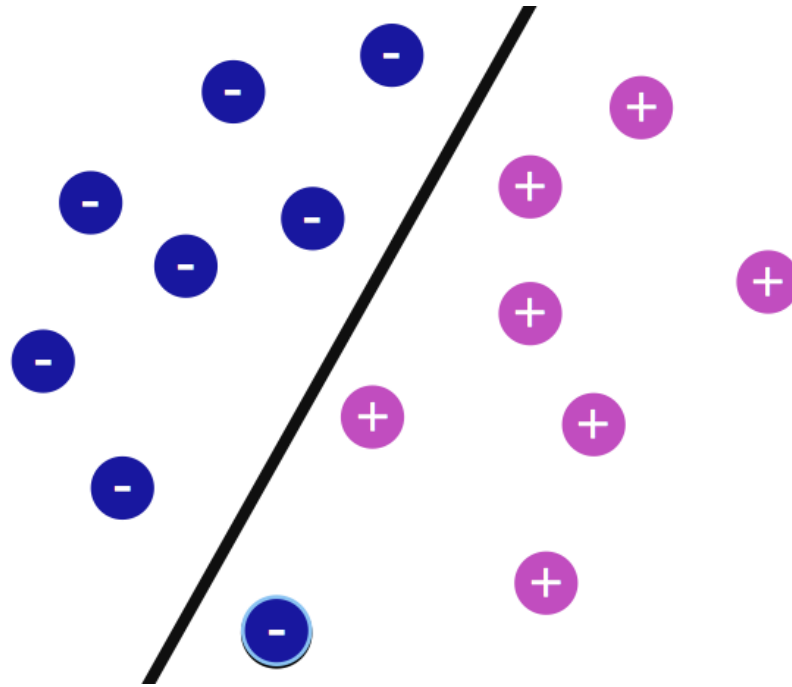


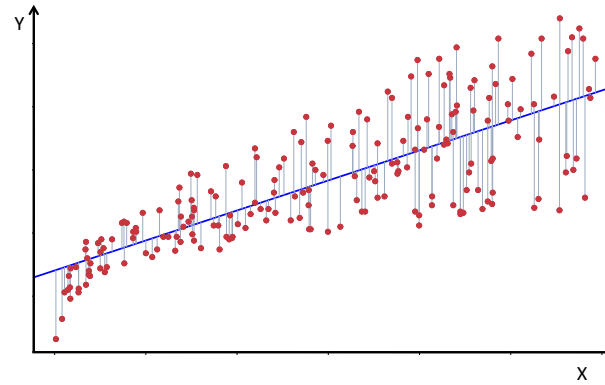
Machine Learning

Linear Models

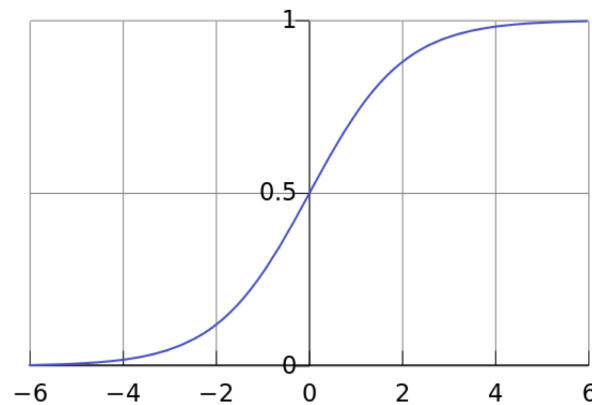


Outline

1. Linear Regression



2. Logistic Regression



Supervised Learning

Training data: “examples” x with “labels” y .

$$(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n) \ / \ \vec{x}_i \in \mathbb{R}^d$$

- **Regression:** y is a real value, $y \in \mathbb{R}$

$$f : \mathbb{R}^d \longrightarrow \mathbb{R}$$

f is called a **regressor**.

- **Classification:** y is discrete. To simplify, $y \in \{-1, +1\}$

$$f : \mathbb{R}^d \longrightarrow \{-1, +1\}$$

f is called a **binary classifier**.

Linear Regression: History

- A very popular technique.
- Rooted in Statistics.
- Method of Least Squares used as early as 1795 by Gauss.
- Re-invented in 1805 by Legendre.
- Frequently applied in **astronomy** to study the large scale of the universe.
- Still a very useful tool today.



Carl Friedrich Gauss

Linear Regression

Given: Training data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ / $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$

Linear Regression

Given: Training data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ / $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$

example $\vec{x}_1 \rightarrow$	x_{11}	x_{12}	\dots	x_{1d}	$y_1 \leftarrow$ label
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_i \rightarrow$	x_{i1}	x_{i2}	\dots	x_{id}	$y_i \leftarrow$ label
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_n \rightarrow$	x_{n1}	x_{n2}	\dots	x_{nd}	$y_n \leftarrow$ label

Linear Regression

Given: Training data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ / $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$

example $\vec{x}_1 \rightarrow$	x_{11}	x_{12}	\dots	x_{1d}	$y_1 \leftarrow$ label
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_i \rightarrow$	x_{i1}	x_{i2}	\dots	x_{id}	$y_i \leftarrow$ label
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_n \rightarrow$	x_{n1}	x_{n2}	\dots	x_{nd}	$y_n \leftarrow$ label

Task: Learn a regression function:

$$f : \mathbb{R}^d \longrightarrow \mathbb{R}$$

$$f(\vec{x}) = y$$

Linear Regression

Given: Training data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ / $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$

example $\vec{x}_1 \rightarrow$	x_{11}	x_{12}	\dots	x_{1d}	$y_1 \leftarrow \text{label}$
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_i \rightarrow$	x_{i1}	x_{i2}	\dots	x_{id}	$y_i \leftarrow \text{label}$
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_n \rightarrow$	x_{n1}	x_{n2}	\dots	x_{nd}	$y_n \leftarrow \text{label}$

Task: Learn a regression function:

$$f : \mathbb{R}^d \longrightarrow \mathbb{R}$$

$$f(\vec{x}) = y$$

Linear Regression: A regression model is said to be linear if it is represented by a linear function.

Linear Regression

Given: Training data: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ / $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$

example $\vec{x}_1 \rightarrow$	x_{11}	x_{12}	\dots	x_{1d}	$y_1 \leftarrow$ label
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_i \rightarrow$	x_{i1}	x_{i2}	\dots	x_{id}	$y_i \leftarrow$ label
\dots	\dots	\dots	\dots	\dots	\dots
example $\vec{x}_n \rightarrow$	x_{n1}	x_{n2}	\dots	x_{nd}	$y_n \leftarrow$ label

Task: Learn a regression function:

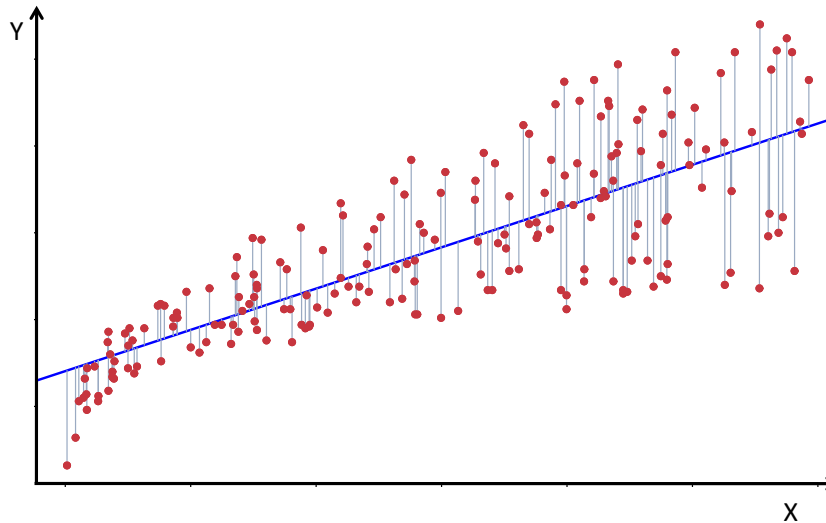
$$f : \mathbb{R}^d \longrightarrow \mathbb{R}$$

$$f(\vec{x}) = y$$

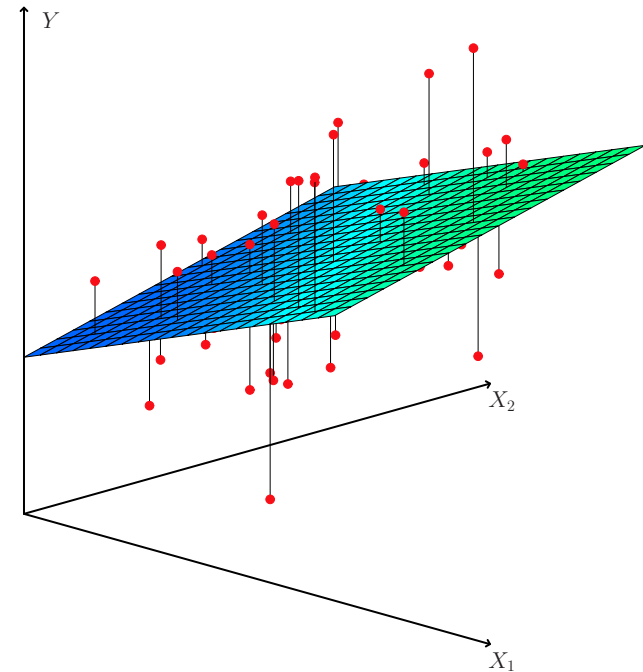
Linear Regression: A regression model is said to be linear if it is represented by a linear function.

Notations: \vec{x}_i or x_i is the i th example. x_j is the j th feature. x_{ij} is the j th feature for the i th example.

Linear Regression



$d = 1$, line in \mathbb{R}^2



$d = 2$, hyperplane in \mathbb{R}^3

Credit: Introduction to Statistical Learning.

Linear Regression

Linear Regression Model:

$$f(\vec{x}_i) = \beta_0 + \sum_{j=1}^d \beta_j x_{ij} \quad \text{with } \beta_j \in \mathbb{R}, \quad j \in \{1, \dots, d\}$$

x_j 's are the features.

β 's are called parameters or coefficients or weights.

Linear Regression

Linear Regression Model:

$$f(\vec{x}_i) = \beta_0 + \sum_{j=1}^d \beta_j x_{ij} \quad \text{with } \beta_j \in \mathbb{R}, \quad j \in \{1, \dots, d\}$$

x_j 's are the features.

β 's are called parameters or coefficients or weights.

Learning the linear model \longrightarrow learning the β 's

Linear Regression

Linear Regression Model:

$$f(\vec{x}_i) = \beta_0 + \sum_{j=1}^d \beta_j x_{ij} \quad \text{with } \beta_j \in \mathbb{R}, \quad j \in \{1, \dots, d\}$$

x_j 's are the features.

β 's are called parameters or coefficients or weights.

Learning the linear model \longrightarrow learning the β 's

Estimation with Least squares:

Use least square loss: $\text{loss}(y_i, f(\vec{x}_i)) = (y_i - f(\vec{x}_i))^2$

Linear Regression

Linear Regression Model:

$$f(\vec{x}_i) = \beta_0 + \sum_{j=1}^d \beta_j x_{ij} \quad \text{with } \beta_j \in \mathbb{R}, \quad j \in \{1, \dots, d\}$$

x_j 's are the features.

β 's are called parameters or coefficients or weights.

Learning the linear model \longrightarrow learning the β 's

Estimation with Least squares:

Use least square loss: $\text{loss}(y_i, f(\vec{x}_i)) = (y_i - f(\vec{x}_i))^2$

We want to minimize the loss over all examples, that is minimize the *risk or cost function* R :

$$R = \frac{1}{2n} \sum_{i=1}^n (y_i - f(\vec{x}_i))^2$$

Linear Regression

A simple case with one feature ($d = 1$):

$$f(x) = \beta_0 + \beta_1 x$$

Linear Regression

A simple case with one feature ($d = 1$):

$$f(x) = \beta_0 + \beta_1 x$$

We want to minimize:

$$R = \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Linear Regression

A simple case with one feature ($d = 1$):

$$f(x) = \beta_0 + \beta_1 x$$

We want to minimize:

$$R = \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$R(\beta) = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Linear Regression

A simple case with one feature ($d = 1$):

$$f(x) = \beta_0 + \beta_1 x$$

We want to minimize:

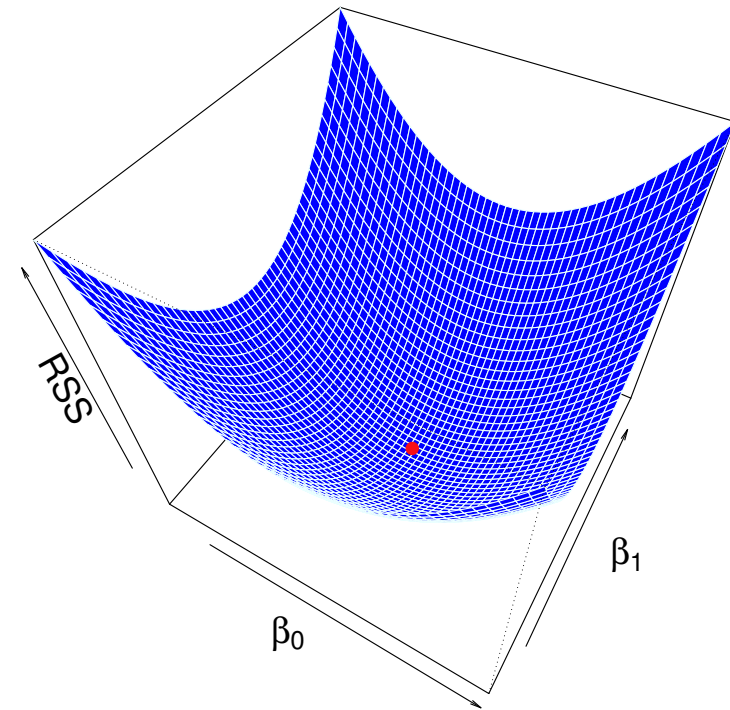
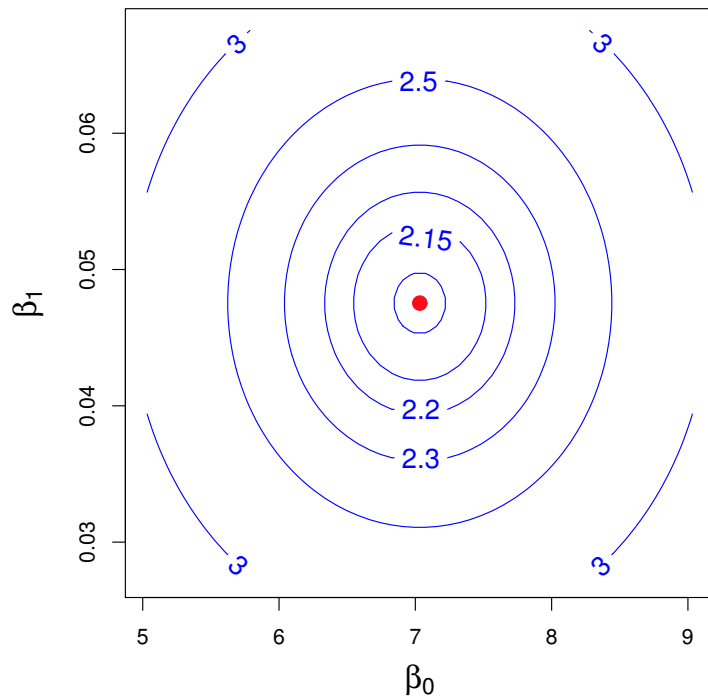
$$R = \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$R(\beta) = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Find β_0 and β_1 that minimize:

$$R(\beta) = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Linear Regression



Credit: Introduction to Statistical Learning.

Linear Regression

Find β_0 and β_1 so that:

$$\operatorname{argmin}_{\beta} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right)$$

Linear Regression

Find β_0 and β_1 so that:

$$\operatorname{argmin}_{\beta} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right)$$

Minimize: $R(\beta_0, \beta_1)$, that is: $\frac{\partial R}{\partial \beta_0} = 0$ $\frac{\partial R}{\partial \beta_1} = 0$

Linear Regression

Find β_0 and β_1 so that:

$$\operatorname{argmin}_{\beta} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right)$$

Minimize: $R(\beta_0, \beta_1)$, that is: $\frac{\partial R}{\partial \beta_0} = 0$ $\frac{\partial R}{\partial \beta_1} = 0$

$$\frac{\partial R}{\partial \beta_0} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_0} (y_i - \beta_0 - \beta_1 x_i)$$

Linear Regression

Find β_0 and β_1 so that:

$$\operatorname{argmin}_{\beta} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right)$$

Minimize: $R(\beta_0, \beta_1)$, that is: $\frac{\partial R}{\partial \beta_0} = 0$ $\frac{\partial R}{\partial \beta_1} = 0$

$$\frac{\partial R}{\partial \beta_0} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_0} (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial R}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-1) = 0$$

Linear Regression

Find β_0 and β_1 so that:

$$\operatorname{argmin}_{\beta} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \right)$$

Minimize: $R(\beta_0, \beta_1)$, that is: $\frac{\partial R}{\partial \beta_0} = 0$ $\frac{\partial R}{\partial \beta_1} = 0$

$$\frac{\partial R}{\partial \beta_0} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_0} (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial R}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-1) = 0$$

$$\beta_0 = \frac{1}{n} \sum_{i=1}^n y_i - \beta_1 \frac{1}{n} \sum_{i=1}^n x_i$$

Linear Regression

$$\frac{\partial R}{\partial \beta_1} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i)$$

Linear Regression

$$\frac{\partial R}{\partial \beta_1} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial R}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-x_i) = 0$$

Linear Regression

$$\frac{\partial R}{\partial \beta_1} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial R}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-x_i) = 0$$

$$\beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - \sum_{i=1}^n \beta_0 x_i$$

Linear Regression

$$\frac{\partial R}{\partial \beta_1} = 2 \times \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial R}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-x_i) = 0$$

$$\beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i - \sum_{i=1}^n \beta_0 x_i$$

Plugging β_0 in β_1 :

$$\beta_1 = \frac{\sum_{i=1}^n y_i x_i - \frac{1}{n} \sum_{i=1}^n y_i \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2 - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n x_i}$$

Linear Regression

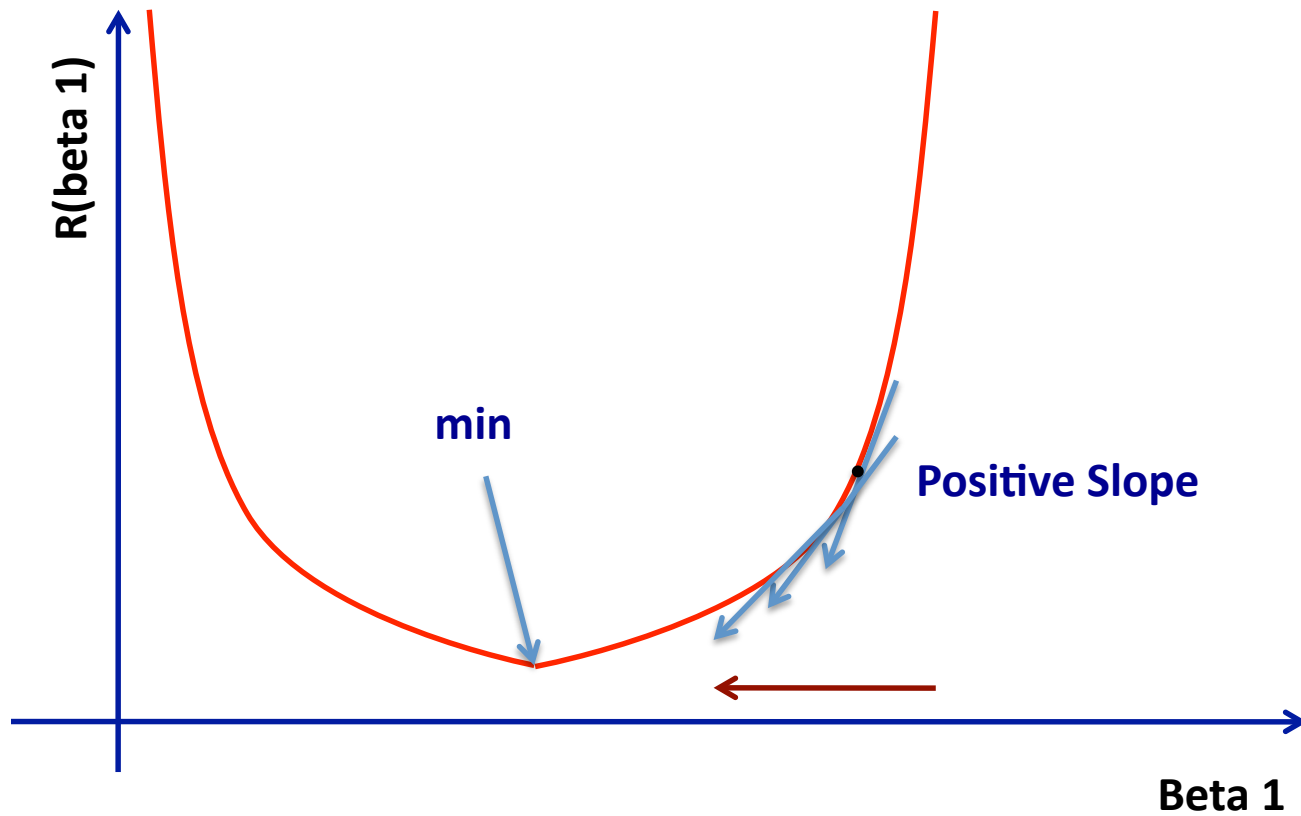
With more than one feature:

$$f(x) = \beta_0 + \sum_{j=1}^d \beta_j x_j$$

Find the β_j that minimize:

$$R = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij})^2$$

Gradient descent



Gradient descent

Gradient Descent is an optimization method.

Repeat until convergence:

Update **simultaneously** all β_j for ($j = 0$ and $j = 1$)

$$\beta_0 := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} R(\beta_0, \beta_1)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} R(\beta_0, \beta_1)$$

Gradient descent

Gradient Descent is an optimization method.

Repeat until convergence:

Update **simultaneously** all β_j for ($j = 0$ and $j = 1$)

$$\beta_0 := \beta_0 - \alpha \frac{\partial}{\partial \beta_0} R(\beta_0, \beta_1)$$

$$\beta_1 := \beta_1 - \alpha \frac{\partial}{\partial \beta_1} R(\beta_0, \beta_1)$$

α is a learning rate.

Gradient descent

In the linear case:

$$\frac{\partial R}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-1) = 0$$

$$\frac{\partial R}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-x_i)$$

Let's generalize it!

Gradient descent

In the linear case:

$$\frac{\partial R}{\partial \beta_0} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-1) = 0$$

$$\frac{\partial R}{\partial \beta_1} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \times (-x_i)$$

Repeat until convergence:

Update **simultaneously** all β_j for ($j = 0$ and $j = 1$)

$$\beta_0 := \beta_0 - \alpha \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i - y_i)$$

$$\beta_1 := \beta_1 - \alpha \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i - y_i)(x_i)$$

Practical considerations

1. **Scaling**: Bring your features to a similar scale.

Practical considerations

1. **Scaling**: Bring your features to a similar scale.

$$x_i := \frac{x_i - \mu_i}{stddev(x_i)}$$

Practical considerations

1. **Scaling**: Bring your features to a similar scale.

$$x_i := \frac{x_i - \mu_i}{stddev(x_i)}$$

2. **Learning rate**: Don't use a rate that is too small or too large.

Practical considerations

1. **Scaling**: Bring your features to a similar scale.

$$x_i := \frac{x_i - \mu_i}{stddev(x_i)}$$

2. **Learning rate**: Don't use a rate that is too small or too large.
3. **R should decrease** after each iteration.

Practical considerations

1. **Scaling**: Bring your features to a similar scale.

$$x_i := \frac{x_i - \mu_i}{stddev(x_i)}$$

2. **Learning rate**: Don't use a rate that is too small or too large.
3. **R should decrease** after each iteration.
4. **Declare convergence** if it start decreasing by less ϵ

Credit

- The elements of statistical learning. Data mining, inference, and prediction. 10th Edition 2009. T. Hastie, R. Tibshirani, J. Friedman.
- Machine Learning 1997. Tom Mitchell.