

Wednesday, February 3, 2021 10:01 AM

Lecture 5

Stochastic Gradient Descent

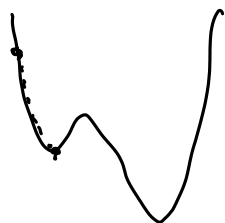
- Vanilla GD is not popular b/c

- full gradient is expensive to compute

$$\mathcal{L} = \sum_{i \in \text{Data}} L(f(x_i; w), y_i)$$

$$\text{compute: } \frac{\partial \mathcal{L}}{\partial w} = \sum_{i \in \text{Data}} \frac{\partial L_i}{\partial w}$$

• have to sum over full dataset for every weight update.



- can easily get stuck in a local minimum.

Wednesday, February 3, 2021 10:31 AM

SGD solves both problems in same way.

IDEA: allow gradients to be noisy

How?

Compute gradients over subset of data.

- Pure SGD: one datapoint

- minibatch SGD: chunk of data.

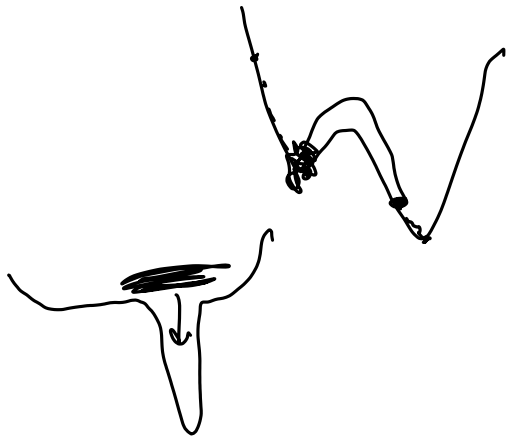
$$\nabla_w L_i$$

$$\downarrow$$

$$\frac{\partial L_i}{\partial w}$$

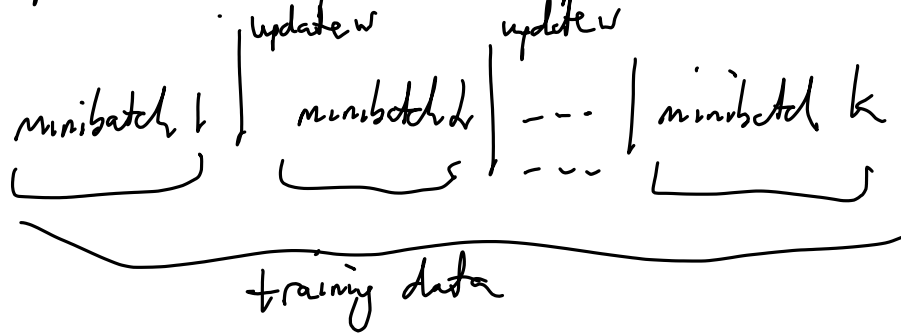
$$\downarrow$$

$$\sum_{i=i_{start}}^{i_{start}+N_{batch}} \frac{\partial L_i}{\partial w}$$



preferred in part b/c computational efficiency

(parallelization over minibatches)

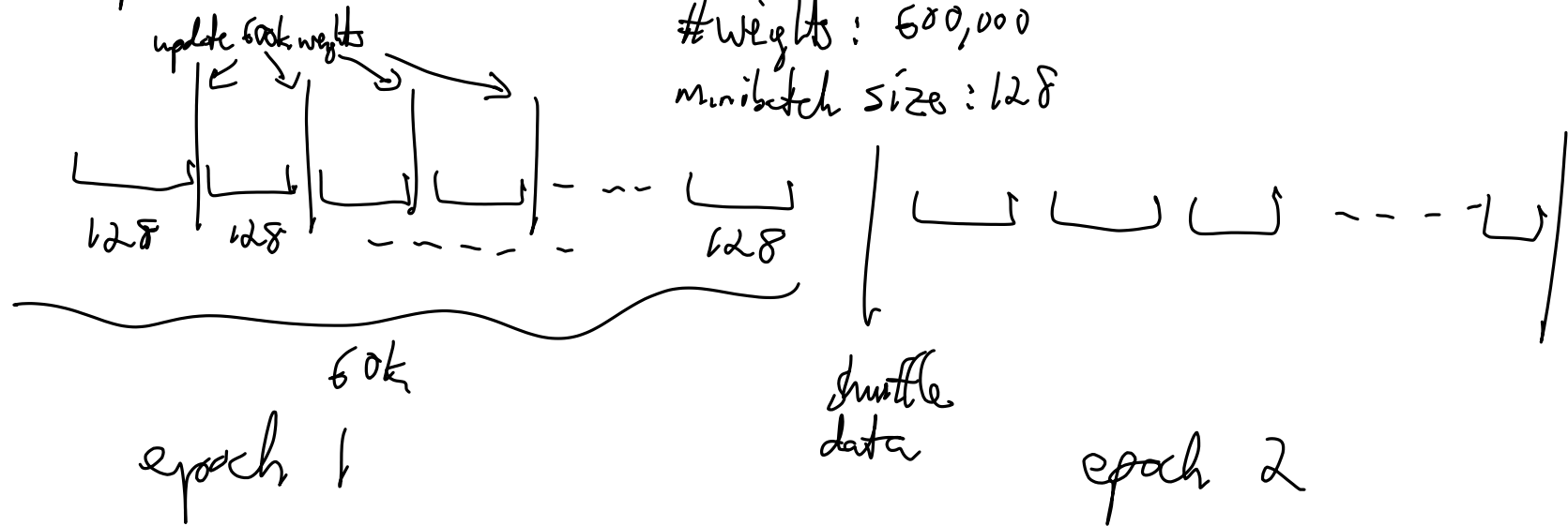


"epoch"
after each epoch,
shuffle data!

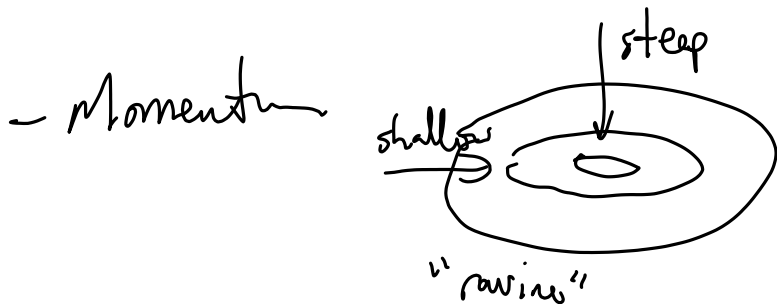
Wednesday, February 3, 2021 10:49 AM

Example: MNIST DNN, Data: 60,000

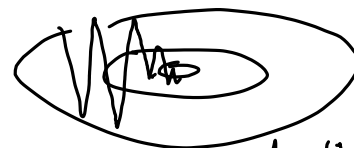
#weights: 600,000
minibatch size: 128



Variations on (S)GD



GD can get in trouble here



idea of momentum: give gradients a push based on their previous direction ^{an acceleration}

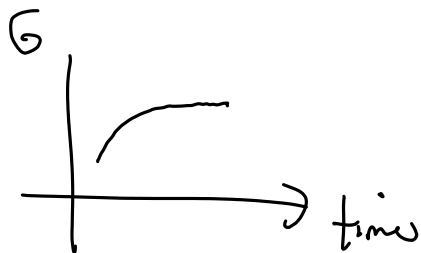
$$GD \text{ w/ momentum: } v_t = \delta v_{t-1} + \gamma \nabla_w L$$

$$w \rightarrow w - v_t$$

- Adaptive learning rates:

Adagrad: $w \rightarrow w - \frac{\gamma}{\sqrt{G}} \nabla_w L$

$$G = \sum_{\text{time steps}} (\nabla_w L)^2$$



large gradients \rightarrow smaller learning rate
 accumulating gradients \rightarrow decrease γ l.r.

★ Adadelta: Adagrad w/ finite time window

$$G = \sum_{\text{past time window}} (\nabla_w L)^2$$

★ Adam (Adaptive Moment estimator): "Adadelta w/ ^{adaptive} moment"

Wednesday, February 3, 2021 11:27 AM

Backpropagation: how gradients are computed efficiently.

FF NN's have a recursive structure, nested

$$f(x) = A(w_h \cdot A(w_{h-1} \cdot A(\dots A(w_1 x))))$$

$$\mathcal{L} = \sum_{i \in \text{data}} L(f(x_i), y)$$

$$\frac{\partial \mathcal{L}}{\partial w_k} \rightarrow \frac{\partial \mathcal{L}}{\partial w_k} \rightarrow \frac{\partial f(x_i; w)}{\partial w_k}$$

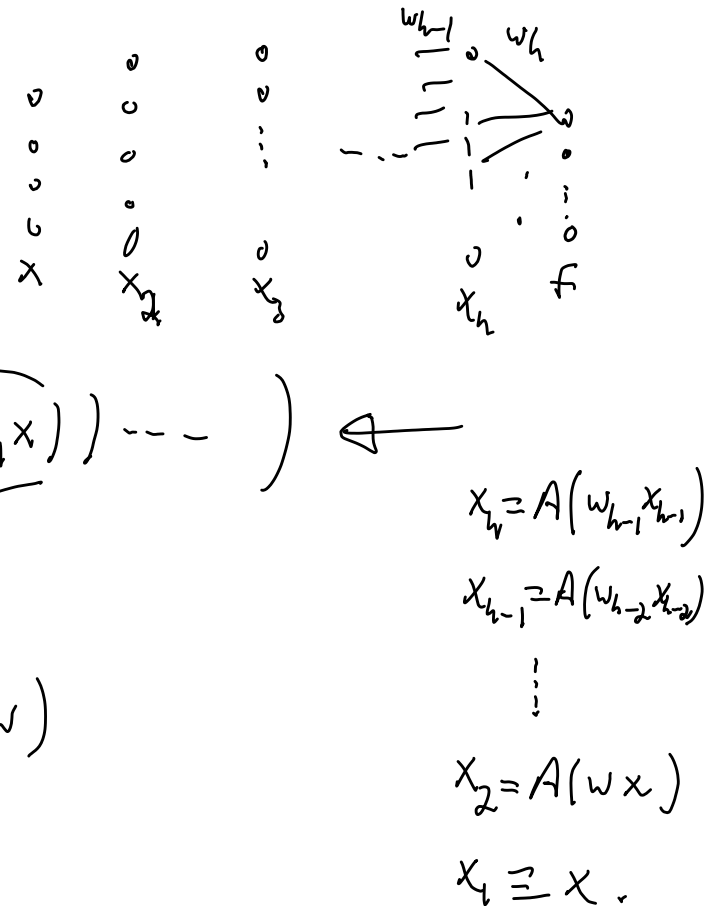
Let's work out a few examples

$$\frac{\partial f}{\partial u_h} = A'(w_h x_h) \cdot x_h \equiv G_h x_h$$

$$\frac{\partial f}{\partial w_{h-1}} = A'(w_h x_h) \cdot w_h A'(w_{h-1} x_{h-1}) x_{h-1} \equiv G_{h-1} x_{h-1}$$

$$\frac{\partial f}{\partial w_{h-2}} = A'(w_h x_h) \cdot v_h A'(w_{h-1} x_{h-1}) w_{h-1} A'(w_{h-2} x_{h-2}) x_{h-2} \equiv G_{h-2} x_{h-2}$$

⋮
 — compute gradients recursively



Backpropagation:

$$G_{k-1} = G_k W_k A (W_{k-1} X_{k-1})$$

- matrix mult. - expensive O^3

- enabled "end-to-end" trans of NNs.

GPU's are key for deep learning.