

# CS391L: Machine Learning WB

## Stochastic Optimization

Inderjit S. Dhillon  
UT Austin

March 30, 2025

# Stochastic Gradient Descent

# Large-scale Problems

- Machine learning: usually minimize the training loss

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \mathbf{x}_n, y_n) \right\} := f(\mathbf{w}) \text{ (linear model)}$$

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(h_{\mathbf{w}}(\mathbf{x}_n), y_n) \right\} := f(\mathbf{w}) \text{ (general hypothesis)}$$

$\ell$ : loss function (e.g.,  $\ell(a, b) = (a - b)^2$ )

- Gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \underbrace{\nabla f(\mathbf{w})}_{\text{main computation}}$$

# Large-scale Problems

- Machine learning: usually minimize the training loss

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \mathbf{x}_n, y_n) \right\} := f(\mathbf{w}) \text{ (linear model)}$$

$$\min_{\mathbf{w}} \left\{ \frac{1}{N} \sum_{n=1}^N \ell(h_{\mathbf{w}}(\mathbf{x}_n), y_n) \right\} := f(\mathbf{w}) \text{ (general hypothesis)}$$

$\ell$ : loss function (e.g.,  $\ell(a, b) = (a - b)^2$ )

- Gradient descent:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \underbrace{\nabla f(\mathbf{w})}_{\text{main computation}}$$

- In general,  $f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N f_n(\mathbf{w})$ ,  
each  $f_n(\mathbf{w})$  only depends on  $(\mathbf{x}_n, y_n)$

# Stochastic gradient

- Gradient:

$$\nabla f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w})$$

- Each gradient computation needs to go through **all training samples**  
slow when millions of samples
- Faster way to compute “**approximate gradient**”?

# Stochastic gradient

- Gradient:

$$\nabla f(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w})$$

- Each gradient computation needs to go through **all training samples** slow when millions of samples
- Faster way to compute “**approximate gradient**”?
- Use **stochastic sampling**:
  - Sample a small subset  $B \subseteq \{1, \dots, N\}$
  - Estimated gradient

$$\nabla f(\mathbf{w}) \approx \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})$$

$|B|$ : batch size

# Stochastic gradient descent

## Stochastic Gradient Descent (SGD)

- Input: training data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Initialize  $\mathbf{w}$  (zero or random)
- For  $t = 1, 2, \dots$ 
  - Sample a **small batch**  $B \subseteq \{1, \dots, N\}$
  - Update parameter

$$\mathbf{w} \leftarrow \mathbf{w} - \eta^t \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})$$

**Extreme case:**  $|B| = 1 \Rightarrow$  **Sample one training data at a time**

# Logistic Regression by SGD

- Logistic regression:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \underbrace{\log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})}_{f_n(\mathbf{w})}$$

## SGD for Logistic Regression

- Input: training data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Initialize  $\mathbf{w}$  (zero or random)
- For  $t = 1, 2, \dots$ 
  - Sample a batch  $B \subseteq \{1, \dots, N\}$
  - Update parameter

$$\mathbf{w} \leftarrow \mathbf{w} - \eta^t \frac{1}{|B|} \sum_{i \in B} \underbrace{\frac{-y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}}}_{\nabla f_n(\mathbf{w})}$$



# Why SGD works?

- Stochastic gradient is an **unbiased estimator** of full gradient:

$$\begin{aligned} E\left[\frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w})\right] &= \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}) \\ &= \nabla f(\mathbf{w}) \end{aligned}$$

- Each iteration updated by

gradient + **zero-mean noise**

# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?

# Stochastic gradient descent

- In gradient descent,  $\eta$  (step size) is a fixed constant
- Can we use fixed step size for SGD?
- SGD with fixed step size **cannot converge to global/local minimizers**
- If  $\mathbf{w}^*$  is the minimizer,  $\nabla f(\mathbf{w}^*) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}^*) = 0$ ,

$$\text{but } \frac{1}{|B|} \sum_{n \in B} \nabla f_n(\mathbf{w}^*) \neq 0 \quad \text{if } B \text{ is a subset}$$

(Even if we got minimizer, SGD will **move away** from it)

# Stochastic gradient descent, step size

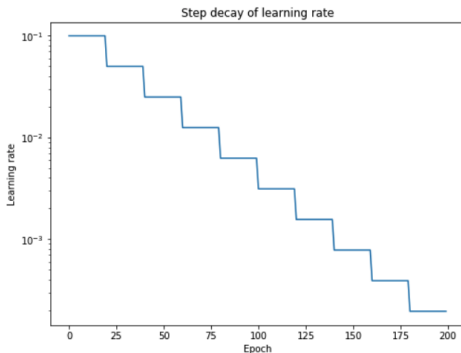
- To make SGD converge:

Step size should decrease to 0

$$\eta^t \rightarrow 0$$

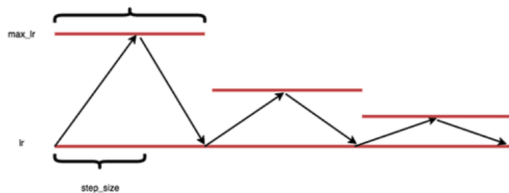
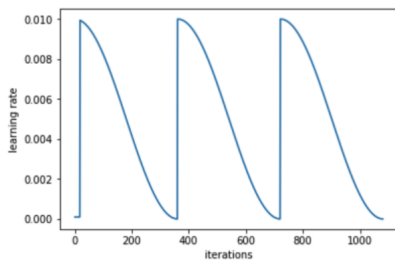
Usually with polynomial rate:  $\eta^t \approx t^{-a}$  with constant  $a > 0$

- Step decay of learning rate



# Learning rate scheduling

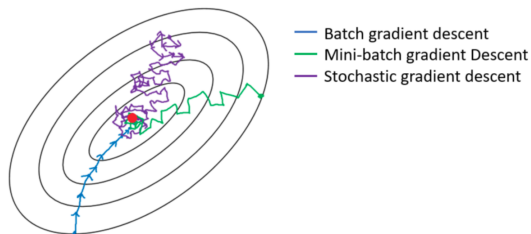
- Other cyclic learning rate scheduling strategies



# Stochastic gradient descent vs Gradient descent

## Stochastic gradient descent:

- pros:
  - cheaper computation per iteration
  - faster convergence in the beginning
- cons:
  - less stable, slower final convergence
  - hard to tune step size



(Figure from <https://medium.com/@ImadPhd/gradient-descent-algorithm-and-its-variants-10f652806a3>)

# Momentum

- Gradient descent: only using current gradient (local information)
- Momentum: use **previous gradient information**
- The momentum update rule:

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \nabla f(\mathbf{w}_t)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \mathbf{v}_t$$

$\beta \in [0, 1)$ : discount factors,  $\alpha$ : step size

- Equivalent to using moving average of gradient:

$$\mathbf{v}_t = (1 - \beta) \nabla f(\mathbf{w}_t) + \beta(1 - \beta) \nabla f(\mathbf{w}_{t-1}) + \beta^2(1 - \beta) \nabla f(\mathbf{w}_{t-2}) + \dots$$

# Momentum gradient descent

## Momentum gradient descent

- Initialize  $\mathbf{w}_0, \mathbf{v}_0 = \mathbf{0}$
- For  $t = 1, 2, \dots$ 
  - Compute  $\mathbf{v}_t \leftarrow \beta \mathbf{v}_{t-1} + (1 - \beta) \nabla f(\mathbf{w}_t)$
  - Update  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \alpha \mathbf{v}_t$

$\alpha$ : learning rate

$\beta$ : discount factor ( $\beta = 0$  means no momentum)



# Momentum stochastic gradient descent

Optimizing  $f(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{w})$

## Momentum stochastic gradient descent

- Initialize  $\mathbf{w}_0, \mathbf{v}_0 = \mathbf{0}$
- For  $t = 1, 2, \dots$ 
  - Sample an  $i \in \{1, \dots, N\}$
  - Compute  $\mathbf{v}_t \leftarrow \beta \mathbf{v}_{t-1} + (1 - \beta) \nabla f_i(\mathbf{w}_t)$
  - Update  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \alpha \mathbf{v}_t$

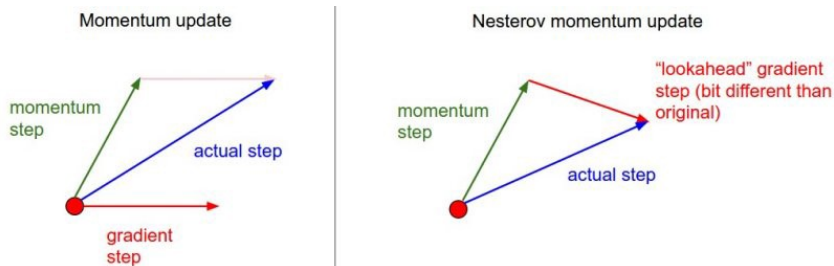
$\alpha$ : learning rate

$\beta$ : discount factor ( $\beta = 0$  means no momentum)

# Nesterov accelerated gradient

- Using the “look-ahead” gradient

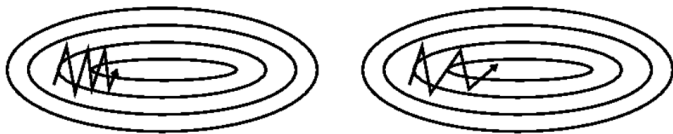
$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + \alpha \nabla f(\mathbf{w}_t - \beta \mathbf{v}_{t-1})$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{v}_t$$



(Figure from <https://towardsdatascience.com>)

# Why momentum works?

- Reduce variance of gradient estimator for SGD
- Even for gradient descent, it's able to speed up convergence in some cases:



Left—SGD without momentum, right—SGD with momentum. (Source: [Genevieve B. Orr](#))

# Adagrad: Adaptive updates (2010)

- SGD update: same step size for all variables
- Adaptive algorithms: each dimension can have a different step size

## Adagrad

- Initialize  $\mathbf{w}_0$
- For  $t = 1, 2, \dots$ 
  - Sample an  $i \in \{1, \dots, N\}$
  - Compute  $\mathbf{g}^t \leftarrow \nabla f_i(\mathbf{w}_t)$
  - $G_j^t \leftarrow G_j^{t-1} + (g_j^t)^2$  for all  $j = 1, \dots, d$
  - Update  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta}{\sqrt{G_i^t + \epsilon}} \mathbf{g}_i^t$

$\eta$ : step size (constant)

$\epsilon$ : small constant to avoid division by 0

# Adagrad

- For each dimension  $i$ , we have observed  $T$  samples  $g_i^1, \dots, g_i^t$
- Standard deviation of  $g_i$ :

$$\sqrt{\frac{\sum_{t'} (g_i^{t'})^2}{t}} = \sqrt{\frac{G_i^t}{t}}$$

- Assume step size is  $\eta/\sqrt{t}$ , then the update becomes

$$w_i^{t+1} \leftarrow w_i^t - \frac{\eta}{\sqrt{t}} \frac{\sqrt{t}}{\sqrt{G_i^t}} g_i^t$$

# Adam: Momentum + Adaptive updates (2015)

## Adam

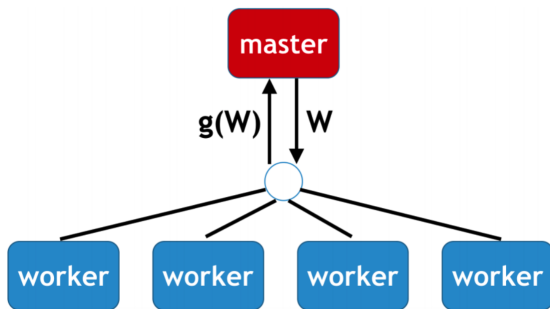
- Initialize  $\mathbf{w}_0$ ,  $\mathbf{m}_0 = \mathbf{0}$ ,  $\mathbf{v}_0 = \mathbf{0}$ ,
- For  $t = 1, 2, \dots$ 
  - Sample an  $i \in \{1, \dots, N\}$
  - Compute  $\mathbf{g}_t \leftarrow \nabla f_i(\mathbf{w}_t)$
  - $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$
  - $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$
  - $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$
  - $\hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$
  - Update  $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \alpha \cdot \hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon)$

# Batch size selection

- Larger batch size  $\Rightarrow$  more computation per update, less noise
- Usually, choose batch size large enough that
  - Fits in (GPU) memory
  - Can fully utilize the computation resource
- For example, 512 batch size for ImageNet training on a standard GPU.

# Large batch training

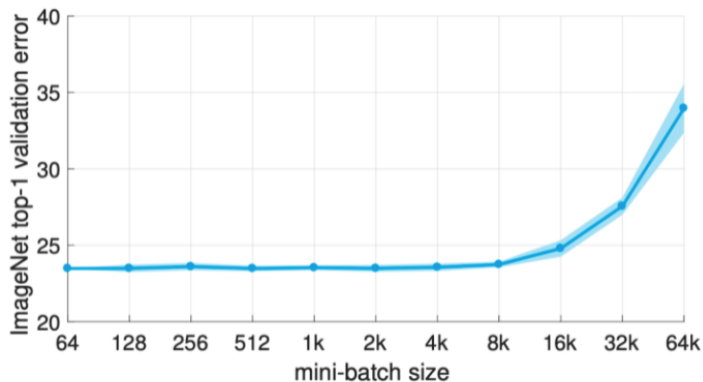
- What if we want to train a model with **hundreds or thousands of GPUs**
- Data parallel distributed computing:  
Increasing the batch size linearly with number of devices.  
⇒ **Large batch training**





# Problem of large batch training

- Tend to converge to models with lower test accuracy when using very large batch size



(Goyal et al., 2017) Training Imagenet in 1 hour

# Sharp vs wide local minimum

- Large-batch SGD/Adam:
  - Usually converge to a **sharp** local minimum (not enough inherent noise in SGD)
  - Harder to generalize to test data

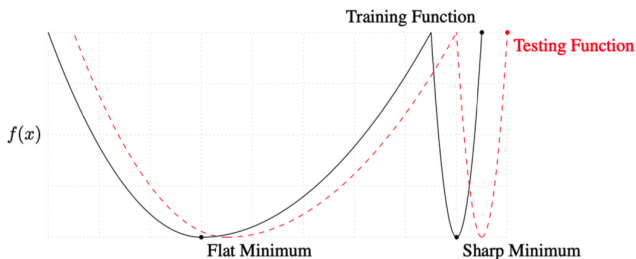


Figure from (Keskar et al., 2017)

## A simple but practical solution: learning rate scaling

- $\text{Var}[\frac{1}{|B|} \sum_{n \in B} \mathbf{g}_n] \approx O(\frac{1}{|B|})$
- Batch size  $\uparrow$ , learning rate  $\uparrow$ 
  - LR scaling: LR as  $O(\sqrt{|B|})$  or  $O(|B|)$

## A simple but practical solution: learning rate scaling

- $\text{Var}[\frac{1}{|B|} \sum_{n \in B} \mathbf{g}_n] \approx O(\frac{1}{|B|})$
- Batch size  $\uparrow$ , learning rate  $\uparrow$ 
  - LR scaling: LR as  $O(\sqrt{|B|})$  or  $O(|B|)$
- However, LR has to be bounded for convergence (as for GD)

$$LR \leq O(1/L) \quad L : \text{Lipchitz constant}$$

- Can't increase LR without limit

# Non-uniform updates between different layers

- Some layer becomes the bottleneck of LR

Layers	$\ w\ _2$	$\ \nabla w\ _2$	$\ w\ _2/\ \nabla w\ _2$
fc8.0	20.24	0.078445	258
fc8.1	0.316	0.006147	51
fc7.0	20.48	0.110949	184
fc7.1	6.400	0.004939	1296
fc6.0	30.72	0.097996	314
fc6.1	6.400	0.001734	<b>3690</b>
conv5.0	6.644	0.034447	193
conv5.1	0.160	0.000961	166
conv4.0	8.149	0.039939	204
conv4.1	0.196	0.000486	403
conv3.0	9.404	0.049182	191
conv3.1	0.196	0.000511	384
conv2.0	5.545	0.057997	96
conv2.1	0.160	0.000649	247
conv1.0	1.866	0.071503	26
conv1.1	0.098	0.004909	<b>20</b>

## Another solution

- Use **Layer-wise Adaptive LR Scaling**

$$\mathbf{w}^{(i)} \leftarrow \mathbf{w}^{(i)} - \eta \frac{\|\mathbf{w}^{(i)}\|}{\|\mathbf{g}^{(i)}\|} \mathbf{g}^{(i)},$$

where  $(\cdot)^{(i)}$  means the  $i$ -th layer of neural network

- LARS:  $\mathbf{g}$  is the stochastic gradient

LAMB:  $\mathbf{g}$  is the Adam update

(“Large Batch Optimization for Deep Learning: Training BERT in 76 minutes, ICLR 2020”)

- Google and Nvidia both use LAMB to train BERT within 1 minute (4096 TPU / 2048 GPU)

# Conclusions

- Stochastic gradient descent
- Variants of SGD used in neural network training.

Questions?