

# Curious behavior of stochastic gradient ①

## Descent

### Setup

data  $x_i$   $y_i$   
↑ ↑  
features labels  $\pm 1$

linear predictor  $w$

loss  $l$ .

so we minimize  $\sum_i l(w_i^T [x_i y_i]) = L(w)$   
by choice of  $w$

### Assume :

• data is separable

• loss is  $e^{-u}$   
(this can be relaxed!)

• for simplicity of notation

all labels are  $\pm 1$   
(so if  $x_i$  has label  $-1$ ,  
use  $-x_i$  instead.)

Notice in this case; there is  
no finite  $w$  that minimizes the  
loss but

$$\lim_{w \rightarrow \infty} L(w) = 0$$

(in the right dirn)

So there is no <sup>finite</sup> critical point of the loss

Q: Where does SGD go?

SGD is

$$\begin{aligned} w_{t+1} &= w_t - \eta \nabla L(w_t) \\ &= w_t - \eta \sum \ell'(\omega^T x_i) x_i \end{aligned}$$

1: (Soudry et al)  
with assumptions,

1)  $\lim_{t \rightarrow \infty} L(w(t)) = 0$

2)  $\lim_{t \rightarrow \infty} \|w_t\| = \infty$

3) for all  $i$ ,  $\lim_{t \rightarrow \infty} \omega_t^T x_i = \infty$

Sketch of proof:

(3)

• Data is separable, so that there is some  $w^*$  st  $w^{*T} x_i > 0$  for all  $i$ .

• This means

$$w^{*T} \nabla L = \sum_i l'(w^{*T} x_i) (w^{*T} x_i) < 0$$

(cause  $l(u) = -e^{-u}$ , so neg; )

• so there aren't any finite critical points.

• But GD goes to a critical point

$$\text{— so } \|w_t\| \rightarrow \infty$$

$$L(w_t) \rightarrow \infty$$

Q: GD goes to  $\infty$ ; but in what dir'n?

we have

$$w_t^T x_i \rightarrow \infty$$

now consider

$$\frac{w_t}{\|w_t\|}$$

if it converges to some limit  $w_\infty$ ,

we have

$$w_t = g(t) \underline{w_\infty} + \underline{p(t)} \leftarrow \text{vector!}$$

where  $g(t) \rightarrow \infty$

$$\lim_{t \rightarrow \infty} \frac{p(t)}{g(t)} = 0$$

$$\begin{aligned} -\nabla L(w) &= \sum_i \exp[-w_t^T x_i] x_i \\ &= \sum_i \exp[-g(t) w_\infty^T x_i] \exp[-p(t)^T x_i] x_i \end{aligned}$$



Now

as  $t \rightarrow \infty$ , consider weight of  $i$ th example (5)

$$\frac{-g(t) w_{\infty}^T x_i}{e} \quad \frac{-p(t)^T x_i}{e}$$

$g(t) \rightarrow \infty$ , so if  $w_{\infty}^T x_i < w_{\infty}^T x_j$   
the  $i$ th example will have much  
higher weight.

$\Rightarrow$  eventually, the only significant  
weights are those associated with  
examples  $s$ , st

$$w_{\infty}^T x_s \leq w_{\infty}^T x_i \text{ for all } i$$

There are several of these (how many?)

$\Rightarrow w_{\infty}$  is a non-neg weighted  
comb of  $x_s$

~~fields~~ Consider

$$\hat{w} = \frac{w_{\infty}}{(\min_i w_{\infty}^T x_i)}$$

$\leftarrow$  scaled so  
it isn't  
infinite

We have

$$\hat{w} = \sum_i \alpha_i x_i$$

AND either  
 $(\alpha_i \geq 0, \hat{w}^T x_i = 1)$

OR  
 $(\alpha_i = 0, \hat{w}^T x_i > 1)$

BUT These are KKT for an SVM  
( $\alpha_i > 0$  for support vectors)

- This should be unexpected — our choice of minimizer has profoundly affected the choice of solution
- Why? — all solns at  $\infty$   
— more than one. (Check!)