

EECS 545: Machine Learning

Lecture 2. Linear Regression

Honglak Lee

1/13/2020

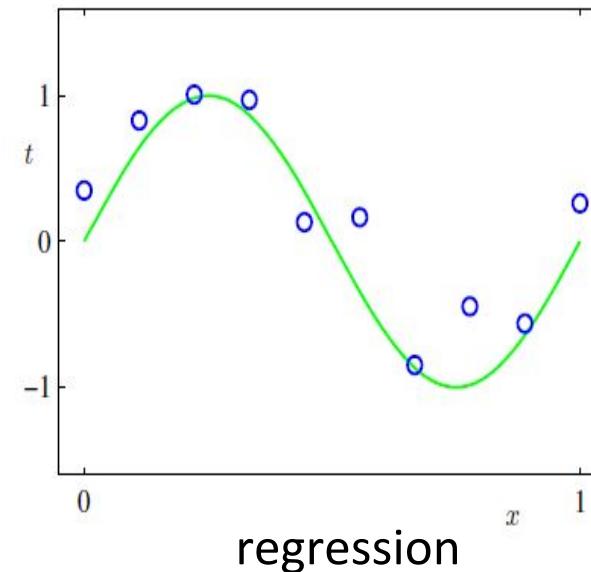
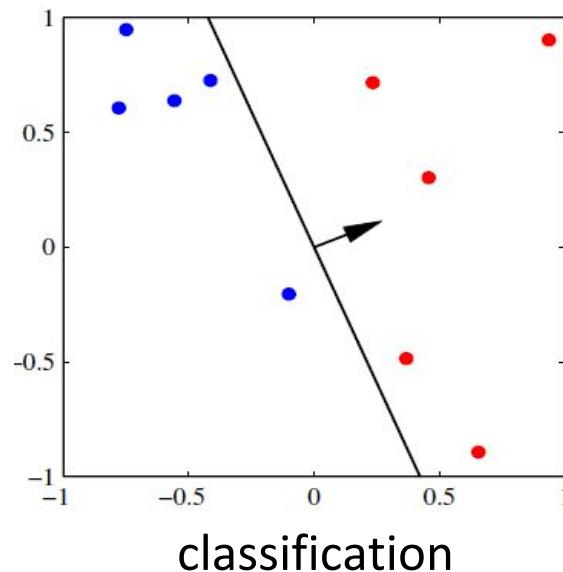


Announcement

- Homework #1 will be released tomorrow (1/14)
 - Due in 2 weeks: 1/28 5pm.
 - Form a study group and start early.
- Honor code
 - Collaboration and discussion is strongly encouraged, but you should write your own solution independently.
 - **Do not** refer to or copy solutions from any other people or other resources. In addition, please do not let other people copy your solution.

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-valued labels: classification
 - Continuous-valued labels: regression (today's topic)



Overview of Topics

- Linear Regression
 - Objective function
 - Vectorization
 - Computing gradient
 - Batch gradient vs. Stochastic Gradient
 - Closed form solution

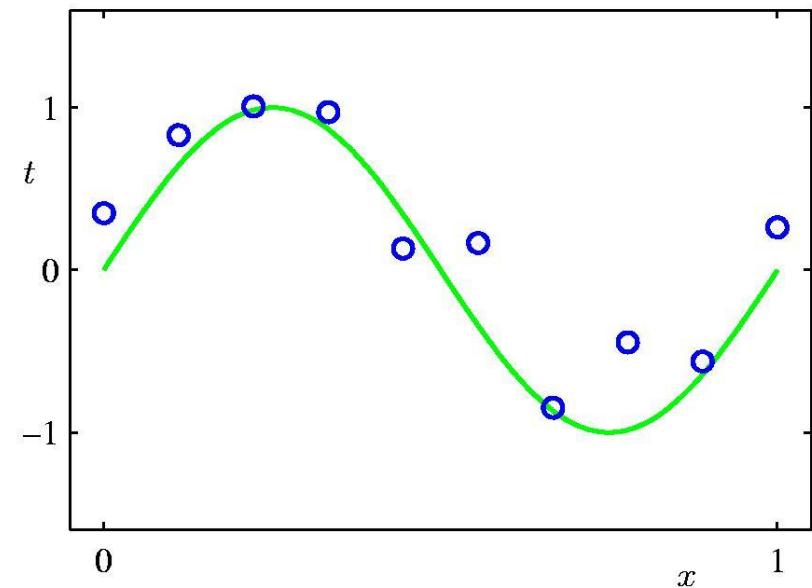
Notation

In this lecture, we will use the following notation:

- $\mathbf{x} \in \mathbb{R}^D$: data (scalar or vector)
- $\phi(\mathbf{x}) \in \mathbb{R}^M$: features for \mathbf{x} (vector)
- $\phi_j(\mathbf{x}) \in \mathbb{R}$: j-th feature for \mathbf{x} (scalar)
- $y \in \mathbb{R}$: continuous-valued label (i.e., target value)
- $\mathbf{x}^{(n)}$: denotes the n-th training example.
- $y^{(n)}$: denotes the n-th training label.

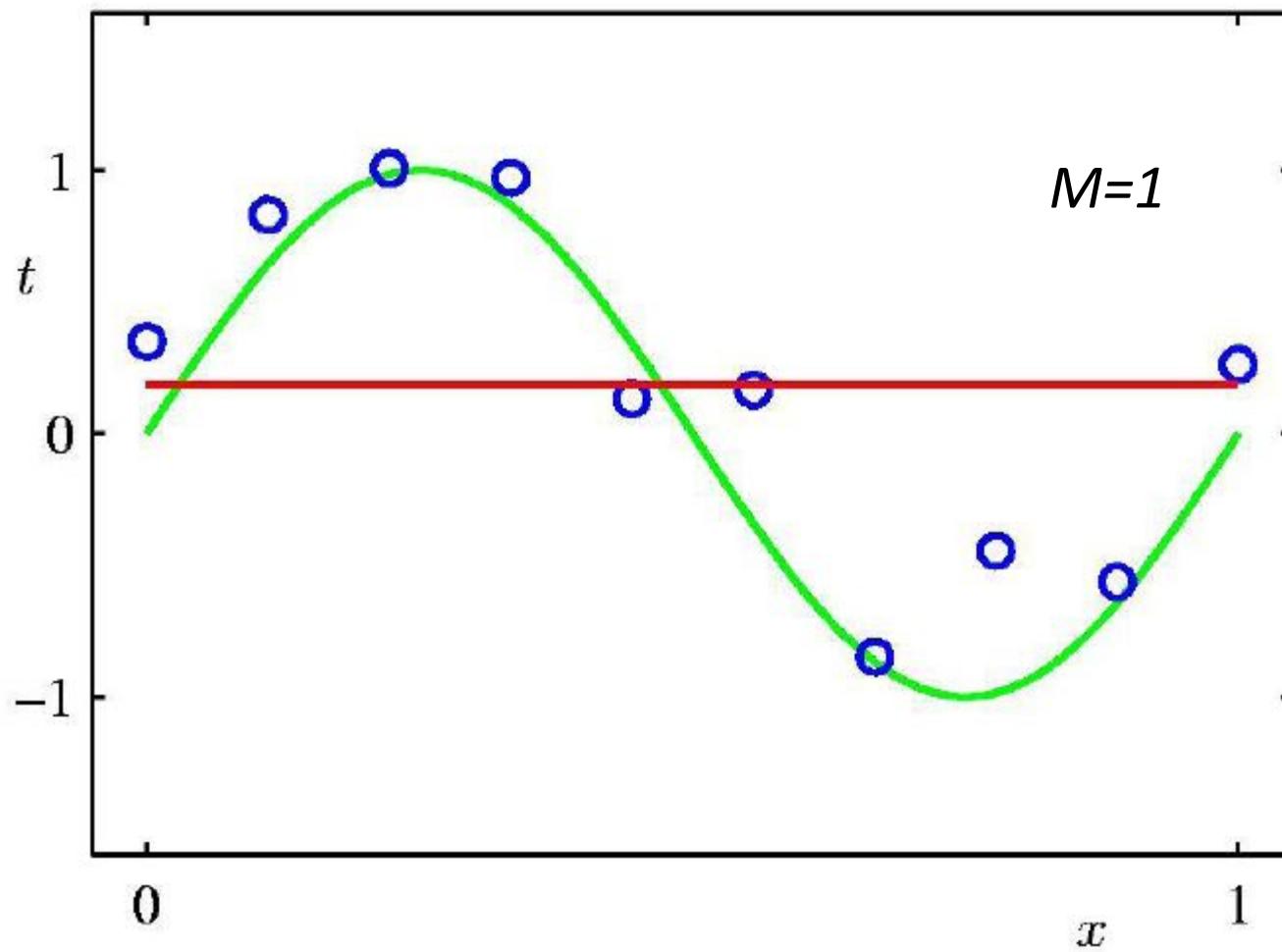
Linear regression (with 1d inputs)

- Consider 1d case (e.g., $D=1$)
- Given a set of observations
 $\{x^{(1)} \dots x^{(N)}\}$
- and corresponding target values: $\{y^{(1)} \dots y^{(N)}\}$
- We want to learn a function $h(x, \mathbf{w}) \approx y$ to predict future values.

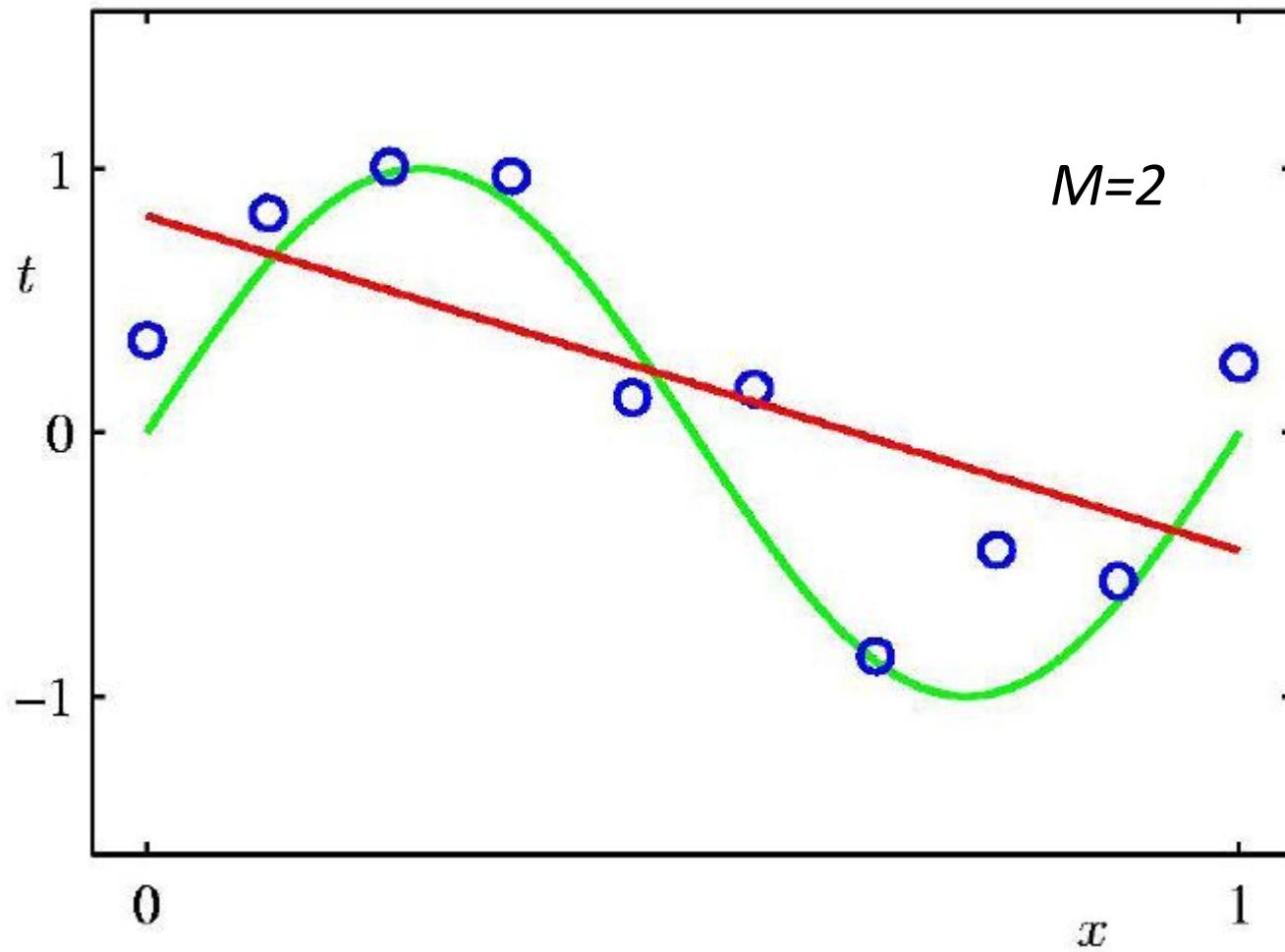


$$h(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_{M-1} x^{M-1} = \sum_{j=0}^{M-1} w_j x^j$$

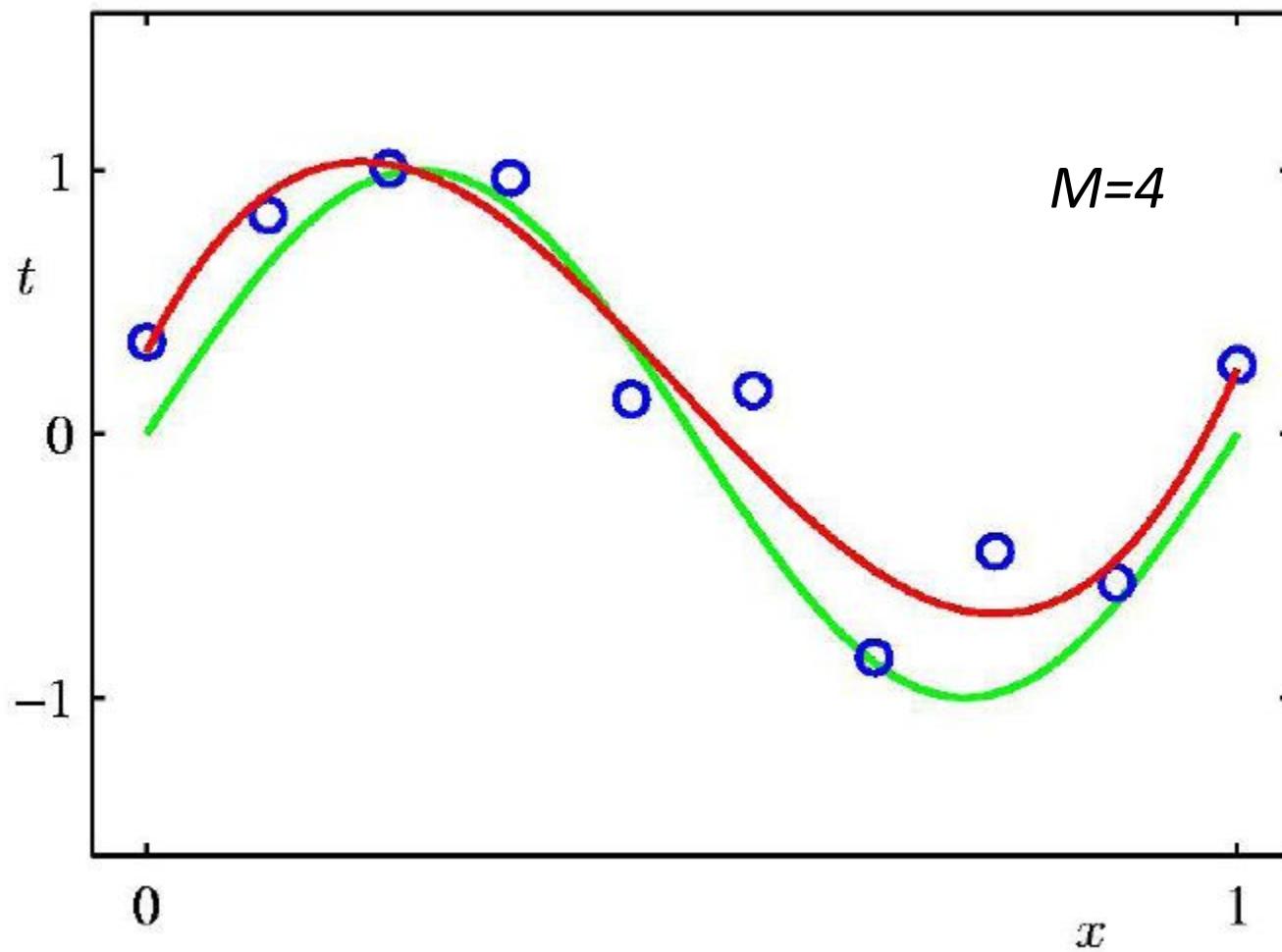
0th Order Polynomial



1st Order Polynomial



3rd Order Polynomial



Linear Regression (general case)

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

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

$$\phi_0(\mathbf{x}) = 1$$
$$h(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

where $\mathbf{w} = (w_0, \dots, w_{M-1}))^T$ (\mathbf{w} and $\phi(\mathbf{x})$ are column vectors)

$$\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$$

Basis Functions

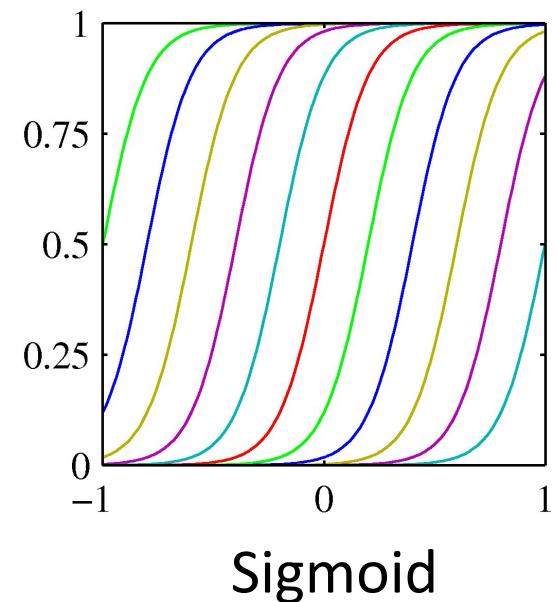
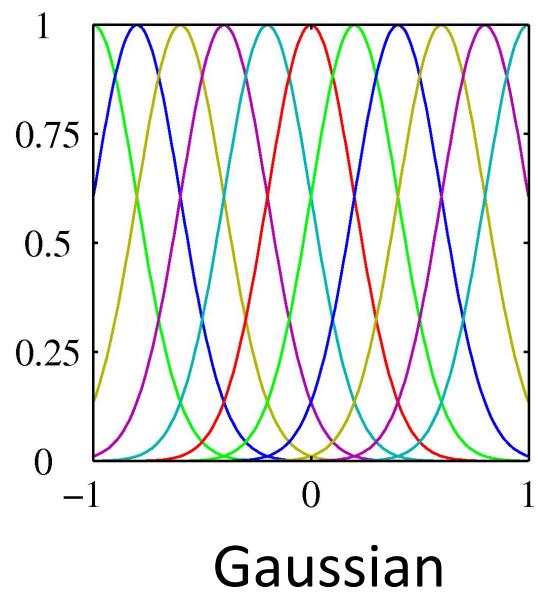
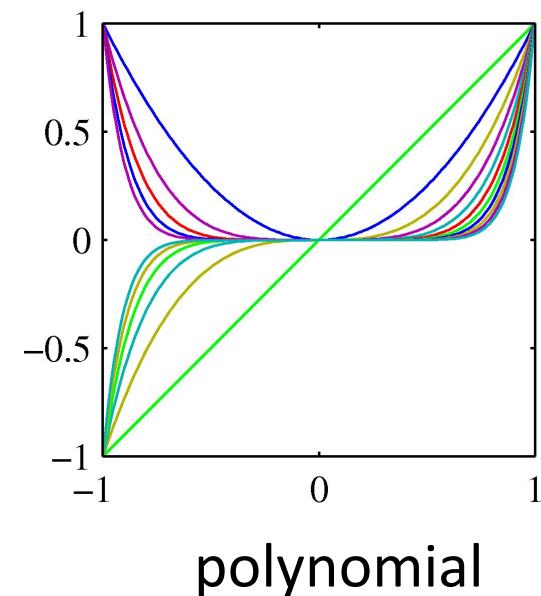
- The basis functions $\phi_j(\mathbf{x})$ need not be linear

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

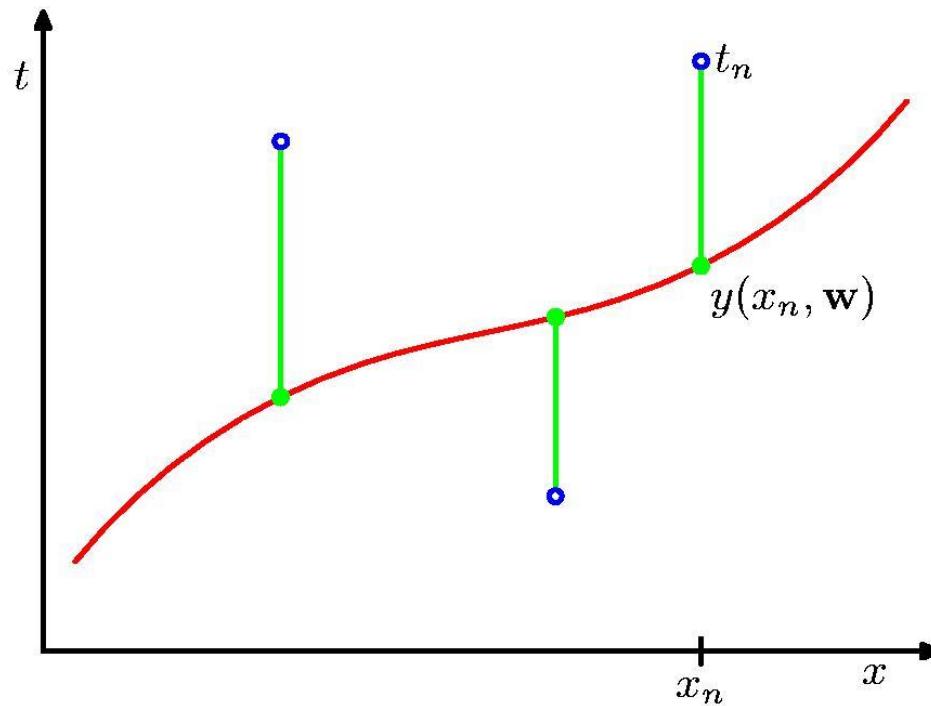
$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$

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

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



Objective: Sum-of-Squares Error Function

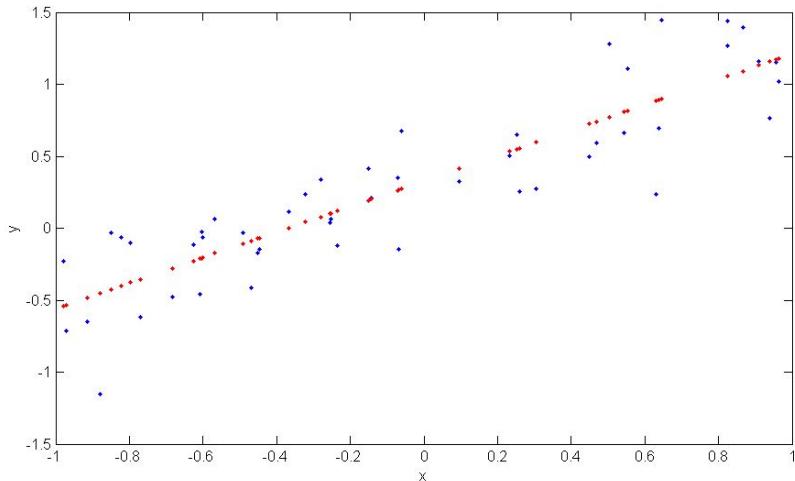


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ h(\mathbf{x}^{(n)}, \mathbf{w}) - y^{(n)} \right\}^2$$

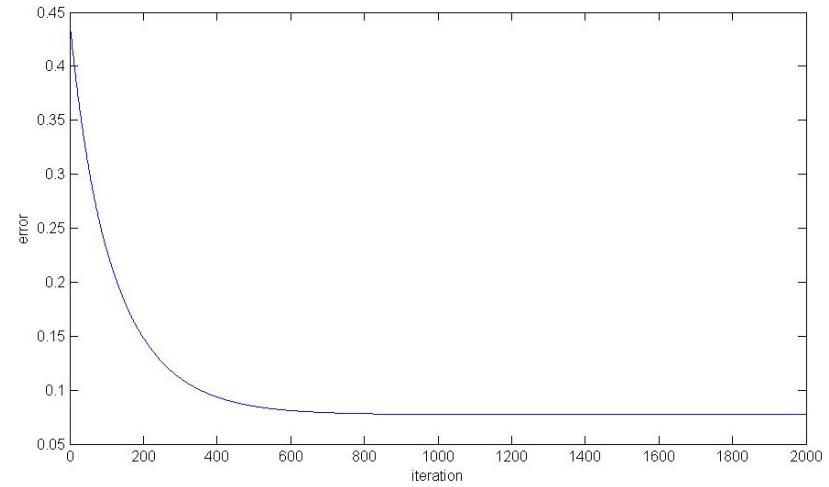
We want to find \mathbf{w} that minimizes $E(\mathbf{w})$ over the training data.

Linear regression via gradient descent (illustration)

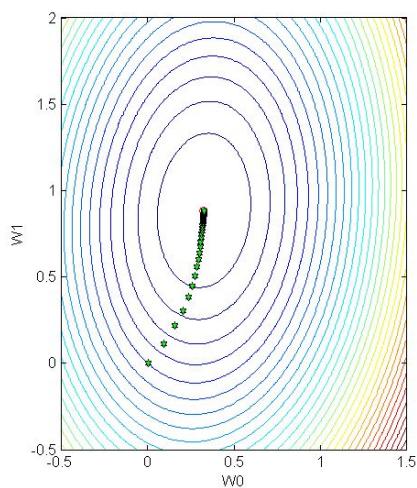
Training data (blue) vs. prediction (red)



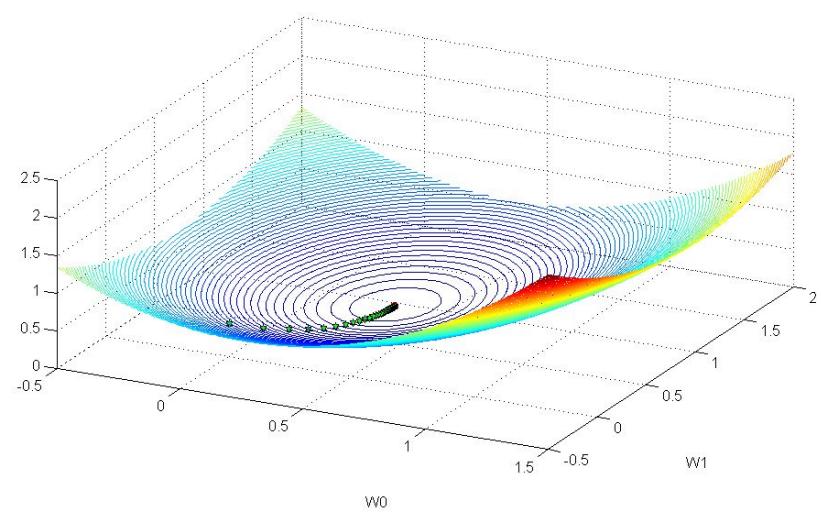
Error curve vs. training epoch



Contour plot of error



Contour plot of error (3d)



Least squares problem

- Objective function

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

- Gradient

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

Least squares problem

- Objective function

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

- Gradient

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \frac{\partial}{\partial w_j} \left(\sum_{j'=0}^{M-1} w_{j'} \phi_{j'}(\mathbf{x}^{(n)}) - y^{(n)} \right) \end{aligned}$$

Least squares problem

- Objective function

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

- Gradient

$$\begin{aligned} \frac{\partial E(\mathbf{w})}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \frac{\partial}{\partial w_j} \left(\sum_{j'=0}^{M-1} w_{j'} \phi_{j'}(\mathbf{x}^{(n)}) - y^{(n)} \right) \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_j(\mathbf{x}^{(n)}) \end{aligned}$$

Concatenate each component of the gradient:

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

We get a vectorized form of the gradient:

$$\nabla_{\mathbf{w}} E(\mathbf{w}) =$$

$$\begin{bmatrix} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{bmatrix}$$

Concatenate each component of the gradient:

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

We get a vectorized form of the gradient:

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \phi(\mathbf{x}^{(n)}) =$$

$$\begin{bmatrix} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{bmatrix} = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \begin{bmatrix} \phi_0(\mathbf{x}^{(n)}) \\ \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \vdots \\ \phi_{M-1}(\mathbf{x}^{(n)}) \end{bmatrix}$$

Concatenate each component of the gradient:

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

We get a vectorized form of the gradient:

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \phi(\mathbf{x}^{(n)}) =$$

$$\begin{bmatrix} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{bmatrix} = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \begin{bmatrix} \phi_0(\mathbf{x}^{(n)}) \\ \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \vdots \\ \phi_{M-1}(\mathbf{x}^{(n)}) \end{bmatrix}$$

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

Concatenate each component of the gradient:

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

We get a vectorized form of the gradient:

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \phi(\mathbf{x}^{(n)}) =$$

$$\begin{bmatrix} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{bmatrix} = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \begin{bmatrix} \phi_0(\mathbf{x}^{(n)}) \\ \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \vdots \\ \phi_{M-1}(\mathbf{x}^{(n)}) \end{bmatrix}$$

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

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

Gradient (compact, vectorized form)

- In summary, we have:

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

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

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \begin{bmatrix} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{bmatrix} \quad \phi(\mathbf{x}^{(n)}) = \begin{bmatrix} \phi_0(\mathbf{x}^{(n)}) \\ \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \vdots \\ \phi_{M-1}(\mathbf{x}^{(n)}) \end{bmatrix}$$

Batch Gradient Descent

- Given data (x, y) , initial w
 - Repeat until convergence

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

where

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

Stochastic Gradient Descent

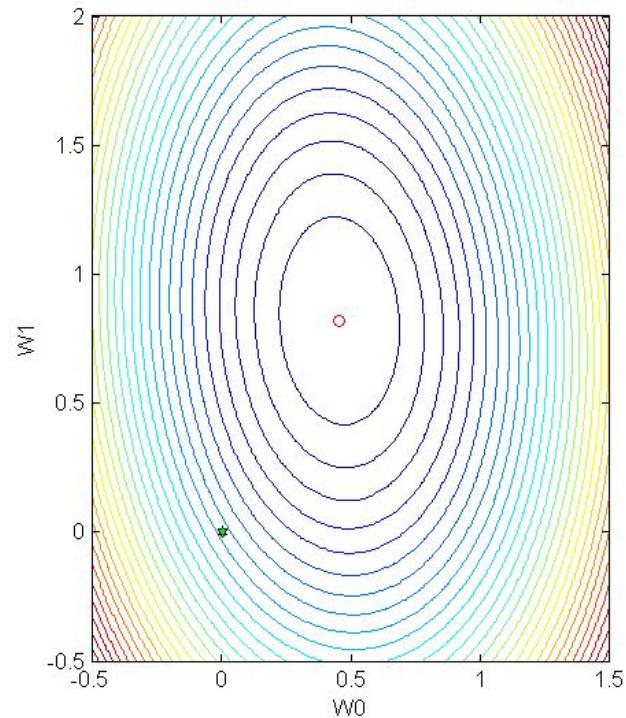
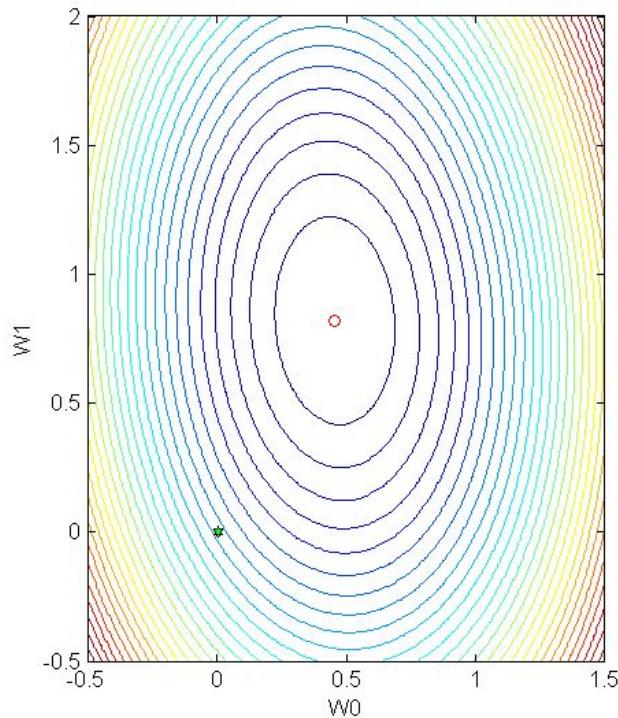
- Main idea: instead of computing batch gradient (over entire training data), just compute gradient for individual example and update
- Repeat until convergence
 - for $n=1,\dots,N$

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

where

$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w} | \mathbf{x}^{(n)}) &= \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi(\mathbf{x}^{(n)}) \\ &= \left(\mathbf{w}^T \phi(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi(\mathbf{x}^{(n)})\end{aligned}$$

Batch gradient vs. 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:

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

- We will derive the gradient from matrix calculus

Closed form solution

- Objective function:

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

Closed form solution

- Objective function:

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

Closed form solution

- Objective function:

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

Closed form solution

- Objective function:

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

- Trick: vectorization (by defining data matrix)

The data matrix

- The design matrix is an NxM matrix, applying
 - the M basis functions (columns)
 - to N data points (rows)

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix}$$

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

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix}$$

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

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix}$$

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

Useful trick: Matrix Calculus

- Idea so far:
 - Compute gradient and set gradient to 0.
(condition for optimal solution)
 - Solve the equation in a closed form using matrix calculus
- Need to compute the first derivative in matrix form

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:

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

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

Matrix calculus: 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 $t \in \mathbb{R}$, $\nabla_x(t f(x)) = t \nabla_x f(x).$

Gradient of Linear Functions

- Linear function

$$f(x) = \sum_{i=1}^n b_i x_i$$

- Gradient

$$\frac{\partial f(x)}{\partial x_k} = \frac{\partial}{\partial x_k} \sum_{i=1}^n b_i x_i = b_k.$$

- Compact form:

$$\nabla_{\mathbf{x}} \mathbf{b}^T \mathbf{x} = \mathbf{b}$$

Gradients of Quadratic Functions

- Quadratic function (A is symmetric):

$$f(\mathbf{x}) = \sum_{i,j=1}^n x_i A_{ij} x_j = \mathbf{x}^T A \mathbf{x}$$

- Gradient:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = 2 \sum_{j=1}^n A_{ij} x_j = 2(A\mathbf{x})_i$$

- Compact form:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = 2A\mathbf{x}$$

Putting together: Solution via matrix calculus

- Compute gradient and set to zero

$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w}) &= \nabla_{\mathbf{w}} \left(\frac{1}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{y} + \frac{1}{2} \mathbf{y}^T \mathbf{y} \right) \\ &= \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{y} \\ &= 0\end{aligned}$$

- Solve the resulting equation (normal equation)

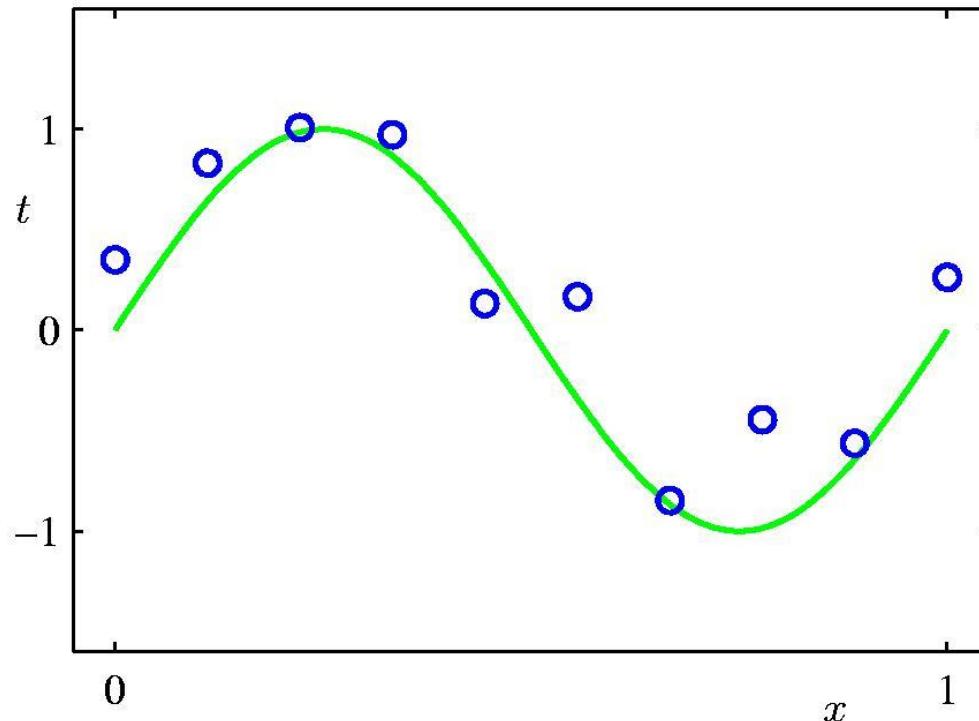
$$\Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{y}$$

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

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

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

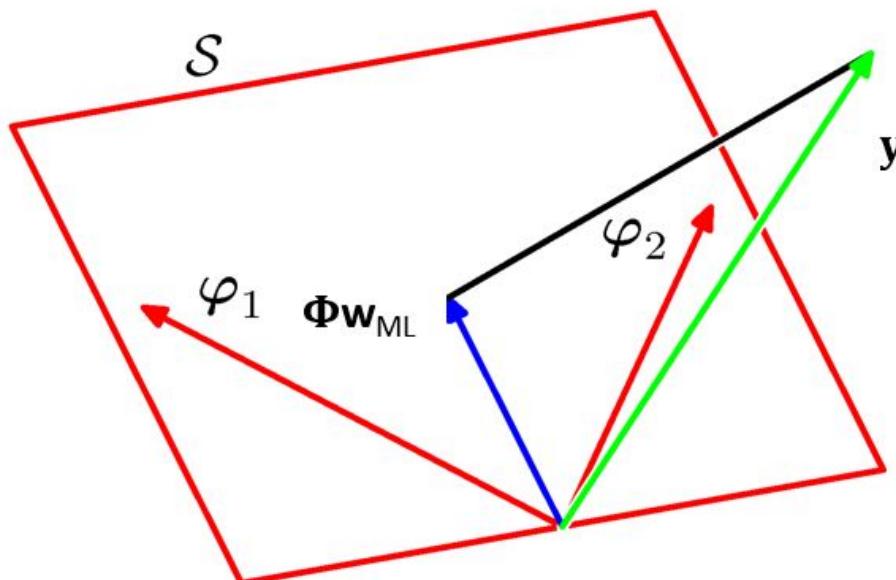
Polynomial Curve Fitting



$$h(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_{M-1} x^{M-1} = \sum_{j=0}^{M-1} w_j x^j$$

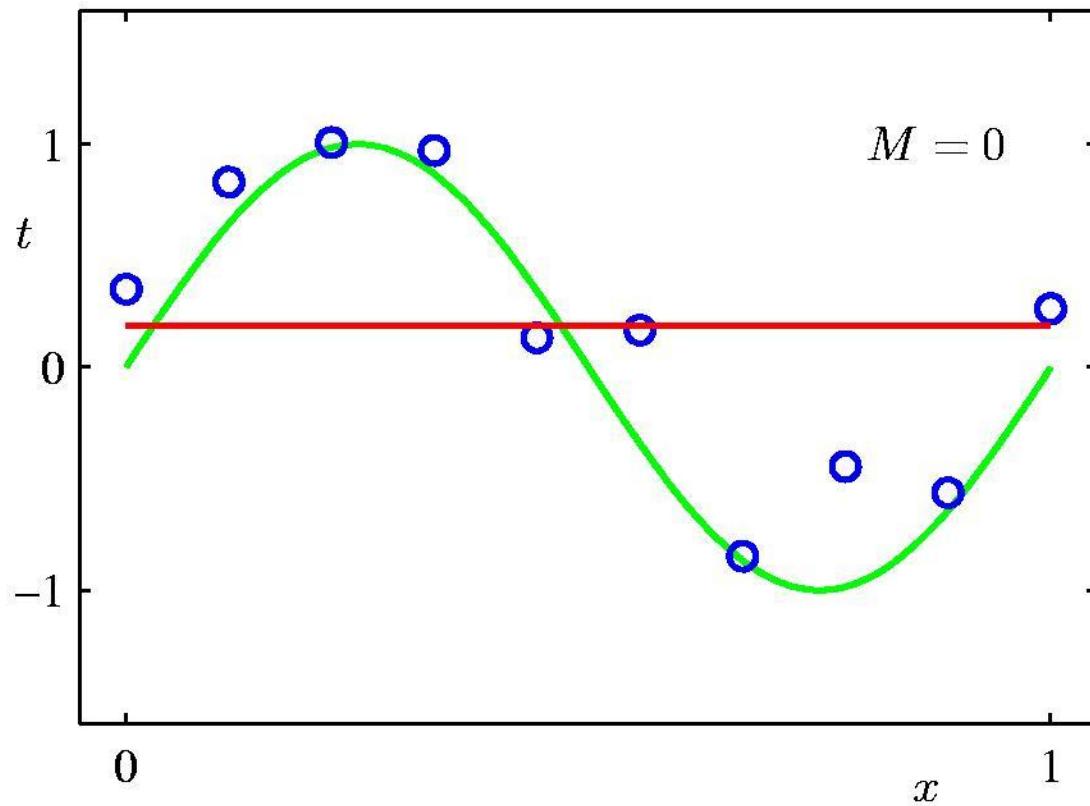
Geometric Interpretation

- Assuming many more observations (N) than the M basis functions $\phi_j(x)$ ($j=0,\dots,M-1$)
- View the observed target values $\mathbf{y} = \{y^{(1)}, \dots, y^{(N)}\}$ as a vector in an N -dim. space.
- The M basis functions $\phi_j(x)$ span the N -dimensional subspace.
 - Where the N -dim vector for ϕ_j is $\{\phi_j(\mathbf{x}^{(1)}), \dots, \phi_j(\mathbf{x}^{(N)})\}$
- $\Phi\mathbf{w}_{ML}$ is the point in the subspace with minimal squared error from \mathbf{y} .
- It's the projection of \mathbf{y} onto that subspace.

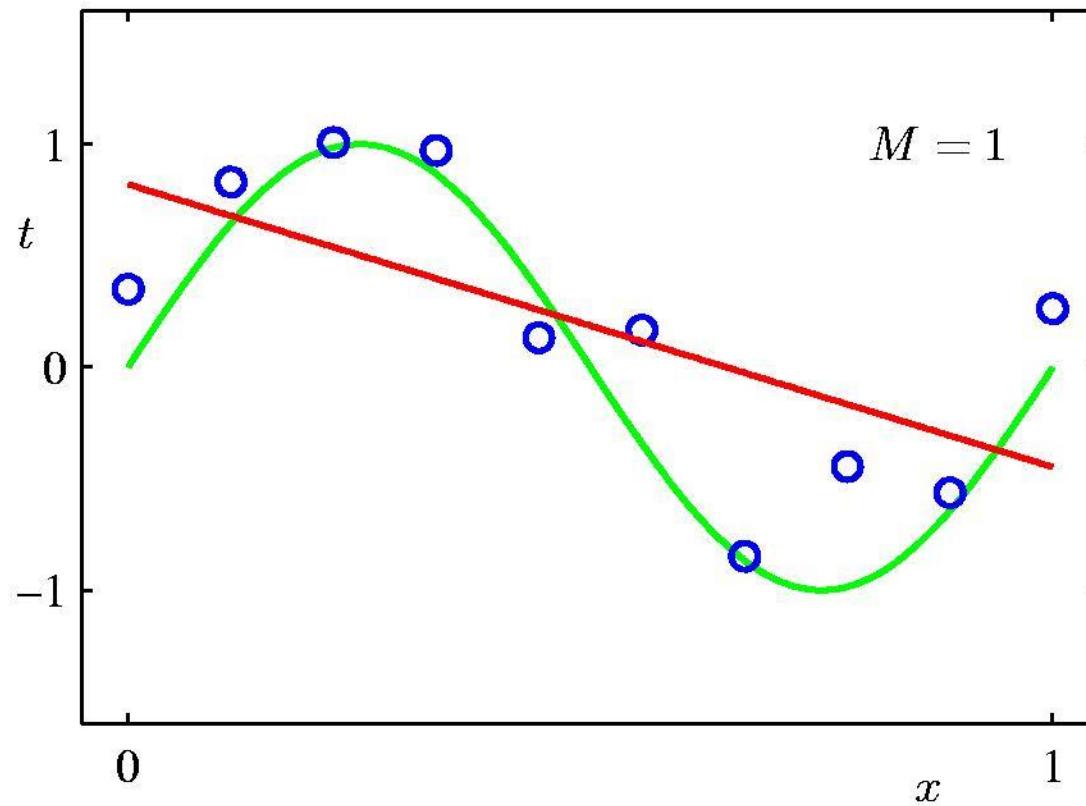


Back to curve-fitting examples

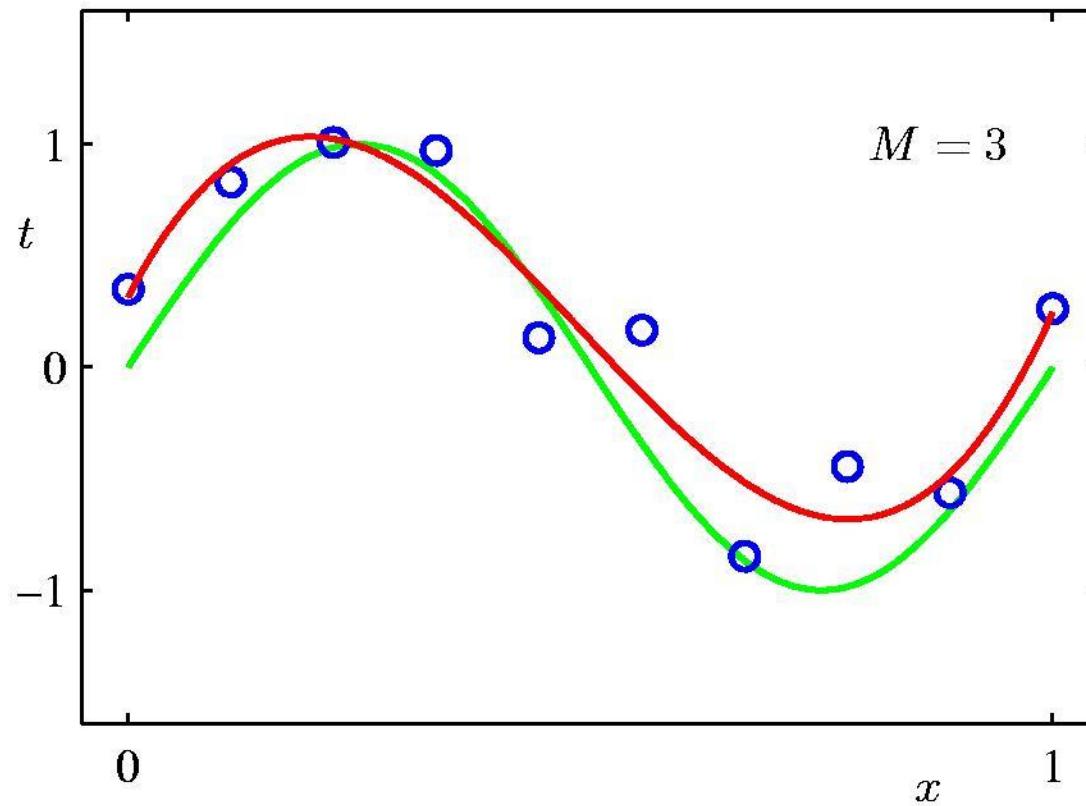
0th Order Polynomial



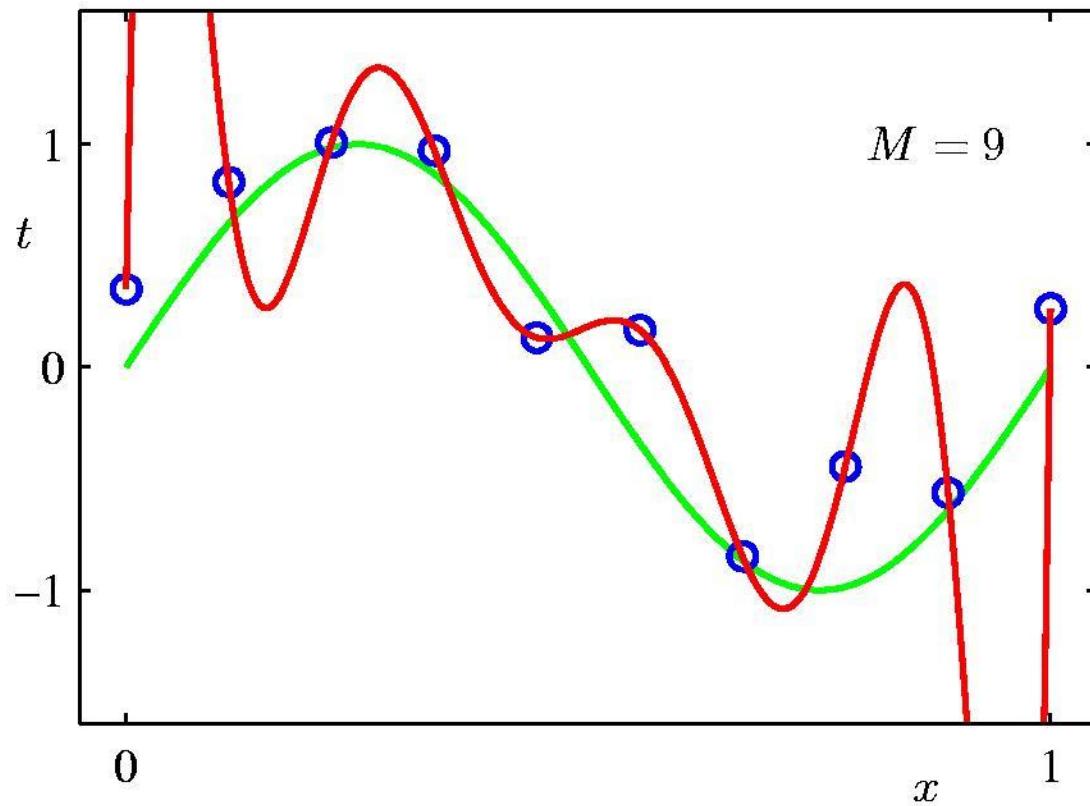
1st Order Polynomial



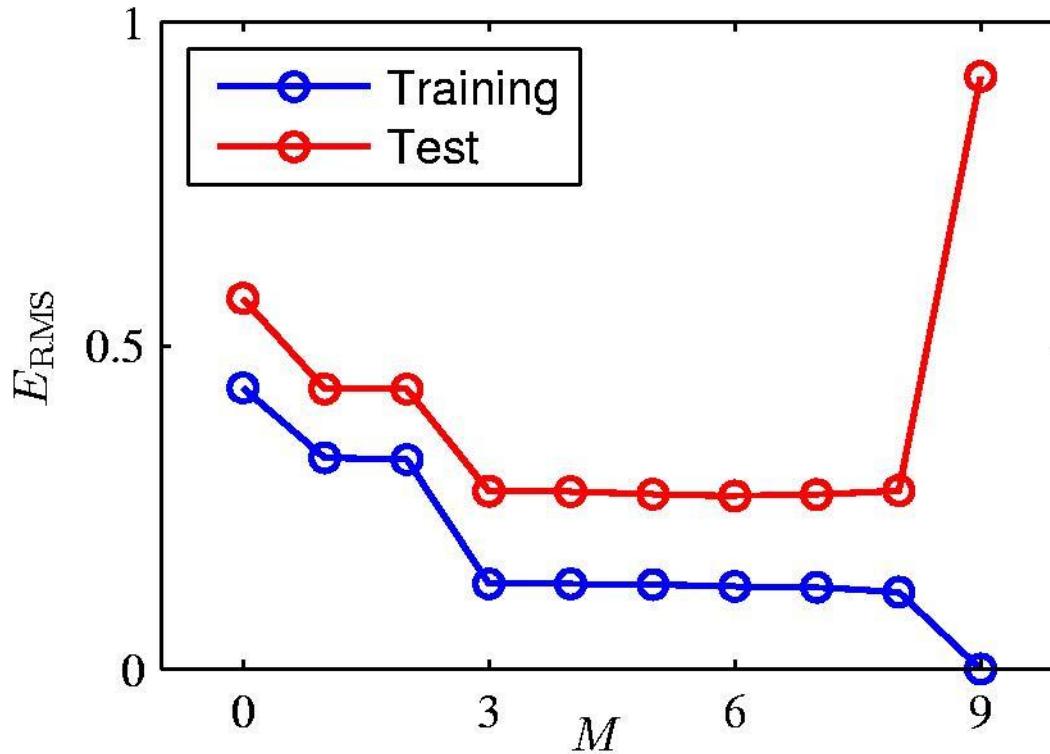
3rd Order Polynomial



9th Order Polynomial



Over-fitting



Root-Mean-Square (RMS) Error:

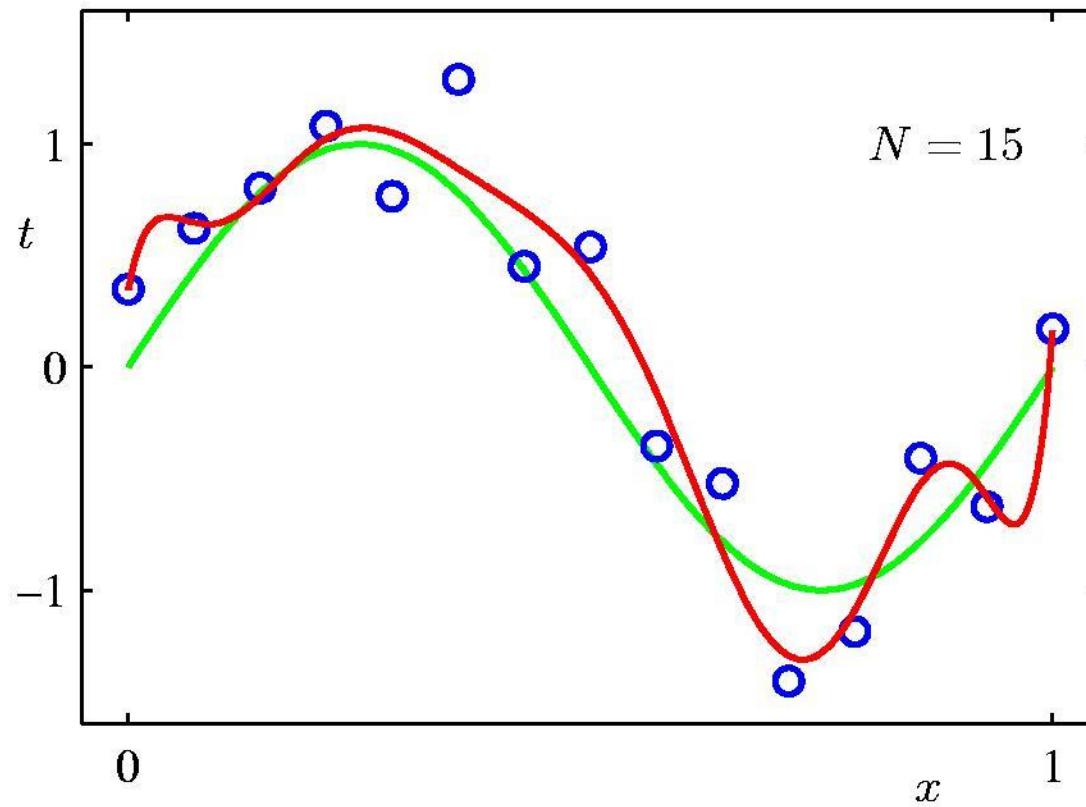
$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

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

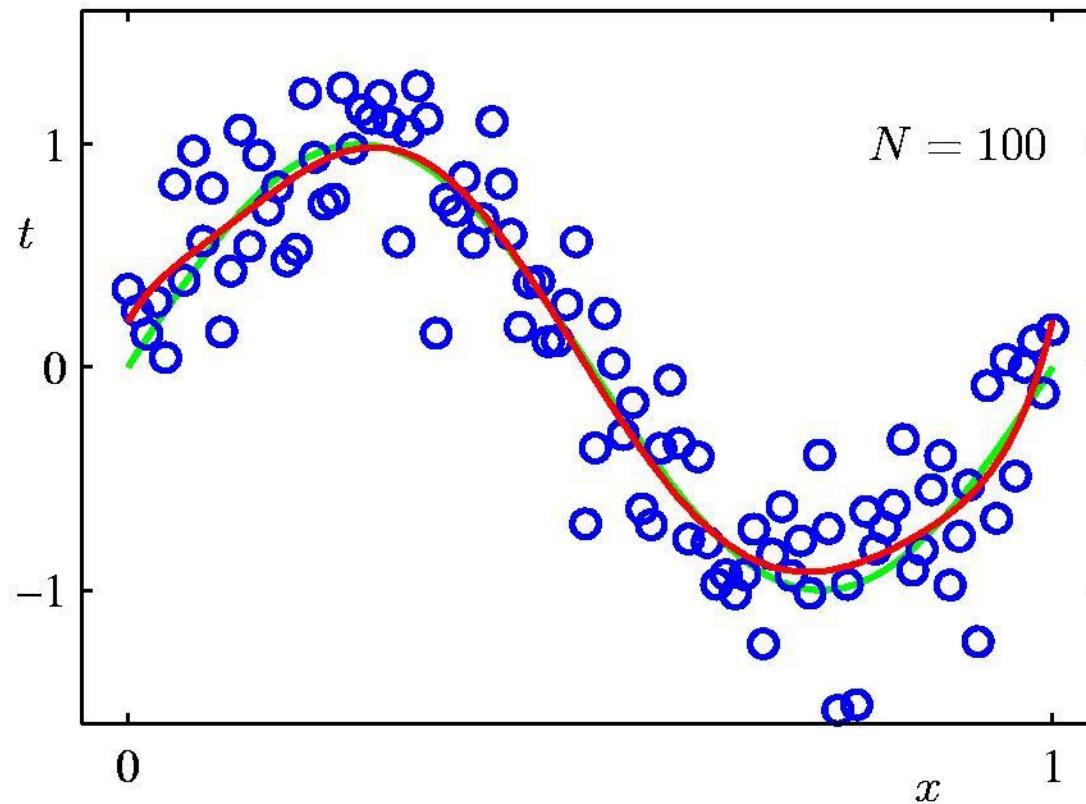
Data Set Size: $N = 15$

9th Order Polynomial



Data Set Size: $N = 100$

9th Order Polynomial



Q. How do we choose the degree of polynomial?

Rule of thumb

- If you have a small number of data points, then you should use low order polynomial (small number of features).
 - Otherwise, your model will overfit
- As you obtain more data points, you can gradually increase the order of the polynomial (more features).
 - However, your model is still limited by the finite amount of the data available (i.e., the optimal model for finite data cannot be infinite dimensional polynomial).
- Controlling model complexity: **regularization**