



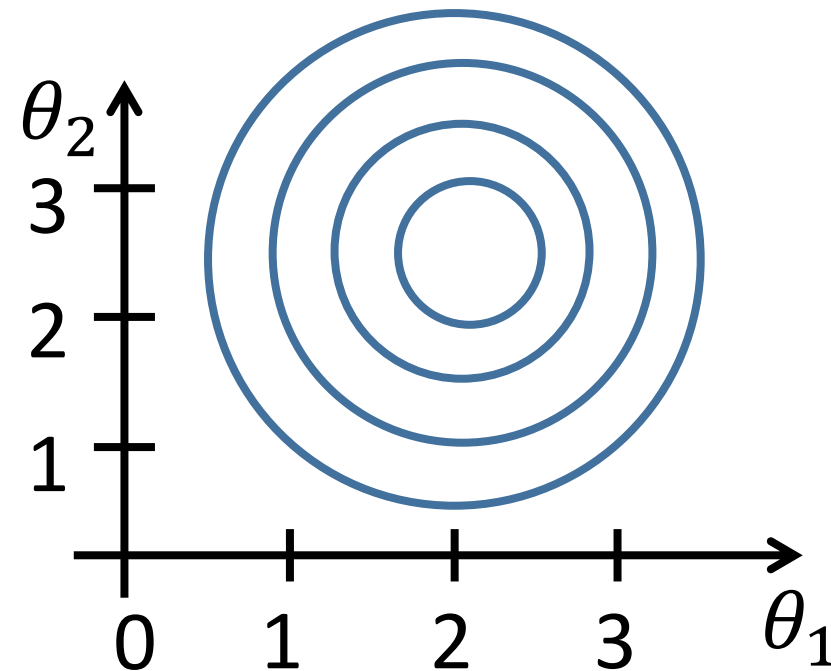
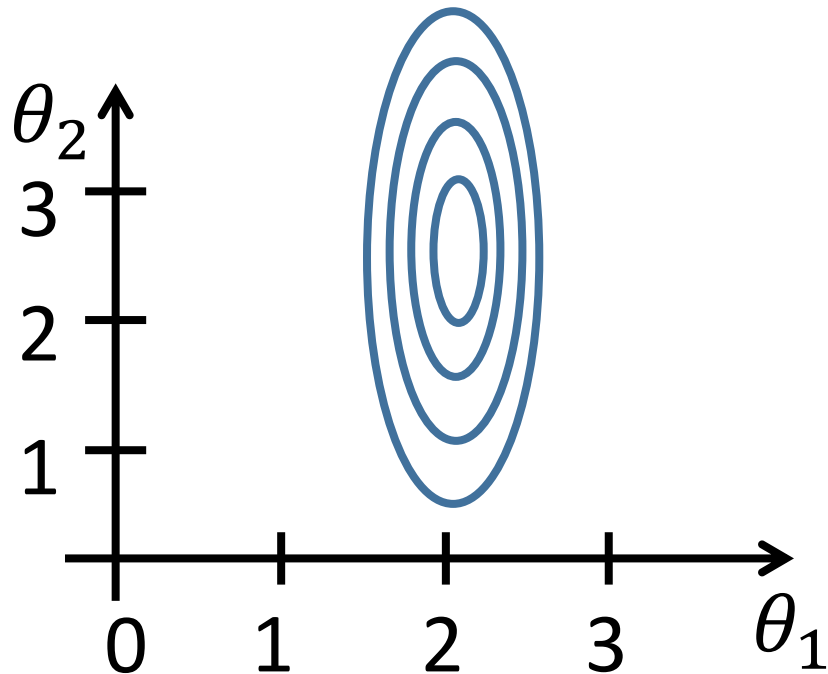
Pushing Regression Further: Polynomial Regression

Alexander G. Ororbia II
Introduction to Machine Learning
CSCI-335
3/16/2026

Gradient Descent in Practice: Feature Scaling

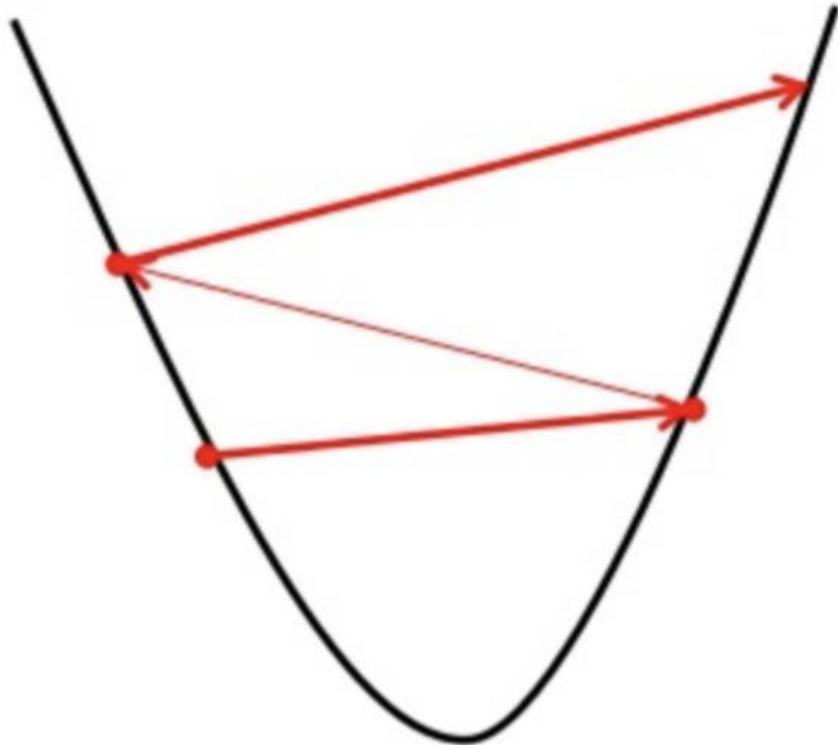
- *Idea:* Make sure features are on a similar scale (e.g., $-1 \leq x_i \leq 1$)
- *Example:* $x_1 = \text{size (0-2000 feat}^2)$
 $x_2 = \text{number of bedrooms (1-5)}$

$$a_i = \frac{v_i - \min v_i}{\max v_i - \min v_i}$$

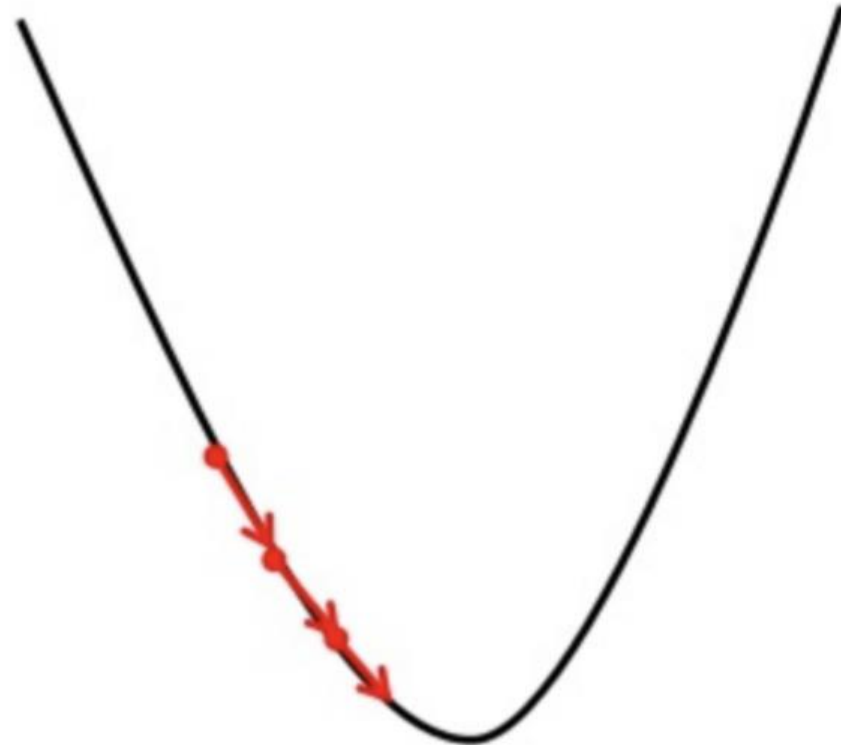


Gradient Descent Learning Rate / Step-Size

Big learning rate



Small learning rate

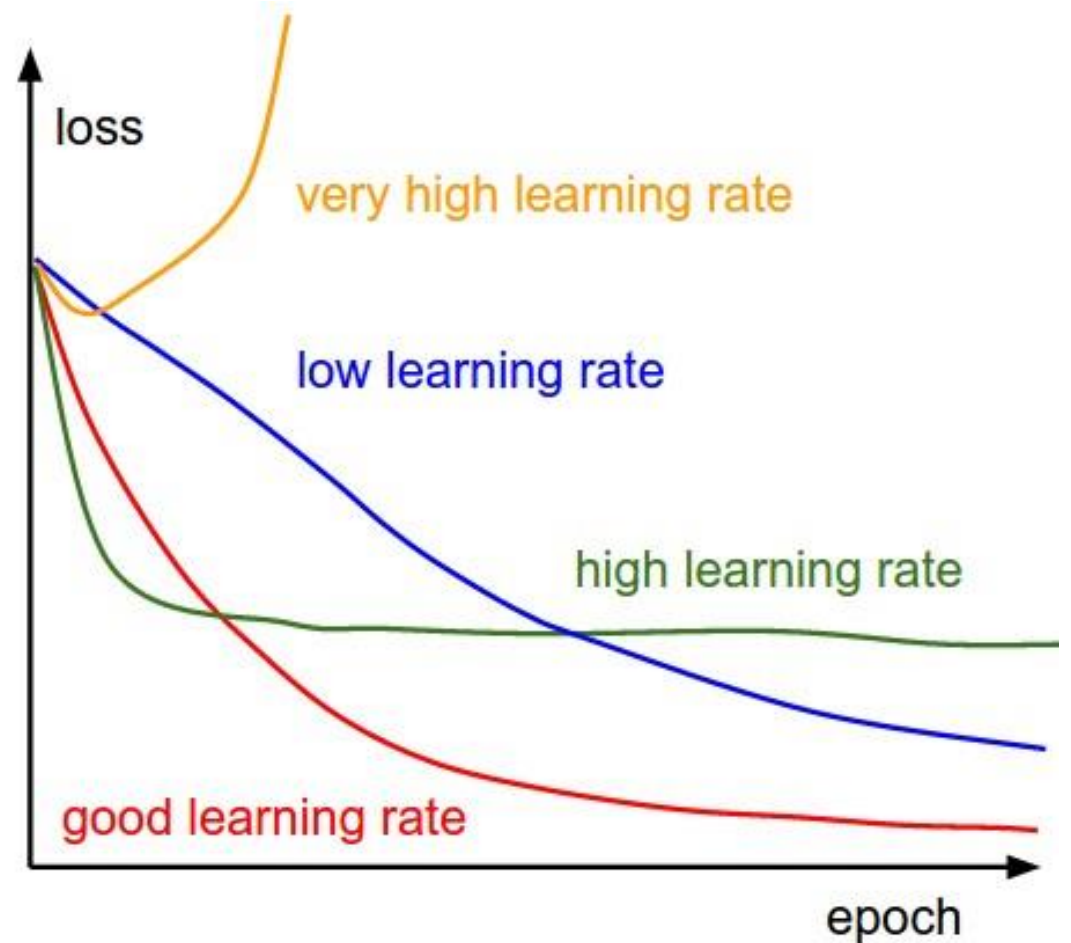


Gradient Descent in Practice: Step Size

- Automatic convergence test
- α too small: slow convergence
- α too large: may not converge

- To choose α , try

0.001, ... 0.01, ..., 0.1, ... , 1



Learning rate is also the “step size” in statistics/mathematics

A Simple Approach to Curve Fitting

Representation

- Fit the data using a *polynomial function*

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

– where M is the order of the polynomial

- Is higher value of M better? We'll see shortly!
- Coefficients w_0, \dots, w_M are collectively denoted by vector \mathbf{w}
- It is a nonlinear function of x , but a linear function of the unknown parameters

(Note: still a linear model; linear in parameter space)

Polynomial Curve Fitting

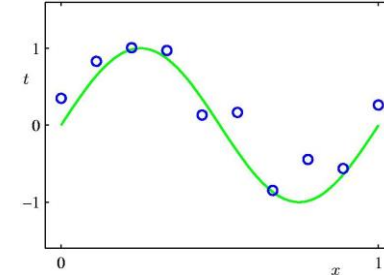
– With a single input variable x

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

M = order of polynomial,

x^j denotes x raised to power j ,

Coefficients w_0, \dots, w_M are collectively denoted by vector \mathbf{w}



Training data set
 $N=10$, Input x , target t

– **Task:** Learn \mathbf{w} from training data $D = \{(x_i, t_i)\}, i = 1, \dots, N$

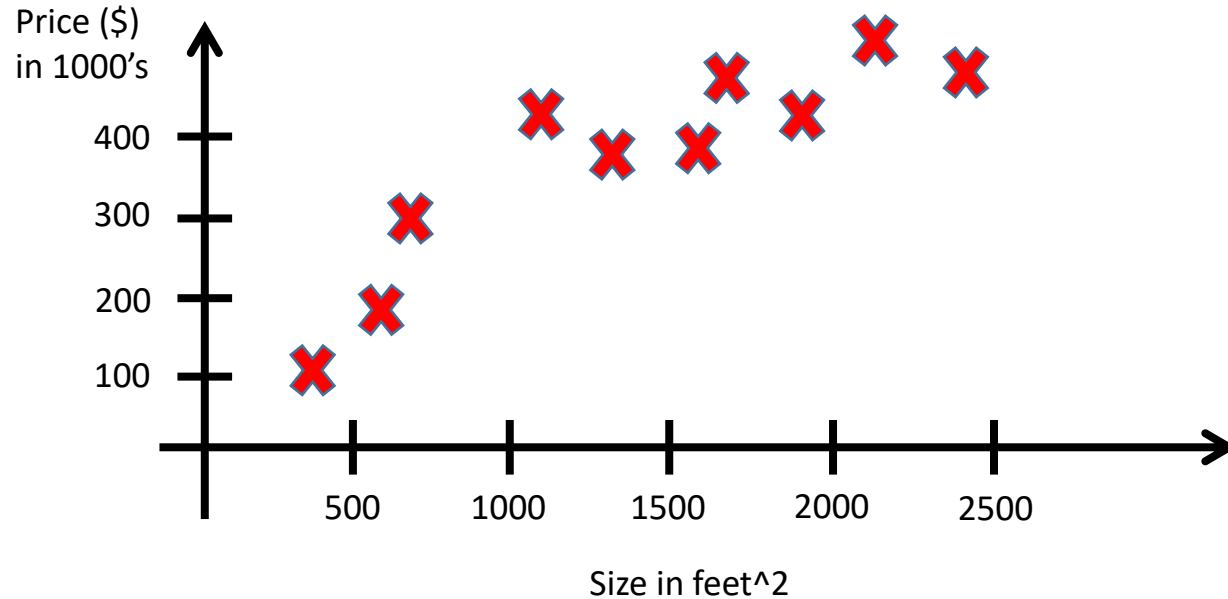
- Can be done by minimizing an error function that minimizes misfit between $y(x, \mathbf{w})$ for any given \mathbf{w} and training data
- One simple choice of error function is sum of squares of error (SSE) between predictions $y(x_n, \mathbf{w})$ for each data point x_n and corresponding target values t_n so that we minimize:

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

- It is zero when function $y(x, \mathbf{w})$ passes exactly through each training data point

A Polynomial Hypothesis

Apply same principles of linear regression derivation for each parameter (just change index $j > 1$)

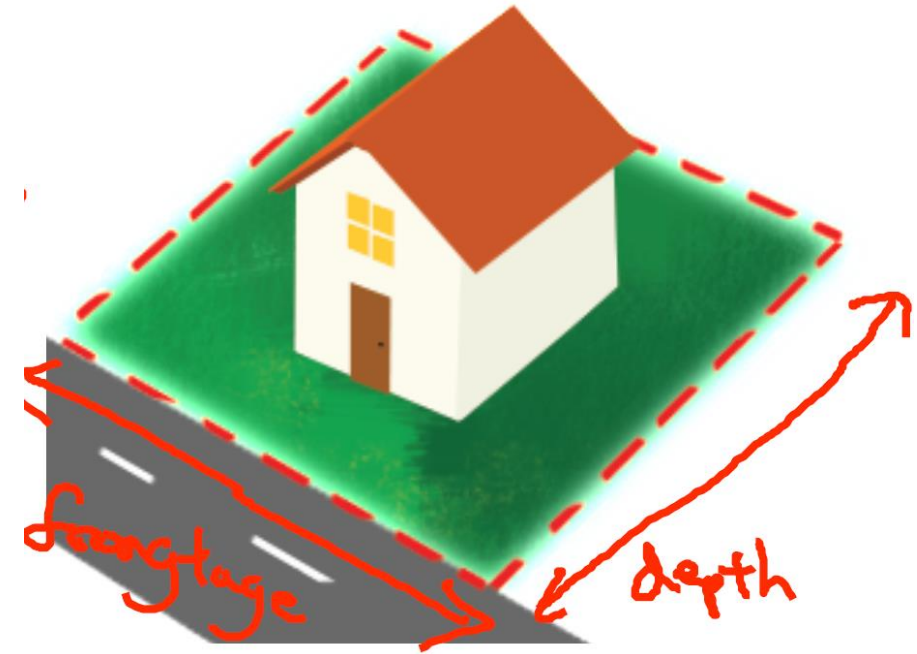


$$\begin{aligned}x_1 &= (\text{size}) \\x_2 &= (\text{size})^2 \\x_3 &= (\text{size})^3\end{aligned}$$

- $$\begin{aligned}h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\&= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3\end{aligned}$$

House Price(s) Prediction

- $h_{\theta}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$
- Area: $xy = \text{frontage (x)} \times \text{depth (y)}$
- Housing model: $h_{\theta}(x) = \theta_0 + \theta_1 xy$
(a multiplicative feature interaction)

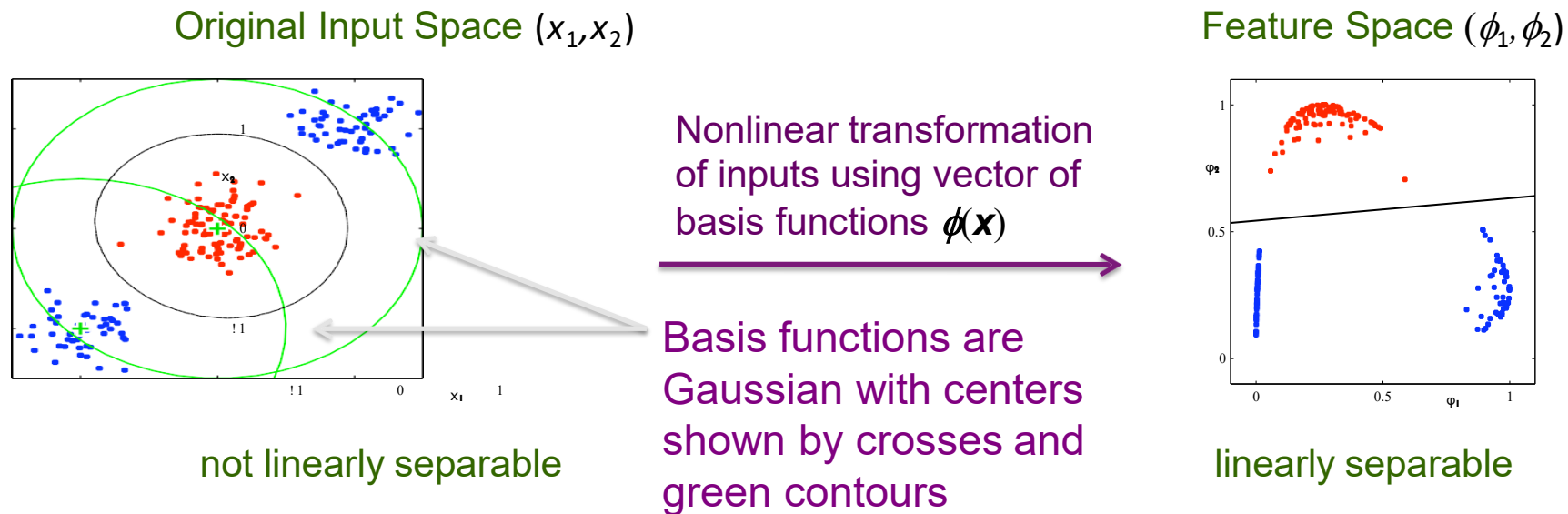


On Basis Functions

- In many applications, we apply some form of fixed-preprocessing, or feature extraction, to the original data variables
- If the original variables comprise the vector \mathbf{x} , then the features can be expressed in terms of basis functions $\{\phi_j(\mathbf{x})\}$
 - By using nonlinear basis functions we allow the function $y(\mathbf{x}, \mathbf{w})$ to be a nonlinear function of the input vector \mathbf{x}
 - They are linear functions of parameters (gives them simple analytical properties), yet are nonlinear wrt input variables

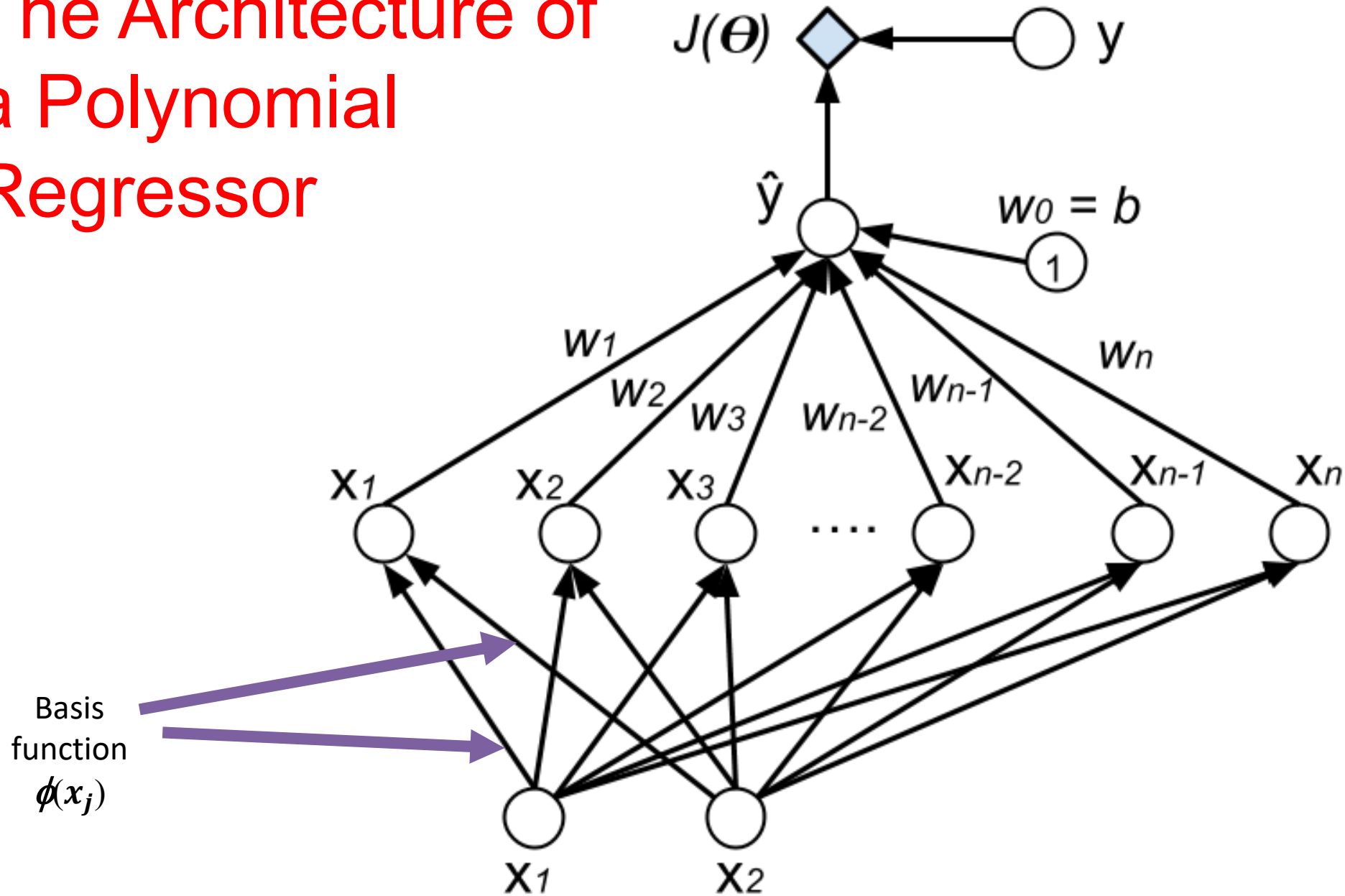
Fixed Basis Functions

Although we use linear (classification) models
Linear-separability in *feature* space *does not* imply
linear-separability in *input* space



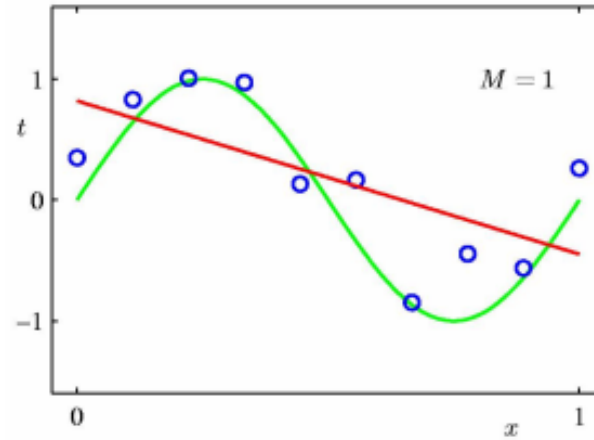
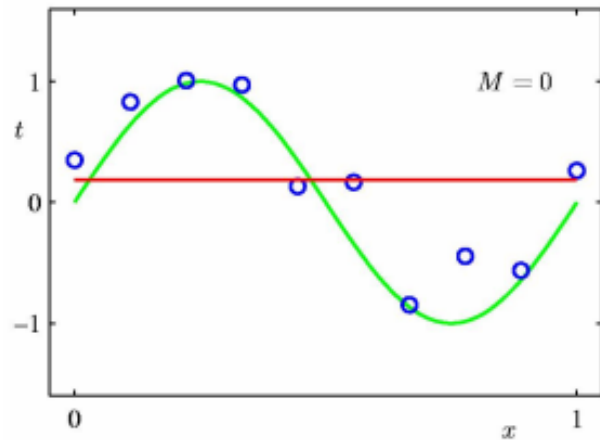
Basis functions with increased dimensionality often used

The Architecture of a Polynomial Regressor

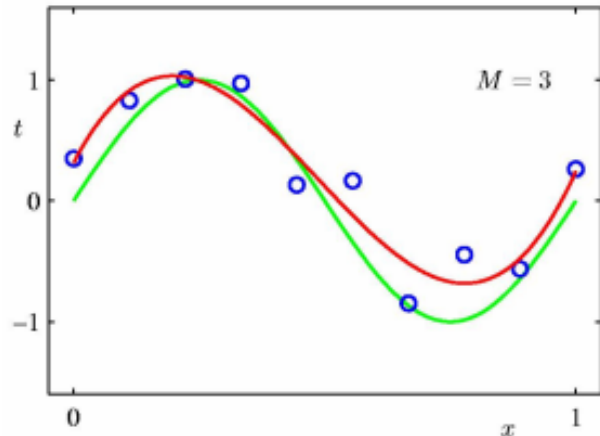


Choosing the Order of M

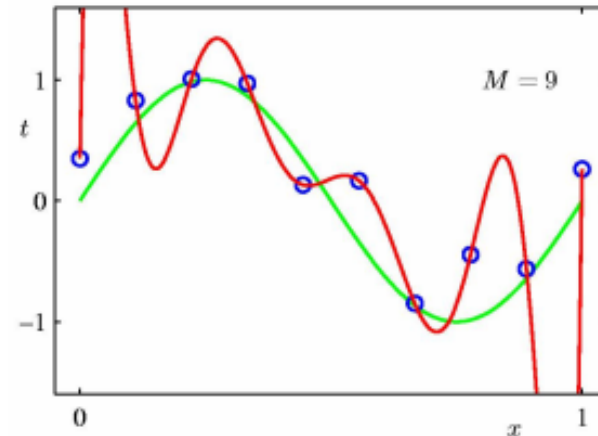
- Model Comparison or Model Selection
- Red lines are best fits with
 - $M = 0, 1, 3, 9$ and $N=10$



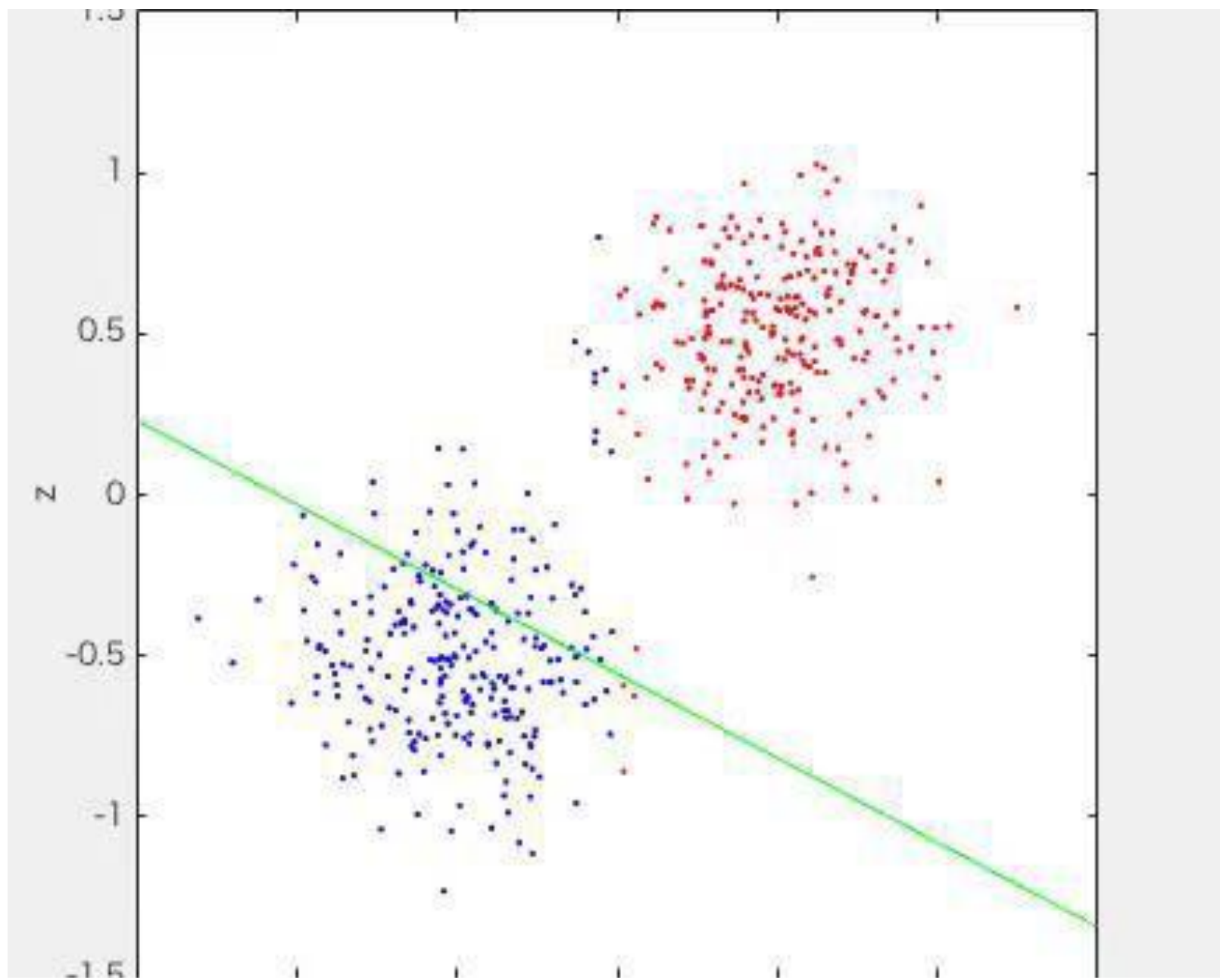
← Poor representations of $\sin(2\pi x)$

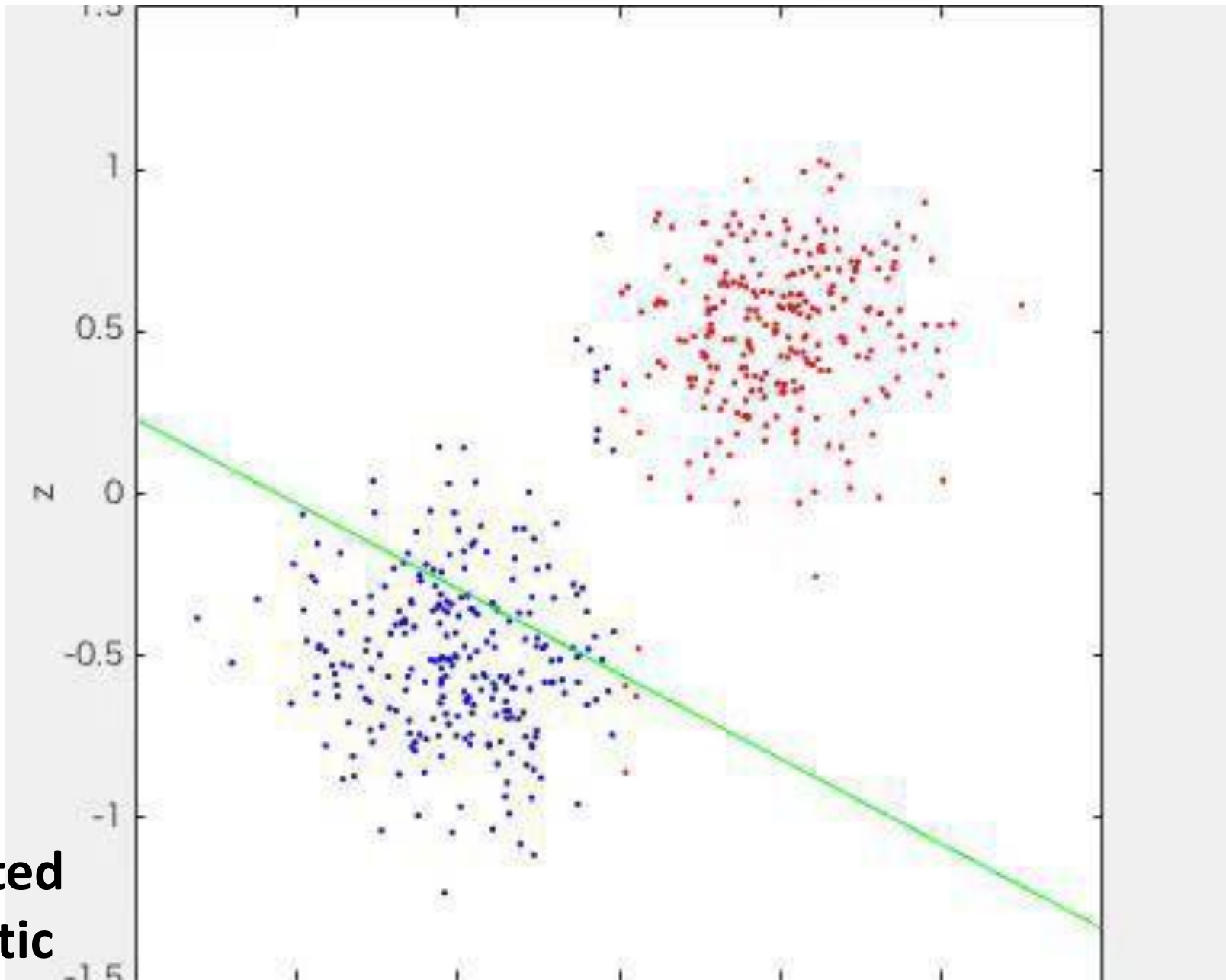


← Best Fit to $\sin(2\pi x)$



Over Fit
Poor representation of $\sin(2\pi x)$





**But if we wanted
some automatic
capacity control?**

Limiting Model Capacity

- Regularization has been used for decades prior to advent of deep learning
- Linear- and logistic-regression allow simple, straightforward and effective regularization strategies
 - Adding a parameter norm penalty $\Omega(\theta)$ to the objective function J :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha\Omega(\theta)$$

- where $\alpha \in [0, \infty)$ is a hyperparameter that weight the relative contribution of the norm penalty term Ω
 - Setting α to 0 results in no regularization. Larger values correspond to more regularization

Gradient of Regularized Objective

- Objective function (with no bias parameter)

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^T w + J(w; X, y)$$

- Corresponding parameter gradient

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y)$$

- To perform single gradient step, perform update:

$$w \leftarrow w - \varepsilon (\alpha w + \nabla_w J(w; X, y))$$

- Written another way, the update is

$$w \leftarrow (1 - \varepsilon \alpha) w - \varepsilon \nabla_w J(w; X, y)$$

- We have modified learning rule to shrink w by constant factor $1 - \varepsilon \alpha$ at each step

Multivariate Regressor Architecture

$$f_{\Theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \leftarrow \text{Your Hypothesis!}$$

Cost:

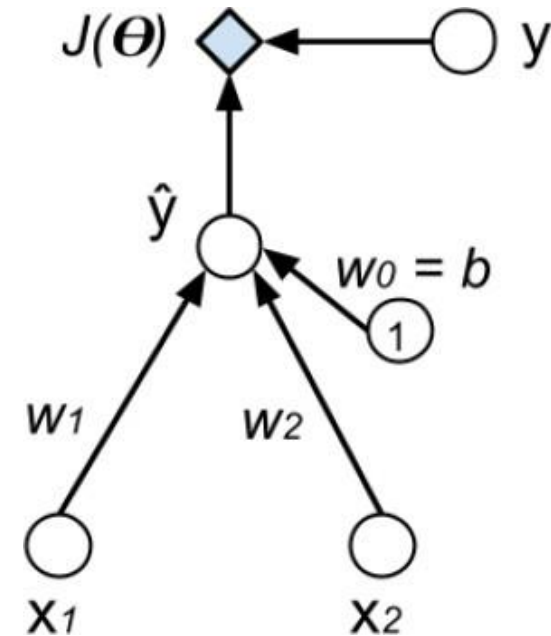
$$\mathcal{J}(\Theta) = \frac{1}{2m} \sum_{i=1}^m (f_{\Theta}(x^i) - y^i)^2 + \frac{\beta}{2m} \sum_{j=1}^n \theta_j^2$$

Derivative/Update:

$$\frac{\partial \mathcal{J}(\Theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (f_{\Theta}(x^i) - y^i) x_j^i + \frac{\beta}{m} \theta_j$$

Optimizer:

$$\theta_j = \theta_j - \alpha \frac{\partial \mathcal{J}(\Theta)}{\partial \theta_j} = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\Theta}(x^i) - y^i) x_j^i, j = 0, 1, 2, \dots, n.$$



Questions?

Deep robots!

Deep questions?!

