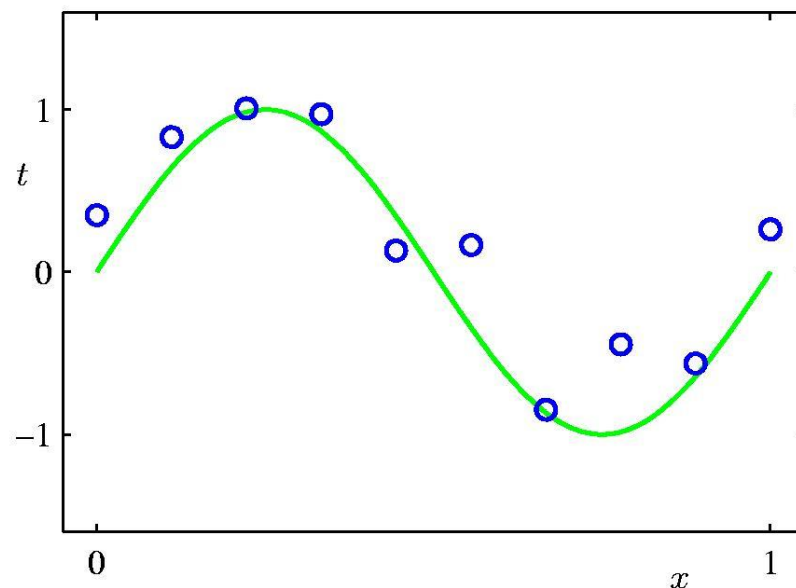
- Recap: ML, Supervised Learning
- **Linear Regression**

# Notation

- In this lecture, we will use
  - $x$ : data (scalar or vector)
  - $\phi(x)$ : features for  $x$
  - $t$  (or  $y$ ): continuous-valued labels (target values)
- We will interchangeably use
  - $x^{(n)} \stackrel{\text{def}}{=} x_n$  to denote  $n$ -th training example.
  - $t^{(n)} \stackrel{\text{def}}{=} t_n$  to denote  $n$ -th target value.

# Regression

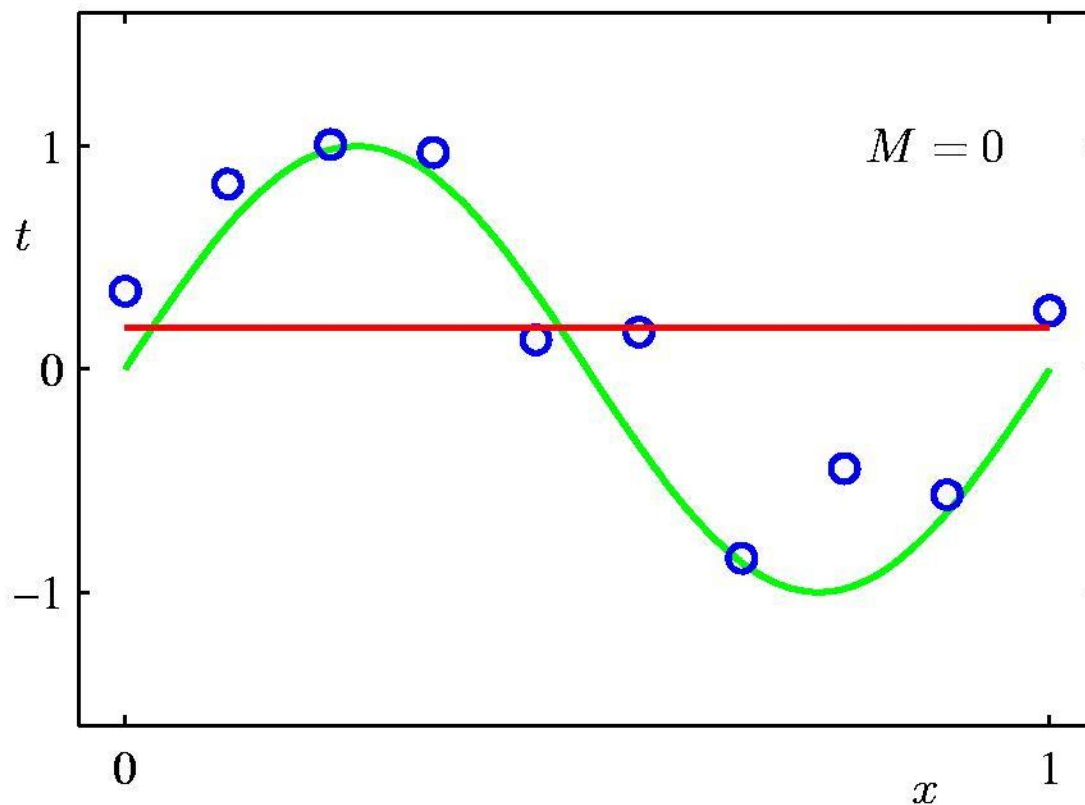
- Given a set of observations
  - $\mathbf{x} = \{x_1 \dots x_N\}$
- And corresponding target values:
  - $\mathbf{t} = \{t_1 \dots t_N\}$
- We want to learn a function  $y(x, \mathbf{w}) = t$  to predict future values.



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

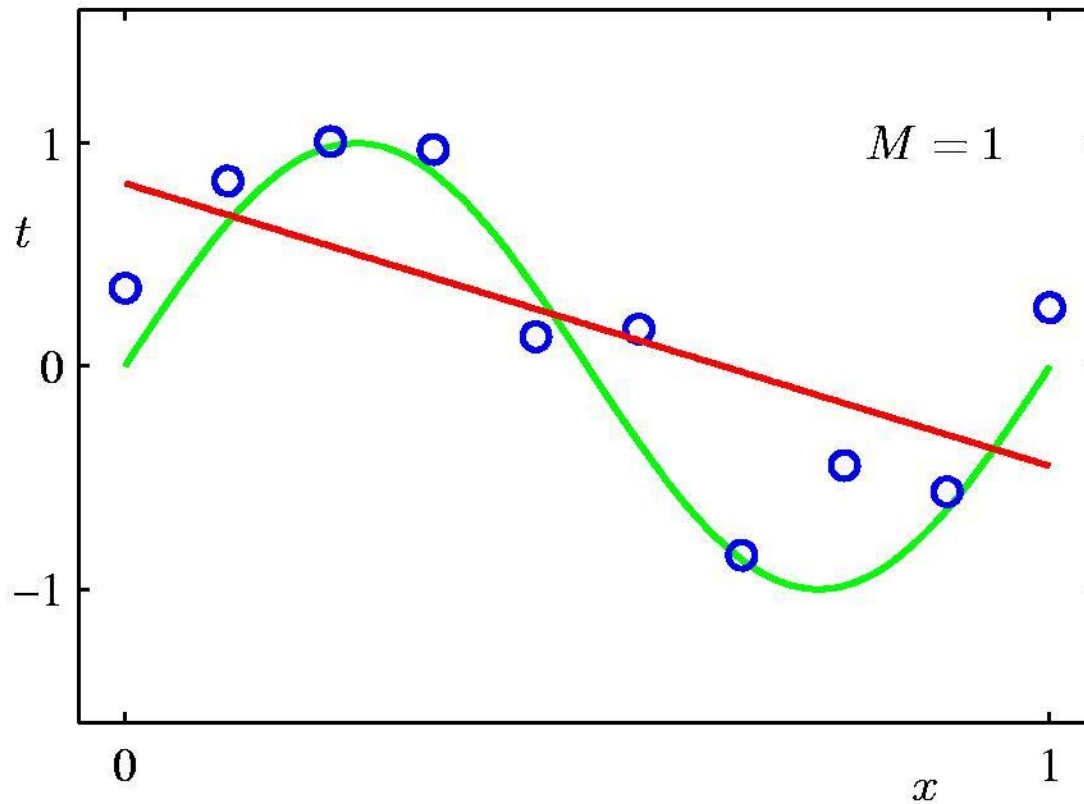
# 0<sup>th</sup> Order Polynomial

$$y = w_0$$



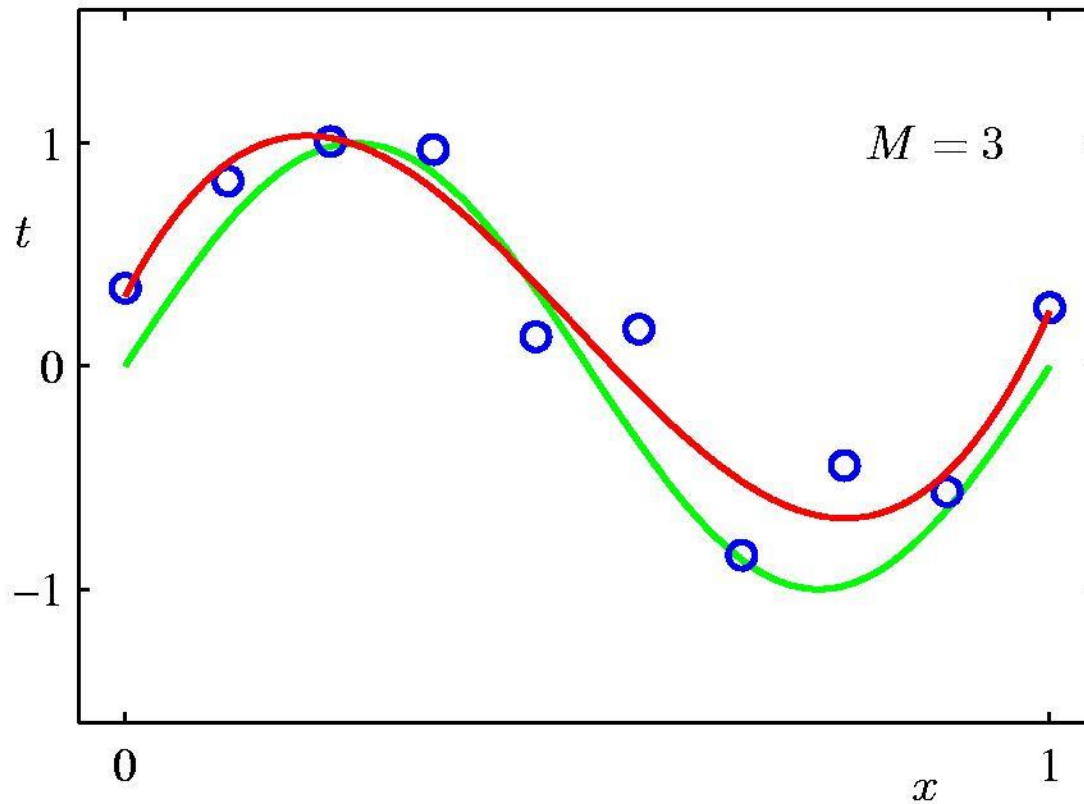
# 1<sup>st</sup> Order Polynomial

$$y = w_0 + w_1 x$$



# 3<sup>rd</sup> Order Polynomial

$$y(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$



# Linear Regression

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

- The function  $y(\mathbf{x}, \mathbf{w})$  is linear in parameters  $\mathbf{w}$ .
  - Goal: find the best value for the weights,  $\mathbf{w}$ .
- For simplicity, add a *bias function*  $\phi_0(\mathbf{x}) = 1$

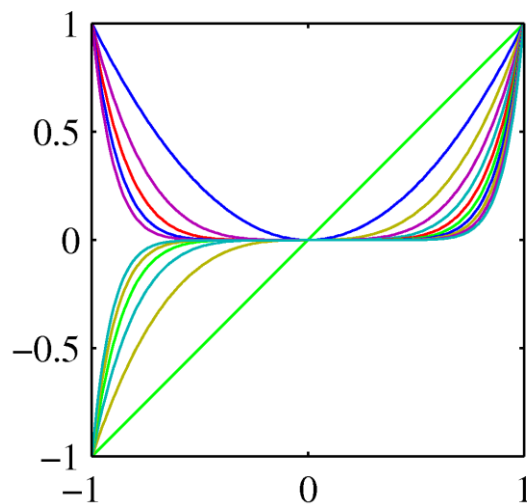
$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

$$\mathbf{w} = (w_0, \dots, w_{M-1})^T \quad \phi = (\phi_0, \dots, \phi_{M-1})^T$$

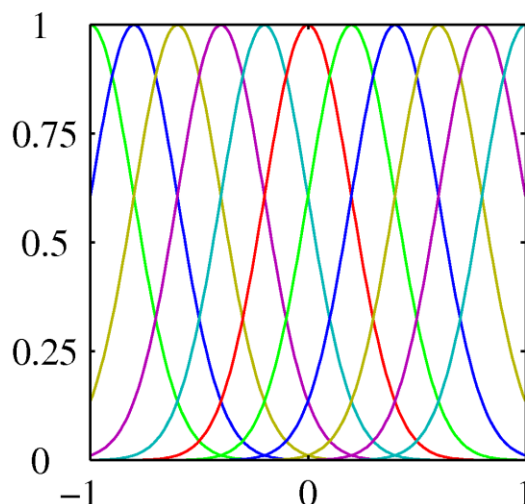
# Basis Functions

- The basis functions  $\phi_j(\mathbf{x})$  need not be linear

$$\phi_j(x) = x^j$$

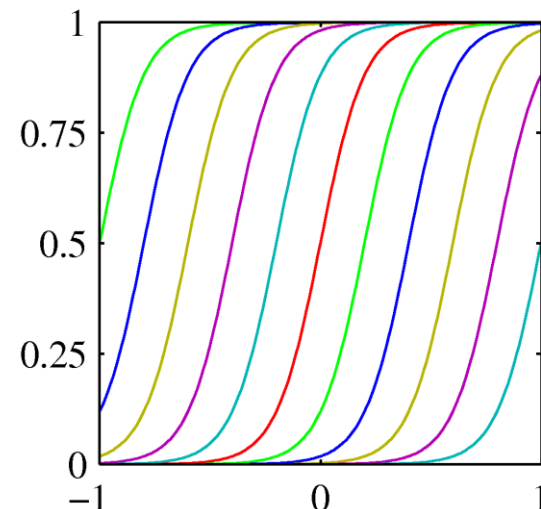


$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$



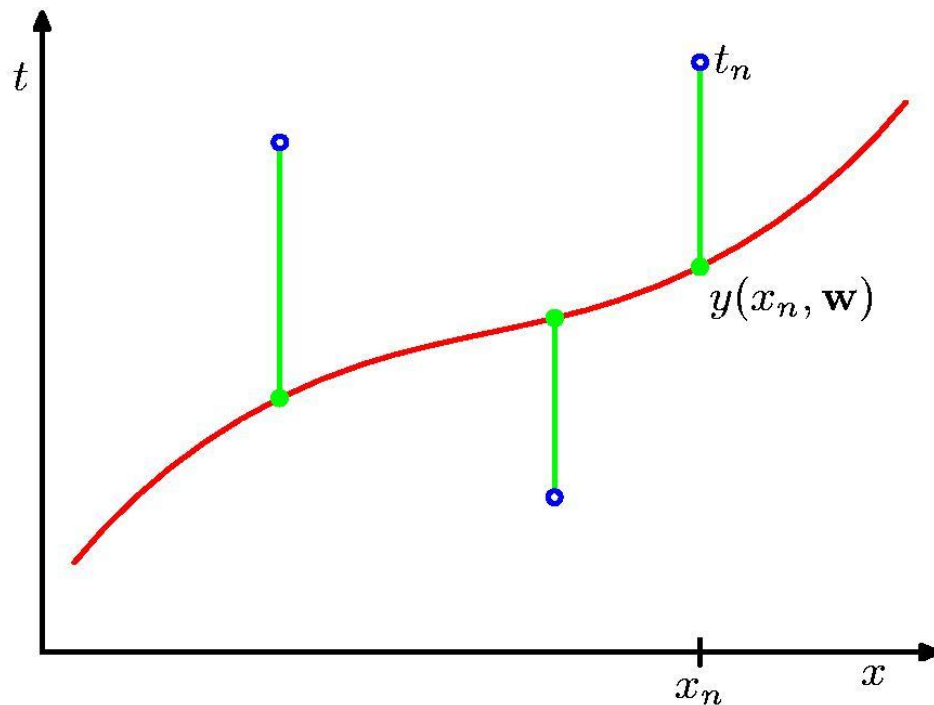
$$\phi_j(x) = \sigma \left( \frac{x - \mu_j}{s} \right)$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$





# Sum-of-Squares Error Function

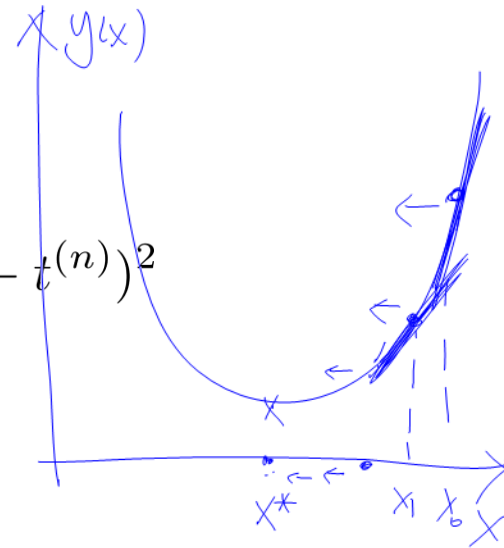


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

# Least squares problem

- Objective function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right)^2$$



- Gradient

$$\frac{\partial E(\mathbf{w})}{\partial w_j} = \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{n=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right)^2$$

$$\begin{bmatrix} \frac{\partial E(\mathbf{w})}{\partial w_1} \\ \frac{\partial E(\mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w_M} \end{bmatrix}$$

$$= \sum_{n=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right) \frac{\partial}{\partial w_j} \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right)$$

$$= \sum_{n=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right) \phi_j(x^{(n)})$$

# Least squares problem

- Gradient (compact, vectorized form)

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \sum_{n=1}^N \left( \sum_{j'=0}^{M-1} w_{j'} \phi_{j'}(x^{(n)}) - t^{(n)} \right) \phi(x^{(n)})$$

def ||

$$\begin{bmatrix} \frac{\partial E(\mathbf{w})}{\partial w_1} \\ \frac{\partial E(\mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w_M} \end{bmatrix}$$

$$= \sum_{n=1}^N (\mathbf{w}^T \phi(x^{(n)}) - t^{(n)}) \phi(x^{(n)})$$

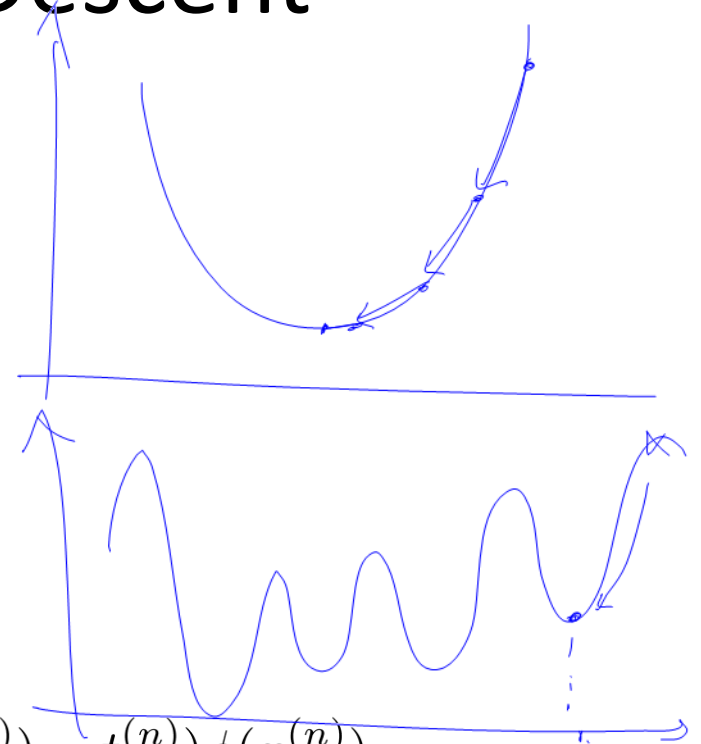
# Batch Gradient Descent

- Given data  $(x, y)$ , initial  $\mathbf{w}$ 
  - Repeat until convergence

$$\mathbf{w} := \mathbf{w} - \eta \nabla_{\mathbf{w}} E(\mathbf{w})$$

where

$$\begin{aligned} \nabla_{\mathbf{w}} E(\mathbf{w}) &= \sum_{n=1}^N \left( \sum_{j'=0}^{M-1} w_{j'} \phi_{j'}(x^{(n)}) - t^{(n)} \right) \phi(x^{(n)}) \\ &= \sum_{n=1}^N (\mathbf{w}^T \phi(x^{(n)}) - t^{(n)}) \phi(x^{(n)}) \end{aligned}$$



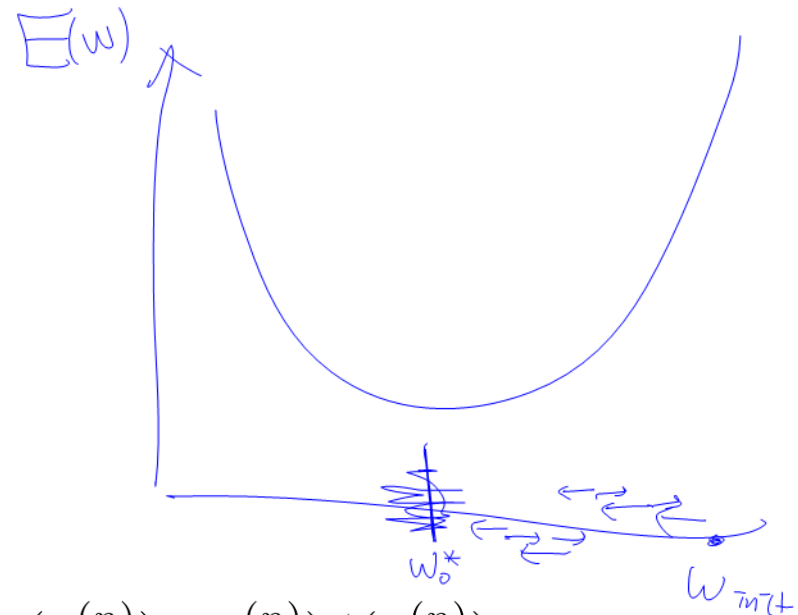
# Stochastic Gradient Descent

- Main idea: instead of computing batch gradient (over entire training data), just compute gradient for individual example and update
- Repeat until convergence
  - for  $n=1, \dots, N$

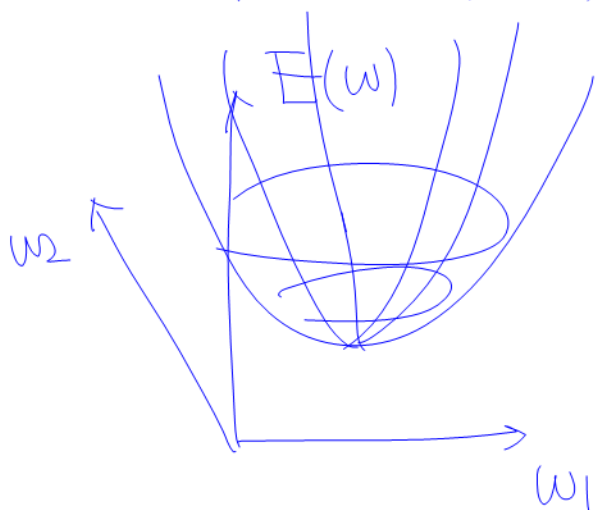
$$\mathbf{w} := \mathbf{w} - \eta \nabla_{\mathbf{w}} E(\mathbf{w} | x^{(n)})$$

where

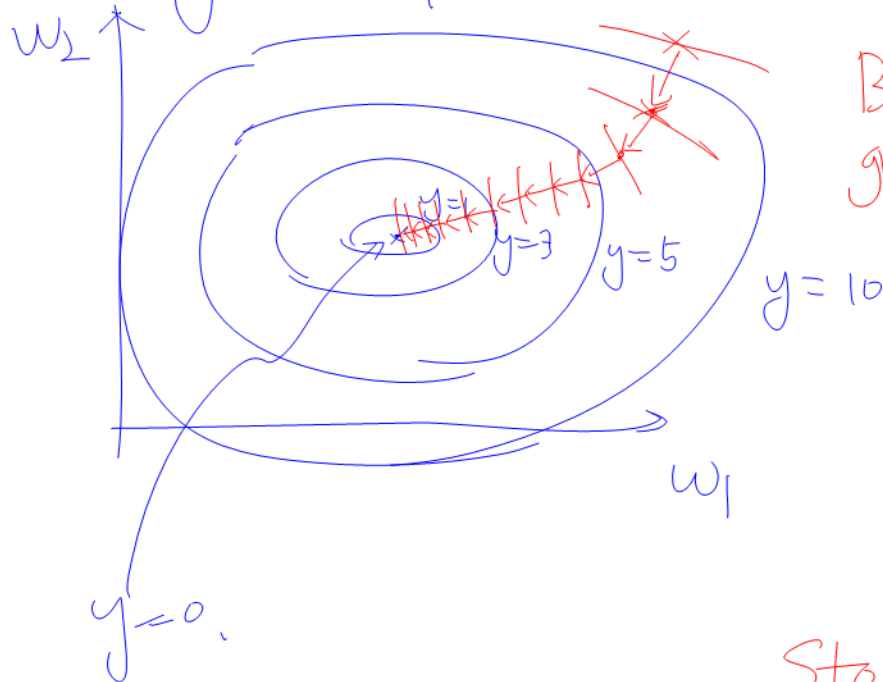
$$\begin{aligned} \nabla_{\mathbf{w}} E(\mathbf{w} | x^{(n)}) &= \left( \sum_{j'=0}^{M-1} w_{j'} \phi_{j'}(x^{(n)}) - t^{(n)} \right) \phi(x^{(n)}) \\ &= (\mathbf{w}^T \phi(x^{(n)}) - t^{(n)}) \phi(x^{(n)}) \end{aligned}$$



$$\phi_1(x), \phi_2(x)$$

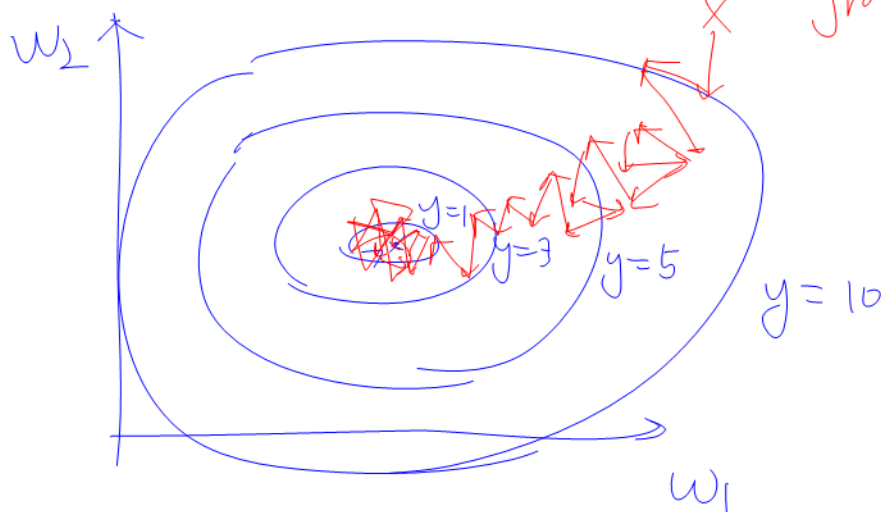


$$y = w_1 \phi_1(x) + w_2 \phi_2(x)$$



Batch  
grad. desc.

$$\eta(t) = \frac{\eta_0}{1 + \frac{t}{t_0}}$$



Stochastic  
gradient

# Closed form solution

- Main idea:
  - Compute gradient and set gradient to 0.  
(condition for optimal solution)
  - Solve the equation in a closed form

- Objective function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right)^2$$

- We will derive the gradient from matrix calculus

# Closed form solution

- Objective function:

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - t^{(n)} \right)^2 \\ &= \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(x^{(n)}) - t^{(n)})^2 \\ &= \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(x^{(n)}))^2 - \sum_{n=1}^N t^{(n)} \mathbf{w}^T \phi(x^{(n)}) + \frac{1}{2} \sum_{n=1}^N t^{(n)2} \\ &= \frac{1}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} \end{aligned}$$

- Recap: matrix calculus (check previous review session)



$$\begin{aligned}
 w^T \phi(x^{(n)}) &= \sum_{j=0}^M w_j \underbrace{\phi_j(x^{(n)})}_{\parallel} \\
 &= \sum_{j=0}^M \Phi_{n,j} w_j
 \end{aligned}$$

$$= \underbrace{\left( \begin{array}{c} \uparrow \\ \Phi \end{array} \right)}_{N \times M} \underbrace{\left( \begin{array}{c} \uparrow \\ w \end{array} \right)}_{M \times 1} \Bigg|_n$$

$$\Phi w = \left[ \begin{array}{c} \\ \\ \\ \end{array} \right]_{N\text{-dim}}$$

$$\sum_{n=1}^N \left( w^T \phi(x^{(n)}) \right)^2 = \sum_{n=1}^N \underbrace{\left( \Phi w \right)_n^2}_{N \times 1}$$

$$= (\Phi w)^T (\Phi w) = w^T \Phi^T \Phi w$$

# The Data

- The design matrix is an  $N \times M$  matrix, applying
  - the  $M$  basis functions (across)
  - to  $N$  data points (down)

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \phi_0(\mathbf{x}_n) & \phi_1(\mathbf{x}_n) & \cdots & \phi_{M-1}(\mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

*n+th example*

$\Phi_{n,j} = \phi_j(\mathbf{x}_n)$

$$\Phi \mathbf{w} \approx \mathbf{t}$$

# Recap on Matrix Calculus

# The Gradient

- Suppose that  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  is a function that takes as input a matrix  $A$  of size  $m \times n$  and returns a real value (scalar). Then the gradient of  $f$  (with respect to  $A \in \mathbb{R}^{m \times n}$ ) is the matrix of partial derivatives, defined as:

$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \cdots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \cdots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \cdots & \frac{\partial f(A)}{\partial A_{mn}} \end{bmatrix}$$

$$(\nabla_A f(A))_{ij} = \frac{\partial f(A)}{\partial A_{ij}}.$$

# The Gradient

Note that the size of  $\nabla_A f(A)$  is always the same as the size of  $A$ . So if, in particular,  $A$  is just a vector  $x \in \mathbb{R}^n$ ,

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}.$$

- $\nabla_x (f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x).$
- For  $t \in \mathbb{R}$ ,  $\nabla_x (t f(x)) = t \nabla_x f(x).$

# Gradients and Hessians of Quadratic and Linear Functions (Recap)

- $\nabla_x b^T x = b$       $\frac{\partial}{\partial x_i} \left( \sum_{i=1}^N b_i x_i \right) = b_i$
- $\nabla_x x^T A x = 2Ax$  (if  $A$  symmetric)
- $\nabla_x^2 x^T A x = 2A$  (if  $A$  symmetric)

$$\frac{\partial}{\partial x_i} \left[ \sum_{j=1}^N x_j A_{ji} x_j \right] = 2 \sum_j A_{ji} x_j$$

# Gradient via matrix calculus

- Compute gradient and set to zero

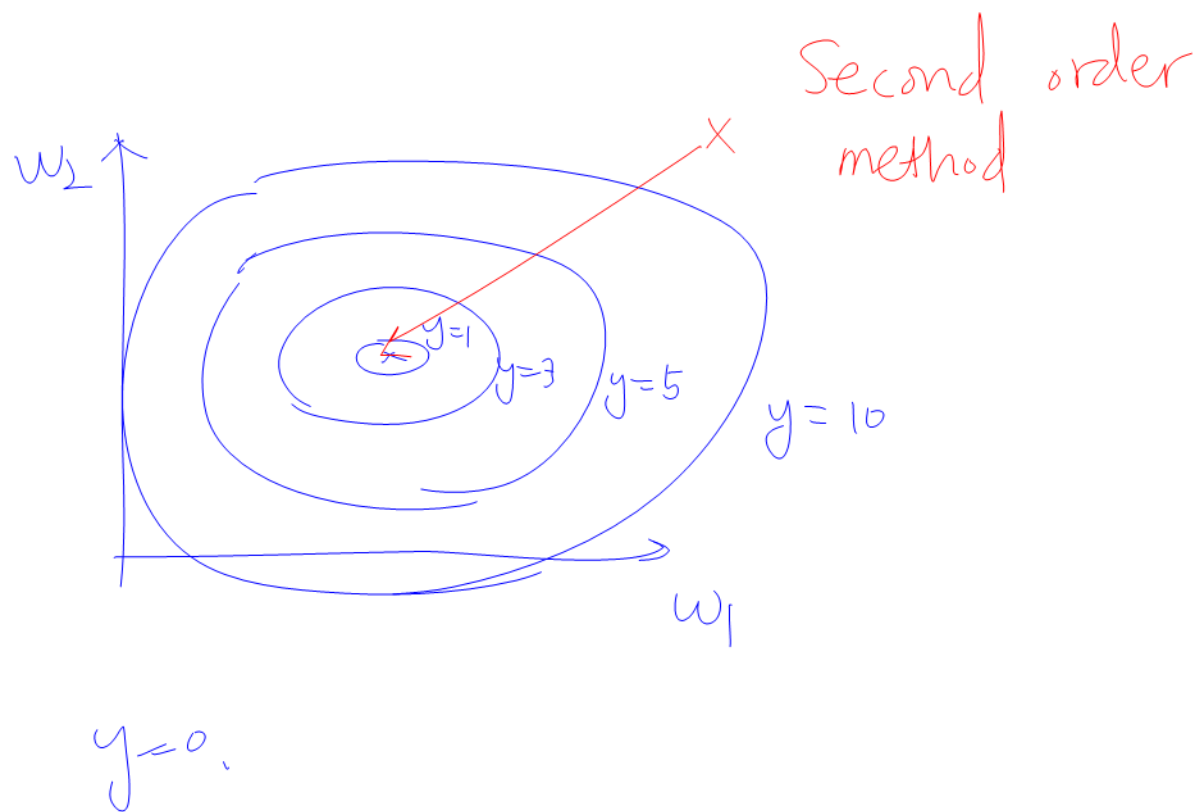
$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w}) &= \nabla_{\mathbf{w}} \left[ \frac{1}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} \right] \\ &= \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}\end{aligned}$$

- Solve the resulting equation (normal equation)

$$\begin{aligned}\Phi^T \Phi \mathbf{w} &= \Phi^T \mathbf{t} \\ \mathbf{w}_{ML} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}\end{aligned}$$

This is the *Moore-Penrose pseudo-inverse*:  $\Phi^\dagger = (\Phi^T \Phi)^{-1} \Phi^T$

applied to:  $\Phi \mathbf{w} \approx \mathbf{t}$



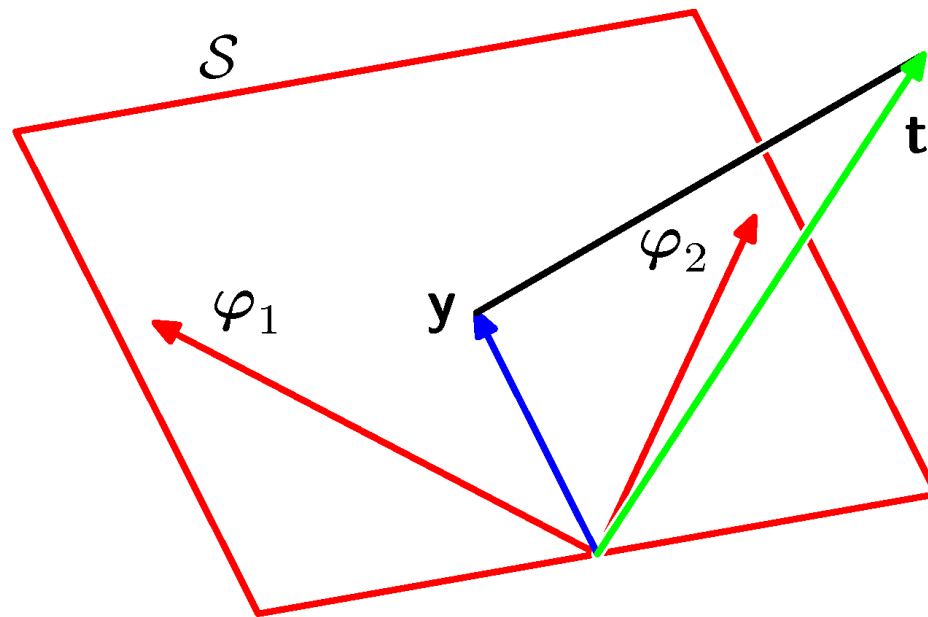


# Geometric Interpretation

- Assuming many more observations ( $N$ ) than the  $M$  basis functions  $\phi_j(\mathbf{x})$
- View the observed target values  $\mathbf{t}=\{t_1 \dots t_N\}$  as a vector in an  $N$ -dimensional space.
- The  $M$  basis functions  $\phi_j(\mathbf{x})$  span an  $M$ -dimensional subspace.
- $y(\mathbf{x}, \mathbf{w}_{ML})$  is the point in the subspace with minimal squared error from  $\mathbf{t}$ .
- It's the projection of  $\mathbf{t}$  onto that subspace.

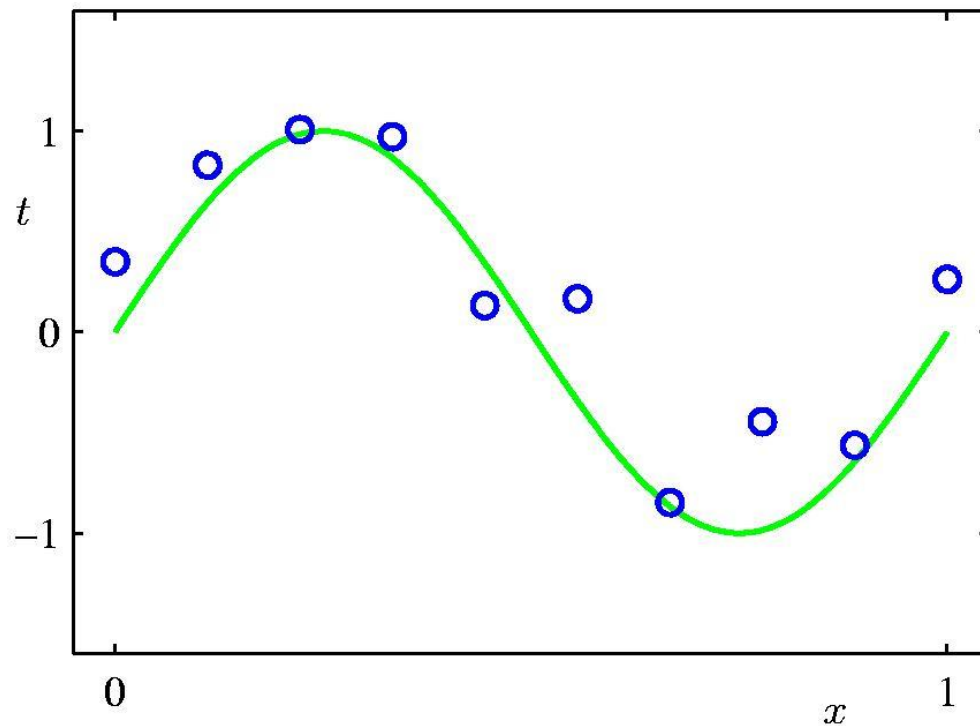
# Geometric Interpretation

- $y(\mathbf{x}, \mathbf{w}_{ML})$  is the projection of  $\mathbf{t}$  onto the subspace spanned by the  $M$  basis functions  $\phi_j(\mathbf{x})$



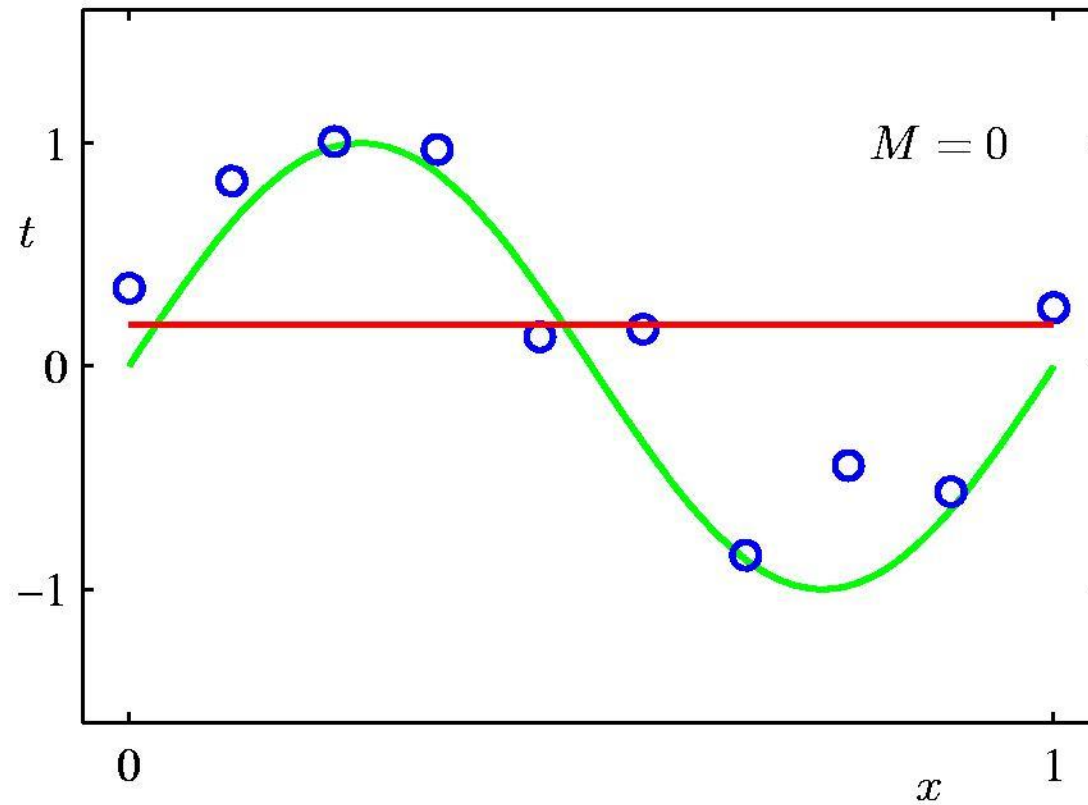
Back to curve-fitting examples

# Polynomial Curve Fitting

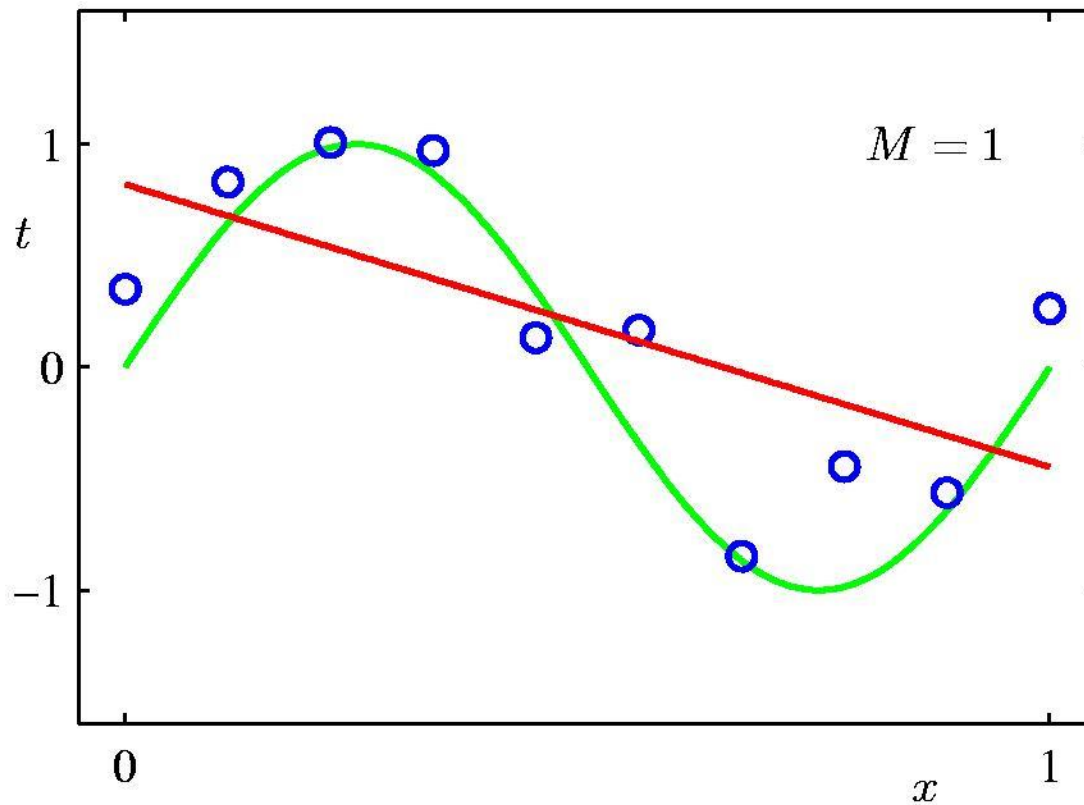


$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

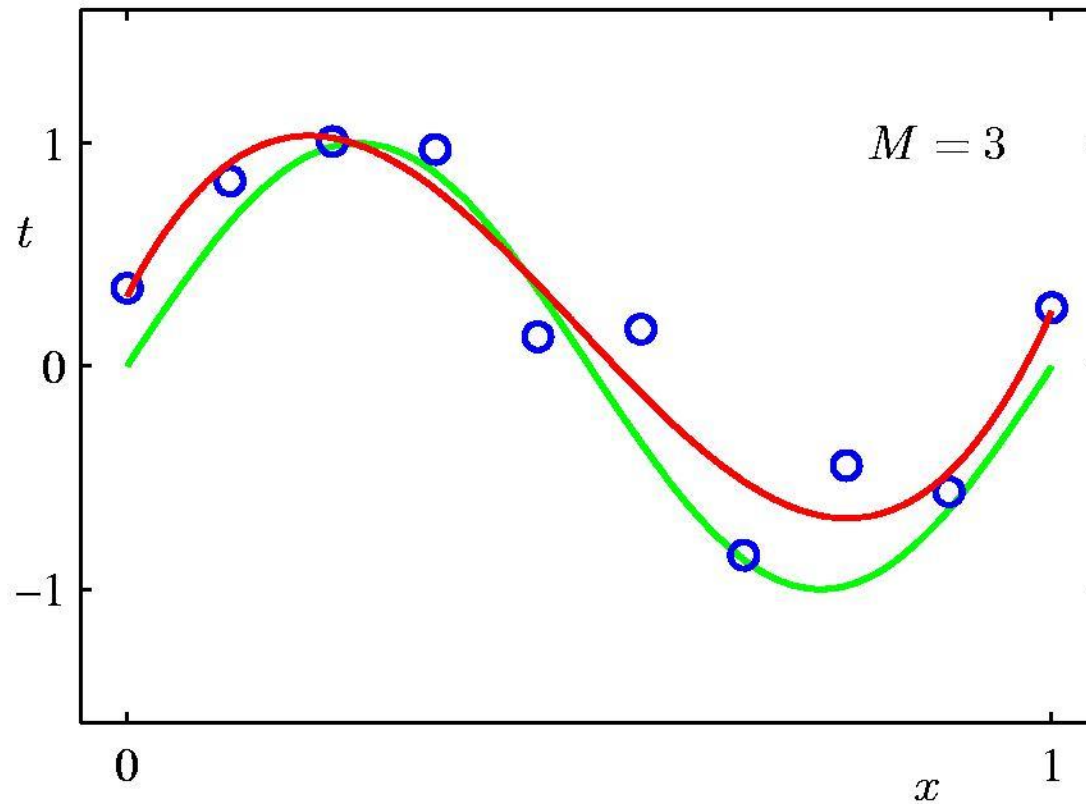
# 0<sup>th</sup> Order Polynomial



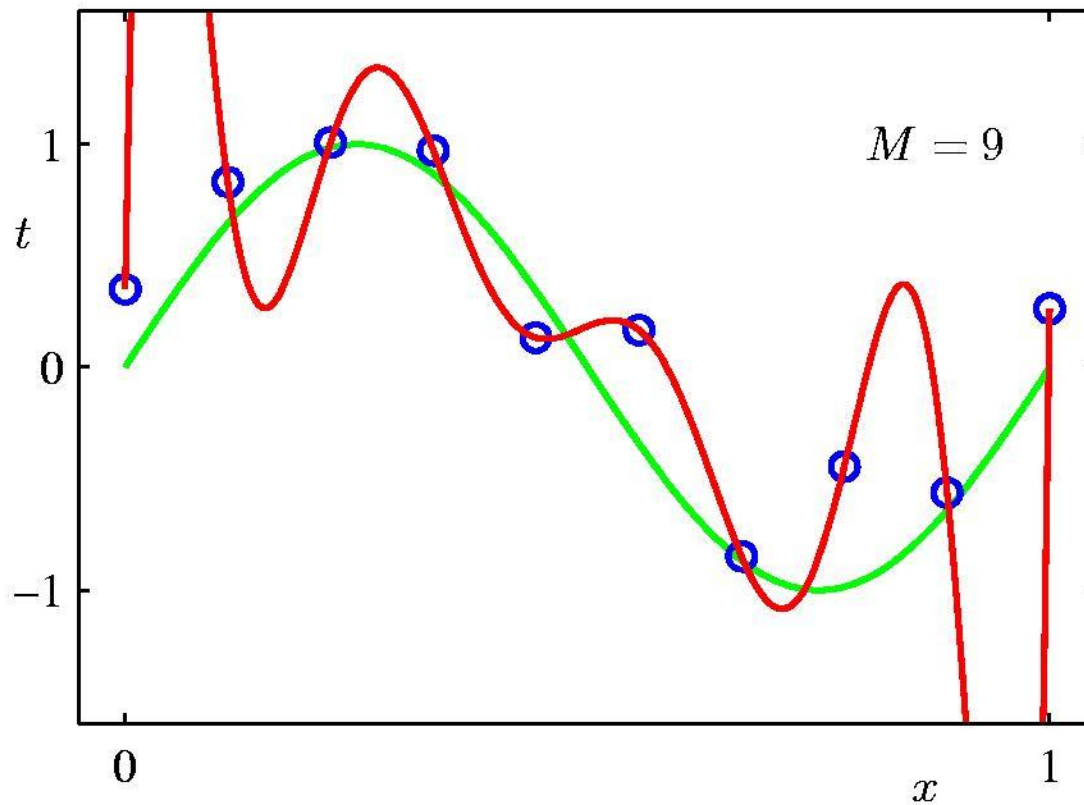
# 1<sup>st</sup> Order Polynomial



# 3<sup>rd</sup> Order Polynomial

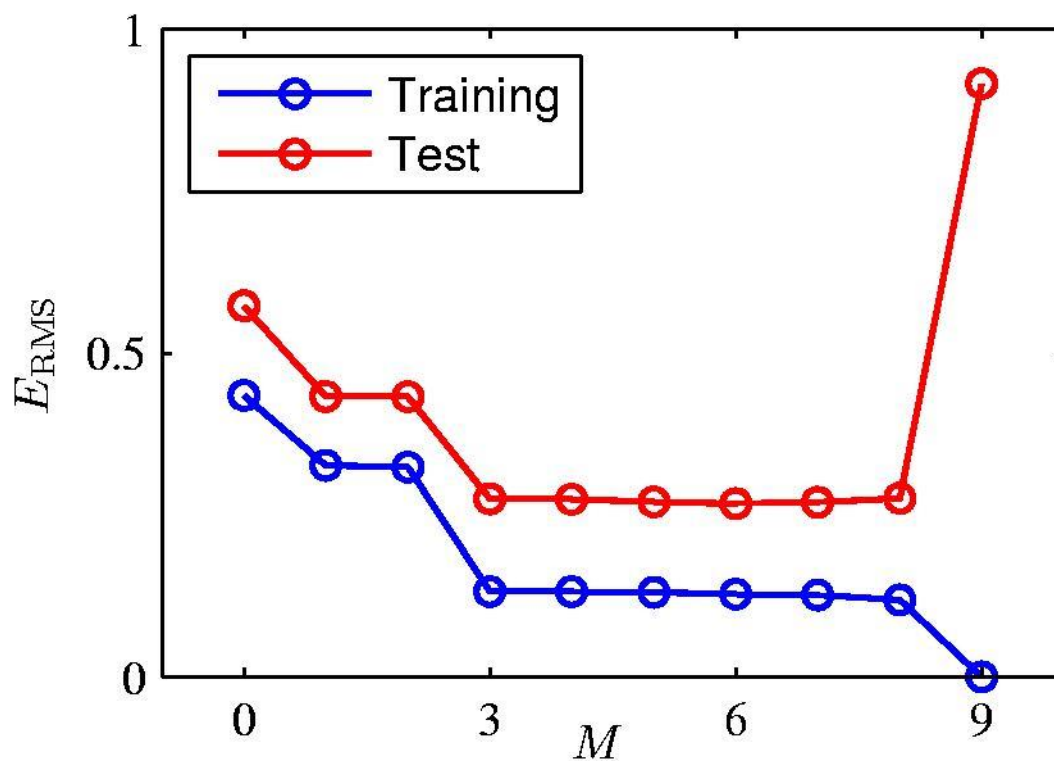


# 9<sup>th</sup> Order Polynomial





# Over-fitting



Root-Mean-Square (RMS) Error:

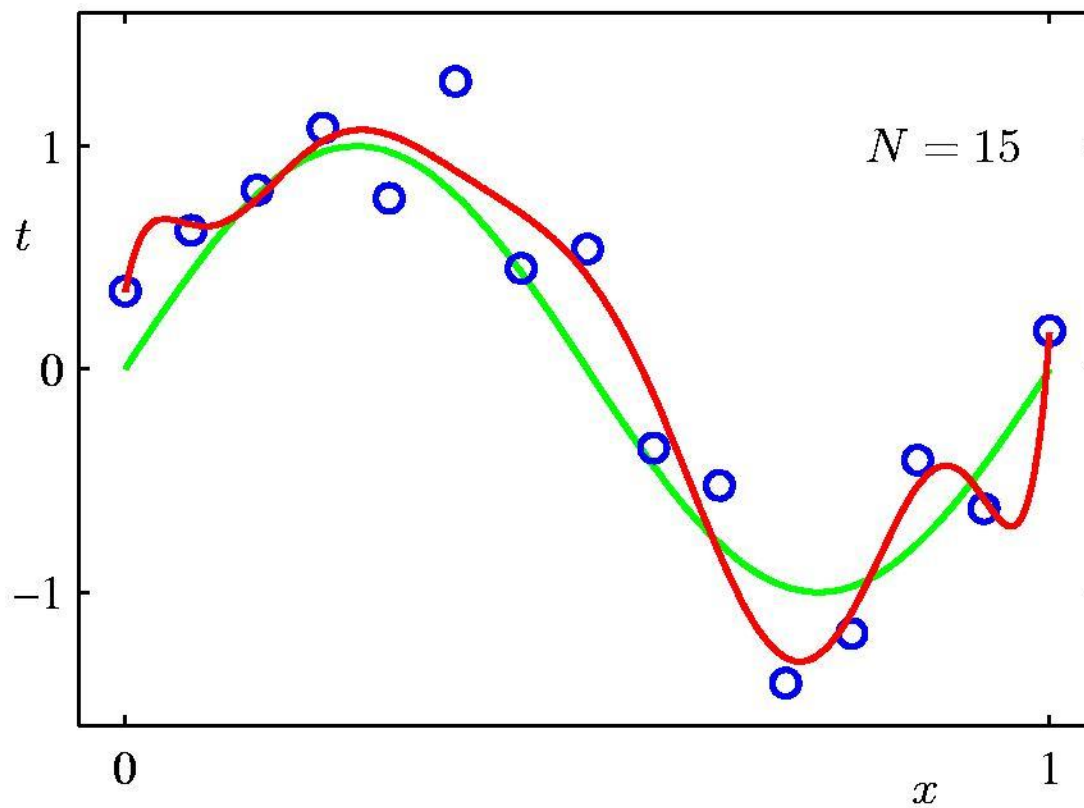
$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

# Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

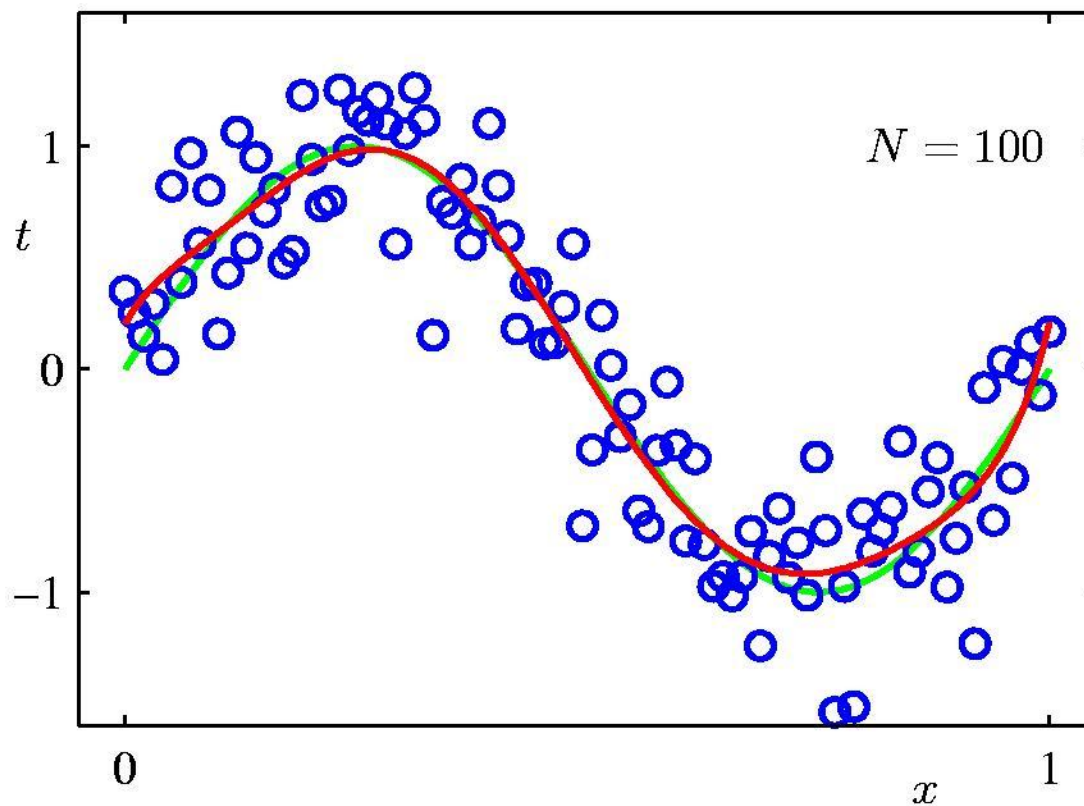
# Data Set Size: $N = 15$

## 9<sup>th</sup> Order Polynomial



Data Set Size:  $N = 100$

9<sup>th</sup> Order Polynomial



# Regularization

- Penalize large coefficient values

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

# Regularized Least Squares

- Add a regularization term to the error function

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

- Usual choices (sums of squares):

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\}^2 \quad E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

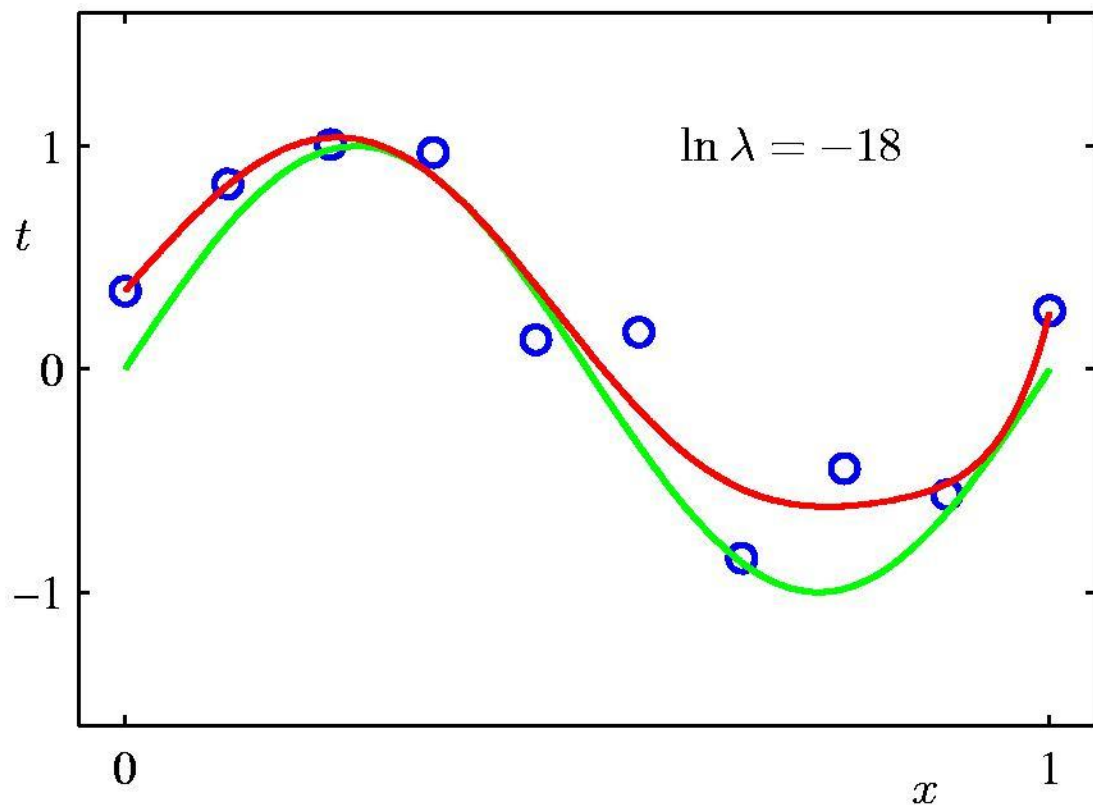
- Total error function becomes:

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

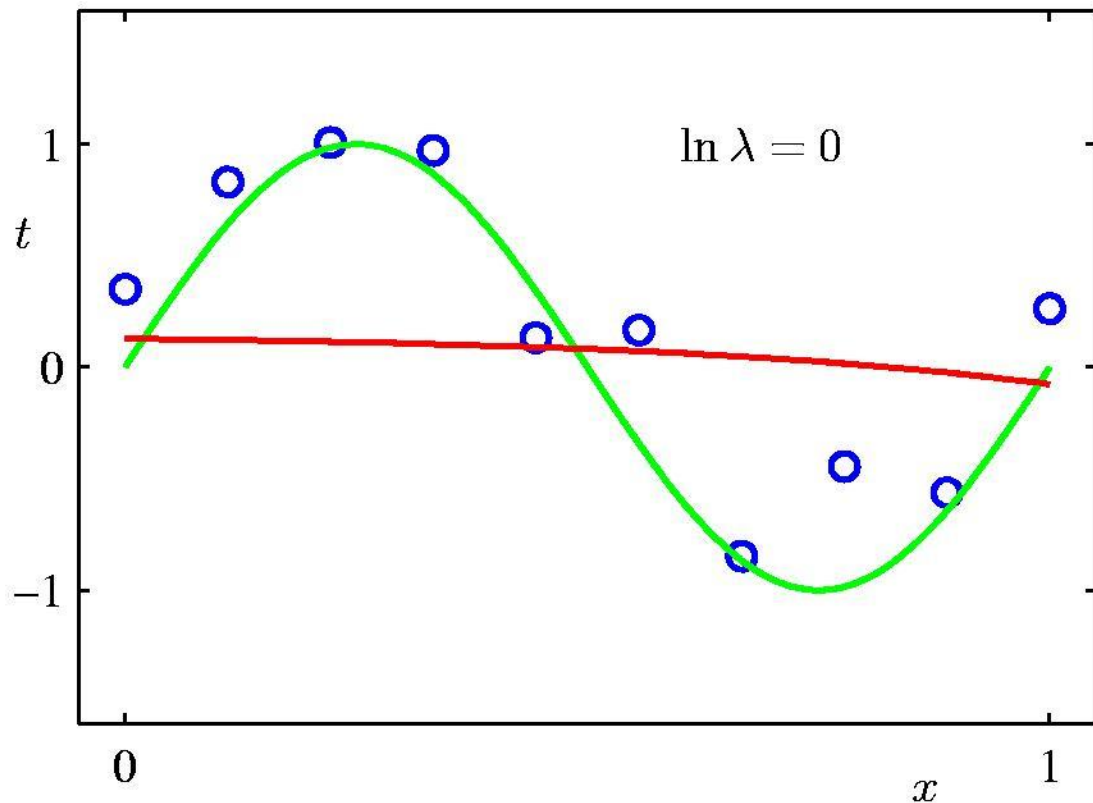
- Can still solve explicitly:

$$\mathbf{w}_{ML} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Regularization:  $\ln \lambda = -18$

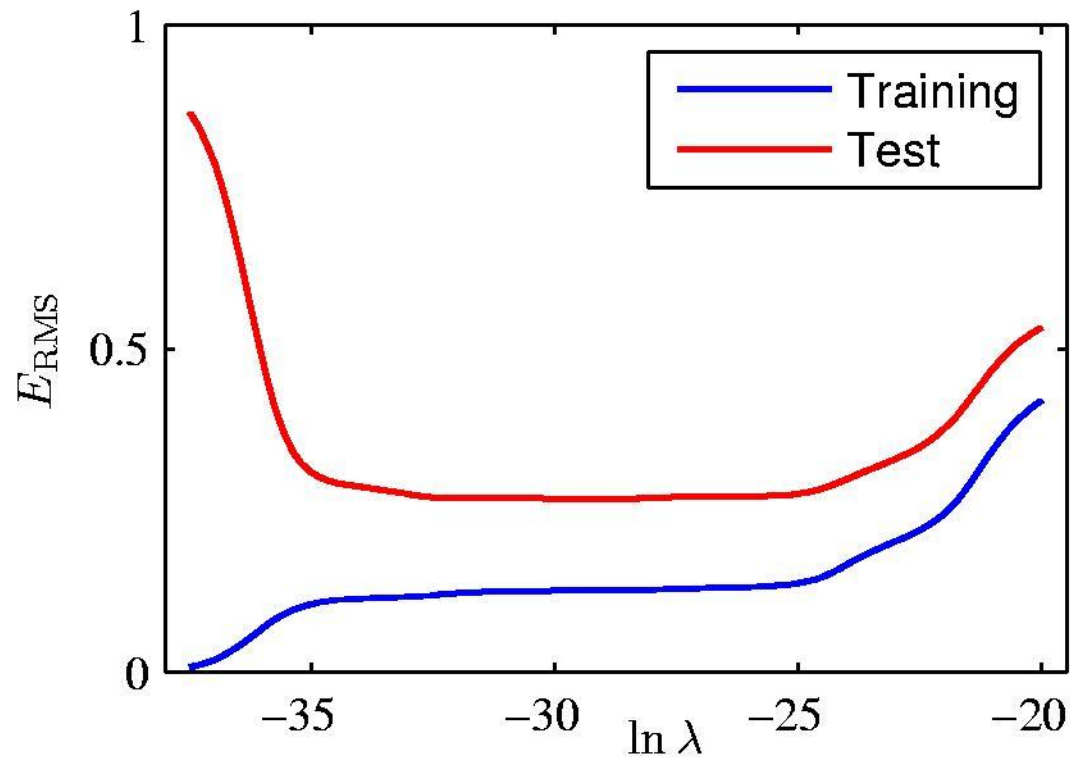


# Regularization: $\ln \lambda = 0$





# Regularization: $E_{\text{RMS}}$ vs. $\ln \lambda$



# Polynomial Coefficients

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01