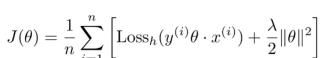
5. Stochastic Gradient Descent **Stochastic Gradient Descent**

Stochastic gradient descent (SGD)



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

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

And that update looks like the perceptor 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.

update looks like.

8:06 / 8:06

▶ 1.25x

X

CC

End of transcript. Skip to the start.

Video Download video file **Transcripts** <u>Download SubRip (.srt) file</u> Download Text (.txt) file



SGD and Hinge Loss

1/1 point (graded)

As we saw in the lecture above,

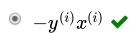
$$J\left(heta, heta_{0}
ight)=rac{1}{n}\sum_{i=1}^{n}\operatorname{Loss}_{h}\left(y^{\left(i
ight)}\left(heta\cdot x^{\left(i
ight)}+ heta_{0}
ight)
ight)+rac{\lambda}{2}\left|\left|\left. heta\left.
ight|
ight|^{2}=rac{1}{n}ig[\sum_{i=1}^{n}\operatorname{Loss}_{h}\left(y^{\left(i
ight)}\left(heta\cdot x^{\left(i
ight)}+ heta_{0}
ight)
ight)+rac{\lambda}{2}\left|\left|\left. heta\left.
ight|
ight|^{2}ig]$$

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

$$heta \leftarrow heta - \eta
abla_{ heta} igl[\operatorname{Loss}_h \left(y^{(i)} \left(heta \cdot x^{(i)} + heta_0
ight)
ight) + rac{\lambda}{2} \mid\mid heta \mid\mid^2 igr]$$

What is $abla_{ heta} \left[\operatorname{Loss}_h \left(y^{(i)} \left(heta \cdot x^{(i)} + heta_0
ight)
ight)
ight]$ if $\operatorname{Loss}_h \left(y^{(i)} \left(heta \cdot x^{(i)} + heta_0
ight)
ight) > 0$?

 $y^{(i)} y^{(i)}$





 $-\lambda\theta$

Solution:

If
$$\operatorname{Loss}_h\left(y^{(i)}\left(heta\cdot x^{(i)}+ heta_0
ight)
ight)>0$$
,

$$\operatorname{Loss}_h\left(y^{(i)}\left(heta\cdot x^{(i)}+ heta_0
ight)
ight)=1-y^{(i)}\left(heta\cdot x^{(i)}+ heta_0
ight)$$

. Thus

$$abla_{ heta} \mathrm{Loss}_h\left(y^{(i)}\left(heta \cdot x^{(i)} + heta_0
ight)
ight) = -y^{(i)}x^{(i)}$$

Submit

You have used 1 of 3 attempts

Answers are displayed within the problem

Comparison with Perceptron

1/1 point (graded)

Observing the update step of SGD,

$$heta \leftarrow heta - \eta
abla_{ heta} igl[\operatorname{Loss}_h \left(y^{(i)} \left(heta \cdot x^{(i)} + heta_0
ight)
ight) + rac{\lambda}{2} \mid\mid heta \mid\mid^2 igr]$$

Which of the following is true?

- ullet As in perceptron, heta is not updated when there is no mistake
- ullet Differently from perceptron, heta is updated even when there is no mistake ullet

Solution:

We can see from

$$heta \leftarrow \left\{ egin{aligned} \left(1 - \lambda \eta
ight) heta ext{ if Loss} = 0 \ \left(1 - \lambda \eta
ight) heta + \eta y^{(i)} x^{(i)} ext{ if Loss} {>} 0 \end{aligned}
ight.$$

that θ 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