

MATH5836: Data and Machine Learning

Week 4: Advanced Techniques for Neural Networks

Sarat Moka

UNSW, Sydney

Key Topics

4.1	Momentum and Nesterov Momentum	4-2
4.2	Stochastic Gradient Descent	4-5
4.3	Adaptive Learning Rate Gradient Descent Methods	4-9
4.4	Dropout Method	4-13
4.5	Weight Decay (L2) regularization	4-15

Books:

- (A) *Mathematical Engineering of Deep Learning* book by Liquet, Moka, and Nazarathy (2024).
Click [here](#) for downloading chapters.
- (B) *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow (2nd Ed.)* by Aurélien Géron (2019).

4.1 Momentum and Nesterov Momentum

- Recall that when $C(\theta)$ is the loss function with θ denoting all the weights and bias in the neural network, the (basic) gradient descent step at the $(t + 1)$ -th iteration is

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla C(\theta^{(t)}),$$

where α is the learning rate.

- For instance, for a weight $w_{i,j}$ in the network,

$$w_{i,j}^{(t+1)} = w_{i,j}^{(t)} - \alpha \frac{\partial C(\theta^{(t)})}{\partial w_{i,j}}.$$

- Observe that the parameter update in the basic gradient descent depends only on the gradient at the current point.
- This implies that in flat regions of the loss landscape, gradient descent essentially stops. Further, at saddle points or local minima, gradient descent stops completely. Such drawbacks may be alleviated with the use of *momentum*.

Momentum

- To get an intuitive feel for the use of momentum in optimization, consider an analogy of rolling a ball downhill on the loss landscape. It gains momentum as it rolls downhill, becoming increasingly faster until it reaches the bottom of a valley. Clearly, the ball keeps memory of the past forces in the form of acceleration. This motivates us to use exponential smoothing on the gradients to obtain an extra step called the *momentum update*.
- It is formulated as follows,

$$v^{(t+1)} = \beta v^{(t)} + (1 - \beta) \nabla C(\theta^{(t)}) \quad (\text{Momentum Update}^1), \quad (4.1)$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha v^{(t+1)} \quad (\text{Parameter Update}), \quad (4.2)$$

for scalar *momentum parameter* $\beta \in [0, 1)$ and *learning rate* $\alpha > 0$, starting with $v^{(0)} = 0$.

- When $\beta = 0$, we have the basic gradient descent method and for larger β , information of previous gradients plays a more significant role.
- The vector $v^{(t)}$ in (4.1) is called the *momentum*² at the t -th update. We can also write that

$$v^{(t+1)} = (1 - \beta) \sum_{\tau=0}^t \beta^{t-\tau} \nabla C(\theta^{(\tau)}), \quad \text{for } t = 0, 1, 2, \dots \quad (4.3)$$

That is, the momentum accumulates all the past (weighted) gradients, providing acceleration to the parameter θ updates on downward surfaces.

²In physics, momentum is defined as the product of the mass of an object and its velocity vector. The object's momentum points in the direction of an object's movement. In contrast, here the momentum vector $v^{(t)}$ points in the direction opposite to the step taken.

- Therefore, for $\beta > 0$, the next step taken via (4.2) is not necessarily taken in the steepest descent, instead the direction is dictated by the (exponential smoothing) average of all the gradients up to the current iteration.
- Figure 4.1 compares the performance of the gradient descent method with momentum for different values of β on the Rosenbrock function presented in a previous lecture. We observe that for large β , the momentum method accelerates as it takes downward steps.

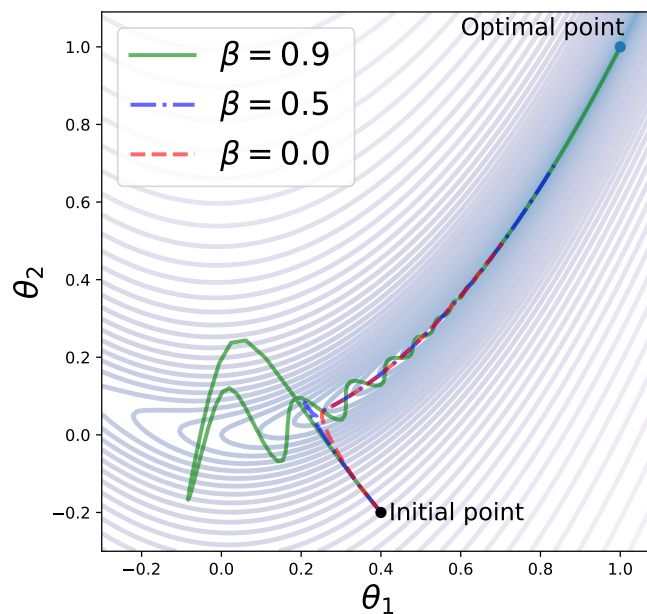


Figure 4.1: Application of momentum to the Rosenbrock function for three different values of β with learning rate $\alpha = 0.001/(1 - \beta)$ and a fixed number of iterations. Note that $\beta = 0$ is basic gradient descent.

Nesterov Momentum

- We have seen that the momentum updates as in (4.1) and (4.2) accelerate like a ball rolling downhill. An issue with this method is that the steps do not slow down after reaching the bottom of a valley, and hence, there is a tendency to overshoot the valley floor. Therefore, if we know approximately where the ball will be after each step, we can slow down the ball before the hill slopes up again.
- The idea of *Nesterov momentum* tries to improve on standard momentum by using the gradient at the predicted future position instead of the gradient of the current position.

- The update equations are,

$$v^{(t+1)} = \beta v^{(t)} + (1 - \beta) \nabla C(\underbrace{\theta^{(t)} - \alpha \beta v^{(t)}}_{\text{predicted next point}}), \quad (4.4)$$

$$\theta^{(t+1)} = \underbrace{\theta^{(t)} - \alpha v^{(t+1)}}_{\text{actual next point}}, \quad (4.5)$$

for constants $\beta \in [0, 1)$ and $\alpha > 0$, with $v^{(1)} = 0$. Compare (4.4) and (4.5) with (4.1) and (4.2), respectively.

- The difference here is that the gradient is computed at a predicted next point $\theta^{(t)} - \alpha \beta v^{(t)}$, which is a proxy for the (unseen) actual next point $\theta^{(t+1)} = \theta^{(t)} - \alpha v^{(t+1)}$ when $\beta \approx 1$.
- Implementing Nesterov momentum via (4.4) and (4.5) requires evaluation of the gradient at the predicted points instead of actual points. This may incur overheads, especially when incorporated as part of other algorithms.
- For example, if we also require $\nabla C(\theta^{(t)})$ at each iteration for purposes such as RMSprop, then the gradient needs to be computed twice instead of once per iteration; once for $\nabla C(\theta^{(t)})$ and once for $\nabla C(\theta^{(t)} - \alpha \beta v^{(t)})$.
- For this reason, and also for simplicity of implementing gradients only at $\theta^{(t)}$, a *look-ahead momentum* mechanism is sometimes used in place of (4.4) and (4.5).
- To see how this mechanism works, first revisit the basic momentum update equations (4.1) and (4.2) and observe that the parameter update can be represented as,

$$\theta^{(t+1)} = \theta^{(t)} - \alpha (\beta v^{(t)} + (1 - \beta) \nabla C(\theta^{(t)})). \quad (4.6)$$

- Observe that $v^{(t)}$ is based on gradients at points $\theta^{(0)}, \dots, \theta^{(t-1)}$, but not on the gradient at $\theta^{(t)}$. Hence a way to incorporate this last gradient is with look-ahead momentum where we replace $v^{(t)}$ of (4.6) by $v^{(t+1)}$. This achieves behaviour similar to Nesterov momentum.
- With such a replacement, we arrive at update equations of the form,

$$v^{(t+1)} = \beta v^{(t)} + (1 - \beta) \nabla C(\theta^{(t)}), \quad (4.7)$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \underbrace{(\beta v^{(t+1)} + (1 - \beta) \nabla C(\theta^{(t)}))}_{\text{look-ahead momentum}}. \quad (4.8)$$

- Using (4.7) and (4.8) aims to provide similar behaviour to the Nesterov momentum equations (4.4) and (4.5), yet only requires evaluation of gradients at $\theta^{(t)}$ as opposed to a shifted point as in (4.4). While look-ahead momentum is not equivalent to Nesterov momentum, the gist of both methods is similar.

4.2 Stochastic Gradient Descent

- In the context of deep learning, the typical form of the loss function is

$$C(\theta) = \frac{1}{n} \sum_{i=1}^n C_i(\theta), \quad (4.9)$$

where n is the number of data samples in a given dataset and $C_i(\theta)$ depends only on the i -th sample in the dataset. For instance, when the loss is mean squared error,

$$C_i(\theta) = (y_i - f_\theta(x_i))^2.$$

- As a consequence, the gradient loss function $C(\theta)$ is given by

$$\nabla C(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla C_i(\theta). \quad (4.10)$$

- The method of *stochastic gradient descent* exploits this structure by computing a noisy gradient using only one randomly selected training sample. That is, it evaluates the gradient for a single $C_i(\theta)$ instead of the gradient for all of $C(\theta)$.
- Computing $\nabla C_i(\theta)$ can be much faster using backpropagation than computing full gradient $\nabla C(\theta)$.
- More precisely, the algorithm operates similarly to gradient descent studied in the previous lecture, yet in the t -th iteration of stochastic gradient descent, an index variable I_t is randomly selected from the set $\{1, \dots, n\}$ and the gradient $\nabla C_{I_t}(\theta)$ is computed only for the I_t -th data sample.
- In particular, the update rule for the decision variable is then,

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla C_{I_t}(\theta^{(t)}). \quad (4.11)$$

- This stochastic algorithm exhibits some interesting properties.
 1. Most importantly, the execution of each iteration can be much faster than that of basic gradient descent because only one sample is used at a time (the gain is of order n).
 2. Further, if the index variable I_t is selected uniformly over $\{1, \dots, n\}$, the parameter update (4.11) is *unbiased* in the sense that the expected descent step for each iteration t is the same as that of gradient descent. That is, in each iteration, we have

$$\mathbb{E}[\nabla C_{I_t}(\theta)] = \frac{1}{n} \sum_{i=1}^n \nabla C_i(\theta) = \nabla C(\theta), \quad \text{for a fixed } \theta. \quad (4.12)$$

and hence the stochastic step (4.11) is “on average correct”.³

³The reader might be tempted to conclude that $\mathbb{E}[\theta_{\text{SGD}}^{(t)}] = \theta_{\text{GD}}^{(t)}$ for all time t , where the subscript SGD is for stochastic gradient descent starting at the same initial condition as GD (“Gradient Descent”). However, in general this is not correct when $\nabla C(\theta)$ is non-linear in θ .

- While being unbiased, the stochastic nature of (4.11) introduces fluctuations in the descent direction.
- Such noisy trajectories may initially appear like a drawback, yet one advantage is that they enable the algorithm to escape flat regions, saddle points, and local minima.
- Thus, while the fluctuations make it difficult to guarantee convergence, in the context of deep learning their effect is often considered desirable in view of highly non-convex loss landscapes.
- Importantly, stochastic gradient descent is generally able to direct the parameters θ towards the region where minima of $C_i(\theta)$ are located.

Illustrative Example

- To get a feel for the last attribute of stochastic gradient descent, we consider another hypothetical example where for $i = 1, \dots, n$,

$$C_i(\theta) = a_i(\theta_1 - u_{i,1})^2 + (\theta_2 - u_{i,2})^2, \quad (4.13)$$

for some constant $a_i > 0$ and $u_i = (u_{i,1}, u_{i,2}) \in \mathbb{R}^2$.

- Here u_i is the unique minimizer of $C_i(\theta)$. These individual loss functions for each observation are convex and hence the total loss function $C(\theta)$ of (4.9) is also convex.⁴
- Let \mathcal{S} be the *convex hull* of u_1, \dots, u_n , that is, \mathcal{S} is the smallest convex set that contains $\{u_1, \dots, u_n\}$. It is now obvious that the global minima of C is also in \mathcal{S} .
- In Figure 4.2 we plot the contours of such a function when $n = 5$ where in Figure 4.2 (a) the points u_1, \dots, u_5 are distinct and in (b) these points are identical. The convex hull is marked by the region bounded in green.
- The figure also plots the evolution of gradient descent and stochastic gradient descent in each of the cases. As is evident, while the stochastic gradient descent trajectory is noisier, outside of the convex hull \mathcal{S} , its trajectory is similar to that of gradient descent.
- Interestingly, in Figure 4.2 (a), once the trajectory hits \mathcal{S} it mostly stays within \mathcal{S} but with a lot of fluctuations.
- The reason for this behavior is that at any point θ not in \mathcal{S} , the descent directions for both gradient descent and stochastic gradient descent are towards the set \mathcal{S} as it contains the minima of $C(\theta)$ and every $C_i(\theta)$.
- To make this a concrete argument, note that the descent step in the t -th iteration of the stochastic gradient descent is

$$-\alpha \nabla C_{I_t}(\theta^{(t)}) = -\alpha \sum_{i=1}^n \mathbf{1}\{I_t = i\} \nabla C_i(\theta^{(t)}).$$

⁴This example does not resemble a general case with multiple local minima since it is convex. Nevertheless, it is useful for understanding some of the behaviour of stochastic gradient descent.

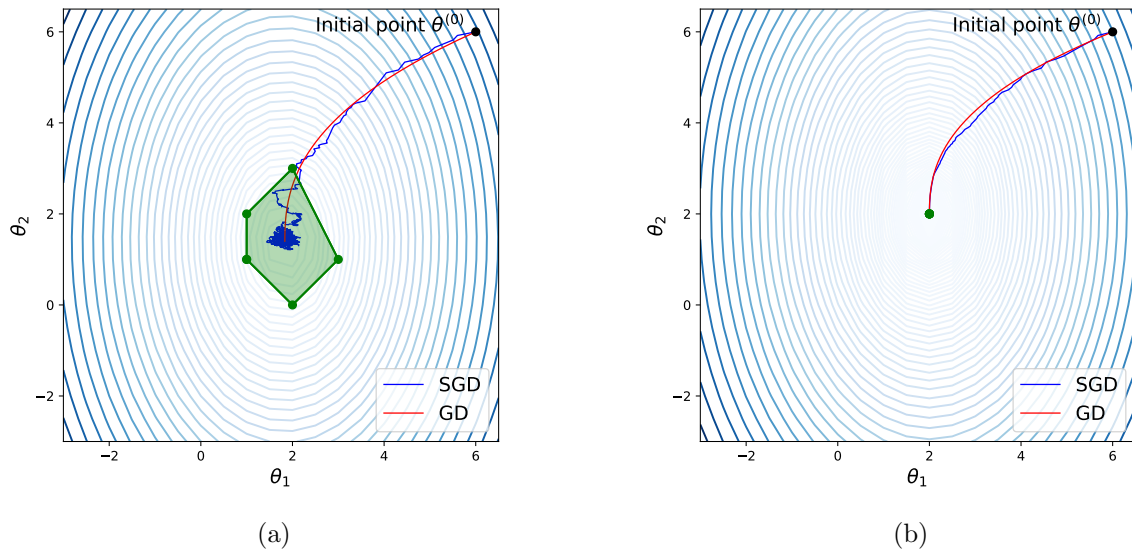


Figure 4.2: An hypothetical example where $C(\theta)$ is composed of individual $C_i(\theta)$ as in (4.13). Comparison of gradient descent (GD) and stochastic gradient descent (SGD). (a) A case where the minimizers u_1, \dots, u_n of $C_i(\theta)$ are different and are each marked by green dots in the plot. (b) A case where the minimizers u_1, \dots, u_n are the same.

- Since every $C_i(\theta)$ has its minimum in the set \mathcal{S} , if $\theta^{(t)} \notin \mathcal{S}$, then every $-\nabla C_i(\theta^{(t)})$, the steepest direction of $C_i(\theta)$ at $\theta^{(t)}$, guides the algorithm towards the set \mathcal{S} . Therefore, irrespective of the choice of I_t , the update moves $\theta^{(t+1)}$ towards \mathcal{S} .

Working with Mini-batches and the Concept of Epochs

- On one extreme, computing the full gradient for basic gradient descent is computationally challenging but the obtained gradients are noiseless. On the other extreme, stochastic gradient descent reduces the computational complexity of each update but the updates are noisy.
- One middle ground approach is to work with mini-batches. Almost any practical deep learning training uses some variant of this approach.
- In particular, a *mini-batch* is simply a small subset of the data points, of size n_b . Mini-batch sizes are typically quite small in comparison to the size of the dataset and are often chosen based on computational capabilities and constraints.
- In practice, a mini-batch size n_b is often selected so the computation of gradients associated with n_b observations can fit on GPU memory. This generally makes computation more efficient in comparison to stochastic gradient descent which cannot take full advantage of parallel computing or GPUs since the gradient evaluation for a single observation index I_t is not easy to parallelize.

- When using mini-batches, for each iteration t , a mini-batch with indices $\mathcal{I}_t \subseteq \{1, \dots, n\}$ is used to approximate the gradient via,

$$\frac{1}{n_b} \sum_{i \in \mathcal{I}_t} \nabla C_i(\theta^{(t)}). \quad (4.14)$$

This *estimated gradient* is then used as part of gradient descent or one of its variants (such as ADAM which we present in the sequel).

- One common practical approach with mini-batches is to shuffle the training data apriori and then use the shuffled indices sequentially. With this process there are n/n_b mini-batches⁵ in the training dataset and each mini-batch is a distinct random subset of the indices.
- In such a case, when using some variant of gradient descent, every pass on the full dataset has n/n_b iterations, with one iteration per mini-batch. Such a pass on all the data via n/n_b iterations is typically called an *epoch*.
- The training process then involves multiple epochs. It is quite common to diagnose and track the training process in terms of epochs.
- Note that one may also randomly shuffle the data (assign new mini-batches) after every epoch to reduce the correlation between the epochs and reduce bias in the optimization process.
- Similarly to (4.12) when using the mini-batch approach, the estimated gradient is unbiased in the sense that

$$\mathbb{E} \left[\frac{1}{n_b} \sum_{j=1}^{n_b} \nabla C_{I_{t,j}}(\theta) \right] = \frac{1}{n_b} \sum_{j=1}^{n_b} \frac{1}{n} \sum_{i=1}^n \nabla C_i(\theta) = \nabla C(\theta), \quad \text{for a fixed } \theta.$$

Here $I_{t,j}$ denotes the j -th element in \mathcal{I}_t and the first equality holds because shuffling makes each $I_{t,j}$ to be uniform over $\{1, 2, \dots, n\}$.

- Importantly, with min-batches, we notice that the noise of the gradients decreases as the mini-batch size n_b increases. To see this, recall that for each t , the index variables $I_{t,1}, \dots, I_{t,n_b}$ are chosen independently. Thus, the variance of the gradient estimate is

$$\text{Var} \left(\frac{1}{n_b} \sum_{j=1}^{n_b} \nabla C_{I_{t,j}}(\theta^{(t)}) \right) = \frac{1}{n_b} \text{Var} \left(\nabla C_{I_{t,1}}(\theta^{(t)}) \right).$$

- Since $\text{Var}(\nabla C_{I_{t,1}}(\theta^{(t)}))$ is the variance of the gradient for an arbitrary data sample, we see that mini-batch gradient descent has n_b times smaller variance than that of stochastic gradient descent.

⁵We assume here that n_b divides n . If it is not the case then there are $\lceil n/n_b \rceil$ mini-batches with the last one having less than n_b samples.

Remark 4.2.1

In practice, for deep neural networks, the stochastic gradient or the mini-batch gradient descent is combined with momentum methods by replacing the gradient $\nabla C(\theta^{(t)})$ in the momentum based methods by with noisy gradients as in stochastic gradient descent or mini-batches.

4.3 Adaptive Learning Rate Gradient Descent Methods

- Observe that the updates of the gradient descent method, with or without momentum, use the same learning rate α for all components of θ .
- However, it is often better to learn the parameters $\theta_1, \dots, \theta_d$ with potentially a specific learning rate for each coordinate θ_i .
- A general approach for overcoming this issue is to scale each coordinate of the descent step with a different factor.
- That is, instead of considering a basic gradient descent step, $-\alpha \nabla C(\theta)$, consider descent steps of the form,

$$\text{step} = -\alpha r \odot \nabla C(\theta), \quad (4.15)$$

where r is some vector of positive entries which recalibrates the descent steps and \odot is the element-wise product of two vectors⁶.

- However, in general, finding a fixed good vector r for (4.15) is difficult, especially when the number of parameters is huge and the actual relationship between parameters and features is not clear. Instead, we use adaptive methods that adjust the descent step during the optimization process.
- We now introduce four such adaptive approaches called *adaptive subgradient* (or *Adagrad*), *root mean square propagation* (or *RMSprop*), *Adadelta*, and *ADAM*.

Adagrad and RMSprop

- In both Adagrad and RMSprop, the update of coordinate i in iteration t can be represented as,

$$\theta_i^{(t+1)} = \theta_i^{(t)} - \frac{\alpha}{\sqrt{s_i^{(t+1)} + \varepsilon}} \frac{\partial C(\theta^{(t)})}{\partial \theta_i}, \quad (4.16)$$

where $s_i^{(0)}, s_i^{(1)}, s_i^{(2)}, \dots$ is a sequence of non-negative values that are updated adaptively and ε taken to be a small value of order 1×10^{-8} to avoid division by zero.

⁶The element-wise product is also called the *Hadamard product* or *Schur product*.

- That is, with this representation, the descent step at time t represented in terms of (4.15) has $r_i = \left(\sqrt{s_i^{(t+1)}} + \varepsilon \right)^{-1}$.
- Generally we aim for s_i to be some smoothed representation of the square of the derivative $\partial C(\theta^{(t)})/\partial \theta_i$.
- Hence, r_i is roughly the inverse of the magnitude of the derivative, and r_i is low when the changes for coordinate i are steep and vice-versa.
- Vector-wise, the parameter update (4.16) can be represented as

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\alpha}{\sqrt{s^{(t+1)}} + \varepsilon} \odot \nabla C(\theta^{(t)}), \quad (4.17)$$

where $s^{(t)} = (s_1^{(t)}, \dots, s_d^{(t)})$, ε is considered a vector, and the addition, division, and square-root operations are all element-wise operations.

- Using the representation (4.16) or (4.17), let us now define how the sequence of $\{s_i^{(t)}\}$ is computed both for Adagrad and RMSprop. Specifically,

$$s_i^{(t+1)} = \begin{cases} \sum_{\tau=0}^t \left(\frac{\partial C(\theta^{(\tau)})}{\partial \theta_i} \right)^2, & \text{for Adagrad,} \\ \gamma s_i^{(t)} + (1 - \gamma) \left(\frac{\partial C(\theta^{(t)})}{\partial \theta_i} \right)^2, & \text{for RMSprop,} \end{cases} \quad (4.18)$$

where the recursive relationship for RMSprop has the initial value at $s_i^{(0)} = 0$.

- RMSprop is also parameterized by a *decay parameter* $\gamma \in [0, 1)$ with typical values at $\gamma = 0.999$ implying that for RMSprop the sequence $\{s_i^{(t)}\}$ is a relatively slow exponential smoothing of the square of the derivative.

Remark 4.3.1

Adagrad may appear simpler than RMSprop since it does not involve any recursive step or any hyper-parameter like γ and it is simply an accumulation of the squares of the derivatives. However, the crucial drawback of Adagrad is that for each i , the sequence $\{s_i^{(t)}\}$ is strictly non-decreasing. As a consequence, the effective learning rate decreases during training, often making it infinitesimally small before convergence. RMSprop was introduced in the context of deep learning after Adagrad, and overcomes this problem via exponential smoothing of the squares of the partial derivatives.

- For RMSprop, the explicit (all-time) representation of the vector $s^{(t+1)}$ is,

$$s^{(t+1)} = (1 - \gamma) \sum_{\tau=0}^t \gamma^{t-\tau} \left(\nabla C(\theta^{(\tau)}) \odot \nabla C(\theta^{(\tau)}) \right), \quad (4.19)$$

while for Adagrad a similar representation holds without the $(1 - \gamma)$ and $\gamma^{t-\tau}$ elements in the formula.

Adadelta

- Recall the computation of $s^{(t+1)}$ for the Adagrad method as in (4.18). When Adagrad's $s^{(t+1)}$ is used in (4.17), it has the problem of monotonically decreasing effective learning rates.
- This is one of the reasons that eventually RMSprop became more popular than Adagrad. However, there are other alternatives that are also very popular.
- One such alternative is the *Adadelta* method which uses exponential smoothing of the squared gradients as in RMSprop, but also uses another exponentially smoothed sequence of the descent step's squares.
- A key motivation for Adadelta is the observation that the update equation (4.17) for RMSprop or Adagrad uses a unit-less quantity as the descent step.
- Specifically, in (4.17) the only unit-full quantity in the coefficient multiplying $\nabla C(\theta^{(t)})$ is the reciprocal of $\sqrt{s^{(t+1)}}$. Hence that coefficient has units which are the inverse of the gradient and these cancel out the units of the gradient implying that the descent step is unit-less.
- With Adadelta, the update equation (4.17) is modified so that the descent step maintains the same units of the gradient, via

$$\theta^{(t+1)} = \theta^{(t)} - \underbrace{\frac{\sqrt{\Delta\theta^{(t)}} + \varepsilon}{\sqrt{s^{(t+1)}} + \varepsilon}}_{\text{Descent step } \tilde{\nabla}C(\theta^{(t)})} \odot \nabla C(\theta^{(t)}), \quad (4.20)$$

where $\Delta\theta^{(t)}$ is adaptively adjusted yet has the same units as $s^{(t+1)}$, making the coefficient of the gradient unit-free.

- Compare (4.20) with (4.17) to observe that $\sqrt{\Delta\theta^{(t)}} + \varepsilon$ replaces the learning rate α . This also means that Adagrad is “learning rate free”. Instead, a potentially more robust parameter $\rho \in [0, 1)$ is used similarly to the γ parameter for RMSprop. This parameter specifies how to exponentially smooth squares of the descent steps.
- Using the descent step of (4.20), the update equation for $\Delta\theta^{(t)}$ is,

$$\Delta\theta^{(t)} = \rho\Delta\theta^{(t-1)} + (1 - \rho) \left(\tilde{\nabla}C(\theta^{(t-1)}) \odot \tilde{\nabla}C(\theta^{(t-1)}) \right),$$

starting with $\Delta\theta^{(0)} = 0$.

- Then at iteration t , updating $\theta^{(t+1)}$ via (4.20), we use both $\Delta\theta^{(t)}$ and $s^{(t+1)}$, where the latter is updated via the RMSprop update equation as in (4.18).

ADAM

- Now that we understand ideas of momentum, and RMSprop, we can piece these together into a single algorithm, namely the *adaptive moment estimation* method, or simply *ADAM*.

- Metaphorically, if the execution of the momentum method is a ball rolling down a slope, the execution of ADAM can be seen as a heavy ball with friction rolling down the slope.
- The key update formula for ADAM is

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \frac{1}{\sqrt{\hat{s}^{(t+1)} + \varepsilon}} \hat{v}^{(t+1)}, \quad (4.21)$$

where all vector operations (division, square root, and addition of ε) are interpreted element wise. Here $\hat{v}^{(t+1)}$ and $\hat{s}^{(t+1)}$ are bias corrected exponential smoothing of the gradient and the squared gradients.

Algorithm 1: ADAM

Input: Dataset $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$,
 objective function $C(\cdot)$, and
 initial parameter vector $\theta^{(0)}$

Output: Approximately *optimal* θ

$t \leftarrow 0$ (Initialize iteration counter)
 $\theta \leftarrow \theta^{(0)}$, $v \leftarrow 0$, and $s \leftarrow 0$ (Initialize state vectors)

repeat

$g \leftarrow \nabla C(\theta)$	(Compute gradient)
$v \leftarrow \beta v + (1 - \beta) g$	(Momentum update)
$s \leftarrow \gamma s + (1 - \gamma) (g \odot g)$	(Second moment update)
$\hat{v} \leftarrow \frac{v}{1 - \beta^{t+1}}$	(Bias correction)
$\hat{s} \leftarrow \frac{s}{1 - \gamma^{t+1}}$	(Bias correction)
$\theta \leftarrow \theta - \alpha \frac{1}{\sqrt{\hat{s} + \varepsilon}} \hat{v}$	(Update parameters)
$t \leftarrow t + 1$	

until *termination condition* is satisfied

return θ

- With ADAM, α is still called the learning rate and is still the most important parameter which one needs to tune.
- The other parameters are $\beta \in [0, 1)$ for the momentum exponential smoothing as used in (4.1) and $\gamma \in [0, 1)$ for the RMSprop exponential smoothing as is used in (4.18).
- ADAM is presented in Algorithm 1 where Step 4 is the gradient computation, steps 5 and 6 are exponential smoothing (momentum and RMSprop), steps 7 and 8 are bias corrections, and finally Step 9 is the descent step.

Remark 4.3.2

Since the introduction of ADAM in 2014, this algorithm became the most widely used algorithm (or “optimizer”) in deep learning frameworks. The common defaults for the hyper-parameters are $\beta = 0.9$ and $\gamma = 0.999$.

4.4 Dropout Method

- The idea of *dropout* is to randomly zero out certain neurons during the training process. This allows training to focus on multiple random subsets of the parameters and yields a form of regularization.
- With dropout, at any backpropagation iteration (forward pass and backward pass) on a mini-batch, only some random subset of the neurons is active.
- Practically neurons in layer ℓ , for $\ell = 0, \dots, L - 1$, have a specified probability $p_{\text{keep}}^{[\ell]} \in (0, 1]$ where if $p_{\text{keep}}^{[\ell]} = 1$ dropout does not affect the layer, and otherwise each neuron i of the layer is “dropped out” with probability $1 - p_{\text{keep}}^{[\ell]}$.
- This is simply a zeroing out of the neuron activation $a_i^{[\ell]}$ as we illustrate in Figure 4.3.

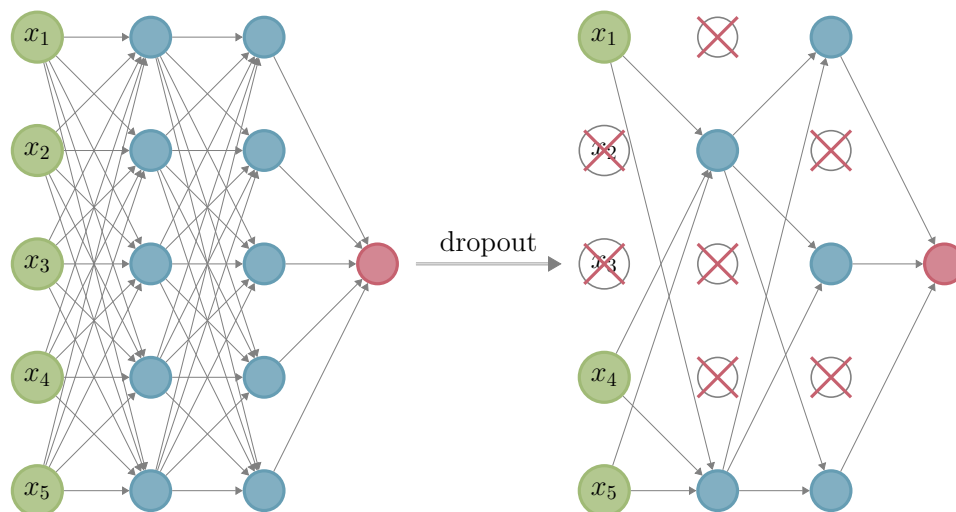


Figure 4.3: An illustration of dropout during training for a network with $L = 3$ layers, $p = 5$ input features, and $q = 1$ output. In each iteration during training, the network is transformed to the network on the right where the dropped out units are randomly selected.

- In the forward pass, when we get to the neurons of layer $\ell + 1$, all the neurons in layer ℓ that were zeroed out do not participate in the computation.
- Specifically, the update for a neuron j in the next layer, assuming a scalar activation function $\sigma(\cdot)$, becomes,

$$a_j^{[\ell+1]} = \sigma(b_j^{[\ell+1]} + \sum_{i \text{ kept}} w_{i,j}^{[\ell+1]} a_i^{[\ell]}). \quad (4.22)$$

- In the backward pass, when neuron i is dropped out in layer ℓ , the weights $w_{i,j}^{[\ell+1]}$ for all neurons $j = 1, \dots, N_{\ell+1}$ are updated based on the gradient $[g_W^{[\ell]}]_{i,j}$ which is set at 0.
- With a pure gradient descent optimizer this means that weights $w_{i,j}^{[\ell+1]}$ are not updated at all during the given iteration, whereas with a momentum based optimizer such as ADAM it means that the descent step for those weights has a smaller magnitude.
- At the end of training, even with dropout implemented, the trained model still has a complete set of weight matrices without any zeroed out elements, similar to the case in which we do not have dropout.
- Hence to account for the fact that neuron i in layer ℓ only took part in a proportion of iterations $p_{\text{keep}}^{[\ell]}$ during the training, when using the model in production (test time), we would like to use a weight matrix $\tilde{W}^{[\ell+1]} = p_{\text{keep}}^{[\ell]} W^{[\ell+1]}$ in place of $W^{[\ell+1]}$.
- The rationale here is to have the production forward pass similar to the training pass. Namely such a production forward pass has,

$$a_j^{[\ell+1]} = \sigma(b_j^{[\ell+1]} + \sum_{i=1}^{N_\ell} p_{\text{keep}}^{[\ell]} w_{i,j}^{[\ell+1]} a_i^{[\ell]}), \quad (4.23)$$

and this serves as an approximation to (4.22).

- In practice, the training–production pair (4.22)–(4.23) is not typically used per-se. The more practical alternative is instead of remembering $p_{\text{keep}}^{[\ell]}$ and deploying it with the production model, the training forward pass is modified to have the reciprocal of $p_{\text{keep}}^{[\ell]}$ as a scaling factor of the weights. Namely, the training forward pass is,

$$a_j^{[\ell+1]} = \sigma(b_j^{[\ell+1]} + \sum_{i \text{ kept}} \frac{1}{p_{\text{keep}}^{[\ell]}} w_{i,j}^{[\ell+1]} a_i^{[\ell]}). \quad (4.24)$$

This form allows to use the resulting model normally in production without having to take dropout into consideration at all. Namely, the production forward pass is,

$$a_j^{[\ell+1]} = \sigma(b_j^{[\ell+1]} + \sum_{i=1}^{N_\ell} w_{i,j}^{[\ell+1]} a_i^{[\ell]}), \quad (4.25)$$

- In practice, this simple and easy to implement idea of dropout has improved performance of deep neural networks in many empirically tested cases. It is now an integral part of deep learning training. We now explore the idea a bit further though the viewpoint of ensemble methods.

4.5 Weight Decay (L2) regularization

- In addition to dropout, addition of a regularization term is another key approach to prevent overfitting and improve generalization performance.
- Augmenting the loss with a regularization term $R_\lambda(\theta)$ restricts the flexibility of the model, and this restriction is sometimes needed to prevent overfitting.
- In the context of deep learning, and especially when ridge regression style regularization is applied, this practice is sometimes called *weight decay* when considering gradient based optimization.
- Take the original loss function $C(\theta)$ and augment it to be $\tilde{C}(\theta) = C(\theta) + R_\lambda(\theta)$. In our discussion here, let us focus on the ridge regression type regularization with parameter $\lambda > 0$, and,

$$R_\lambda(\theta) = \frac{\lambda}{2} R(\theta), \quad \text{with} \quad R(\theta) = \|\theta\|^2 = \theta_1^2 + \dots + \theta_d^2.$$

- Further, we may even restrict regularization to certain layers and not others.
- Now assume we execute basic gradient descent steps. With a learning rate $\alpha > 0$, the update at iteration t is,

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla \tilde{C}(\theta^{(t)}).$$

- In our ridge regression style penalty case we have $\nabla \tilde{C}(\theta) = \nabla C(\theta) + \lambda \theta$, and hence the gradient descent update can be represented as

$$\theta^{(t+1)} = (1 - \alpha\lambda)\theta^{(t)} - \alpha \nabla C(\theta^{(t)}). \quad (4.26)$$

- Now the interesting aspect of (4.26), assuming that $\alpha\lambda < 2$, is that it involves shrinkage or weight decay directly on the parameters in addition to gradient based learning. That is, independently of the value of the gradient $\nabla C(\theta^{(t)})$, in every iteration, (4.26) continues to decay the parameters, each time multiplying the previous parameter by a factor $1 - \alpha\lambda$.
- This weight decay phenomena can then be extended algorithmically to enforce regularization not directly via addition of a regularization term, but rather simply by augmenting the gradient descent updates to include weight decay.
- For example we may consider popular gradient based algorithms such as ADAM, and in each case add an additional step which incorporates multiplying the weights by a constant less than but close to unity.