

Thursday, Oct. 10

Our goal: Let  $c(w; x, y)$  be our objective for one sample

$$\min_{w \in \mathbb{R}^d} E[c(w; X, Y)]$$

$= \int_P(x, y) c(w; x, y) dx dy$

be our objective for one sample.

Empirical estimate of true obj.:  $c(w) = \frac{1}{n} \sum c_i(w)$

$$c_i(w) = c(w; x_i, y_i)$$

Batch grad. descent

$$\nabla c(w) = \nabla \left( \frac{1}{n} \sum c_i(w) \right)$$
$$= \frac{1}{n} \sum \nabla c_i(w)$$

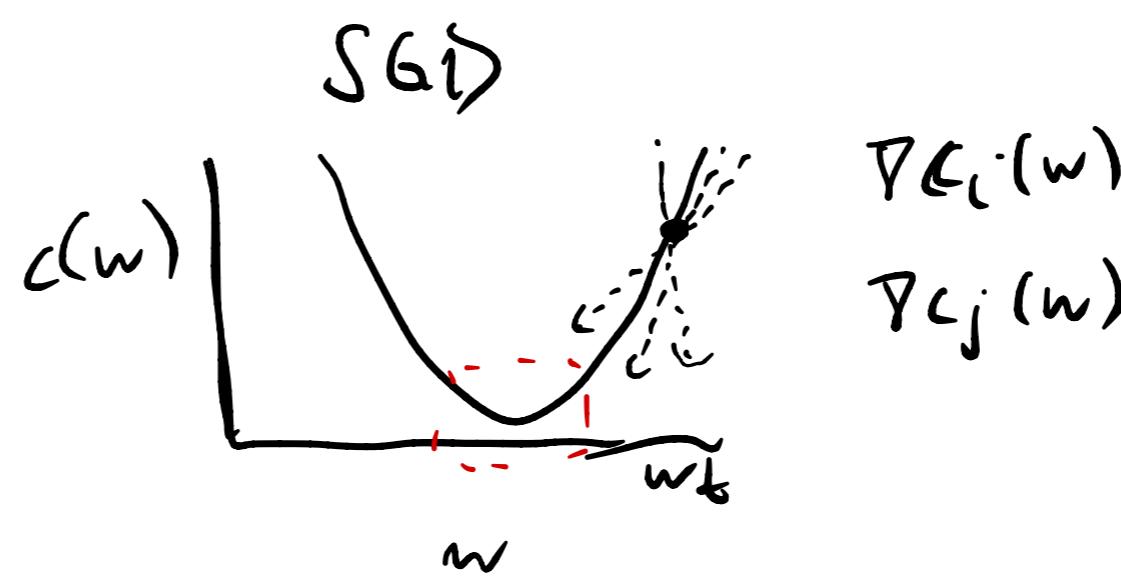
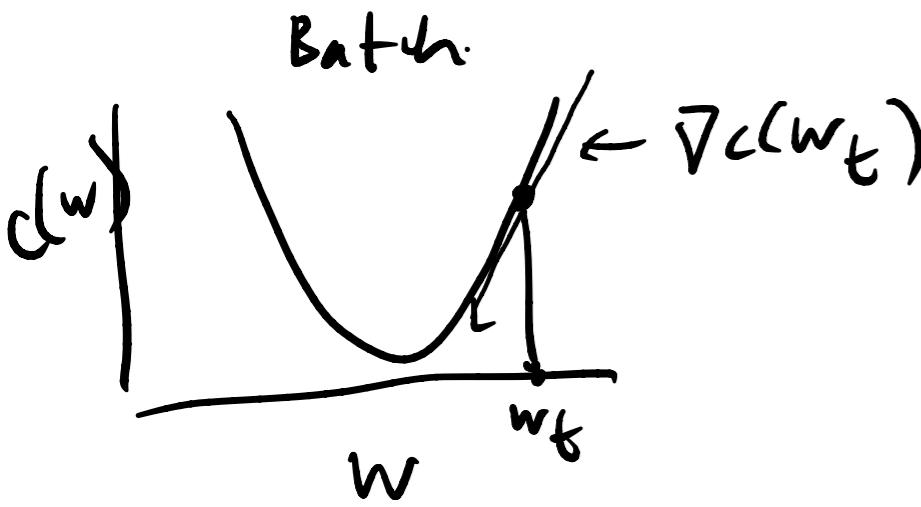
$O(dn)$  per iteration

Stochastic Grad Descent

$$\nabla c(w) \approx \nabla c_i(w)$$

for random  $i \in \{1, \dots, n\}$ .

old) per iteration



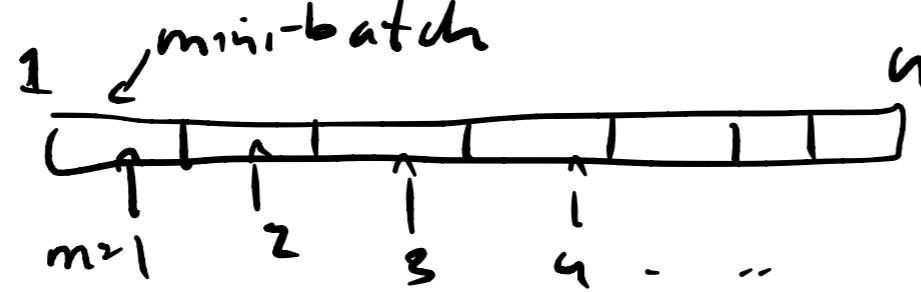
$i$  is randomly sample from  $1, \dots, n$

$$E[\nabla c_i(w)] = \sum_{i=1}^n p(i) \nabla c_i(w) = \frac{1}{n} \sum_{i=1}^n \nabla c_i(w) = \nabla c(w)$$

In between : mini-batch SGD  $\nabla c(w) \approx \frac{1}{b} \sum_{j=1}^b \nabla c_{i_j}(w)$

Cost per iteration :  $O(db)$   $i_j$  uniform random in  $\{1, \dots, n\}$ .

Extremes :  $b = n$  or  $b = 1$  (very noisy)  
 Usually better to be more moderate, (e.g.  $b = 16$ )

Implementation :  Epoch = one processing of the data

mini-batch SGD:

for  $k = 1$ : # of epochs (e.g. 100 epochs)

→ Shuffle data  $\{c_1, \dots, c_n\}$

for  $m = 1 : n/b$   $b$  is mini-batch size.

$$w = w - \eta \text{ (gradient-on-mini-batch)}$$

Batch-grad. = First-order grad. descent

Goal for  $\eta$ : Find  $\eta = \underset{\eta \geq 0}{\operatorname{argmin}} C(w_t - \eta \nabla C(w_t))$

Approximate line search, start at  $\eta = 1.0$ ,  $0 < \gamma < 1$   
reduce by  $\gamma$  if overshoot,  $\eta \in \mathbb{Z}_\eta$

$$w_t = 2.0, g = 5.0, \eta = 1.0, \gamma = 0.5$$

$$\hat{w} = w_t - 1.0 \cdot g = -3.0 \quad \hat{C}(w)$$

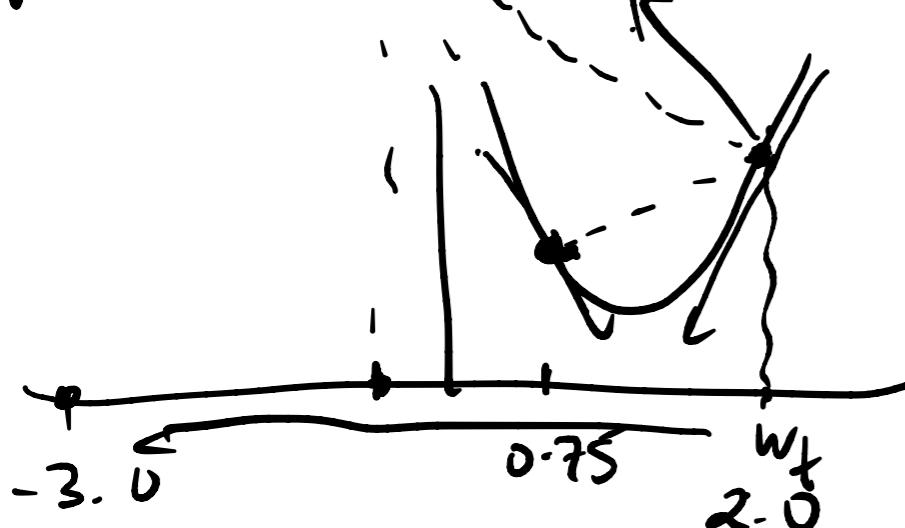
$$\hat{w} = w_t - 0.5 \cdot g = -0.5 \quad \hat{C}(w)$$

$$\hat{w} = w_t - 0.25 \cdot g = 0.75 \quad \hat{C}(w)$$

Accept

$$w_{t+1} = w_t - 0.25g. \quad X$$

Should we use this for SGD:  $\underset{\eta \geq 0}{\operatorname{argmin}} c_i(w_t - \eta \nabla c_i(w))$   
not our goal



Heuristic approach for  $\eta$  for SGD:

e.g. Adagrad, use vector of stepsizes.  $\vec{\eta} \in \mathbb{R}^d$ .

$$b_t = b_t + \nabla c_i(w_t)^2 \leftarrow \text{element-wise squaring}$$
$$\eta_t = \frac{\eta_0}{\sqrt{b_t}} \leftarrow \text{element-wise}$$
$$g^2 = \begin{bmatrix} g_1^2 \\ g_2^2 \\ \vdots \\ g_d^2 \end{bmatrix}$$
$$\eta = \begin{bmatrix} 1.0 \\ 0.01 \end{bmatrix}$$

statistics on magnitudes of gradients

what is the batch GD and SGD updates for  $l_2$  regularized linear regression?

$$c(w) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (x_i^\top w - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2 = \frac{1}{n} \sum_{i=1}^n c_i(w)$$

$$c_i(w) = \frac{1}{2} (x_i^\top w - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2 \rightarrow \frac{1}{n} \sum_{i=1}^n \frac{\lambda}{2} \|w\|_2^2$$

Update: for random  $k$

$$w_{t+1} = w_t - \eta_t [(x_k^\top w - y_k) x_k + \lambda w]$$

gradient  
 $= \lambda w$

Batch update for  $l_1$ , i.e.  $\|w\|_1 = \sum_{j=1}^d |w_j|$

$$\min_w \frac{1}{2n} \|Xw - y\|_2^2 + \lambda \|w\|_1,$$

compute gradient

$$\frac{\partial \|w\|_1}{\partial w_j} = \frac{\partial \sum_{k=1}^d |w_k|}{\partial w_j} = \frac{\partial |w_j|}{\partial w_j}$$

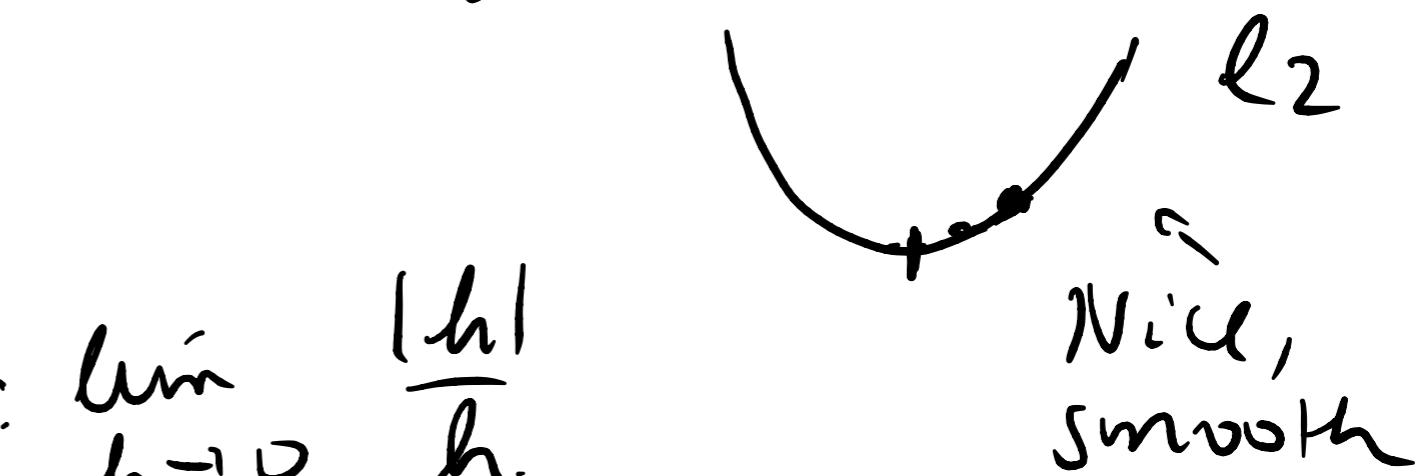
$$\frac{df(w)}{dw} = \lim_{h \rightarrow 0} \frac{f(w+h) - f(w)}{h}$$

$h = -10^{-7}v$   
from left : -1

$h$  from right = +1

$$h = +10^{-7}v$$

limit at  $w=0$ :  $\lim_{h \rightarrow 0} \frac{|h|}{h}$ .



Alternative : descend and project

$$c(w) = \frac{1}{2n} \|Xw - y\|_2^2$$

(smooth)

$$R(w) = \lambda \|w\|_1$$

(non-smooth)

$$\text{Descend: } \tilde{w}^{(t+1)} = w^{(t)} - \eta \nabla c(w^{(t)}) \quad \begin{matrix} \text{makes} \\ \text{error} \\ \text{smaller} \end{matrix}$$

Project:  
 (back onto  
 space where  
 $R(w)$  is  
 reasonably small)

$$w^{(t+1)} = \text{prox}(\tilde{w}^{(t+1)})$$

$$= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \frac{\eta \lambda \|w\|_1 + \frac{1}{2} \|\tilde{w}^{(t+1)} - w\|_2^2}{\eta R(w)}$$

$$w_j^{(t+1)} = \begin{cases} \tilde{w}_j^{(t+1)} - \eta \lambda & \text{if } \tilde{w}_j^{(t+1)} > \eta \lambda \\ \tilde{w}_j^{(t+1)} + \eta \lambda & \text{if } \tilde{w}_j^{(t+1)} < -\eta \lambda \\ 0 & \text{else} \end{cases}$$

big +  
big -