Backpropagation, SGD, and Adam

Deep Learning (ENGG*6600*01)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Summer 2023 **Gradient Descent**



- Gradient descent is one of the fundamental first-order methods.
- It was first suggested by Cauchy in <u>1874</u> [1] and <u>Hadamard</u> in <u>1908</u> [2] and its convergence was later analyzed in [3].
- Onconstrained optimization:

$$\begin{array}{c} \underset{x}{\text{minimize}} \quad f(x). \\ \end{array} \tag{1}$$

 In numerical optimization for unconstrained optimization, we start with a random feasible initial point and iteratively update it by step Δx:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \Delta \mathbf{x}.$$
(2)

- Continue until we converge to (or get sufficiently close to) the desired optimal point x*.
- The step Δx is also denoted by **p** in the literature, i.e., $\mathbf{p} := \Delta x$.



Gradient descent: update



Assume the gradient of function f(x) is L-smooth where L is the Lipschitz constant. In gradient descent, the update at every iteration is (see my "Optimization Techniques" course for proof):

• The problem is that often we either do not know the Lipschitz constant \underline{L} or it is hard to compute. Hence, rather than $\Delta x = -\frac{1}{L} \nabla f(\mathbf{x}^{(k)})$, we use:

$$\Delta \boldsymbol{x} = -\eta \nabla f(\boldsymbol{x}^{(k)}), \text{ i.e. } \boldsymbol{x}^{(k+1)} := \boldsymbol{x}^{(k)} - \eta \nabla f(\boldsymbol{x}^{(k)}), \tag{4}$$

where $\eta > 0$ is the step size, also called the learning rate in data science literature.

- If the optimization problem is maximization rather than minimization, the step should be $(\Delta x = \eta \nabla f(x^{(k)}))$ rather than Eq. (4). In that case, the name of method is gradient ascent.
- The learning rate can be found by **line search** (see my "Optimization Techniques" course for more information), which is used often in optimization and not in deep learning.



Gradient descent: series of solutions

• For a convex function, the series of solutions converges to the optimal solution while the function value decreases iteratively until the local minimum:

$$\frac{\{\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots\} \rightarrow \mathbf{x}^{*},}{f(\mathbf{x}^{(0)}) \ge f(\mathbf{x}^{(1)}) \ge f(\mathbf{x}^{(2)}) \ge \dots \ge f(\mathbf{x}^{*}).}$$

 If the optimization problem is a convex problem, the solution is the global solution; otherwise, the solution is local.

Gradient descent: cost versus iterations



Convergence criterion

Convergence criteria

- For all numerical optimization methods including gradient descent, there exist several methods for convergence criterion to stop updating the solution and terminate optimization. 7f(x) = 0
- Some of them are:
 - Small norm of gradient:

$$\|\nabla f(\boldsymbol{x}^{(k+1)})\|_2 \leq \boldsymbol{\epsilon},$$

where ϵ is a small positive number.

-

- ★ The reason for this criterion is the first-order optimality condition (recall that at the local optimum, we have $\|\nabla f(\mathbf{x}^*)\|_2 = 0$).
- * If the function is not convex, this criterion has the risk of stopping at a saddle point.
- Small change of <u>cost function</u>:

Small change of gradient of function:

$$\left| \nabla f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)}) \right| \leq \epsilon.$$

Reaching maximum desired number of iterations, denoted by

$$k < \max_k$$
.

Line-Search

Line-search

- We saw the step size of gradient descent requires knowledge of the Lipschitz constant for the smoothness of gradient. However, we may not know the exact Lipschitz constant. Hence, we can find the suitable step size η by a search which is named the line-search.
- In line-search of every optimization iteration, we start with $(\underline{\eta} = 1)$ and if it does not satisfy:

with step
$$\Delta \mathbf{x} = -\eta \nabla f(\mathbf{x}^{(k)})$$
: $\rightarrow \mathbf{x}^{(k)} + \Delta \mathbf{x} - \eta \nabla f(\mathbf{x}^{(k)}) < 0,$ (5)
 $f(\mathbf{x}^{(k)} + \Delta \mathbf{x}) < f(\mathbf{x}^{(k)}) \Rightarrow f(\mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)})) < f(\mathbf{x}^{(k)}),$ (6)
we have it $n \leftarrow n/2$

we halve it, $\eta \leftarrow \eta/2$.

• This halving step size is repeated until this equation is satisfied, i.e., until we have a decrease in the objective function. Note that this decrease will happen when the step size becomes small enough to satisfy (see my "Optimization Techniques" course for proof):

$$\bigstar \boxed{\eta < \frac{1}{L}}$$
(7)

A more sophisticated line-search method is the <u>Armijo line-search</u> [4], also called the <u>backtracking line-search</u>. Another more sophisticated line-search is <u>Wolfe conditions</u> [5]. We will learn it later in the course. See my "Optimization Techniques" course for more information about these.

Gradient descent with line-search

The algorithm of gradient descent with line-search:



As this algorithm shows, line-search has its own internal iterations inside every iteration of gradient descent.

Momentum

Gradient descent with momentum

- Gradient descent and other first-order methods can have a momentum term. Momentum, proposed in [6], makes the change of solution Δx a little similar to the previous change of solution.
- Hence, the change adds a history of previous change to Eq. (4):

$$(\Delta \mathbf{x})^{(k)} := \alpha (\Delta \mathbf{x})^{(k-1)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}),$$
(8)

where $\alpha > 0$ is the momentum parameter which weights the importance of history compared to the descent direction.

• We use this $(\Delta x)^{(k)}$ in Eq. (2) for updating the solution:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + (\Delta \mathbf{x})^{(k)}.$$

 Because of faithfulness to the track of previous updates, momentum reduces the amount of oscillation of updates in gradient descent optimization.



Steepest Descent

Steepest Descent

- Steepest descent is similar to gradient descent but there is a difference between them.
- In steepest descent, we move toward the negative gradient as much as possible to reach the smallest function value which can be achieved at every iteration.
- Hence, the step size at iteration k of steepest descent is calculated as [7]:

$$(\eta^{(k)}) := \arg\min_{\eta} f(\mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)}))), \qquad (9)$$

and then, the solution is updated using Eq. (4) as in gradient descent:

$$\begin{aligned} \mathbf{x}^{(k+1)} &:= \mathbf{x}^{(k)} - \eta \nabla f(\mathbf{x}^{(k)}). \end{aligned}$$

Backpropagation

Neural network

Neural network:



• Every neuron in neural network:



Let x_{ji} denote the weight connecting neuron i to neuron j. Let a_i and z_i be the output of neuron i before and after applying its activation function σ_i(.): ℝ → ℝ, respectively.

$$\mathbf{A}_{i} = \sum_{\ell=1}^{m} \mathbf{x}_{i\ell} \mathbf{x}_{\ell},$$



Backpropagation



- If layer *i* is the last layer, δ_i can be computed by derivative of error (loss function) w.r.t. the output.
- However, if i is one of the hidden layers, δ_i is computed by chain rule as:

• The term
$$\partial a_j / \partial a_i$$
 is calculated by chain rule as:

$$\underbrace{\left(\begin{array}{c} \partial e \\ \partial a_i \end{array}\right)}_{i} = \underbrace{\left(\begin{array}{c} \partial e \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial e \\ \partial a_i \end{array}\right)}_{j} = \sum_j \left(\begin{array}{c} \partial e \\ \partial a_i \end{array}\right) = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_i \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} = \underbrace{\left(\begin{array}{c} \partial a_j \\ \partial a_j \end{array}\right)}_{j} =$$

where (a) is because $a_i = \sum_i x_{ji} z_i$ and $z_i = \sigma(a_i)$ and $\sigma'(.)$ denotes the derivative of activation function. Parting Eq. (12) in Eq. (11) gives:

$$\chi_{j|z_{l}} \in \mathcal{N}_{j_{l}} := \chi_{j_{l}}$$

Backpropagation

We found:

1

• Putting this equation in Eq. (10),
$$\frac{\partial e}{\partial x_{i\ell}} = \delta_i \times (z_\ell) \text{ gives:}$$

$$\underbrace{\partial e}{\partial x_{i\ell}} = (z_\ell) \sigma'(a_i) \sum_j (\delta_j \times_{ji}). \quad (13)$$

• **Backpropagation** uses the gradient in Eq. (13) for updating the weight $x_{i\ell}, \forall i, \ell$ by gradient descent:

$$x_{i\ell}^{(k+1)} := x_{i\ell}^{(k)} - \eta^{(k)} \frac{\partial e}{\partial x_{i\ell}}.$$

- This tunes the weights from last layer to the first layer for every iteration of optimization.
- Therefore, <u>backpropagation</u>, proposed in <u>1986</u> [6], is actually gradient descent with <u>chain</u> rule in derivatives because of <u>having layers of parameters</u>. It is the most well-known optimization method used for training neural networks.

Stochastic gradient methods

Stochastic gradient descent

- Assume we have a dataset of *n* data points, $\{a_i \in \mathbb{R}^d\}_{i=1}^n$ and their labels $\{l_i \in \mathbb{R}\}_{i=1}^n$.
- Let the cost function f(.) be decomposed into summation of n terms $\{f_i(x)\}_{i=1}^n$. Some well-known examples for the cost function terms are:

Least squares error:
$$f_i(\mathbf{x}) = 0.5 (\mathbf{a}_i^{\top}(\mathbf{x}) - l_i)^2$$
, **K**

► Logistic loss (for
$$l_i \in \{-1, 1\}$$
): $\log(\frac{1}{1 + \exp(-l_i \boldsymbol{a}_i^\top \boldsymbol{x})})$.

• The optimization problem becomes:

$$\underset{\mathbf{x}}{\underset{\mathbf{x}}{\text{minimize}}} \left(\frac{1}{n} \sum_{i=1}^{n} f_i(\mathbf{x}). \right)$$
(14)

In this case, the full gradient is the average gradient, i.e:

$$\chi \overset{(k+1)}{\leftarrow} \chi \overset{k}{\leftarrow} \Delta \chi \overset{(k)}{\leftarrow} \Delta x \overset{(k)}{\leftarrow} \nabla f(x) = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x), \qquad (15)$$
so $\Delta x = -\eta \nabla f(x^{(k)}),$ becomes $\Delta x = -(\eta/n) \sum_{i=1}^{n} \nabla f_i(x^{(k)}).$ This is what gradient descent uses for updating the solution at every iteration.

Stochastic gradient descent

$$\bigstar \quad \nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\mathbf{x}),$$

- Calculation of this <u>full gradient</u> is time-consuming and <u>inefficient</u> for <u>large values of n</u>, especially as it needs to be recalculated at every iteration.
- Stochastic Gradient Descent (SGD), also called <u>stochastic gradient method</u>, approximates gradient descent <u>stochastically</u> and <u>samples</u> (i.e. <u>bootstraps</u>) one of the points at every iteration for updating the solution. Hence, it uses:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - (\eta^{(k)}) \nabla f_i(\mathbf{x}^{(k)}), \qquad (16)$$

rather than Eq. (4), $\underline{x^{(k+1)}} := \underline{x^{(k)}} - \eta \nabla f(\underline{x^{(k)}})$.

- The idea of stochastic approximation was first proposed in 1951 [8]. It was first used for machine learning in 1998 [9].
- As Eq. (16) states, SGD often uses an <u>adaptive step</u> size which changes in every iteration. The step size can be decreasing because in initial iterations, where we are far away from the optimal solution, the step size can be large; however, it should be small in the last iterations which is supposed to be close to the optimal solution. Some well-known adaptations for the step size are:

$$\eta^{(k)} := \frac{1}{k}, \quad \eta^{(k)} := \frac{1}{\sqrt{k}}, \quad \eta^{(k)} := \eta.$$
(17)



- Consider a <u>convex</u> and <u>differentiable</u> function f(.), with <u>domain</u> \mathcal{D} , whose gradient is <u>L-smooth</u>. Let f^* be the <u>minimum of cost function</u> and x^* be the minimizer. Starting from the <u>initial point</u> $x^{(0)}$, after <u>t</u> iterations of the optimization algorithm, we will have the following.
- The convergence rate of gradient descent:

$$f(\mathbf{x}^{(t+1)}) - f^{*} \leq \frac{2L \|\mathbf{x}^{(0)} - \mathbf{x}^{*}\|_{2}^{2}}{(t+1)} = O(\frac{1}{t}).$$

$$f^{*} \leq f(\mathbf{x}^{(t+1)})$$

$$f_{(\mathbf{x}^{*})=f^{*} \dots f^{*}$$

$$\mathbf{x}^{*}$$
(18)

Convergence Rate of Stochastic Gradient Descent

• Consider a function $f(\mathbf{x}) = \sum_{i=1}^{n} f_i(\mathbf{x})$ and which is bounded below and each f_i is differentiable. Let the domain of function f(.) be $\underline{\mathcal{D}}$ and its gradient be \underline{L} -smooth. Assume $\mathbb{E}[||\nabla f_i(\mathbf{x}_k)||_2^2 | \mathbf{x}_k] \leq \beta^2$ where $\underline{\beta}$ is a constant. Assume $\frac{\mathbb{E}[||\nabla f_i(\mathbf{x}_k)||_2^2 | \mathbf{x}_k] \leq \beta^2}{||\nabla f_i(\mathbf{x}_k)||_2^2 | \mathbf{x}_k] \leq \beta^2}$

• Depending on the step size, the convergence rate of SGD is:

$$\int f(\mathbf{x}^{(t+1)}) - f^* \leq \boxed{\mathcal{O}(\frac{1}{\log t})} \text{ if } \eta^{(\tau)} = \frac{1}{\tau},$$
(19)

$$f(\mathbf{x}^{(t+1)}) - f^* \leq \mathcal{O}(\frac{\log t}{\sqrt{t}}) \quad \text{if} \quad \eta^{(\tau)} = \frac{1}{\sqrt{\tau}}, \tag{20}$$

$$\bigstar \quad f(\mathbf{x}^{(t+1)}) - f^* \leq \mathcal{O}(\frac{1}{t} + \eta) \quad \text{if} \quad \underline{\eta}^{(\tau)} = \eta, \tag{21}$$

where τ denotes the iteration index.

• If the functions f_i 's are μ -strongly convex, then the convergence rate of SGD is:

$$f(\mathbf{x}^{(t+1)}) - f^* \leq \mathcal{O}(\frac{1}{t}) \quad \text{if} \quad \eta^{(\tau)} = \frac{1}{\mu\tau}, \tag{22}$$

$$f(\mathbf{x}^{(t+1)}) - f^* \leq \mathcal{O}((1 - \frac{\mu}{L})^t + \eta) \quad \text{if} \quad \eta^{(\tau)} = \eta.$$
(23)

Analysis of convergence rates

• Recall Eqs. (21) and (23):

nd (23): convex or non-convex: $\mathcal{O}(\underbrace{\ell}_{t}^{t} + \eta)$ if $\underbrace{\eta^{(\tau)} = \eta}_{t}$, strongly convex: $\mathcal{O}((1 - \frac{\mu}{L}) + \eta)$ if $\underbrace{\eta^{(\tau)} = \eta}_{t}$,

つ(り)

0(Ł)

- These equations show that with a fixed step size η, SGD converges sublinearly for a non-convex function and exponentially for a strongly convex function in the initial iterations.
- However, in the late iterations, it <u>stagnates</u> to a neighborhood around the optimal point and never reaches it. Hence, SGD has less accuracy than gradient descent (whose convergence rate is $O(\frac{1}{t})$ as in Eq. (18)).
- The <u>advantage of SGD over gradient descent</u> is that its every iteration is much faster than every iteration of gradient descent because of less computations for gradient. This faster pacing of every iteration shows off more when <u>n is huge</u>.
- In summary, SGD has fast convergence to a low accurate optimal solution.
- It is noteworthy that the full gradient is not available in SGD to use for checking convergence, as discussed before. One can use other criteria or merely check the norm of gradient for the sampled point.
- SGD can be used with the line-search methods, too. SGD can also use a momentum term.

- Gradient descent uses the <u>entire n data points</u> and <u>SGD</u> uses one randomly sampled point at every iteration. For large datasets, <u>gradient descent is very slow and intractable</u> in every iteration while <u>SGD will need a significant number of iterations to roughly cover all</u> data. Besides, SGD has low accuracy in convergence to the optimal solution.
- We can have a <u>middle case where we use a batch of b randomly sampled points at every</u> iteration. This method is named the <u>mini-batch SGD</u> or the <u>hybrid</u> <u>deterministic-stochastic gradient</u> method. This batch-wise approach is wise for large datasets.
- Usually, before start of optimization, the <u>n data points</u> are randomly divided into $\left(\frac{n}{b}\right)$ batches of size <u>b</u>. This is equivalent to simple random sampling for sampling points into batches without replacement. We denote the dataset by $\underline{\mathcal{D}}$ (where $|\underline{\mathcal{D}}| = n$) and the *i*-th batch by \mathcal{B}_i (where $|\mathcal{B}_i| = b$). The batches are disjoint:

$$\left[\bigcup_{i=1}^{\lfloor n/b \rfloor} \mathcal{B}_i = \mathcal{D},\right]$$
(24)

$$\mathcal{B}_{i} \cap \mathcal{B}_{j} = \varnothing, \quad \forall i, j \in \{1, \dots, \lfloor n/b \rfloor\}, \ i \neq j.$$
(25)

• Another less-used approach for making batches is to sample points for a batch during optimization. This is equivalent to <u>bootstrapping</u> for sampling points into batches with replacement. In this case, the batches are not <u>disjoint</u> anymore and Eqs. (24) and (25) do not hold.

Definition (Epoch)

In mini-batch SGD, when all $\lfloor n/b \rfloor$ batches of data are used for optimization once, an **epoch** is completed. After completion of an epoch, the next epoch is started and epochs are repeated until convergence of optimization.

In mini-batch SGD, if the k-th iteration of optimization is using the <u>k'-th batch</u>, the update of solution is done as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \overset{\mathbf{k}}{\not b} \sum_{i \in \mathcal{B}_{k'}} \nabla f_i(\mathbf{x}^{(k)}).$$

- The scale factor 1/b is sometimes dropped for simplicity.
- Mini-batch SGD is used significantly in machine learning, especially in <u>neural networks</u> [9, 10].
- Because of dividing data into batches, mini-batch SGD can be solved on parallel servers as a distributed optimization method.

(26)

Theorem (Convergence rates for mini-batch SGD)

Consider a function $f(\mathbf{x}) = \sum_{i=1}^{n} f_i(\mathbf{x})$ which is <u>bounded below</u> and each f_i is differentiable. Let the domain of function f(.) be <u>D</u> and its gradient <u>be L-smooth</u> and assume $\eta^{(k)} = \eta$ is fixed. The <u>batch-wise gradient</u> is an <u>approximation to the full gradient</u> with some error e_t for the t-th iteration:

The convergence rate of mini-batch SGD for non-convex and convex functions are:

$$\mathcal{O}\big(\frac{1}{t} + \|\boldsymbol{e}_t\|_2^2\big), \tag{28}$$

where t denotes the iteration index. If the functions f_i 's are μ -strongly convex, then the convergence rate of mini-batch SGD is:

$$\mathcal{O}ig((1-rac{\mu}{L})^t+\|e_t\|_2^2ig).$$

Therefore, the convergence rate of mini-batch gets closer to that of gradient descent, $\mathcal{O}(1/t)$ if the batch size increases.

(29)

• If we sample the batches without replacement (i.e., sampling batches by simple random sampling before start of optimization) or with replacement (i.e., bootstrapping during optimization), the expected error is [11, Proposition 3]:

respectively, where σ^2 is the variance of whole dataset.

- According to Eqs. (30) and (31), the accuracy of SGD by sampling without and with replacement increases by $\underline{b \rightarrow n}$ and $\underline{b \rightarrow \infty}$, respectively.
- However, this increase makes every iteration slower so there is trade-off between accuracy and speed.

Adaptive Learning Rate

Adaptive Gradient (AdaGrad)

- We can <u>adapt the learning rate</u> in <u>stochastic gradient methods</u>. Three most well-known methods for adapting the learning rate are AdaGrad, RMSProp, and <u>Adam</u>.
- Adaptive Gradient (AdaGrad) method, proposed in 2011 [12], updates the solution iteratively as:

$$\bigstar \quad (\mathbf{x}^{(k+1)}) := \mathbf{x}^{(k)} - \eta^{(k)} \mathbf{C}^{-1} \nabla f_i(\mathbf{x}^{(k)}), \qquad (32)$$

where **G** is a $(d \times d)$ diagonal matrix whose (j, j)-th element is:

$$\mathbf{\mathcal{K}} \quad \mathbf{\mathcal{G}}(j,j) := \sqrt{\varepsilon} + \sum_{\tau=0}^{k} (\nabla_{j} (\mathbf{x}^{(\tau)}))^{2}, \qquad (33)$$

where $\varepsilon \ge 0$ is for stability (making <u>G</u> full rank), i_{τ} is the randomly sampled point (from $\{1, \ldots, n\}$) at iteration τ , and $\nabla_j f_{i_{\tau}}(.)$ is the partial derivative of $f_{i_{\tau}}(.)$ w.r.t. its j-th element (note that $f_{i_{\tau}}(.)$ is d-dimensional). • Putting Eq. (33) in Eq. (32) can simplify AdaGrad to:

$$\mathbf{x}_{j}^{(k+1)} := \mathbf{x}_{j}^{(k)} - \left[\underbrace{\eta^{(k)}}_{\sqrt{\varepsilon + \sum_{\tau=0}^{k} (\nabla_{j} f_{i_{\tau}}(\mathbf{x}^{(\tau)}))}} \nabla f_{j}(\mathbf{x}_{j}^{(k)}) \right]$$

$$(34)$$

 AdaGrad keeps a history of the sampled points and it takes derivative for them to use. During the iterations so far, if a dimension has changed significantly, it dampens the learning rate for that dimension (see the inverse in Eq. (32)); hence, it gives more weight for changing the dimensions which have not changed noticeably. In this way, all dimensions will have a fair chance to change.

Root Mean Square Propagation (RMSProp)

- <u>Root Mean Square Propagation</u> (RMSProp) was first proposed in 2012 [13] which is unpublished.
- It is an improved version of <u>Rprop (resilient backpropagation)</u>, proposed in <u>1992</u> [14], which uses the sign of gradient in optimization.
- Inspired by momentum in Eq. (8):

$$\bigstar \bigstar \left((\Delta \mathbf{x})^{(k)} := \alpha(\Delta \mathbf{x})^{(k-1)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}), \right)$$

it updates a scalar variable v as [15]:

$$\mathbf{x} \quad \mathbf{v}^{(k+1)} := \gamma \mathbf{v}^{(k)} + (1-\gamma) \|\nabla f_i(\mathbf{x}^{(k)})\|_2^2, \tag{35}$$

where $\gamma \in [0, 1]$ is the forgetting factor (e.g., $\gamma = 0.9$). Then, it uses this v to weight the learning rate:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \underbrace{\frac{\eta^{(k)}}{\sqrt{\varepsilon + (\mathbf{v}^{(k+1)})}} \nabla f_j(\mathbf{x}_j^{(k)}), \qquad (36)$$

where $\epsilon \ge 0$ is for stability not to have division by zero. Comparing Eqs. (34) and (36) shows that RMSProp has a similar form to AdaGrad.

Adaptive Moment Estimation (Adam)

- Adam (Adaptive Moment Estimation) optimizer [16] improves over RMSProp by adding a momentum term.
- It updates the scalar v and the vector $m \in \mathbb{R}^d$ as:

$$\underbrace{ \mathbf{m}^{(k+1)} := \gamma_1 \mathbf{m}^{(k)} + (1 - \gamma_1) \nabla f_i(\mathbf{x}^{(k)}), }_{\mathbf{v}^{(k+1)} := \gamma_2 \mathbf{v}^{(k)} + (1 - \gamma_2) \| \nabla f_i(\mathbf{x}^{(k)}) \|_2^2, }$$

$$(37)$$

$$(37)$$

$$(37)$$

$$(38)$$

where $\gamma_1, \gamma_2 \in [0, 1]$. It normalizes these variables as:

$$\widehat{m{m}}^{(k+1)} := rac{1}{1-\gamma_1^k} m{m}^{(k+1)}, \ \widehat{m{arphi}^{(k+1)}} := rac{1}{1-\gamma_2^k} m{v}^{(k+1)}.$$

Then, it updates the solution as:

$$\boldsymbol{x}^{(k+1)} := \boldsymbol{x}^{(k)} - \underbrace{\boldsymbol{y}^{(k)}}_{\sqrt{\varepsilon + \widehat{\boldsymbol{y}}^{(k+1)}}} \widehat{\boldsymbol{m}}^{(k+1)},$$

(39)

which is stochastic gradient descent with momentum while using RMSProp.

• The Adam optimizer is one of the mostly used optimizers in neural networks.

Coding a Neural Network

Neural network: importing packages

Importing packages

```
[188] # installation in Google Colab's Jupyter notebook:
pip install torch
```

Looking in indexes: <u>https://pyi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/</u> Requirement already satisfied: torch in /usr/local/lib/python3.8/dist-packages (1.13.1+cu116) Requirement already satisfied: typing-extensions in /usr/local/lib/python3.8/dist-packages (from torch) (4.4.0)

```
[320] # importing packages (libraries):
import torch
import torch.nn as nn
from torch.utils.data import DataLoader, Dataset
import matplotlib.pyplot as plt
import numpy as np
from tqdm import tqdm
```

Neural network: defining the network

- Defining the network

```
[216] device = torch.device("cuda:0" if torch.cuda.is available() else "cpu")
[339] # defining the structure of neural network:
        class NeuralNetwork(nn.Module):
           def __init__(self):
                super(NeuralNetwork, self). init ()
                self.laver1 = nn.Linear(1, 10)
                self.layer2 = nn.Linear(10, 20)
                self.laver3 = nn.Linear(20, 10)
               self.layer4 = nn.Linear(10, 1)
               self.relu1 = nn.ReLU()
                self.relu2 = nn.ReLU()
                self.relu3 = nn.ReLU()
           def forward(self, x):
               x = self.relu1(self.layer1(x))
               x = self.relu2(self.laver2(x))
               x = self.relu3(self.layer3(x))
               x = self.laver4(x)
               return x
[340] # instantiate the class of neural network:
        net = NeuralNetwork()
        print(net)
       NeuralNetwork(
         (layer1): Linear(in features=1, out features=10, bias=True)
         (layer2): Linear(in_features=10, out_features=20, bias=True)
         (layer3): Linear(in_features=20, out_features=10, bias=True)
         (laver4): Linear(in features=10, out features=1, bias=True)
         (relu1): ReLU()
         (relu2): ReLU()
         (relu3): ReLU()
```

Neural network: optimizer

- Optimizer

```
>> [341] # define optimizer:
optimizer = 'Adam'
if optimizer == 'SGD':
optimizer = torch.optim.SGD(net.parameters(), lr=0.02)
elif optimizer == 'Adam':
optimizer = torch.optim.Adam(net.parameters(), lr=0.02)
# define the loss function:
loss_func = torch.nn.MSELoss()
```

Neural network: data loader

Data loader

```
[342] class Data(Dataset):

            def init (self, x, y):
                self.data = torch.from numpv(x.reshape(-1,1)).float()
                self.label = torch.from numpy(y.reshape(-1,1)).float()
            def __len__(self):
                return len(self.data)
            def getitem (self, item):
                data point = self.data[item]
               label point = self.label[item]
                return data point, label point
/ [343] batch_size = 16
       def load dataset(x train, y train, x test, y test):
            # data loader for training data:
           train ds = Data(x train, y train)
            train loader = DataLoader(train ds, batch size=batch size, shuffle=True)
            # data loader for test data:
            test ds = Data(x test, v test)
            test loader = DataLoader(test ds, batch size=batch size, shuffle=False)
            return train loader, test loader
```

Neural network: dataset

```
dataset type = 'nonlinear'
    if dataset_type == 'linear':
        # almost linear dataset:
        x train = np.random.rand(100)
        y train = np.sin(x train) * (x train**3) + 3*x train + np.random.rand(100)*0.8
        x test = np.random.rand(100)
        y_test = np.sin(x_test) * (x_test**3) + 3*x_test + np.random.rand(100)*0.8
    elif dataset_type == 'nonlinear':
        # dataset:
        x train = np.random.rand(100) * 10
        y_train = np.sin(x_train) + np.random.rand(100)*0.8
        x test = np.random.rand(100) * 10
        v test = np.sin(x test) + np.random.rand(100)*0.8
    # reshape to have samples in rows:
    x train = x train.reshape((-1, 1))
    x_test = x_test.reshape((-1, 1))
    # visualize data:
    plt.scatter(x_train, y_train, c='r', label='train')
    plt.scatter(x_test, y_test, c='b', label='test')
    plt.xlabel('x')
    plt.ylabel('y')
    plt.legend()
    plt.show()
```



Neural network: training

Training neural network

```
/ [345] # load dataset:
        train_loader, test_loader = load_dataset(x_train, y_train, x_test, y_test)
[346] n_epochs = 1000
        loss list = []
        for epoch in tadm(range(n epochs), desc='epochs');
            loss list in epoch = []
            for step, (data_batch, label_batch) in enumerate(train_loader):
                data_batch, label_batch = data_batch.to(device), label_batch.to(device)
                prediction = net(data_batch)
                loss = loss func(prediction, label batch)
                loss list in epoch.append(loss.cpu().detach().item())
                optimizer.zero grad()
                loss.backward()
                optimizer.step()
            loss list.append(np.mean(loss list in epoch))
```

epochs: 100%

Neural network: test (evaluation)

- Test (evaluation) phase

```
prediction_list = []
with torch.no_grad():
    for step, (data_batch, label_batch) in enumerate(test_loader):
        prediction = net(data_batch)
        prediction_list.extend(prediction)
```

```
% [350] # visualize the predicted and actual data:
plt.scatter(x_test, y_test, c='b', label='test')
plt.scatter(x_test, prediction_list, c='g', label='prediction')
plt.ylabel('x')
plt.ylabel('y')
plt.legend()
plt.show()
```



Acknowledgement

- Some slides of this slide deck are inspired by the lectures of Prof. Kimon Fountoulakis at the University of Waterloo.
- Some slides of this slide deck are inspired by the lectures of Prof. Stephen Boyd at the Stanford University.
- Our tutorial also has the materials of this slide deck: [17]
- See my "Optimization Techniques" course on my YouTube channel for more information about first-order optimization including these methods.

References

- C. Lemaréchal, "Cauchy and the gradient method," *Doc Math Extra*, vol. 251, no. 254, p. 10, 2012.
- J. Hadamard, Mémoire sur le problème d'analyse relatif à l'équilibre des plaques élastiques encastrées, vol. 33.
 Imprimerie nationale, 1908.
- [3] H. B. Curry, "The method of steepest descent for non-linear minimization problems," *Quarterly of Applied Mathematics*, vol. 2, no. 3, pp. 258–261, 1944.
- [4] L. Armijo, "Minimization of functions having Lipschitz continuous first partial derivatives," *Pacific Journal of mathematics*, vol. 16, no. 1, pp. 1–3, 1966.
- [5] P. Wolfe, "Convergence conditions for ascent methods," SIAM review, vol. 11, no. 2, pp. 226–235, 1969.
- [6] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [7] E. K. Chong and S. H. Zak, An introduction to optimization. John Wiley & Sons, 2004.
- [8] H. Robbins and S. Monro, "A stochastic approximation method," The annals of mathematical statistics, pp. 400–407, 1951.

References (cont.)

- [9] L. Bottou *et al.*, "Online learning and stochastic approximations," *On-line learning in neural networks*, vol. 17, no. 9, p. 142, 1998.
- [10] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [11] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karray, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," arXiv preprint arXiv:2011.00901, 2020.
- [12] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of machine learning research*, vol. 12, no. 7, 2011.
- [13] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude," COURSERA: Neural networks for machine learning, vol. 4, no. 2, pp. 26–31, 2012.
- [14] M. Riedmiller and H. Braun, "Rprop-a fast adaptive learning algorithm," in *Proceedings of the International Symposium on Computer and Information Science VII*, 1992.
- [15] G. Hinton, N. Srivastava, and K. Swersky, "Neural networks for machine learning lecture 6a overview of mini-batch gradient descent," tech. rep., Department of Computer Science, University of Toronto, 2012.
- [16] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014.

References (cont.)

[17] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "KKT conditions, first-order and second-order optimization, and distributed optimization: Tutorial and survey," arXiv preprint arXiv:2110.01858, 2021.