

Stochastic gradient descent

Many learning problems look like

$$\text{arg min}_w \quad L(w, x_i, y_i) + \lambda \|w\|_?$$

\uparrow
loss on examples

\uparrow some regularization
 typically 2-norm
 or 1-norm.

eg. SVM.

loss is $\sum_i \text{hinge loss}(w, \text{example } i)$

eg. logistic regression

loss is $\sum_i \text{exp loss}$.

In these cases, two important obs.

1) We don't so much care about

$E = \sum_{i \in \text{examples}} L(w, \text{example}_i)$ as about

for empirical loss

$\sum_{i \in \text{everything}} L(w, \text{example}_i) = J$ for true loss

this means that the min of the optimization problem w^* , may not be best w (which is w^{opt})

$w^{\text{opt}} = \text{argmin } J$

$w^* = \text{argmin } E$

error incurred by exact optimization

⇒ This means inexact opt might be OK.

2) Exact optimization could be
huge expensive

(3)

(millions of examples; millions of feats)

Idea: Pick example (feature) at
random, compute gradient for that
alone, take small step.

STOCHASTIC GRADIENT DESCENT

• How far?

step lengths S_i should have property

$$\bullet S_N \rightarrow 0, \quad N \rightarrow \infty$$

$$\bullet \sum_{i=1}^R S_i \rightarrow \infty \quad R \rightarrow \infty$$

(ie. path ~~can~~ ^{must} be infinitely long, but
steps must become infinitely small)

eg. Logistic regression with 2 Norm (4)

$$\min_w \sum_i \left[\mathcal{L}_e(w, y, x) + \frac{\lambda}{N} \|w\|_2^2 \right] = E + R$$

for each step choose 1 example, UAR. Say k 'th

$$g^{(n)} = \left[\begin{array}{c} \frac{x_k}{1 + e^{w^{(n)T} x_k}} - \left(\frac{y_k + 1}{2} \right) x_k \\ + \frac{2\lambda}{N} w^{(n)} \end{array} \right]$$

$$w^{(n+1)} = w^{(n)} - \delta_i \cdot g^{(n)}$$

Notice, because we choose UAR,

$$E(g^{(n)}) = \nabla_w [E + R]$$

typically, step lengths look like

$$\delta_i = \frac{1}{i}$$

(this satisfies constraints).

We get (with work)

$$\frac{1}{(w_t - w^*)^2}$$

grows linearly in t

i.e. early steps really help, late steps
don't do much.

Notice: algorithm is on-line

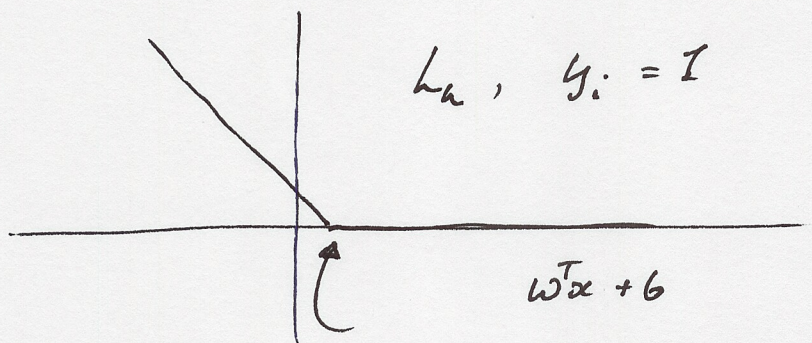
What about SVM?

(6)

$$\sum_i \left[L_h(w, x_i, y_i) + \frac{\lambda}{N} w^T w \right]$$

which isn't differentiable.

$$L_h = \max(0, 1 - y_i(\omega^T x_i + b))$$



however, this is convex.

- We invoke the sub-gradient

consider the graph of a function

$$(x_1, x_2, x_3, \dots, x_n, f(x_1, \dots, x_n))$$

this is a surface.

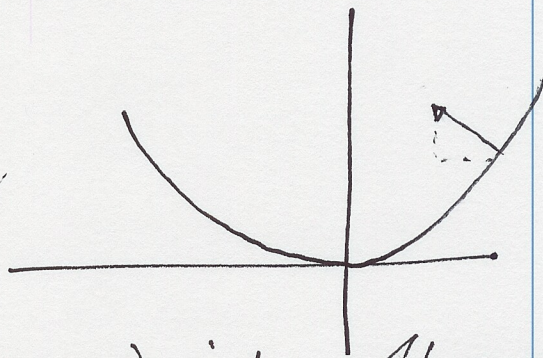
• assume ~~the~~ f is differentiable

— surface has a normal

$$N = k \left(-\frac{\partial f}{\partial x_1}, -\frac{\partial f}{\partial x_2}, \dots, +1 \right)$$

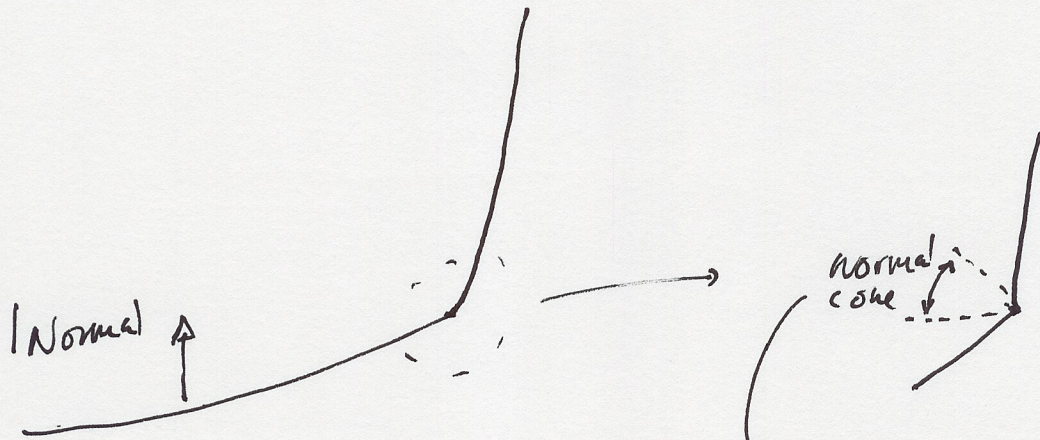
$$= k (-\nabla f, +1)$$

eg curve in the plane
is graph of fn of 1 var



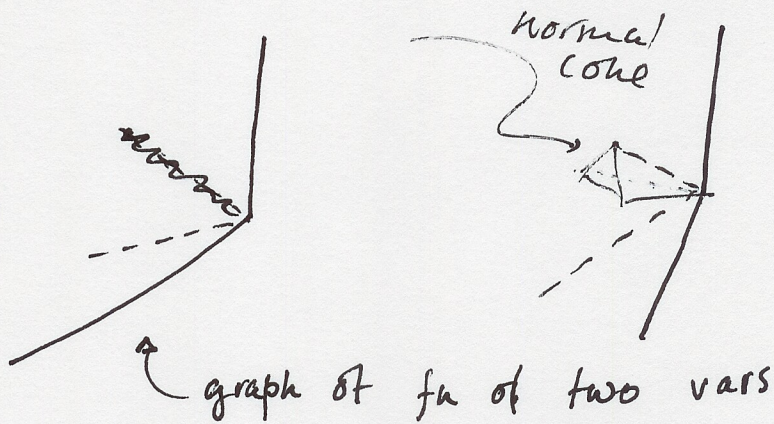
Notice we can read gradient off normal

Now, if ~~curve~~ function is not diff, we ⑧
 can come up w/ normal cone



all of these vectors
 are normals at
 this pt.

(Notice, f needs to be convex so we know what inferior is for this construction)



⑨

- read off gradient corresp to any element of normal cone
- this is subgradient ($\tilde{\nabla}$)
- moving backward along subgradient will guarantee descent for small enough step
- subgradient of diff. fn = gradient

SVM.

$$\tilde{\nabla}_w L_h(w) = \begin{cases} 0 & \text{if } L_h(w) = 0 \\ -y_i x_i & \text{otherwise} \end{cases}$$

So ~~subgr~~ Stochastic Subgradient Descent ⁽¹⁰⁾
is

- choose k 'th example at random
- if right, $w^{(n+1)} = w^{(n)} \left(1 - \frac{\lambda \delta_i}{2N}\right)$
- if wrong, $w^{(n+1)} = w^{(n)} - \delta_i \begin{bmatrix} -y_i x_i \\ y_i \end{bmatrix} + \frac{\lambda w^{(n)}}{2N}$

This is amazingly effective.

L_1 regularization

(11)

$$\|w\|_1 = \sum_i |w_i|$$

$$\|w\|_2^2 = \sum_i w_i^2$$

Notice: in l_2 norm, small w_i have little effect hence, not much advantage in driving to zero; in l_1 norm, effect of small w_i is substantial; they tend to go to zero, resulting in sparsity (desirable)

eg h_1 neg logistic regression

$$\min_w \sum_i h_e(w, y_i, x_i) + \lambda \|w\|_1$$

- not differentiable
- subgradient unlikely to enforce sparsity

alternative (equivalent)

$$\min_w \sum_i h_e(w, y_i, x_i) + \sum_e h_e$$

$$-h_k \leq w_k \leq h_k$$

(Notice this is a box problem;
 Solve with interior point method;
 excellent version due to
 Koh et al)