

CS260: Machine Learning Algorithms

Lecture 4: Stochastic Gradient Descent

Cho-Jui Hsieh
UCLA

Jan 16, 2019

Large-scale Problems

- Machine learning: usually minimizing the training loss

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

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

ℓ : 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 minimizing the training loss

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

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

ℓ : 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

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

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:

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

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, η (step size) is a fixed constant
- Can we use fixed step size for SGD?

Stochastic gradient descent

- In gradient descent, η (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**

Stochastic gradient descent

- In gradient descent, η (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, $\nabla f(\mathbf{w}^*) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}^*) = 0$,

Stochastic gradient descent

- In gradient descent, η (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, $\nabla f(\mathbf{w}^*) = \frac{1}{N} \sum_{n=1}^N \nabla f_n(\mathbf{w}^*) = 0$,

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

Stochastic gradient descent

- In gradient descent, η (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

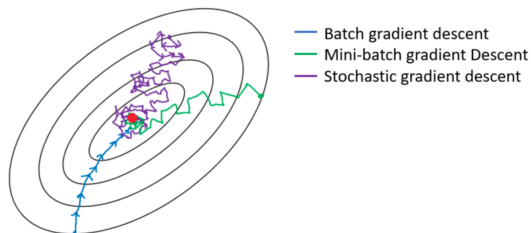
$$\eta^t \rightarrow 0$$

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

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)

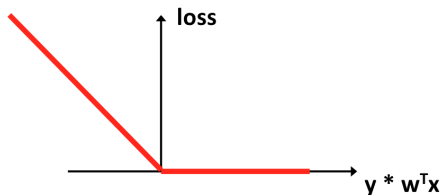
Revisit perceptron Learning Algorithm

- Given a classification data $\{\mathbf{x}_n, y_n\}_{n=1}^N$
- Learning a linear model:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \mathbf{x}_n, y_n)$$

- Consider the loss:

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$



What's the gradient?

Revisit perceptron Learning Algorithm

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$

Consider two cases:

- Case I: $y_n \mathbf{w}^T \mathbf{x}_n > 0$ (prediction **correct**)
 - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
 - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$

Revisit perceptron Learning Algorithm

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$

Consider two cases:

- Case I: $y_n \mathbf{w}^T \mathbf{x}_n > 0$ (prediction **correct**)
 - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
 - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
- Case II: $y_n \mathbf{w}^T \mathbf{x}_n < 0$ (prediction **wrong**)
 - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{w}^T \mathbf{x}_n$
 - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{x}_n$

Revisit perceptron Learning Algorithm

$$\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = \max(0, -y_n \mathbf{w}^T \mathbf{x}_n)$$

Consider two cases:

- Case I: $y_n \mathbf{w}^T \mathbf{x}_n > 0$ (prediction **correct**)
 - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
 - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = 0$
- Case II: $y_n \mathbf{w}^T \mathbf{x}_n < 0$ (prediction **wrong**)
 - $\ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{w}^T \mathbf{x}_n$
 - $\frac{\partial}{\partial \mathbf{w}} \ell(\mathbf{w}^T \mathbf{x}_n, y_n) = -y_n \mathbf{x}_n$

SGD update rule: Sample an index n

$$\mathbf{w}^{t+1} \leftarrow \begin{cases} \mathbf{w}^t & \text{if } y_n \mathbf{w}^T \mathbf{x}_n \geq 0 \text{ (predict correct)} \\ \mathbf{w}^t + \eta^t y_n \mathbf{x}_n & \text{if } y_n \mathbf{w}^T \mathbf{x}_n < 0 \text{ (predict wrong)} \end{cases}$$

Equivalent to Perceptron Learning Algorithm when $\eta^t = 1$

Momentum

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

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, α : step size

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, α : step size

- Equivalent to using moving average of gradient:

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

Momentum

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

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

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

- Equivalent to using moving average of gradient:

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

- Another equivalent form:

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

Momentum gradient descent

Momentum gradient descent

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

α : learning rate

β : 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 $\mathbf{w}_0, \mathbf{v}_0 = 0$
- 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 $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \alpha \mathbf{v}_t$

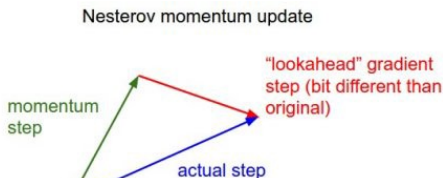
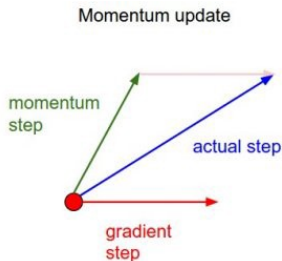
α : learning rate

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

Nesterov accelerated gradient

- Using the “look-ahead” gradient

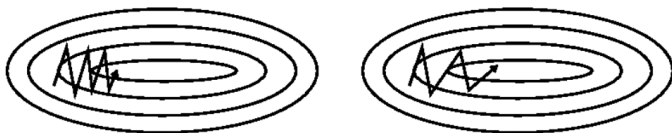
$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + \alpha \nabla f(\mathbf{w}_t - \beta \mathbf{v}_{t-1})$$
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{v}_t$$



(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: 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 $\mathbf{g}^t \leftarrow \nabla f_i(\mathbf{w}_t)$
 - $G_i^t \leftarrow G_i^{t-1} + (g_i^t)^2$
 - Update $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \frac{\eta}{\sqrt{G_i^t + \epsilon}} \mathbf{g}_i^t$

η : step size (constant)

ϵ : 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 :

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

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

$$w_i^{t+1} \leftarrow w_i^t - \frac{\eta}{\sqrt{t}} \frac{\sqrt{t}}{\sqrt{(G_i^t)^2}} g_i^t$$

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 $\mathbf{w}_t \leftarrow \mathbf{w}_t - \alpha \cdot \hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon)$

Conclusions

- Stochastic gradient descent
- Momentum & adaptive updates

Questions?