

EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the [Privacy Policy](#).



# 5. Stochastic Gradient Descent

## Stochastic Gradient Descent



### Stochastic gradient descent (SGD)

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \left[ \text{Loss}_h(y^{(i)} \theta \cdot x^{(i)}) + \frac{\lambda}{2} \|\theta\|^2 \right]$$

Select  $i \in \{1, \dots, n\}$  at random

$$\theta \leftarrow \theta - \eta_t \nabla_{\theta} \left[ \text{Loss}_h(y^{(i)} \theta \cdot x^{(i)}) + \frac{\lambda}{2} \|\theta\|^2 \right]$$



derivative

from the last terms, if the loss is non-zero.  
And that update looks like the perceptron update, but it is actually made even if we correctly classify the example.  
If the example is within the margin boundaries, you would get a non-zero loss.  
So here, we have just a better way of writing what that stochastic gradient descent update or SGD **update looks like.**

▶ 8:06 / 8:06

▶ Speed 1.50x

🔊

⌂

CC

“

[End of transcript. Skip to the start.](#)

Video  
[Download video file](#)

Transcripts  
[Download SubRip\(.srt\) file](#)  
[Download Text\(.txt\) file](#)

## SGD and Hinge Loss

1/1 point (graded)

As we saw in the lecture above,

$$J(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2 = \frac{1}{n} \left[ \sum_{i=1}^n \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2 \right]$$

With stochastic gradient descent, we choose  $i \in \{1, \dots, n\}$  at random and update  $\theta$  such that

$$\theta \leftarrow \theta - \eta \nabla_{\theta} [\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2]$$

What is  $\nabla_{\theta} [\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0))]$  if  $\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) > 0$ ?

☐  $y^{(i)} x^{(i)}$

☒  $-y^{(i)} x^{(i)}$  ✓

☐ 0

☐  $\lambda \theta$

☐  $-\lambda \theta$

**Solution:**

If  $\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) > 0$

$$\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) = 1 - y^{(i)}(\theta \cdot x^{(i)} + \theta_0)$$

. Thus

$$\nabla_{\theta} \text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) = -y^{(i)} x^{(i)}$$

Submit

You have used 1 of 3 attempts

**i** Answers are displayed within the problem

## Comparison with Perceptron

1/1 point (graded)

Observing the update step of SGD,

$$\theta \leftarrow \theta - \eta \nabla_{\theta} [\text{Loss}_h(y^{(i)}(\theta \cdot x^{(i)} + \theta_0)) + \frac{\lambda}{2} \|\theta\|^2]$$

Which of the following is true?

☐ As in perceptron,  $\theta$  is not updated when there is no mistake

☒ Differently from perceptron,  $\theta$  is updated even when there is no mistake ✓

### Solution:

We can see from

$$\theta \leftarrow \begin{cases} (1 - \lambda\eta) \theta & \text{if Loss}=0 \\ (1 - \lambda\eta) \theta + \eta y^{(i)} x^{(i)} & \text{if Loss}>0 \end{cases}$$

that  $\theta$  is updated even when the sum of losses is 0. This is different from perceptron.

Submit

You have used 1 of 1 attempt


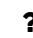


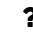
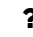
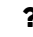


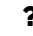
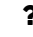
 Answers are displayed within the problem

Discussion

Hide Discussion

**Topic:** Unit 1 Linear Classifiers and Generalizations (2 weeks):Lecture 4. Linear Classification and Generalization / 5. Stochastic Gradient Descent

Add a Post

Show all posts	by recent activity
 <a href="#">[Staff] Formula in SGD and Hinge Loss</a>	3
 <a href="#">Constraints on Learning Parameter <math>\eta</math> (eta)</a> I am unable to understand the reasons and requirements behind constraints mentioned by Professor on the learning rate parameter $\eta$ (starting at 3:18 in the video). What is ...	2
 <a href="#">Missing norm(...) in regularization term at 6:37?</a> Shouldn't there be the norm of theta instead of just theta when we took the derivative and outlined the case when the loss is not zero?	5
 <a href="#">How to check if the shape of objective function</a> Of perception, we want to minimize objective, but how to know the shape of it ? or all objective functions are assumed to have convex or concave shape ?	5
 <a href="#">Updating offset parameter</a>	4
 <a href="#">regularzation term when loss = 0</a> In the lecture the prof. only note the regularzation term when loss >0, should we also include this term when loss = 0?	5
 <a href="#">difference between gradient descent and SGD</a> I am a bit confused with this lecture, can I say the difference between both of them is SGD converge faster than GD? because the step size of SGD is a function with time (e.g.,...	3
 <a href="#">Doubt about gradient of regularization parameter term.</a>	2
 <a href="#">Three differences with Perceptron</a> The professor mentioned three differences with perceptron. I could see two. Can anyone help me understand what the third difference is?	7
 <a href="#">x is missing superscript in several places on this page</a> .	2
 <a href="#">Comparison with Perceptron Confusion</a> [STAFF] I'm a bit confused on the Comparison with Perceptron question. It seems at odd with the lecture, and the updating algo in the solution is not what Prof wrote in his sli...	2

Learn About Verified Certificates