

# Estimation of the Gradient I

- Recall the function is

$$f(\boldsymbol{\theta}) = \frac{1}{2C} \boldsymbol{\theta}^T \boldsymbol{\theta} + \frac{1}{l} \sum_{i=1}^l \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i})$$

- The gradient is

$$\begin{aligned} & \frac{\boldsymbol{\theta}}{C} + \frac{1}{l} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^l \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i}) \\ &= \frac{\boldsymbol{\theta}}{C} + \frac{1}{l} \sum_{i=1}^l \nabla_{\boldsymbol{\theta}} \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i}) \end{aligned}$$

# Estimation of the Gradient II

- Going over all data is time consuming
- If data are from the same distribution

$$E(\nabla_{\theta} \xi(\theta; \mathbf{y}, Z^1)) = \frac{1}{l} \nabla_{\theta} \sum_{i=1}^l \xi(\theta; \mathbf{y}^i, Z^{1,i})$$

then we may just use a **subset**  $S$

$$\frac{\theta}{C} + \frac{1}{|S|} \nabla_{\theta} \sum_{i:i \in S} \xi(\theta; \mathbf{y}^i, Z^{1,i})$$

# Stochastic Gradient Algorithm I

- 1: Given an initial learning rate  $\eta$ .
- 2: **while do**
- 3:     Choose  $S \subset \{1, \dots, l\}$ .
- 4:     Calculate

$$\theta \leftarrow \theta - \eta \left( \frac{\theta}{C} + \frac{1}{|S|} \nabla_{\theta} \sum_{i:i \in S} \xi(\theta; \mathbf{y}^i, Z^{1,i}) \right)$$

- 5:     May adjust the learning rate  $\eta$
  - 6: **end while**
- It's known that deciding a suitable learning rate is difficult

# Stochastic Gradient Algorithm II

- Too small learning rate: very slow convergence
- Too large learning rate: the procedure may diverge

# Stochastic Gradient “Descent” I

- In comparison with gradient descent you see that we don't do line search
- Indeed we cannot. Without the full gradient, the sufficient decrease condition may never hold.

$$f(\boldsymbol{\theta} + \alpha\Delta\boldsymbol{\theta}) < f(\boldsymbol{\theta}) + \nu\nabla f(\boldsymbol{\theta})^T(\alpha\Delta\boldsymbol{\theta})$$

- Therefore, we don't have a “descent” algorithm here
- It's possible that

$$f(\boldsymbol{\theta}^{\text{next}}) > f(\boldsymbol{\theta})$$

- Though people frequently use “SGD,” it's unclear if “D” is suitable in the name of this method

# Momentum I

- This is a method to improve the convergence speed
- A new vector  $\mathbf{v}$  and a parameter  $\alpha \in [0, 1)$  are introduced

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \eta \left( \frac{\boldsymbol{\theta}}{C} + \frac{1}{|S|} \nabla_{\boldsymbol{\theta}} \sum_{i:i \in S} \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i}) \right)$$
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v}$$

# Momentum II

- Essentially what we do is

$$\begin{aligned}\theta \leftarrow & \theta - \eta(\text{current sub-gradient}) \\ & - \alpha\eta(\text{prev. sub-gradient}) \\ & - \alpha^2\eta(\text{prev. prev. sub-gradient}) - \dots\end{aligned}$$

- There are some reasons why doing so can improve the convergence speed, though details are not discussed here

# AdaGrad I

- Scaling learning rates inversely proportional to the square root of sum of past gradient squares (Duchi et al., 2011)
- Update rule:

$$\begin{aligned}\mathbf{g} &\leftarrow \frac{\boldsymbol{\theta}}{C} + \frac{1}{|S|} \nabla_{\boldsymbol{\theta}} \sum_{i:i \in S} \xi(\boldsymbol{\theta}; \mathbf{y}^i, Z^{1,i}) \\ \mathbf{r} &\leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g} \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\mathbf{r}} + \delta} \odot \mathbf{g}\end{aligned}$$

- $\mathbf{r}$ : sum of past gradient squares



# AdaGrad II

$\epsilon$  and  $\delta$  are given constants

- $\odot$ : Hadamard product (element-wise product of two vectors/matrices)
- A large  $\mathbf{g}$  component  
 $\Rightarrow$  a larger  $\mathbf{r}$  component  
 $\Rightarrow$  fast decrease of the learning rate
- Conceptual explanation from Duchi et al. (2011):
  - frequently occurring features  $\Rightarrow$  low learning rates
  - infrequent features  $\Rightarrow$  high learning rates

# AdaGrad III

“the intuition is that each time an infrequent feature is seen, the learner should **take notice**.”

- But how is this explanation related to  $\mathbf{g}$  components?
- Let's consider **linear** classification. Recall our optimization problem is

$$\frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^l \xi(\mathbf{w}; y_i, \mathbf{x}_i)$$

# AdaGrad IV

- For methods such as SVM or logistic regression, the loss function can be written as a function of  $\mathbf{w}^T \mathbf{x}$

$$\xi(\mathbf{w}; y, \mathbf{x}) = \hat{\epsilon}(\mathbf{w}^T \mathbf{x})$$

Then the gradient is

$$\mathbf{w} + C \sum_{i=1}^I \hat{\epsilon}'(\mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i$$

- Thus the gradient is related to the **density of features**

# AdaGrad V

- The above analysis is for linear classification
- But now we have a **non-convex** neural network!
- **Empirically**, people find that the sum of squared gradient since the beginning causes **too fast decrease of the learning rate**

# RMSPprop I

- The original reference seems to be the lecture slides at [https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture\\_slides\\_lec6.pdf](https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf)
- Idea: they think AdaGrad's learning rate may be too small before reaching a locally convex region
- That is, OK to sum all past gradient squares in convex, but not non-convex
- Thus they do “exponentially weighted moving average”

# RMSPprop II

- Update rule

$$\begin{aligned} \mathbf{r} &\leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g} \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g} \end{aligned}$$

- AdaGrad:

$$\begin{aligned} \mathbf{r} &\leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g} \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\sqrt{\mathbf{r} + \delta}} \odot \mathbf{g} \end{aligned}$$

# RMSPprop III

- Somehow the setting is a bit heuristic and the reason behind the change (from AdaGrad to RMSPprop) is not really that strong

# References I

- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.