

# REGRESSION



# Regression

Relationships among several quantities

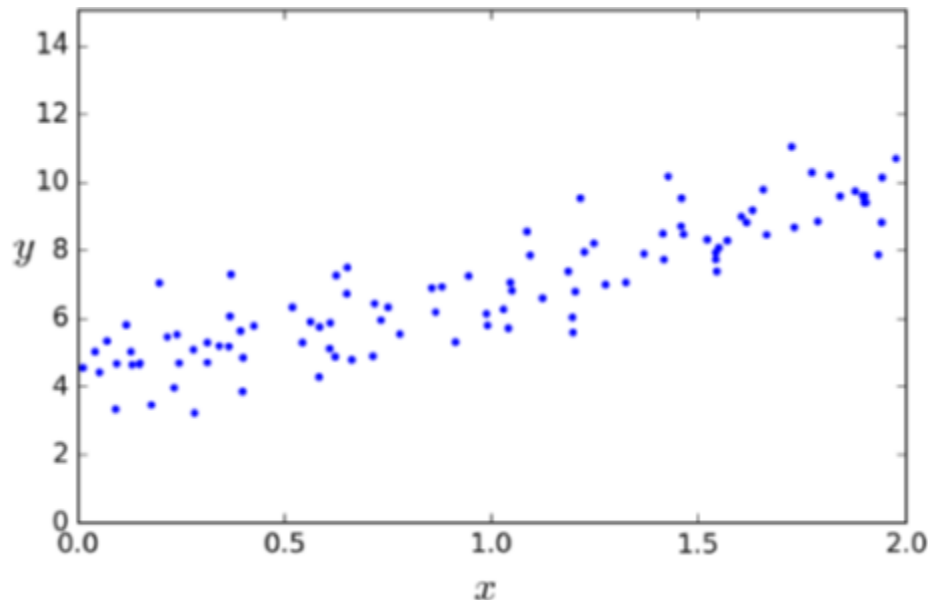
build a model that predicts the value of one variable as a function of other variables

simplest relation between two variables  $x$  and  $y$  is the linear equation

Regression  
coefficients

$$y = \beta_0 + \beta_1 x$$

$$(x_1, y_1), \dots, (x_n, y_n)$$



Regression of  $y$  on  $x$

# Regression

$$y = \beta_0 + \beta_1 x$$

that's the Linear Regression model—but how do we train it?

training a model means setting its parameters so  
that the model best fits the training set

# Regression


If the data points were on the line, the parameters would satisfy the equations

Predicted <i>y</i> -value		Observed <i>y</i> -value
$\beta_0 + \beta_1 x_1$	=	$y_1$
$\beta_0 + \beta_1 x_2$	=	$y_2$
$\vdots$		$\vdots$
$\beta_0 + \beta_1 x_n$	=	$y_n$

if the data points  
don't lie on a line

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$X\beta = y$$

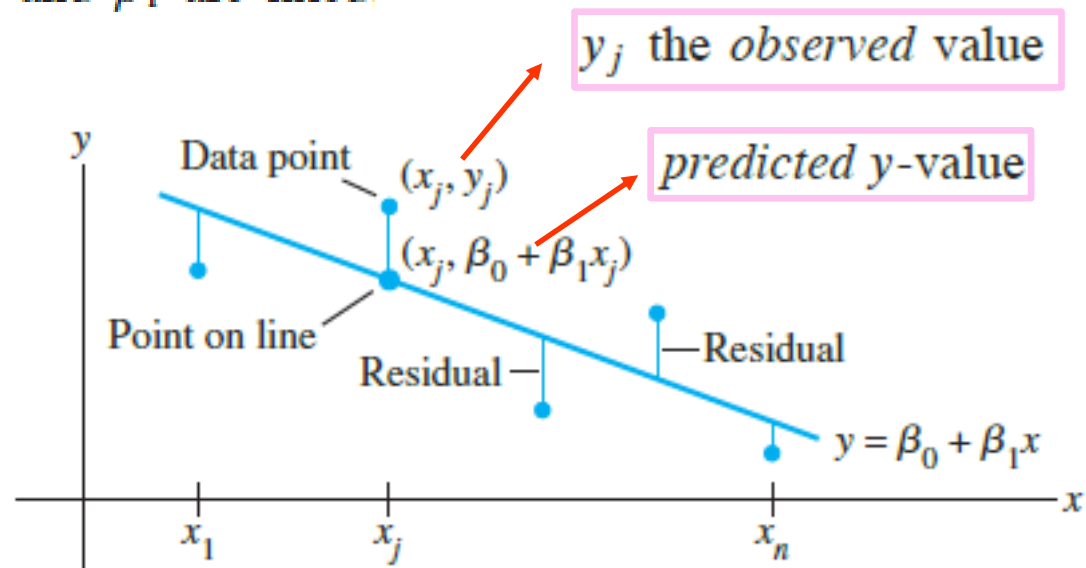


training a model means setting its parameters so  
that the model best fits the training set

we first need a measure of how well (or poorly) the  
model fits the training data.

# Regression

Suppose  $\beta_0$  and  $\beta_1$  are fixed.



# Regression

There are several ways to measure how “close” the line is to the data

The usual choice is to add the squares of the residuals

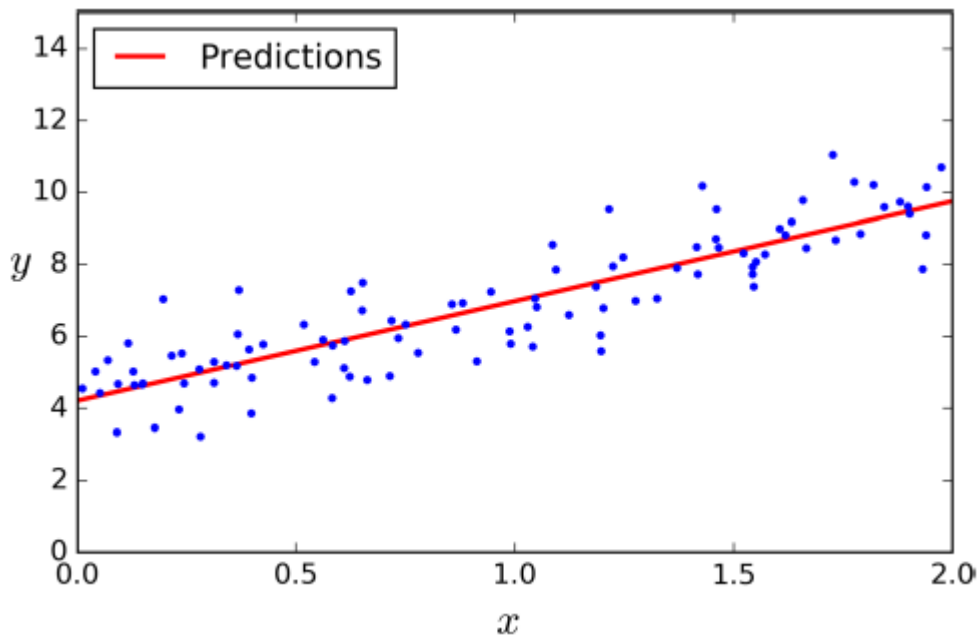
**least-squares line** is the that minimizes the sum of the squares of the residuals

$$\text{residual} = \epsilon = y - X\beta$$

$$\min \|X\beta - y\|_2^2$$

**Cost  
function**

MSE  
RMSE

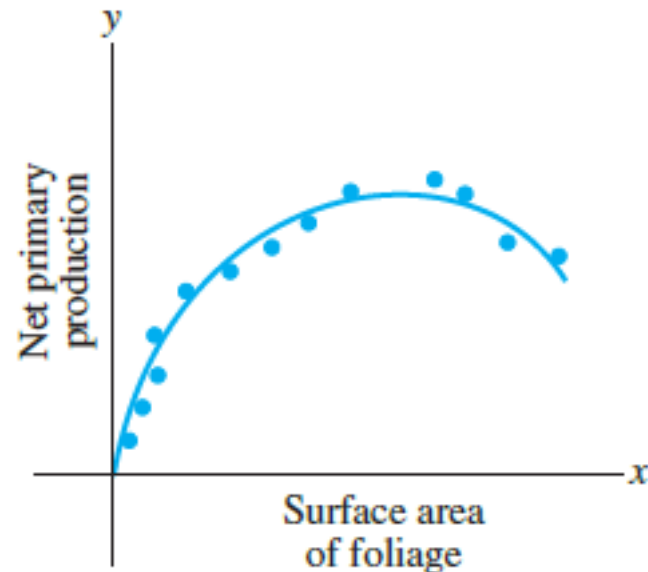
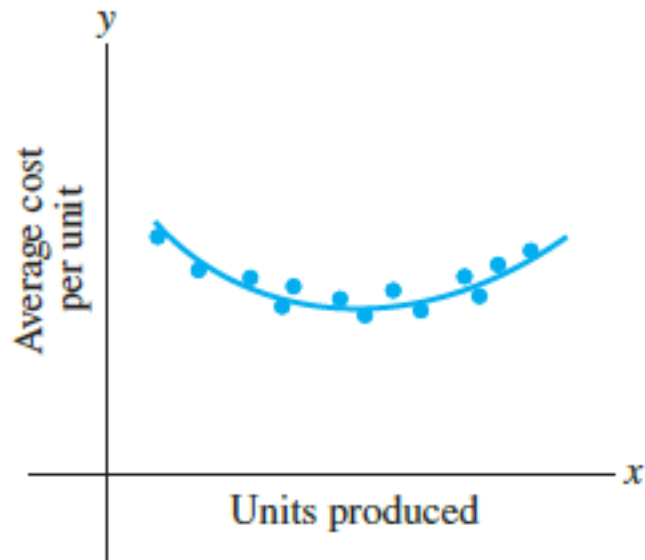


# Regression(Curve fitting )

data points  $(x_1, y_1), \dots, (x_n, y_n)$  on a scatter plot do not lie close to any line,

some other functional relationship between  $x$  and  $y$

$$y = \beta_0 f_0(x) + \beta_1 f_1(x) + \dots + \beta_k f_k(x)$$





# Regression(Curve fitting )

## Example

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

$$(x_1, y_1), \dots, (x_n, y_n)$$

$$\begin{aligned} y_1 &= \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \epsilon_1 \\ y_2 &= \beta_0 + \beta_1 x_2 + \beta_2 x_2^2 + \epsilon_2 \\ &\vdots \\ y_n &= \beta_0 + \beta_1 x_n + \beta_2 x_n^2 + \epsilon_n \end{aligned} \quad \Rightarrow \quad \begin{matrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \\ \mathbf{y} \end{matrix} = \begin{matrix} \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \\ \mathbf{X} \end{matrix} \begin{matrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \\ \boldsymbol{\beta} \end{matrix} + \begin{matrix} \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \\ \boldsymbol{\epsilon} \end{matrix}$$

$$\text{residual} = \epsilon = y - X\beta$$

$$\min \|X\beta - y\|_2^2$$



# Multiple Regression

# Multiple Regression

We have  $n$  features and we want to predict  $y$  based on them

$$x_1, x_2, \dots, x_m \quad \longrightarrow \quad y$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$$

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_m x_m$$

$$y = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_m x_m$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_m^{(1)} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^{(n)} & \dots & x_m^{(n)} \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_m \end{bmatrix}$$

$$\text{residual} = \epsilon = y - X\beta$$

$$\min \|X\beta - y\|_2^2$$

# Multiple Regression

$$y = \beta_0 + \beta_1 f_1(x_1) + \beta_2 f_2(x_2) + \cdots + \beta_m f_m(x_m)$$

$$X = \begin{bmatrix} 1 & f_1(x_1^{(1)}) & \cdots & f_m(x_m^{(1)}) \\ \vdots & \vdots & & \vdots \\ 1 & f_1(x_1^{(n)}) & \cdots & f_m(x_m^{(n)}) \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_m \end{bmatrix}$$

**Least Square  
Problem**

$$\text{residual} = \epsilon = y - X\beta$$

$$\min \|X\beta - y\|_2^2$$

# Solving Least Square Problem

least-squares solution is a solution of the normal equations

$$\min \|X\beta - y\|_2^2$$



$$X^T X \beta = X^T y$$

Optimization  
Methods:  
Gradient  
Descent

High  
computational  
complexity

# Optimization algorithms

$x_0$



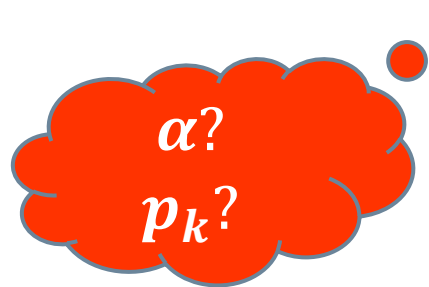
generate a sequence of iterates  $\{x_k\}_{k=0}^{\infty}$



terminate : no more progress or a solution point with sufficient accuracy

$$x_{k+1} = x_k + \alpha p_k$$

direction



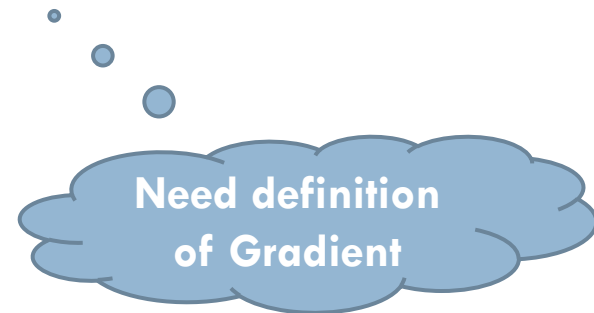
Step length or Learning Rate

# Optimization algorithms

direction

Descent methods

✓ any descent direction is guaranteed to produce a decrease in  $f$ , provided that the step length is sufficiently small



# Gradient

$$f : \mathbf{R}^n \rightarrow \mathbf{R}$$

**f is  
real-valued**

$$\nabla f(x) \quad \rightarrow \quad \nabla f(x)_i = \frac{\partial f(x)}{\partial x_i}, \quad i = 1, \dots, n.$$



# Gradient: example

## Example:

quadratic function

$$f : \mathbf{R}^n \rightarrow \mathbf{R}$$

$$f(x) = (1/2)x^T P x + q^T x + r$$

$P \in \mathbf{S}^n$ ,  $q \in \mathbf{R}^n$ , and  $r \in \mathbf{R}$

$$\nabla f(x) = P x + q$$

# Directional Derivative

$$\nabla_p f(x) = \langle \nabla f(x), p \rangle$$

$$\nabla_p f(x) < 0$$



P is a descent direction

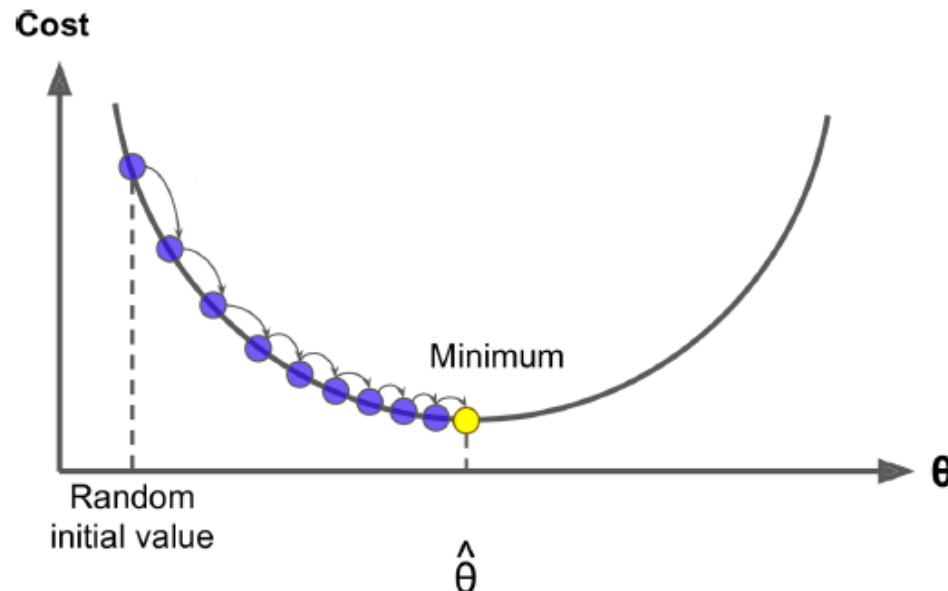
# Gradient Descent

## steepest descent direction

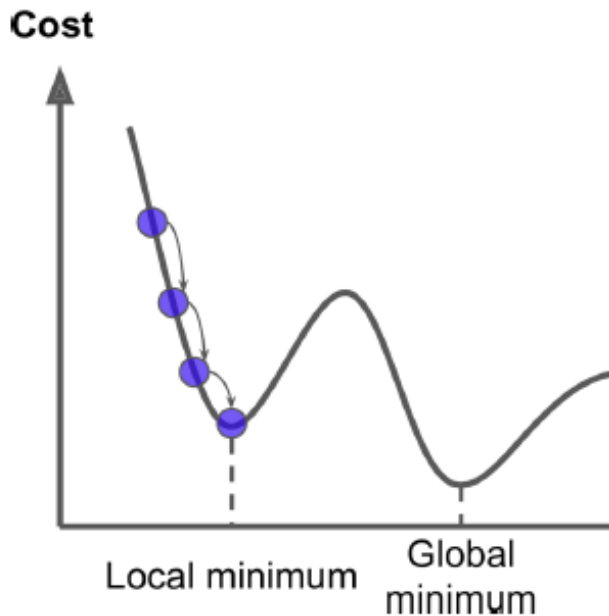
✓ steepest descent direction  $-\nabla f_k$  is the most obvious choice for search direction for a line search method.

✓ choose the step length  $\alpha$  in a variety of ways

$$x_{k+1} = x_k + \alpha_k (-\nabla f(x_k))$$



# Gradient Descent



$$f(\beta) = \|X\beta - y\|_2^2$$

Convex function

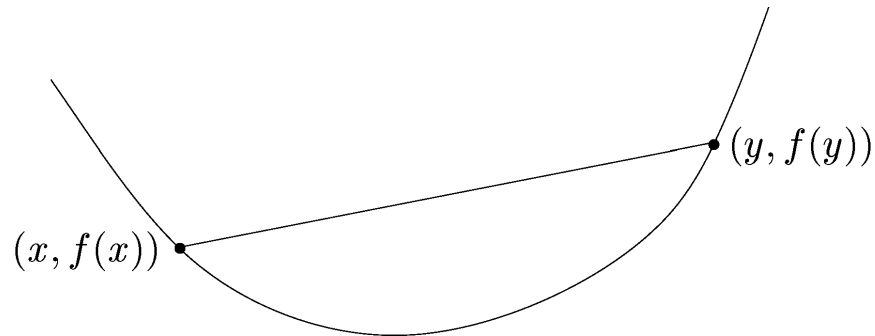
$$\nabla f(\beta) = X^T X \beta - X^T y$$

# Definition

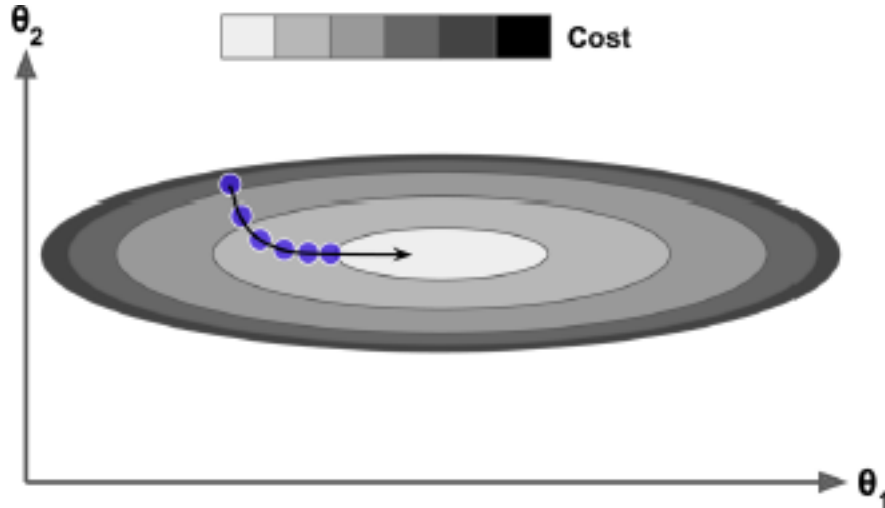
## Convex function

A function  $f : \mathbf{R}^n \rightarrow \mathbf{R}$  is *convex* if  $\mathbf{dom} f$  is a convex set and if for all  $x, y \in \mathbf{dom} f$ , and  $\theta$  with  $0 \leq \theta \leq 1$ , we have

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$



$$f(\boldsymbol{\beta}) = \|X\boldsymbol{\beta} - \mathbf{y}\|_2^2$$

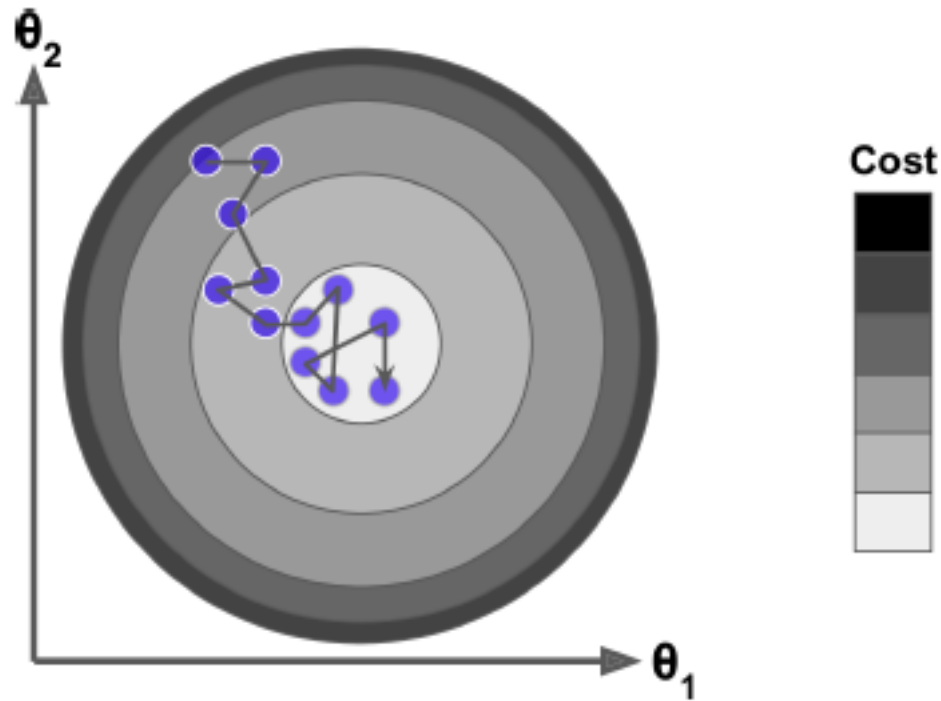


- ✓ training a model means searching for a combination of model parameters minimizes a cost function (over the training set)
- ✓ a search in the model's *parameter space*

# Gradient Descent

- ✓ more parameters a model has, more dimensions this space has, and the harder search
- ✓ **Batch Gradient Descent** uses the whole training set to compute the gradients at every step, which makes it very slow when the training set is large.
- ✓ **Stochastic Gradient Descent** : picks a random instance in the training set at every step and computes the gradients based only on that single instance.
- ✓ is much less regular than Batch Gradient Descent
- ✓ instead of gently decreasing until it reaches the minimum, the cost function will bounce up and down, decreasing only on average.
- ✓ **Mini- Batch Gradient Descent**
- ✓ Mini-batch GD computes the gradients on small random sets of instances called mini-batches

# Stochastic Gradient Descent



final parameter values are good, but not optimal

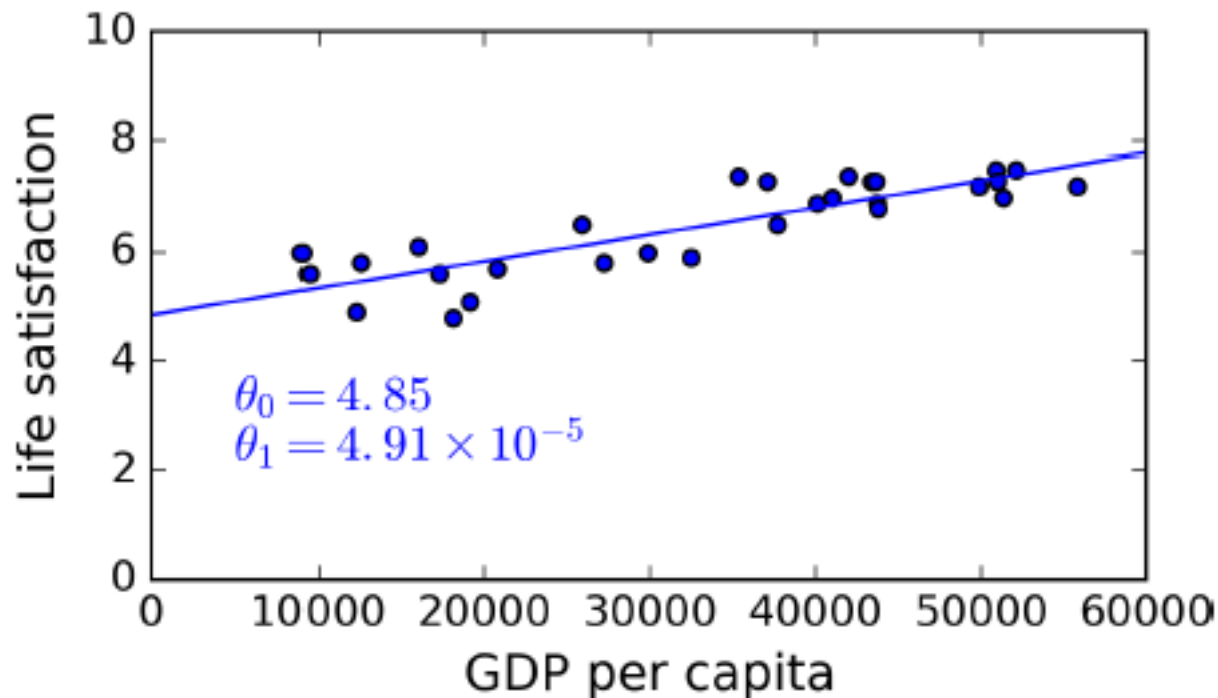




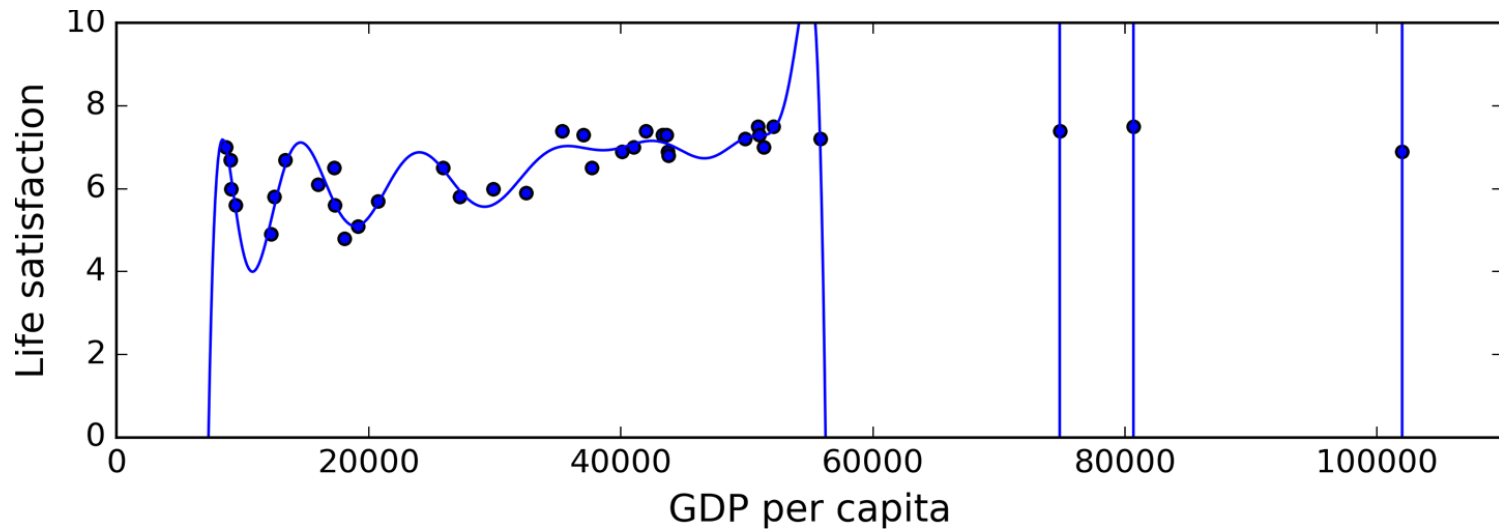
# Overfitting and Underfitting

# Overfitting

- ❖ Overgeneralizing is something that we humans do all too often
- ❖ machines can fall into the same trap
- ❖ In Machine Learning this is called overfitting
- ❖ model performs well on the training data, but it does not generalize well



# Overfitting



Even though it performs much better on the training data than the simple linear model, would you really trust its predictions?

# Overfitting

Overfitting happens when the model is too complex relative to the amount and noisiness of the training data

- ❖ simplify the model
- ❖ gather more training data
- ❖ reduce the noise in the training data

## Underfitting

- ❖ underfitting is the opposite of overfitting
- ❖ it occurs when your model is too simple to learn the underlying structure of the data
- ❖ more powerful model

# Bias/Variance Tradeoff

## **Bias**

- ❖ due to wrong assumptions
- ❖ Assuming that the data is linear when it is actually quadratic.
- ❖ A high-bias model is most likely to underfit the training data.

## **Variance**

- ❖ model's excessive sensitivity to small variations in the training data.
- ❖ A model with many degrees of freedom is likely to have high variance
- ❖ overfit the training data.

# Regularization

- ❖ Constraining a model to make it simpler and reduce the risk of overfitting is called regularization
  - ❖ degrees of freedom
  - ❖ find the right balance between fitting the data perfectly and keeping the model simple enough to ensure that it will generalize well
  - ❖ A simple way to regularize a polynomial model is to reduce the number of polynomial degrees.
- 
- ✓ the fewer degrees of freedom model has, the harder it will be for it to overfit the data
  - ✓ For a linear model, regularization is typically achieved by constraining the weights of the model.

# Regularized Regression

Ridge Regression : Tikhonov regularization

*regularization term*

$$\|\beta\|^2$$

*Cost Function*

$$f(\beta) = \|X\beta - y\|_2^2 + \alpha\|\beta\|_2^2$$

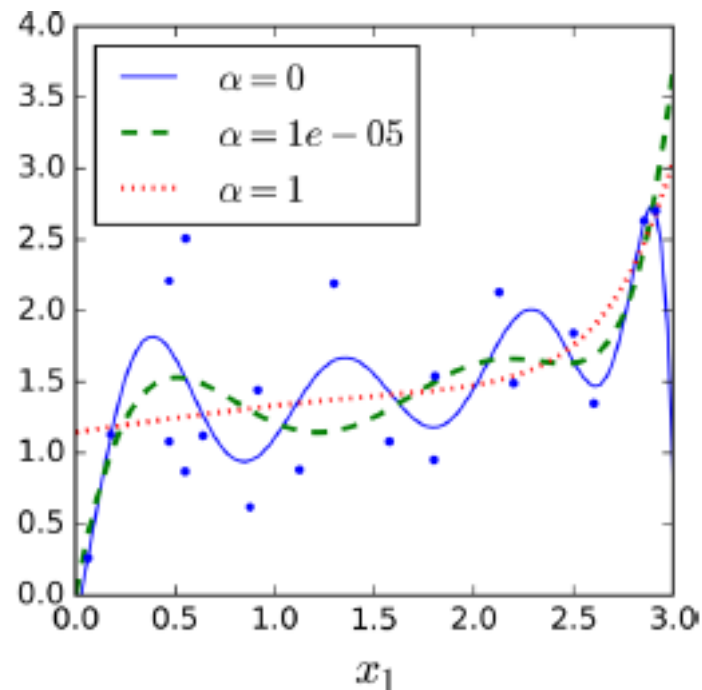


Hyperparameter

Once the model is trained, you want to use the unregularized performance measure to evaluate the model's performance.

It is important to scale the data before performing Ridge Regression, as it is sensitive to the scale of the input features. This is true of most regularized models.

# Tikhonov regularization



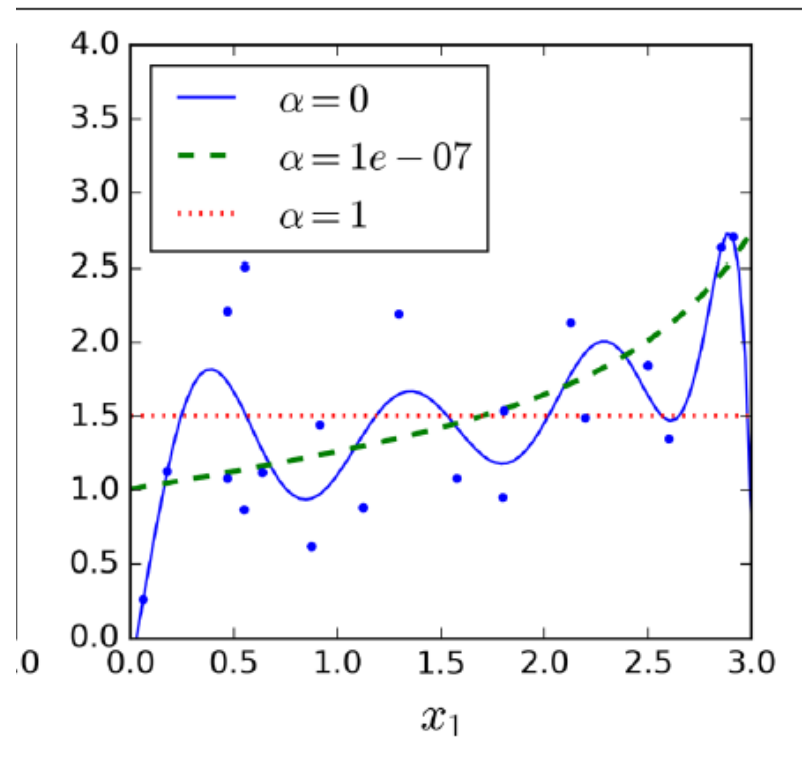
As with Linear Regression, we can perform Ridge Regression either by computing a closed-form equation or by performing Gradient Descent



# Lasso Regression

*Least Absolute Shrinkage  
Selection Operator Regression*

$$f(\beta) = \|X\beta - y\|_2^2 + \alpha\|\beta\|_1$$



# Lasso Regression

Lasso Regression automatically performs  
feature selection and outputs a  
*sparse model*

Lasso cost function is not differentiable  
Gradient Descent  
Subgradient

$$\alpha \|\beta\|_1 \quad \rightarrow \quad \alpha \begin{bmatrix} \text{sign}(\beta_1) \\ \text{sign}(\beta_2) \\ \vdots \\ \text{sign}(\beta_m) \end{bmatrix}$$

**Elastic Net**

a middle ground between Ridge Regression  
and Lasso Regression

## تمرین (سری اول)

□ الف) در منظم سازی **tikhonov** معادله نرمال مربوطه را به دست آورید. ب) رابطه هر گام روش **GD** را برای این روش بنویسید ج) با استفاده از تعریف تابع محدب نشان دهید ترم منظم ساز اضافه شده در این روش یک ترم محدب است

$$f(\beta) = \|X\beta - y\|_2^2 + \alpha\|\beta\|_2^2$$