

Linear Regression and SGD

EE599 Deep Learning

Kuan-Wen (James) Huang

Spring 2020

Linear Regression and SGD

- Setting:
 - Input: $\mathbf{x} \in R^D$ (features, observations, inputs)
 - Output: $\hat{y} \in R$ (targets, responses, outputs)
 - Training Data: $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - Model: $\hat{y} = f(\mathbf{x}) = b + \mathbf{w}^T \mathbf{x}$

We call $\mathbf{w} = [w_1, w_2, \dots, w_D]$ **weights** or **parameter vector** and b is the **bias**. Sometimes, we use the augmented weights $\tilde{\mathbf{w}} = [w_1, w_2, \dots, w_D, b]^T$ with the augmented features $\tilde{\mathbf{x}} = [x_1, x_2, \dots, x_D, 1]^T$ so that $\hat{y} = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$.

Linear Regression and SGD

- Objective

- We want to minimize the sum of squared error, i.e.

$$e(\mathbf{w}) = \sum_{i=1}^N (y^{(i)} - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^{(i)})^2.$$

- It can be also written as

$$e(\mathbf{w}) = \|\mathbf{y} - \tilde{X}\tilde{\mathbf{w}}\|^2$$

where

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}, \tilde{X} = \begin{bmatrix} \tilde{\mathbf{x}}^{(1)T} & 1 \\ \vdots & \vdots \\ \tilde{\mathbf{x}}^{(N)T} & 1 \end{bmatrix}$$

Linear Regression and SGD

- The optimal solution which minimizes $e(\mathbf{w})$ is

$$\mathbf{w}^* = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \mathbf{y}$$

if $\tilde{X}^T \tilde{X}$ is invertible.

- In general, not knowing whether the objective is convex quadratic, we solve it via **stochastic gradient descent**.

Linear Regression and SGD

- Let's do some coding!