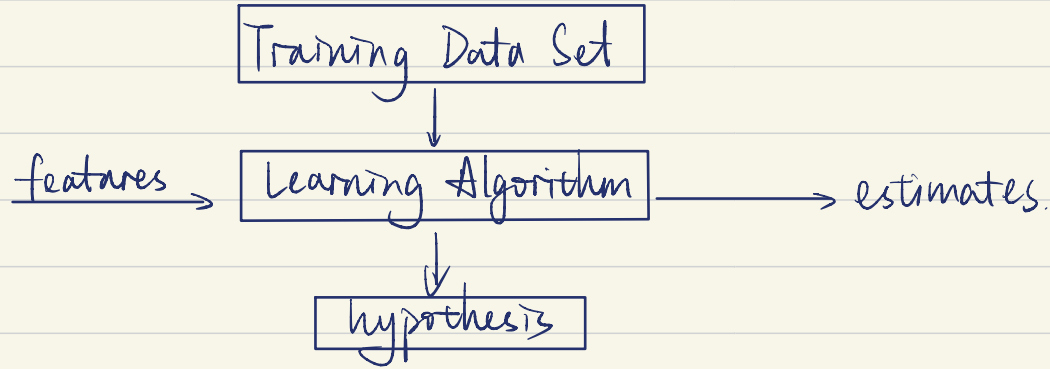


Lecture 2 Linear Regression and Gradient Descent

1.



hypothesis in linear regression represents as: $h(x) = \theta_0 + \theta_1 x_1 + \dots$
equals: $h(x) = \sum_{j=0}^n \theta_j x_j$, where $x_0 = 0$

Notations:

m : # training example

n : # features

$x^{(i)}$: the i^{th} training example

2. Choose θ st. $h(x) \approx y$ for training examples

equals: minimize $\sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 = J(\theta) \leftarrow$ cost function

use gradient descent:

① start with some θ (say $\theta = \vec{0}$)

② keep changing θ to reduce $J(\theta)$

$$\theta_j := \theta_j - \boxed{\alpha} \cdot \frac{\partial}{\partial \theta_j} J(\theta), j = 0, 1, 2, \dots$$

↘ learning rate

Batch Gradient Descent (use the whole data set)

Disadvantage: Computational Expensive

<proof>: when there's only one example:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2 &= \frac{1}{2} \cdot 2 (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h(x) - y) \\ &= (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h(x) - y) = (h(x) - y) \cdot x_j \end{aligned}$$

⇒ Repeat until convergence:

$$\theta := \theta + \alpha \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x^{(i)}$$

3. Stochastic Gradient Descent (more computational efficient)

Repeat {

For $j=1$ to m {

$$\theta := \theta - \alpha (h(x^{(i)}) - y^{(i)}) x^{(i)}$$

}

}

a slightly noisier and more random path
will never quite converge, oscillating around the minimum

4. Normal Equation

Notation: $\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \dots & \frac{\partial f}{\partial A_{1d}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{n1}} & \dots & \frac{\partial f}{\partial A_{nd}} \end{bmatrix}$ for a $f: \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$

Derivations: If A is a square matrix ($A \in \mathbb{R}^{n \times n}$)

$\text{tr } A = \text{sum of diagonal entries} = \sum_i^n A_{ii}$

$$\text{tr } A = \text{tr } A^T$$

$$\text{tr } AB = \text{tr } BA$$

$$\text{If } f(A) = \text{tr } AB$$

$$\text{tr } ABC = \text{tr } CAB$$

$$\Rightarrow \nabla_A f(A) = B^T$$

$$\nabla_A \text{tr } AA^T C = CA + C^T A \quad \left(\frac{d}{dx} a^2 c = 2ac \right)$$

For least squares:

$$X_{m \times d} = \begin{bmatrix} -(x^{(1)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\therefore h_\theta(x) - \vec{y} = X\theta - \vec{y} = \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(m)})^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\therefore J(\theta) = \frac{1}{2} (h_\theta(x) - \vec{y})^2 = \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \quad (\vec{z}^T \vec{z} = \sum_i \vec{z}^2)$$

<proof> Normal Equation:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (\mathbf{x}\theta - \vec{y})^T (\mathbf{x}\theta - \vec{y}) \\&= \frac{1}{2} \nabla_{\theta} (\mathbf{x}\theta)^T \mathbf{x}\theta - (\mathbf{x}\theta)^T \vec{y} - \vec{y}^T (\mathbf{x}\theta) + \vec{y}^T \vec{y}) \\&= \frac{1}{2} \nabla_{\theta} (\theta^T (\mathbf{x}^T \mathbf{x}) \theta - \vec{y}^T (\mathbf{x}\theta) - \vec{y} (\mathbf{x}\theta)) \\&= \frac{1}{2} \nabla_{\theta} (\theta^T (\mathbf{x}^T \mathbf{x}) \theta - 2(\mathbf{x}^T \vec{y})^T \theta) \\&= \frac{1}{2} (2\mathbf{x}^T \mathbf{x} \theta - 2\mathbf{x}^T \vec{y}) \\&= \mathbf{x}^T \mathbf{x} \theta - \mathbf{x}^T \vec{y}\end{aligned}$$