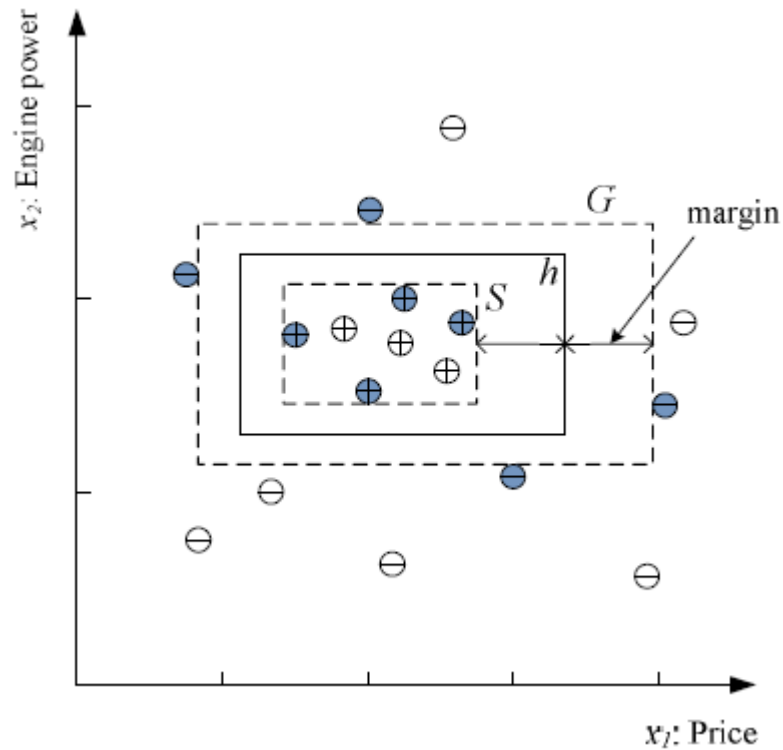


# CS60020: Foundations of Algorithm Design and Machine Learning

Sourangshu Bhattacharya

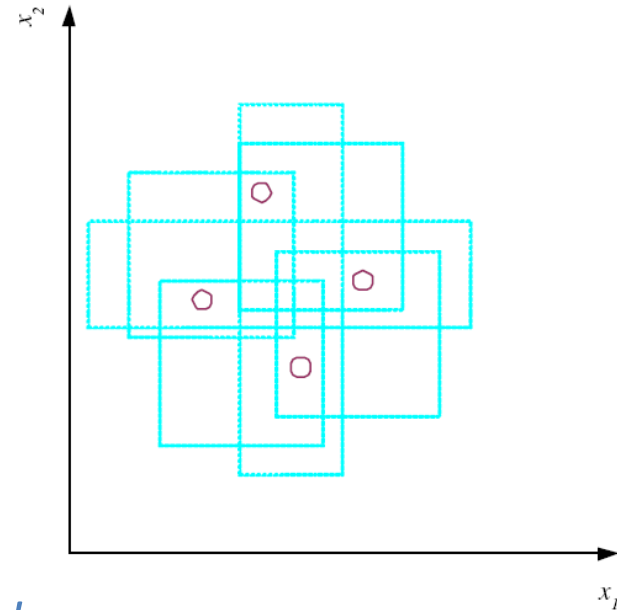
# Margin

- Choose  $h$  with largest margin



# VC Dimension

- $N$  points can be labeled in  $2^N$  ways as  $+/-$
- $\mathcal{H}$  shatters  $N$  if there exists  $h \in \mathcal{H}$  consistent for any of these:  
$$VC(\mathcal{H}) = N$$

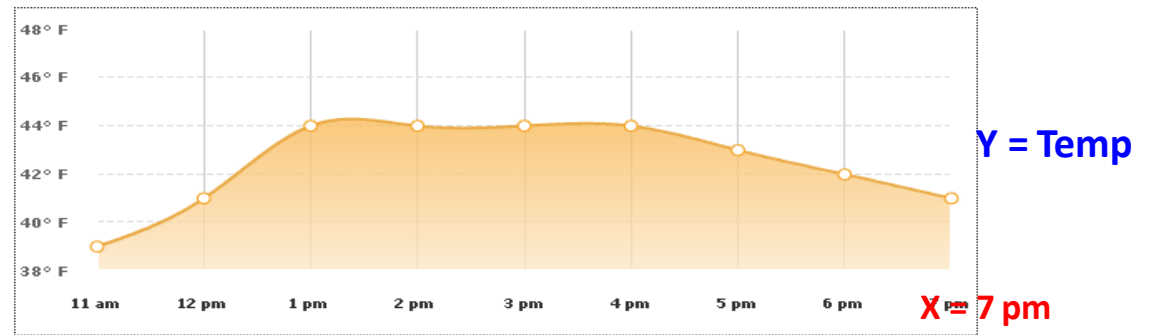
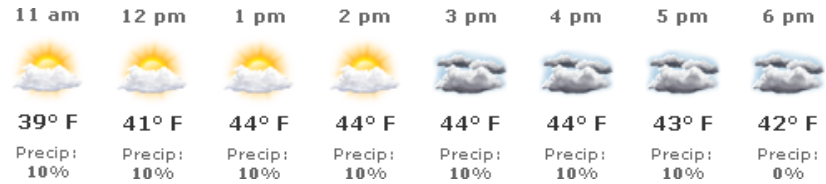


*An axis-aligned rectangle shatters 4 points only !*

# **STATISTICAL MACHINE LEARNING**

# Supervised Learning Tasks - Regression

## Weather Prediction



## Estimating Contamination



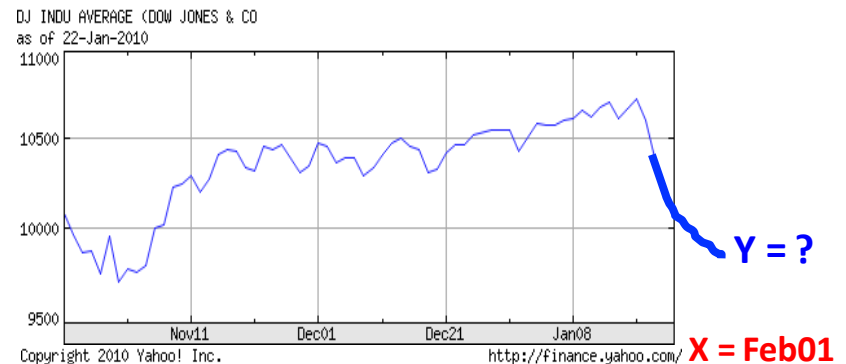
X = new location  
Y = sensor reading

# Supervised Learning

**Goal:** Construct a **predictor**  $f : X \rightarrow Y$  to minimize loss function (performance measure)



Sports  
Science  
News



**Classification:**

$$P(f(X) \neq Y)$$

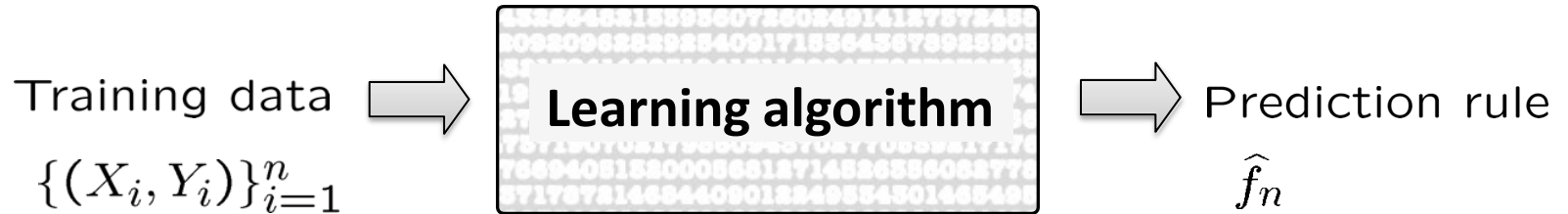
**Probability of Error**

**Regression:**

$$\mathbb{E}[(f(X) - Y)^2]$$

**Mean Squared Error**

# Regression algorithms



Linear Regression

# Replace Expectation with Empirical Mean

Optimal predictor:  $f^* = \arg \min_f \mathbb{E}[(f(X) - Y)^2]$

Empirical Minimizer:  $\hat{f}_n = \arg \min_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \right)$

**Empirical mean**

Law of Large Numbers:

$$\frac{1}{n} \sum_{i=1}^n [\text{loss}(Y_i, f(X_i))] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$



# Restrict class of predictors

Optimal predictor:

$$f^* = \arg \min_f \mathbb{E}[(f(X) - Y)^2]$$

Empirical Minimizer:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

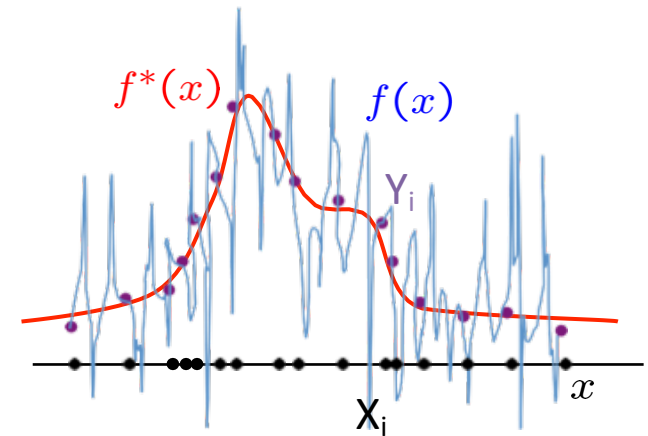
Class of predictors

Why?

Overfitting!

Empirical loss minimized by any function of the form

$$f(x) = \begin{cases} Y_i, & x = X_i \text{ for } i = 1, \dots, n \\ \text{any value,} & \text{otherwise} \end{cases}$$



# Restrict class of predictors

Optimal predictor:  $f^* = \arg \min_f \mathbb{E}[(f(X) - Y)^2]$

Empirical Minimizer:  $\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$

**Class of predictors**

- Class of Linear functions
- Class of Polynomial functions
- Class of nonlinear functions

# Linear Regression

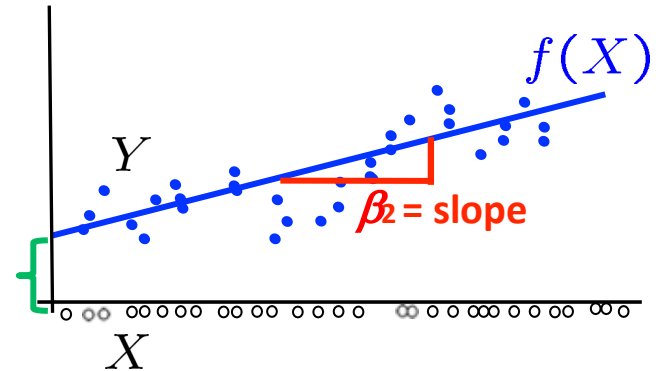
$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad \text{Least Squares Estimator}$$

$\mathcal{F}_L$  -Class of Linear functions

Uni-variate case:

$$f(X) = \beta_1 + \beta_2 X$$

$\beta_1$  -intercept



Multi-variate case:

$$f(X) = f(X^{(1)}, \dots, X^{(p)}) = \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$$

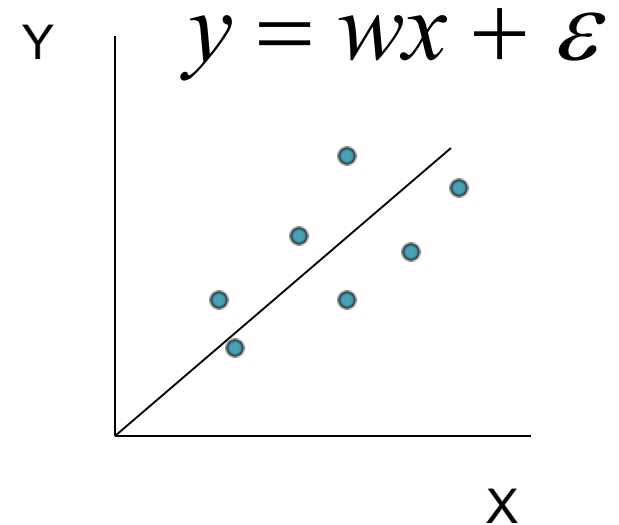
$$= X\beta \quad \text{where} \quad X = [X^{(1)} \dots X^{(p)}], \quad \beta = [\beta_1 \dots \beta_p]^T$$

# Linear regression in 1D

- Our goal is to estimate  $w$  from a training data of  $\langle x_i, y_i \rangle$  pairs
- Optimization goal: minimize squared error (least squares):

$$\arg \min_w \sum_i (y_i - wx_i)^2$$

- Why least squares?
  - minimizes squared distance between measurements and predicted line
  - the math is pretty



# Solving Linear Regression in 1D

- To optimize – closed form:
- We just take the derivative w.r.t. to  $w$  and set to 0:

$$\frac{\partial}{\partial w} \sum_i (y_i - wx_i)^2 = 2 \sum_i -x_i (y_i - wx_i) \Rightarrow$$

$$2 \sum_i x_i (y_i - wx_i) = 0 \Rightarrow 2 \sum_i x_i y_i - 2 \sum_i wx_i x_i = 0$$

$$\sum_i x_i y_i = \sum_i wx_i^2 \Rightarrow$$

$$w = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

# Least Squares Estimator

$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2$$

$$f(X_i) = X_i \beta$$



$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (X_i \beta - Y_i)^2$$

$$\hat{f}_n^L(X) = X \hat{\beta}$$

$$= \arg \min_{\beta} \frac{1}{n} (\mathbf{A} \beta - \mathbf{Y})^T (\mathbf{A} \beta - \mathbf{Y})$$

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

# Least Squares Estimator

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

$$J(\beta) = (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

$$\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\hat{\beta}} = 0$$

# Normal Equations

$$\underbrace{(\mathbf{A}^T \mathbf{A})}_{p \times p} \underbrace{\hat{\beta}}_{p \times 1} = \underbrace{\mathbf{A}^T \mathbf{Y}}_{p \times 1}$$

If  $(\mathbf{A}^T \mathbf{A})$  is invertible,

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \qquad \hat{f}_n^L(X) = X \hat{\beta}$$

When is  $(\mathbf{A}^T \mathbf{A})$  invertible ?

Recall: **Full rank matrices are invertible.**

What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?



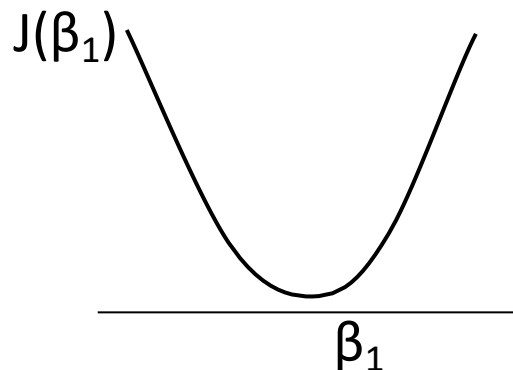
# Gradient Descent

Even when  $(\mathbf{A}^T \mathbf{A})$  is invertible, might be computationally expensive if  $\mathbf{A}$  is huge.

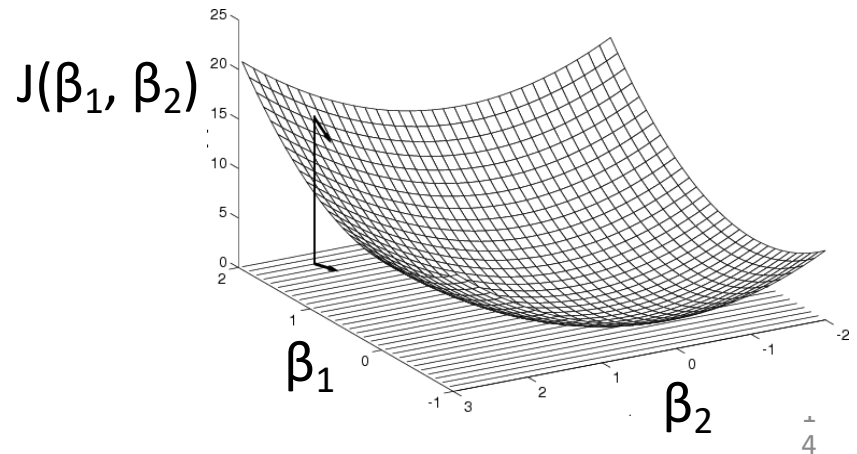
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Treat as optimization problem

Observation:  $J(\beta)$  is convex in  $\beta$ .



**How to find the minimizer?**



# Gradient Descent

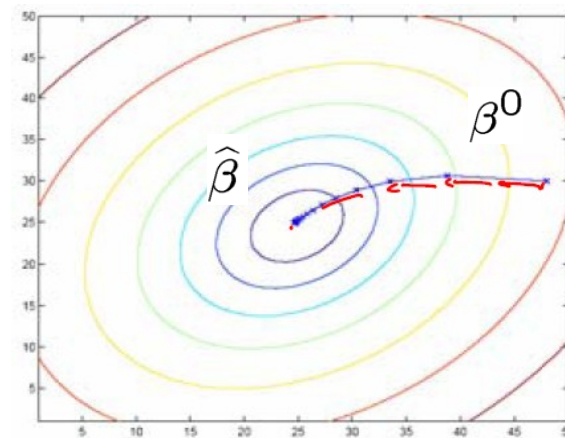
Even when  $(\mathbf{A}^T \mathbf{A})$  is invertible, might be computationally expensive if  $\mathbf{A}$  is huge.

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Since  $J(\beta)$  is convex, move along negative of gradient

Initialize:  $\beta^0$

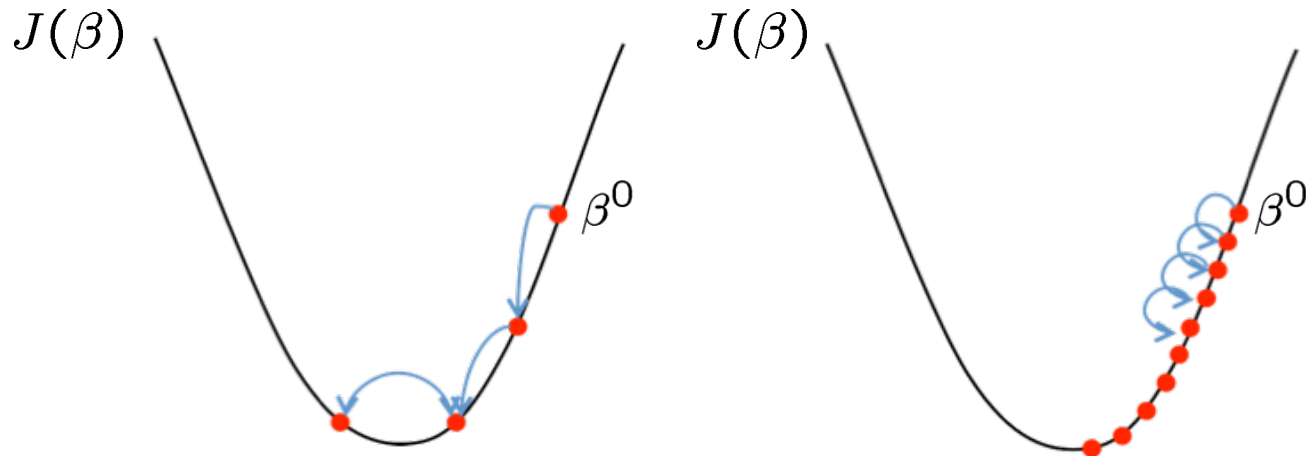
$$\begin{aligned} \text{Update: } \beta^{t+1} &= \beta^t - \overset{\text{step size}}{\alpha} \left. \frac{\partial J(\beta)}{\partial \beta} \right|_t \\ &= \beta^t - \alpha \mathbf{A}^T (\mathbf{A}\beta^t - \mathbf{Y}) \\ &\quad \underbrace{\hspace{10em}}_{0 \text{ if } \hat{\beta} = \beta^t} \end{aligned}$$



Stop: when some criterion met e.g. fixed # iterations, or  $\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\beta^t} < \epsilon$ .

# Effect of step--size

$\alpha$



Large  $\alpha \Rightarrow$  Fast convergence but larger residual error  
Also possible oscillations

Small  $\alpha \Rightarrow$  Slow convergence but small residual error

# Stochastic Gradient Descent

- Gradient descent (also known as Batch Gradient Descent) computes the gradient using the whole dataset
- Stochastic Gradient Descent computes the gradient using a single sample (or a mini-batch).

# Non-Linear basis function

- So far we only used the observed values  $x_1, x_2, \dots$
- However, linear regression can be applied in the same way to **functions** of these values
  - Eg: to add a term  $w x_1 x_2$  add a new variable  $z = x_1 x_2$  so each example becomes:  $x_1, x_2, \dots, z$
- As long as these functions can be directly computed from the observed values the parameters are still linear in the data and the problem remains a multi-variate linear regression problem

$$y = w_0 + w_1 x_1^2 + \dots + w_k x_k^2 + \varepsilon$$

# Non-linear basis functions

- What type of functions can we use?
- A few common examples:

- Polynomial:  $\phi_j(x) = x^j$  for  $j=0 \dots n$

- Gaussian:  $\phi_j(x) = \frac{(x - \mu_j)}{2\sigma_j^2}$

- Sigmoid:  $\phi_j(x) = \frac{1}{1 + \exp(-s_j x)}$

- Logs:  $\phi_j(x) = \log(x + 1)$

Any function of the input values can be used. The solution for the parameters of the regression remains the same.

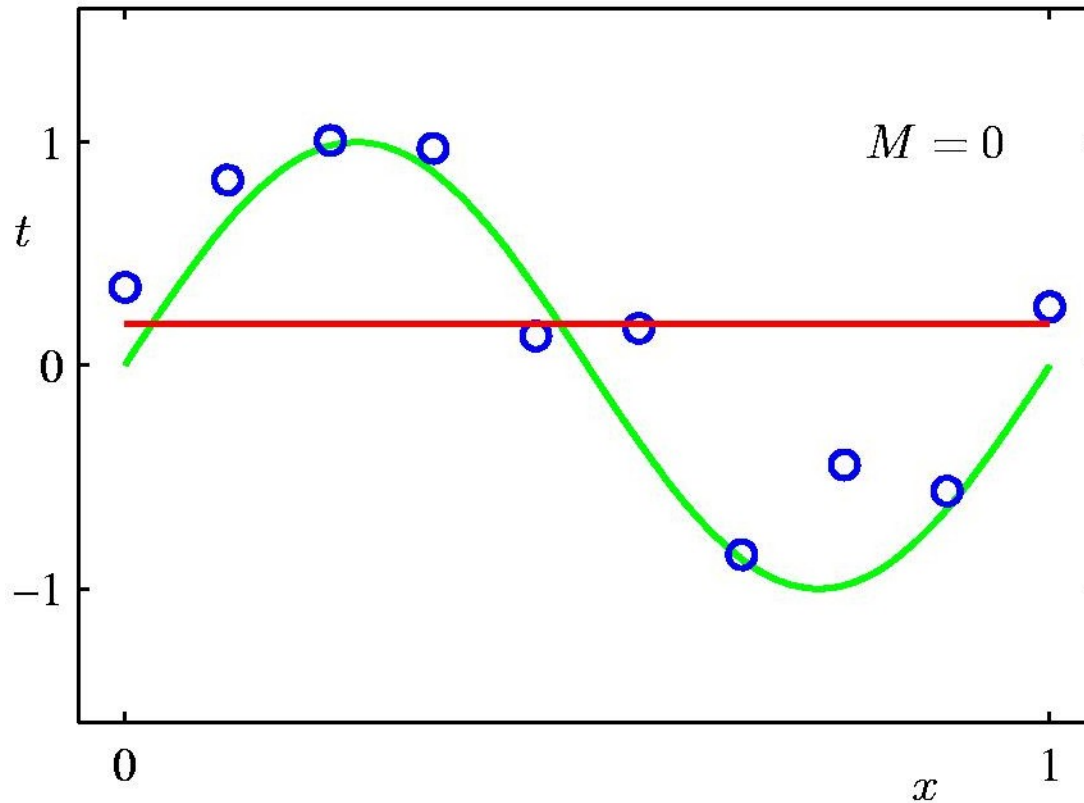
# General linear regression problem

- Using our new notations for the basis function linear regression can be written as

$$y = \sum_{j=0}^n w_j \phi_j(x)$$

- Where  $\phi_j(\mathbf{x})$  can be either  $x_j$  for multivariate regression or one of the non-linear basis functions we defined
- ... and  $\phi_0(\mathbf{x})=1$  for the intercept term

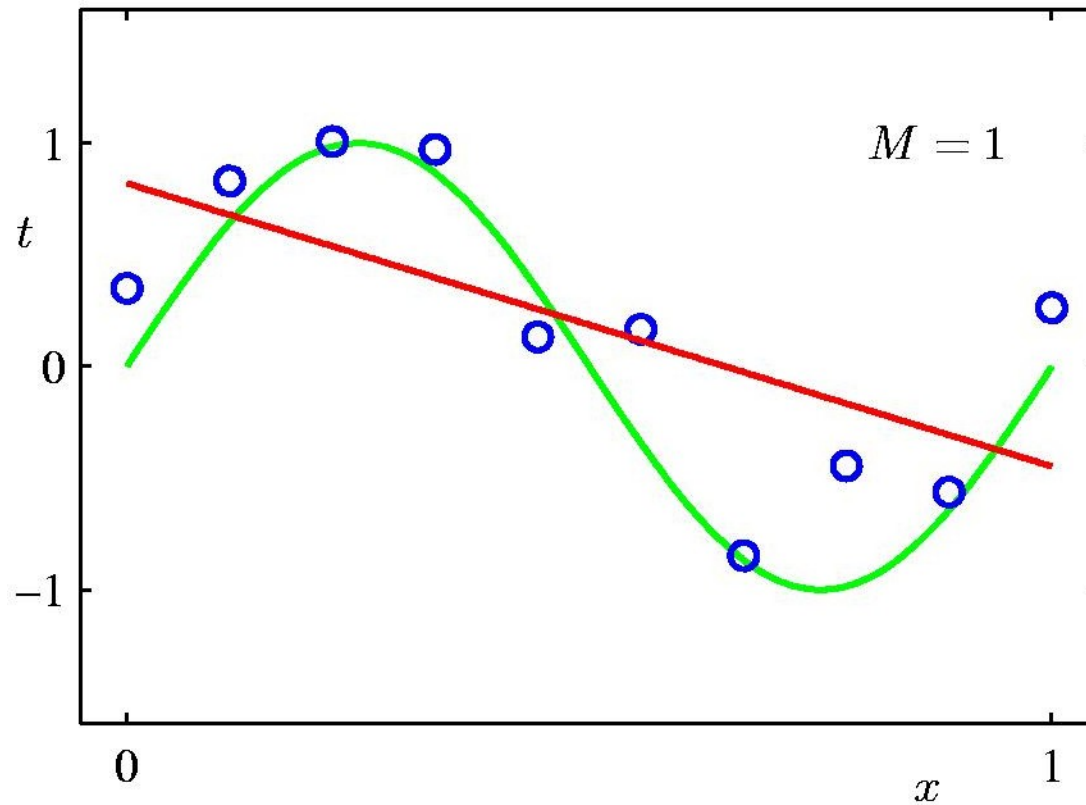
# 0<sup>th</sup> Order Polynomial



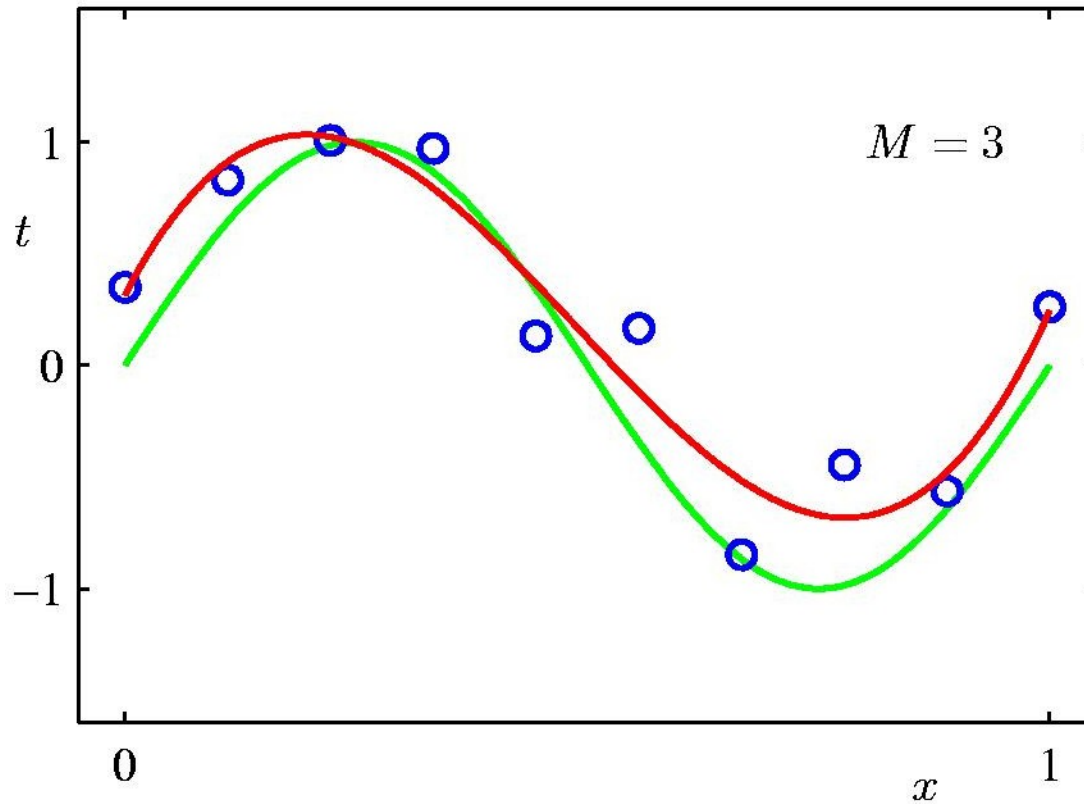
$n=10$



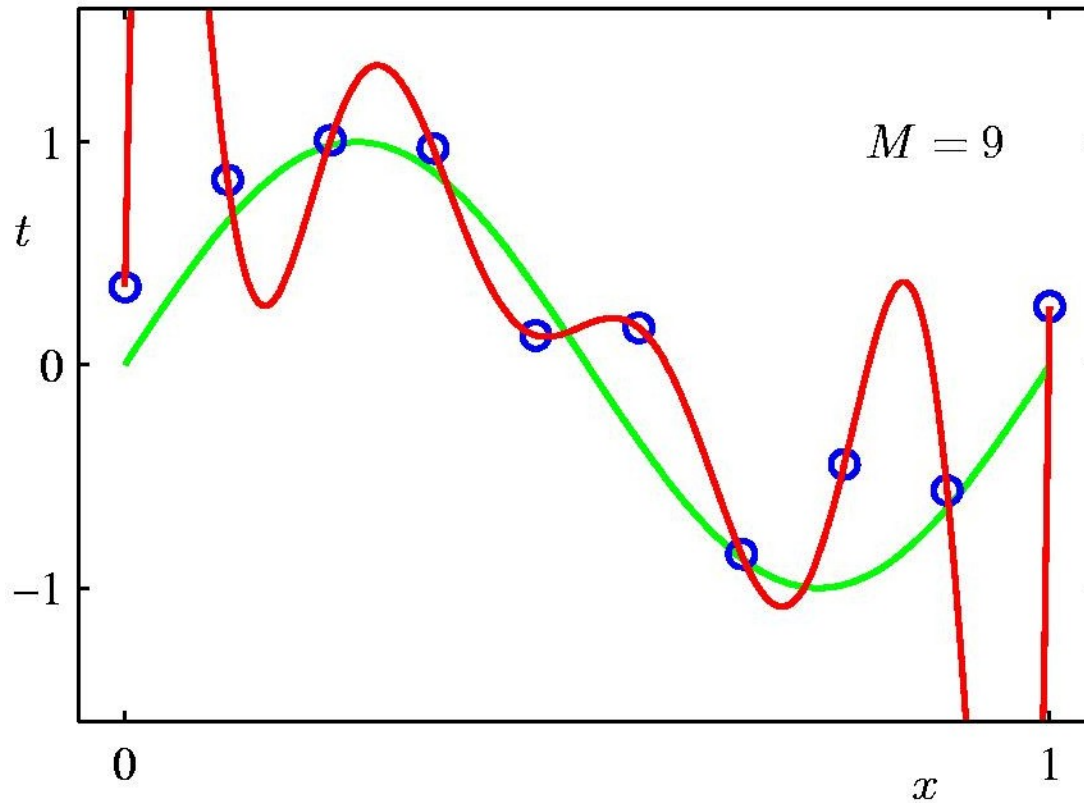
# 1<sup>st</sup> Order Polynomial



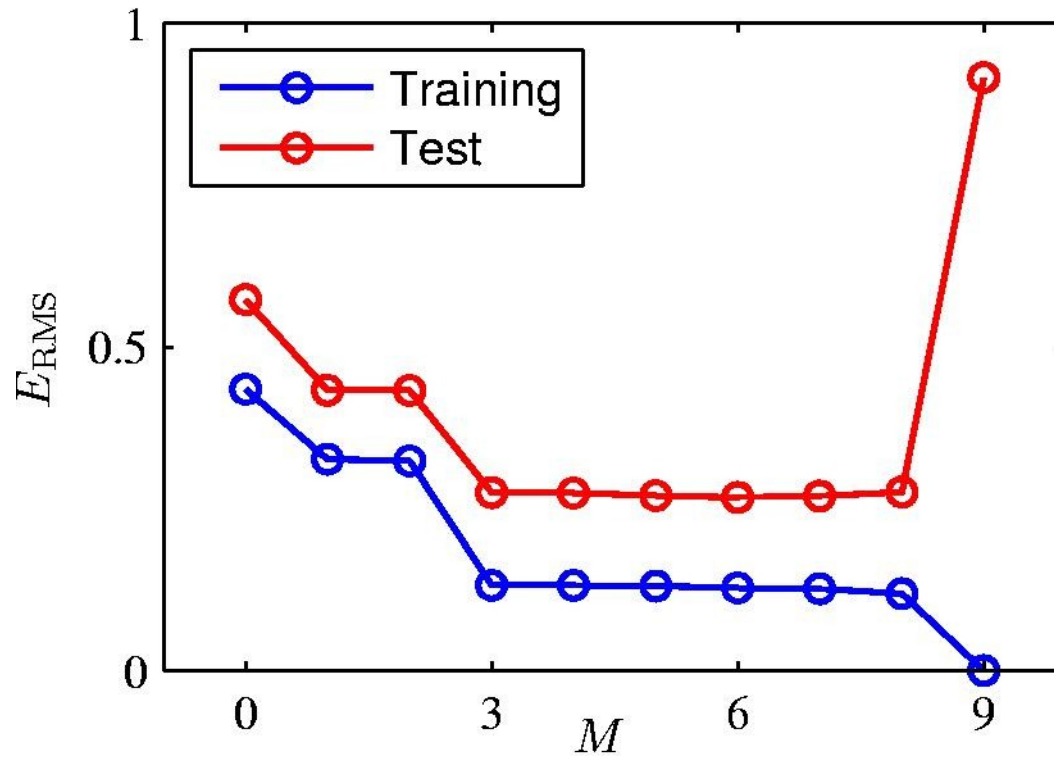
# 3<sup>rd</sup> Order Polynomial



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



# Over-fitting



Root-Mean-Square (RMS) Error

# 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

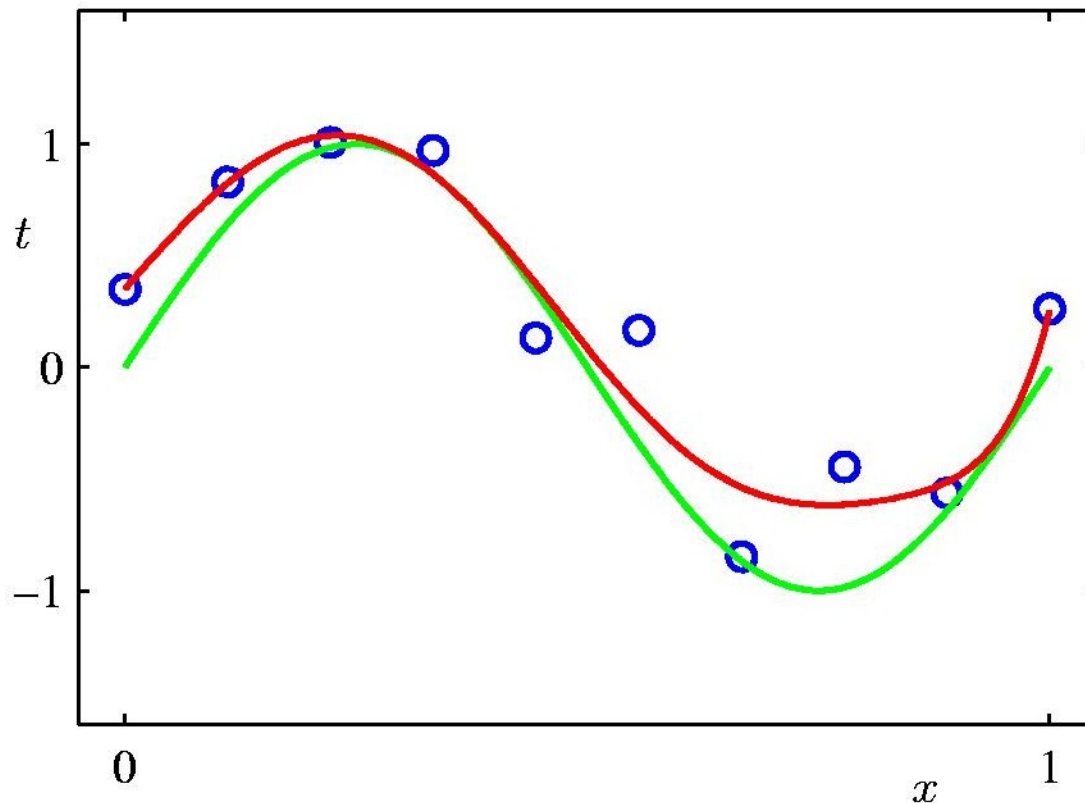
# Regularization

Penalize large coefficient values

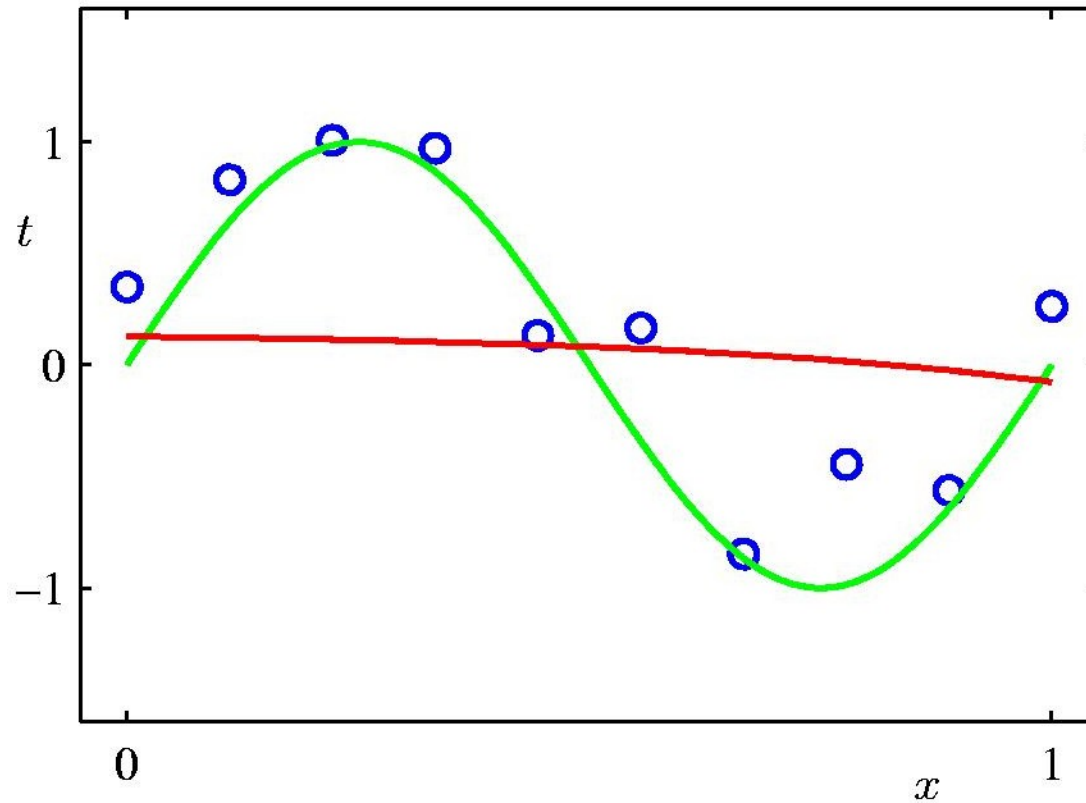
$$J_{\mathbf{x},\mathbf{y}}(\mathbf{w}) = \frac{1}{2} \sum_i \left( y^i - \sum_j w_j \phi_j(\mathbf{x}^i) \right)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

# Regularization:

$$\ln \lambda = -18$$



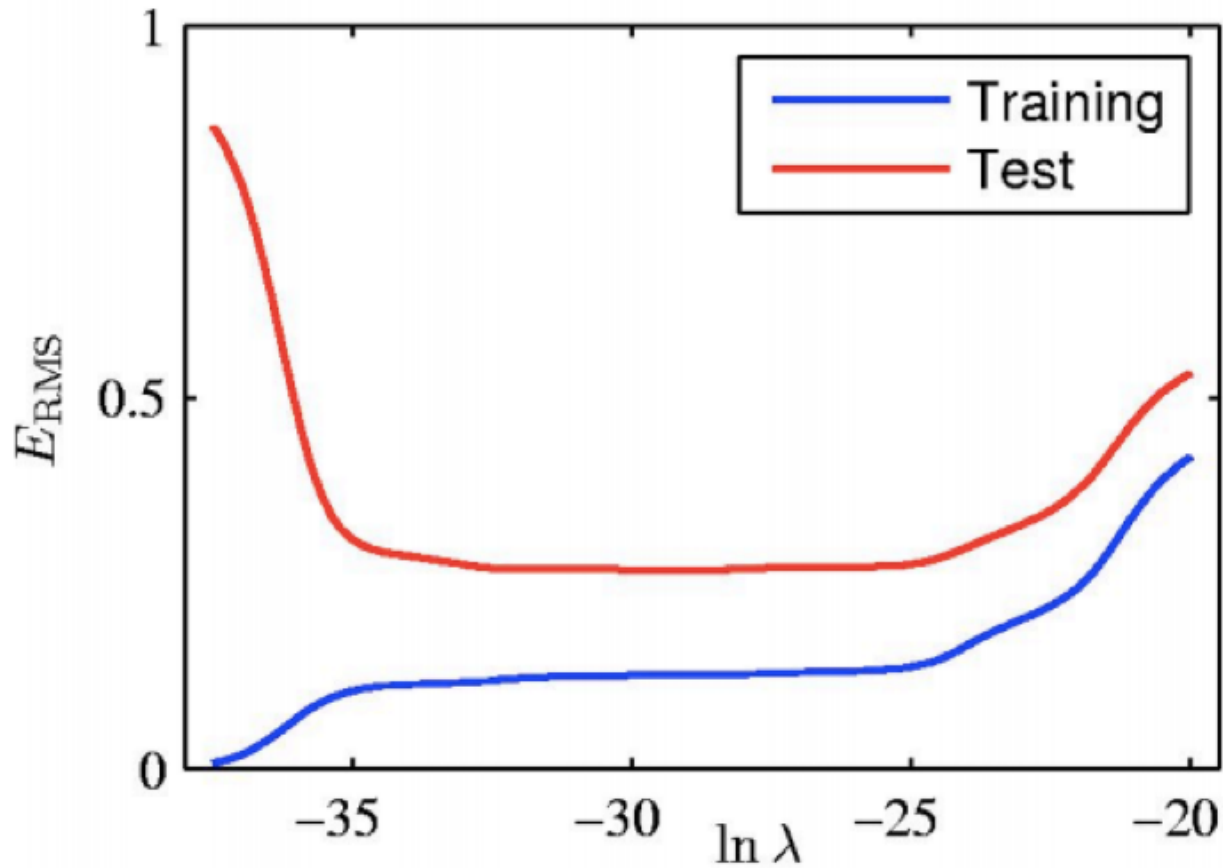
# Over Regularization





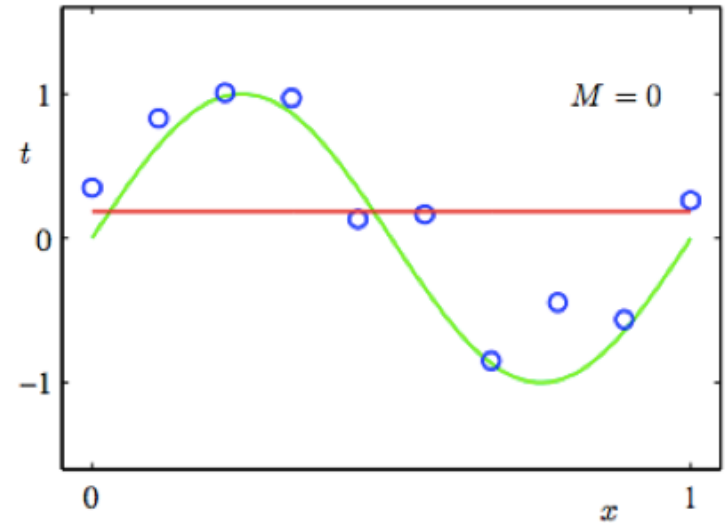
# Regularization

9<sup>th</sup> Order Polynomial

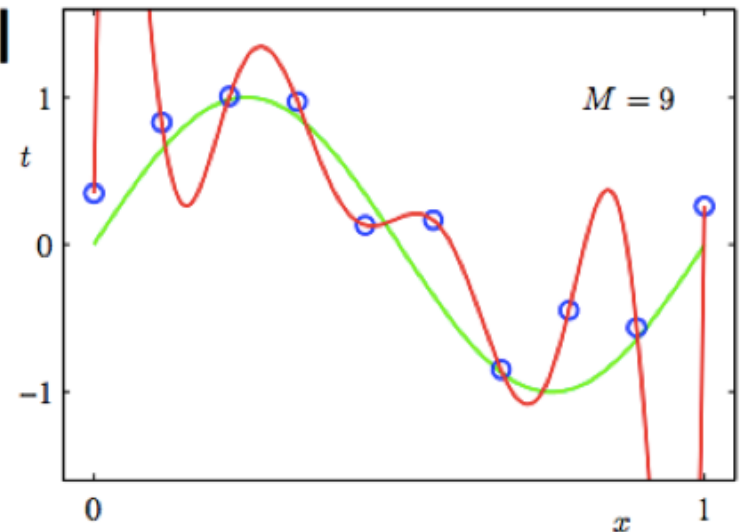


# Bias-Variance Tradeoff

- **Model too simple:** does not fit the data well
  - A *biased* solution



- **Model too complex:** small changes to the data, solution changes a lot
  - A *high-variance* solution



# Effect of Model Complexity

- If we allow very complicated predictors, we could overfit the training data.

