



Linear Regression

Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Andrew Ng (Stanford), Eric Eaton (UPenn), David Kauchak (Pomona), and the many others who made their course materials freely available online.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

Linear Regression Setup

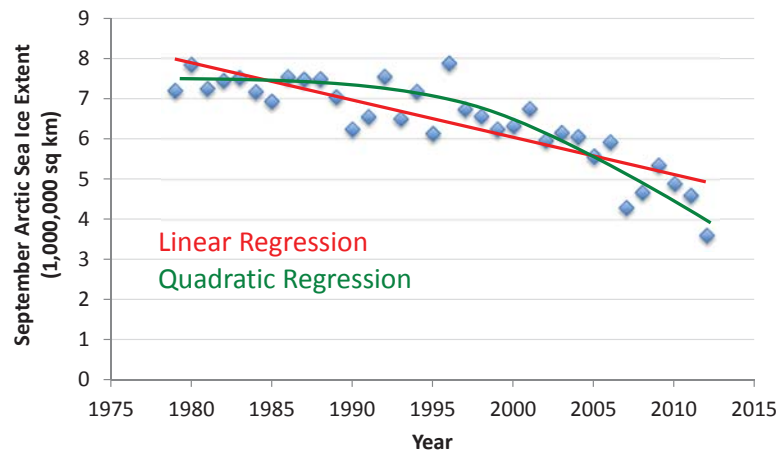
Learning Goals

- Describe how regression differs from classification
- Describe the goal of linear regression

Regression

Given:

- Data $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ where $\mathbf{x}^{(i)} \in \mathbb{R}^d$
- Corresponding labels $\mathbf{y} = \{y^{(1)}, \dots, y^{(n)}\}$ where $y^{(i)} \in \mathbb{R}$



Based on slide by Eric Eaton

[Data from G. Witt. Journal of Statistics Education, Volume 21, Number 1 (2013)]

Linear Regression

Hypothesis

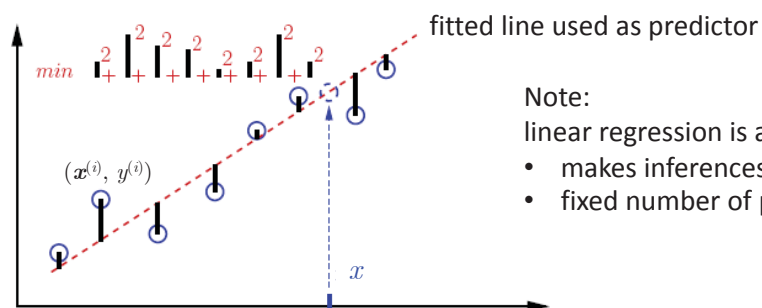
$$y = h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$$

← assume $x_0 = 1$

$$= \sum_{j=0}^d \theta_j x_j = \boldsymbol{\theta}^T \mathbf{x}$$

θ_j 's are parameters (weights)

Fit model (find $\boldsymbol{\theta}$) by **minimizing sum of squared errors**



Note:

linear regression is a **parametric** method

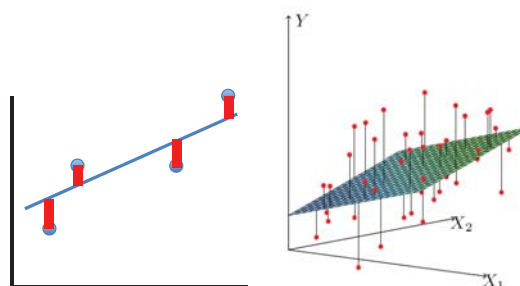
- makes inferences about model
- fixed number of parameters ($d+1$)

Based on slide by Eric Eaton [Figures courtesy of Greg Shakhnarovich]

Cost Function

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

Fit by solving $\min_{\theta} J(\theta)$



More generally...

$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n \text{Loss} \left(y^{(i)}, h_{\theta}(\mathbf{x}^{(i)}) \right) \quad \text{empirical risk / training error}$$

$$z = y^{(i)} - h_{\theta}(\mathbf{x}^{(i)})$$

$$\text{Loss}(z) = \frac{z^2}{2}$$

squared error

different people use $1/(2n)$ vs $1/2$

Q1: Why use squared error?

Q2: How do we optimize θ ?

Based on slide by Eric Eaton

(This slide intentionally left blank.)

Probabilistic Interpretation

Learning Goals

- Describe why we minimize squared error

Probabilistic Interpretation

Assume $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$

What is $\epsilon^{(i)}$?

- error term
- captures un-modeled effects (e.g. missing features)
- captures random noise

How can we model $\epsilon^{(i)}$?



Probabilistic Interpretation



(This slide intentionally left blank.)

Maximum Likelihood Estimate (MLE)



(This slide intentionally left blank.)

Solving Linear Regression

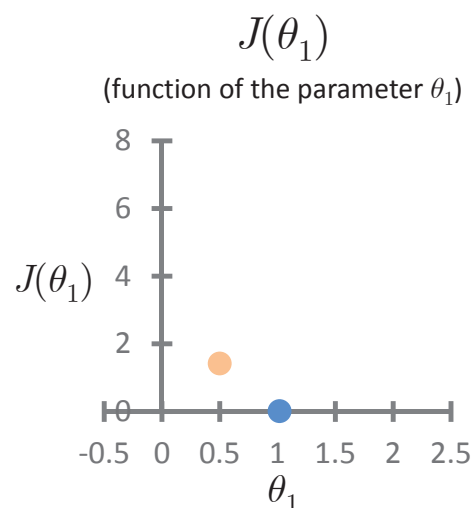
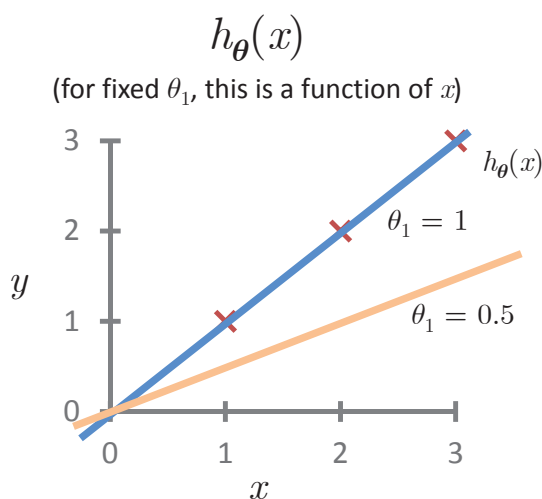
Learning Goals

- Describe shape of $J(\theta)$
- Describe two approaches for optimizing θ
 - Gradient Descent (stochastic and batch version)
 - Normal Equations
- Compare tradeoffs of GD vs Normal Equations

Intuition behind Cost Function

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

For insight on $J(\theta)$,
assume $x^{(i)} \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]^T$.
Also fix $\theta_0 = 0$.

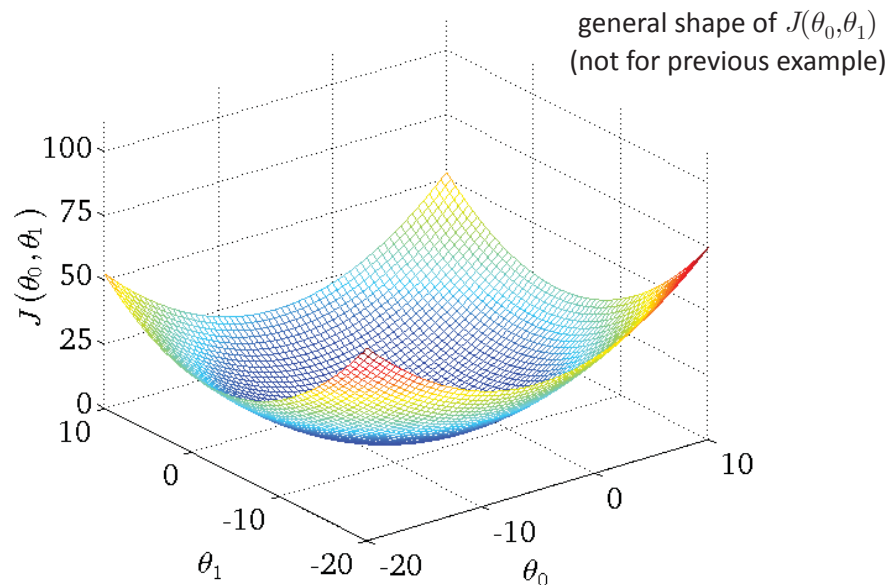


Based on slide by Eric Eaton
[Example by Andrew Ng]

$$J(0.5) = \frac{1}{2} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] = 1.75$$

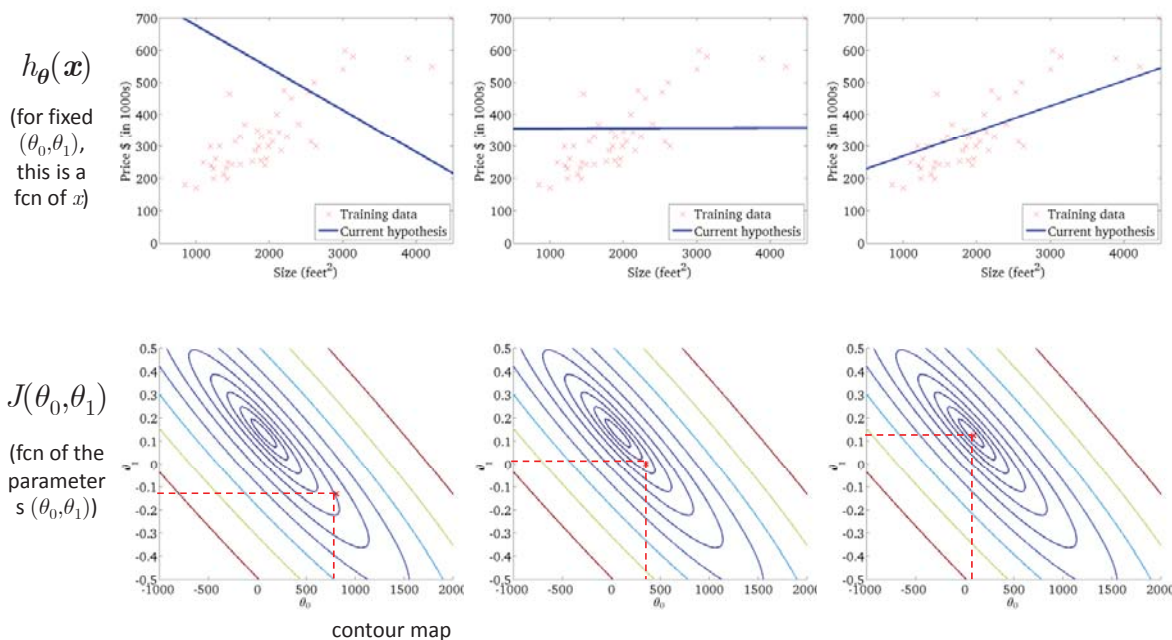


Intuition behind Cost Function



Based on slide by Andrew Ng

Intuition behind Cost Function



Based on slide by Andrew Ng

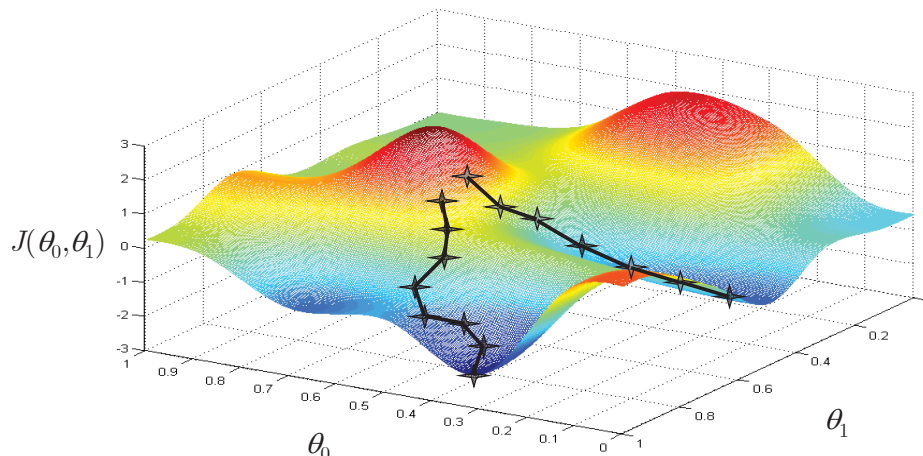
Solving Linear Regression

Learning Goals

- ✓ Describe shape of $J(\theta)$
- Describe two approaches for optimizing θ
 - Gradient Descent (stochastic and batch version)
 - Normal Equations
- Compare tradeoffs of GD vs Normal Equations

Basic Search Procedure

- Idea
- Choose initial value for θ
 - Until we reach minimum
 - Choose new value for θ to reduce $J(\theta)$



since the least squares objective function is convex,
we do not need to worry about local minima

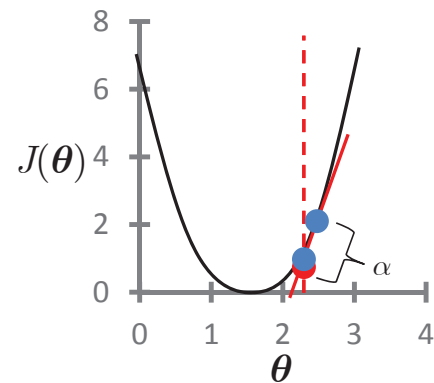
Approach 1: Gradient Descent

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

learning rate
(small, e.g., $\alpha = 0.05$)

For one training example (x, y)

update simultaneously for all $j = 0, \dots, d$



Based on slide by Eric Eaton



(This slide intentionally left blank.)

Batch versus Stochastic Mode

Batch Gradient Descent (BGD)

repeat until convergence {

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} \quad \text{for every } j \text{ simultaneously}$$

}

converged when
 $\|\theta_{\text{new}} - \theta_{\text{old}}\|_2 < \epsilon$

To achieve simultaneous update:

- At start of each GD iteration, compute $h_{\theta}(\mathbf{x}^{(i)})$
- Use this stored value in update step

Stochastic Gradient Descent (SGD)

repeat until convergence {

for $i = 1$ to n { ← often will randomly shuffle points

$$\theta_j \leftarrow \theta_j - \alpha \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} \quad \text{for every } j \text{ simultaneously}$$

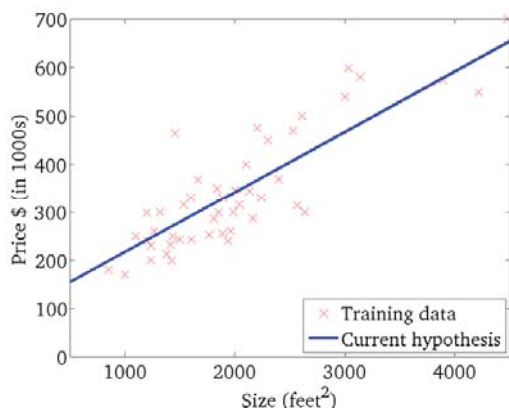
}

}

Gradient Descent Example

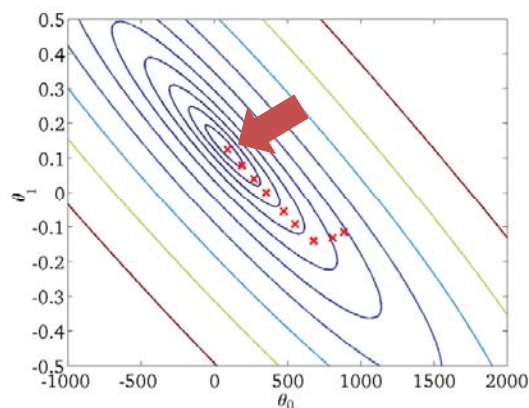
$$h_{\theta}(x)$$

(for fixed (θ_0, θ_1) , this is a function of x)



$$J(\theta_0, \theta_1)$$

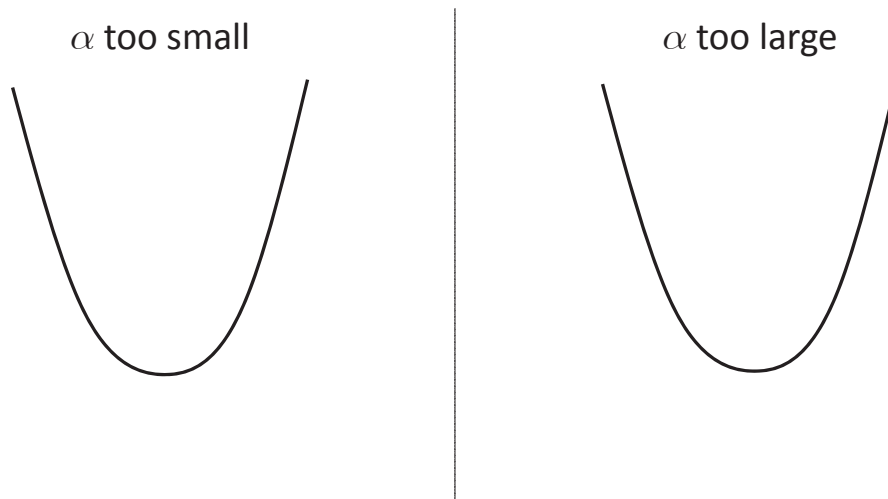
(function of the parameter (θ_0, θ_1))



Batch versus Stochastic GD

- Batch has to scan through entire training set before taking single step (costly if n is large)
- Stochastic makes progress “right away”
 - often θ gets “close” to minimum much faster
 - but θ may never converge if it oscillates around $\arg \min J(\theta)$
 - in practice, most values near minimum will be good enough approximation to true minimum
 - alternatively, we can guarantee convergence if we allow step size α to decrease, e.g. $\alpha_k = \frac{1}{1+k}$, where k = outer loop iteration

Choosing step size α



To check if GD is working, print out $J(\theta)$ each iteration

- value should decrease each iteration
- if it does not, adjust α



Solving Linear Regression

Learning Goals

- ✓ Describe shape of $J(\theta)$
- Describe two approaches for optimizing θ
 - ✓ Gradient Descent (stochastic and batch version)
 - Normal Equations
- Compare tradeoffs of GD vs Normal Equations

Approach 2: Normal Equations

(closed form solution)

Vectorization

- more compact equations
- faster code (using optimized matrix libraries)

Let us consider our model

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

$$\text{Let } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \text{ and } \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}. \text{ Then } h_{\theta}(\mathbf{x}) = \theta^T \mathbf{x}.$$

Matrix-Vector Form

For n instances:
$$h_{\theta}(\mathbf{x}^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)}$$

Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} - & (\mathbf{x}^{(1)})^T & - \\ & \vdots & \\ - & (\mathbf{x}^{(i)})^T & - \\ & \vdots & \\ - & (\mathbf{x}^{(n)})^T & - \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}$$

Then

$$h_{\theta}(\mathbf{X}) = \mathbf{X}\boldsymbol{\theta} \quad \text{predictions for all instances}$$

$\downarrow \quad \swarrow \quad \searrow$
 $\mathbb{R}^{n \times 1} \quad \mathbb{R}^{n \times (d+1)} \quad \mathbb{R}^{(d+1) \times 1}$

Matrix-Vector Form of Cost Function

Let

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$



Solve for θ analytically

$$\theta^* = \arg \min_{\theta} J(\theta) \text{ attained when } \nabla_{\theta} J(\theta) = 0$$

Expanding $J(\theta)$:



Solve for θ analytically

Take gradient, set equal to 0, then solve for θ :



Closed-Form Solution

Can obtain θ by plugging in \mathbf{X} and \mathbf{y} into

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$
$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- Cost of inverting $\mathbf{A} = \mathbf{X}^T \mathbf{X} \sim O(d^3)$
 - expensive for large d
- If \mathbf{A} is not invertible (i.e. \mathbf{A} is singular), may need to
 - use pseudo-inverse (`np.linalg.pinv`)
 - Cholesky decomposition better
 - remove redundant (not linearly independent) features
 - remove extra features to ensure $d \leq n$

Based on slide by Eric Eaton

Solving Linear Regression

Learning Goals

- ✓ Describe shape of $J(\theta)$
- ✓ Describe two approaches for optimizing θ
 - ✓ Gradient Descent (stochastic and batch version)
 - ✓ Normal Equations
- Compare tradeoffs of GD vs Normal Equations

Gradient Descent vs Normal Equation

Gradient Descent

- requires multiple iterations
- need to choose α
- works well when d is large
- can support online learning

Normal Equations

- non-iterative
- no need for α
- slow if d is large
 - matrix inversion is $O(d^3)$

Based on slide by Eric Eaton

What Should You Be Able To Do?

- Name cost function for linear regression and provide rationale
- Given model, find Maximum Likelihood Estimator (MLE)
- Optimize functions by taking derivative / gradient
- State tradeoffs of gradient descent vs normal equations
- Homework
 - Extend linear regression to weighted linear regression
 - Implement (regularized polynomial) regression using SGD and normal equations

(Extra Slides)

Linear Regression Extensions

Learning Goals

- Describe how to extend linear regression to more complex models using basis functions
- Describe why we might scale features

Basis Functions

So far,

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

transform
features by ϕ

Generally,

$$h_{\theta}(\mathbf{x}) = \sum_{j=0}^{d'} \theta_j \underbrace{\phi_j(\mathbf{x})}_{\text{basis function}}$$

d' ← new dimension

Typically, $\phi_0(\mathbf{x}) = 1$ so that θ_0 acts as bias.

In the simplest case, we use linear basis functions:

$$\phi_j(\mathbf{x}) = x_j$$

More complex basis functions allow use of linear regression techniques to **fit non-linear datasets**.

Extending Linear Regression to More Complex Models

The inputs \mathbf{X} for linear regression can be ...

- original quantitative inputs
- transformation of quantitative inputs
 - e.g. log, exp, square root, square, etc
- polynomial transformation
 - e.g. $y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$
- dummy coding of categorical inputs
- interactions between variables
 - e.g. $x_3 = x_1 x_2$

Based on slide by Eric Eaton

Polynomial Regression

Let $x \in \mathbb{R}$.

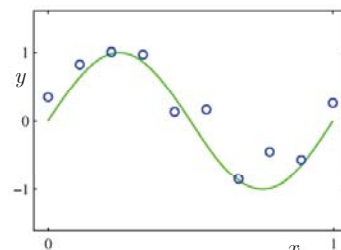
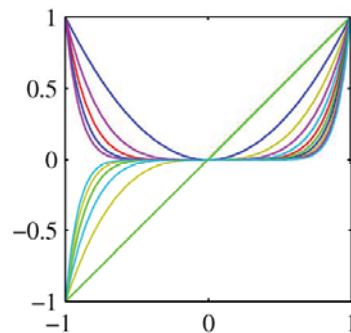
Polynomial basis functions:

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

These are “global” (a small change in x affects all basis functions).

Fit a polynomial curve with a linear model

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_p x^p = \sum_{j=0}^p \theta_j x^j$$



Based on slide by Eric Eaton
(originally by Christopher Bishop [PRML])

Summary

Basic Linear Model

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

Generalized Linear Model

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{j=0}^{d'} \theta_j \phi_j(\mathbf{x})$$

Once we have replaced the data by the outputs of the basis functions, **fitting the generalized linear model is same problem as fitting the basic model**, e.g.

$$\Phi = \begin{bmatrix} - & (\boldsymbol{\phi}(\mathbf{x}^{(1)}))^T & - \\ & \vdots & \\ - & (\boldsymbol{\phi}(\mathbf{x}^{(i)}))^T & - \\ & \vdots & \\ - & (\boldsymbol{\phi}(\mathbf{x}^{(n)}))^T & - \end{bmatrix}$$

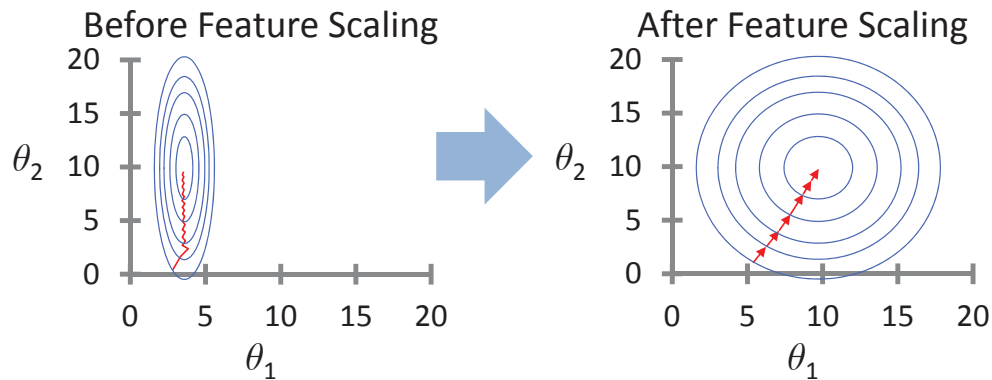
$$\boldsymbol{\theta}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Based on slide by Eric Eaton
(originally by Geoff Hinton)

(This slide intentionally left blank.)

Feature Scaling

Ensure features have similar scales



Why?

- Make GD converge much faster
- Avoid numerical issues in matrix-vector products

Based on slide by Eric Eaton

Feature Standardization

Rescale features to have zero-mean, unit-variance

Let

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)} \quad \sigma_j = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_j^{(i)} - \mu_j)^2}$$

Replace each value with

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad \text{for } j = 1, \dots, d$$

do not rescale x_0

Remember to apply same transformation to both training and prediction

- learn parameters (μ, σ) using only training set
- see `sklearn.preprocessing`
- outliers can cause problems

Based on slide by Eric Eaton