



Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**



Machine Learning: Yoshinari Fujinuma

University of Colorado Boulder

LECTURE 10

Slides adapted from Chenhao Tan, Jordan Boyd-Graber, Chris Ketelsen

Logistics

- HW2 available on Github
- Small edits made on the slides from the last lecture
- Final project team formulation due on March 1st

Learning objectives

- Understand stochastic gradient descent

Outline

Stochastic Gradient Descent (SGD)

Problem with Gradient Descent

- $(\mathbf{x}, y) = \{(\mathbf{x}_0, y_0), \dots, (\mathbf{x}_N, y_N)\}$: Training data with N examples. x_{ij} is the j th feature of i th example.
- $\beta = (\beta_0, \dots, \beta_d)$: Parameters of logistic regression.
- η : Learning rate (step size)

Updating parameters by gradient descent is

$$\beta'_j \leftarrow \beta_j - \eta \frac{\partial \mathcal{L}}{\partial \beta_j} \quad (1)$$

where $\frac{\partial \mathcal{L}}{\partial \beta_j}$ for logistic regression is

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = \sum_i^N - (y_i - \sigma_i) \mathbf{x}_{ij} \quad (2)$$

Problem: \sum_i^N indicates that you need to go through **all** training examples

Approximating the Gradient

- Training datasets are big these days (to fit into memory)
- What if we compute an update just from one observation?

Intuition of SGD: Analogy to Getting to Union Station

Pretend it's a pre-smartphone world and you want to get to Union Station



Stochastic Gradient Descent

- $(\mathbf{x}, y) = \{(\mathbf{x}_0, y_0), \dots, (\mathbf{x}_N, y_N)\}$: Training data with N examples. x_{ij} is the j th feature of i th example.
- $\beta = (\beta_0, \dots, \beta_n)$: Parameters of logistic regression.
- η : Learning rate (step size)

Updating the parameters by stochastic gradient descent

$$\beta'_j \leftarrow \beta_j - \eta \frac{\partial \mathcal{L}}{\partial \beta_j} \quad (3)$$

where $\frac{\partial \mathcal{L}}{\partial \beta_j}$ for logistic regression is

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = -(y_i - \sigma_i) \mathbf{x}_{ij} \quad (4)$$

We now compute $\frac{\partial \mathcal{L}}{\partial \beta_j}$ without \sum_i^N i.e., only from **one** training example

Gradient Descent vs. Stochastic Gradient Descent

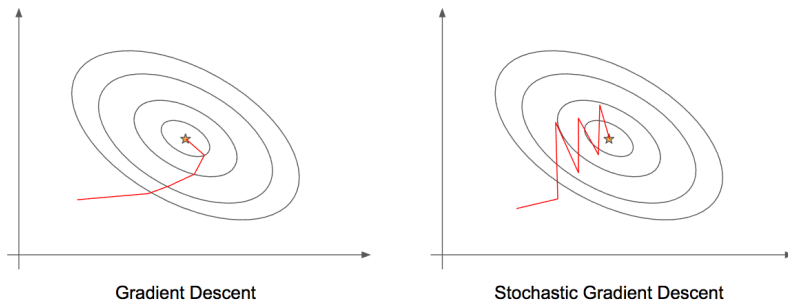


Image from <https://pythonmachinelearning.pro/complete-guide-to-deep-neural-networks-part-2/>

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle \beta_0 = 0, \beta_1 = 0, \beta_2 = 0, \beta_3 = 0, \beta_4 = 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

You first see the positive example. First, compute σ_1

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

You first see the positive example. First, compute σ_1

$$\sigma_1 = \Pr(y_1 = 1 | \mathbf{x}_1) = \frac{1}{1 + \exp -\beta^T \mathbf{x}_1} =$$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

You first see the positive example. First, compute σ_1

$$\sigma_1 = \Pr(y_1 = 1 | \mathbf{x}_1) = \frac{1}{1 + \exp -\beta^T \mathbf{x}_1} = \frac{1}{1 + \exp 0} = 0.5$$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

$\sigma_1 = 0.5$ What's the updated β'_0 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_0 ? $\beta'_0 = \beta_0 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,0} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 1.0$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_0 ? $\beta'_0 = \beta_0 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,0} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 1.0 = 0.5$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_1 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_1 ? $\beta'_1 = \beta_1 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,1} = 0.0 + 1.0 \cdot (1 - 0.5) \cdot 4$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_1 ? $\beta'_1 = \beta_1 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,1} = 0.0 + 1.0 \cdot (1 - 0.5) \cdot 4 = 2.0$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_2 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_2 ? $\beta'_2 = \beta_2 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,2} = 0.0 + 1.0 \cdot (1 - 0.5) \cdot 3$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_2 ? $\beta'_2 = \beta_2 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,2} = 0.0 + 1.0 \cdot (1 - 0.5) \cdot 3 = 1.5$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_3 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_3 ? $\beta'_3 = \beta_3 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,3} = 0.0 + 1.0 \cdot (1 - 0.5) \cdot 1$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_3 ? $\beta'_3 = \beta_3 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,3} = 0.0 + 1.0 \cdot (1 - 0.5) \cdot 1 = 0.5$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_4 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_4 ? $\beta'_4 = \beta_4 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,4} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 0$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0, 0, 0, 0, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_4 ? $\beta'_4 = \beta_4 + \eta \cdot (y_1 - \sigma_1) \cdot \mathbf{x}_{1,4} = 0.0 + 1.0 \cdot (1.0 - 0.5) \cdot 0 = 0.0$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

Now you see the negative example. What's σ_2 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

Now you see the negative example. What's σ_2 ?

$$\sigma_2 = \Pr(y_2 = 0 | \vec{x}_2) = \frac{\exp -\beta^T \mathbf{x}_i}{1 + \exp -\beta^T \mathbf{x}_i} = \frac{\exp\{-(.5+1.5+1.5+0)\}}{1 + \exp\{-(.5+1.5+1.5+0)\}} =$$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

Now you see the negative example. What's σ_2 ?

$$\sigma_2 = \Pr(y_2 = 0 | \vec{x}_2) = \frac{\exp -\beta^T \mathbf{x}_i}{1 + \exp -\beta^T \mathbf{x}_i} = \frac{\exp\{-(.5+1.5+1.5+0)\}}{1 + \exp\{-(.5+1.5+1.5+0)\}} = 0.97$$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

Now you see the negative example. What's σ_2 ?

$$\sigma_2 = 0.97$$

What's the updated β'_0 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_0 ? $\beta'_0 = \beta_0 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,0} = 0.5 + 1.0 \cdot (0 - 0.97) \cdot 1$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_0 ? $\beta'_0 = \beta_0 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,0} = 0.5 + 1.0 \cdot (0 - 0.97) \cdot 1 = -0.47$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_1 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_1 ? $\beta'_1 = \beta_1 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,1} = 2.0 + 1.0 \cdot (0 - 0.97) \cdot 0$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_1 ? $\beta'_1 = \beta_1 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,1} = 2.0 + 1.0 \cdot (0 - 0.97) \cdot 0 = 2.0$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_2 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_2 ? $\beta'_2 = \beta_2 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,2} = 1.5 + 1.0 \cdot (0 - 0.97) \cdot 1$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_2 ? $\beta'_2 = \beta_2 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,2} = 1.5 + 1.0 \cdot (0 - 0.97) \cdot 1 = 0.53$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$
$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_3 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_3 ? $\beta'_3 = \beta_3 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,3} = 0.5 + 1.0 \cdot (0 - 0.97) \cdot 3$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_3 ? $\beta'_3 = \beta_3 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,3} = 0.5 + 1.0 \cdot (0 - 0.97) \cdot 3 = -2.41$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$
$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_4 ?

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\begin{aligned}\beta'_j &= \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij} \\ \vec{\beta} &= \langle 0.5, 2, 1.5, 0.5, 0 \rangle\end{aligned}$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_4 ? $\beta'_4 = \beta_4 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,4} = 0 + 1.0 \cdot (0 - 0.97) \cdot 4$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle 0.5, 2, 1.5, 0.5, 0 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

What's the updated β'_4 ? $\beta'_4 = \beta_4 + \eta \cdot (y_2 - \sigma_2) \cdot \mathbf{x}_{2,4} = 0 + 1.0 \cdot (0 - 0.97) \cdot 4 = -3.88$

SGD Example

Note that $-\eta \frac{\partial \mathcal{L}}{\partial \beta_j} = \eta(y_i - \sigma_i)\mathbf{x}_{ij}$

$$\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$$

$$\vec{\beta} = \langle -0.47, 2, 0.53, -2.41, -3.88 \rangle$$

$$y_1 = 1$$

$$\mathbf{x}_1 = (1, 4, 3, 1, 0)$$

(Assume step size $\eta = 1.0$.)

$$y_2 = 0$$

$$\mathbf{x}_2 = (1, 0, 1, 3, 4)$$

Overview of Optimizing β using SGD

1. Initialize a vector β to be all zeros
2. For $t = 1, \dots, T$ (i.e. number of epochs)
 - For each example \mathbf{x}_i, y_i and each feature j :
 - Compute $\sigma_i = \Pr(y_i | \mathbf{x}_i)$
 - Set $\beta'_j = \beta_j + \eta(y_i - \sigma_i)\mathbf{x}_{ij}$
3. Output the parameters β_0, \dots, β_d .

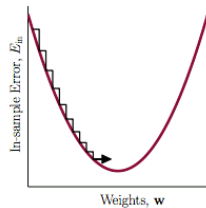
Overview of Optimizing β using SGD

1. Initialize a vector β to be all zeros
2. For $t = 1, \dots, T$ (i.e. number of epochs)
 - For each example \mathbf{x}_i, y_i and each feature j :
 - Compute $\sigma_i = \Pr(y_i | \mathbf{x}_i)$
 - Set $\beta'_j = \beta_j + \eta(y_i - \sigma_i)x_{ij}$
3. Output the parameters β_0, \dots, β_d .

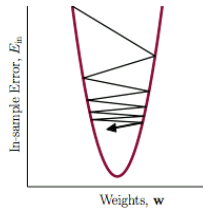
How to decide η ?

Choosing learning rate

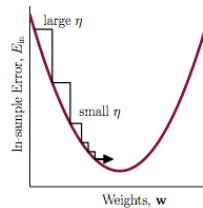
η too small



η too large



variable η_t – just right



Learning rate decay

Decay schedule can be seen as a hyperparameter too.

- Decay after each epoch (e.g., $\frac{\eta_0}{t^2}$)

Advanced stochastic gradient descent:

<http://ruder.io/optimizing-gradient-descent/>

(Bonus) Mini-Batch Stochastic Gradient Descent

- $(\mathbf{x}, y) = \{(\mathbf{x}_0, y_0), \dots, (\mathbf{x}_N, y_N)\}$: Training data with N examples. x_{ij} is the j th feature of i th example.
- $\beta = (\beta_0, \dots, \beta_n)$: Parameters of logistic regression.
- η : Learning rate (step size)

Updating the parameters by stochastic gradient descent

$$\beta'_j \leftarrow \beta_j - \eta \frac{\partial \mathcal{L}}{\partial \beta_j} \quad (5)$$

where $\frac{\partial \mathcal{L}}{\partial \beta_j}$ for logistic regression is

$$\frac{\partial \mathcal{L}}{\partial \beta_j} = \sum_i^M - (y_i - \sigma_i) \mathbf{x}_{ij} \quad (6)$$

We now compute $\frac{\partial \mathcal{L}}{\partial \beta_j}$ from M training examples (less noisy)