

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2024

Contents

B	Optimizing Neural Networks	2
B.1	Strategies towards adaptive step-size	2
B.1.1	Momentum	2
B.1.2	Adadelta	10
B.1.3	Adagrad	14
B.1.4	Adam	15
B.2	Batch Normalization Details	19
B.2.1	Applying batch normalization to backprop	19

APPENDIX B

Optimizing Neural Networks

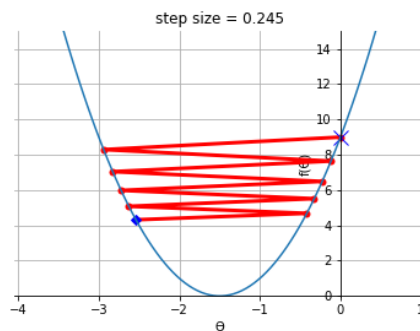
These notes carry over from the [Neural Networks](#) chapter. Here, we include topics that we previously skimmed over.

B.1 Strategies towards adaptive step-size

B.1.1 Momentum

B.1.1.1 Solving Oscillation

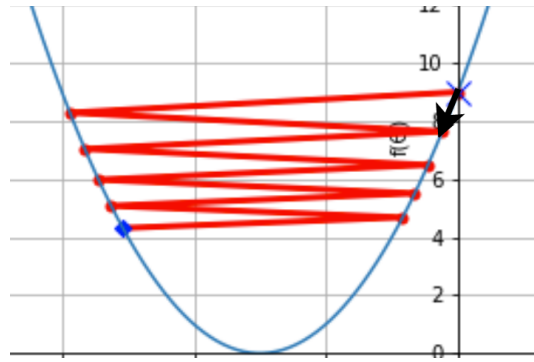
Let's look at one common problem we have with gradient descent: **oscillation**.



We overshoot our target, and then have to take another step that **undoes** most of what happened in the previous step. So, we waste a lot of time correcting the last step.

This can significantly **slow** down how quickly we converge.

For example, our first two steps land us in almost the same place we started!



The black arrow shows the combined effect of our first two steps: almost nothing!

We don't want to waste time, so we want to remove the "part" of the gradient that is likely to **cancel** out.

The **next** gradient cancels out some of the previous. Our first two steps add up, or "**average** out" to a small improvement.

If our steps are effectively "averaging", we'll speed up that process: we'll average together the gradients *before* taking our step!

Concept 1

Since our gradient descent steps **combine** to give us our new model, we can think of them as adding, or "**averaging**" to a more accurate improvement.

This means we can take a bigger step in a direction we won't have to cancel!

When our function **oscillates**, we get the same pattern **multiple** times: past steps indicate the sort of pattern we'll see in **future** steps.

So, we'll average our current gradient with past gradients: that way, the **component** that gets cancelled out is **removed**, and we won't have to undo our mistake over and over.

The only real difference between adding and averaging is whether we divide by the number of terms.

Concept 2

Oscillation causes us to move back and forth over the same region **multiple** times, where each step mostly **cancels** out the last.

One solution is to **average out** multiple of our gradients: the part that is "**cancelled** out" should be eliminated by the average.

So, we average our **past** gradients (past oscillation) with our **current** gradient, so we move in a more **efficient** direction, speeding up our algorithm.

Another way to think about it: when we **average** out our current and previous gradient, we're cancelling out what they "**disagree**" on, and keeping what they **agree** on.

So, we're taking a step in the direction that multiple gradients agree will **improve** our model!

B.1.1.2 Weighted averages

We could naively average all of our gradients **equally**. But, this would be a bad idea:

- It doesn't give you as much control of the algorithm: what if we care more about the **present** gradient, than the previous one?
- Gradients further in the **past** are **less likely** to matter: we've moved further away from those positions.
 - We also need to **scale** down past terms, so they don't take up most of the average.

The first problem is easy to solve: we'll **weigh** each of our terms differently.

If you're averaging 100 terms, and you add one more... it's not going to change much.

Concept 3

A **weighted average** is used when we want some terms to affect our **average** more than others.

We represent this with **weights**: each weight represents the **proportion** of our average from that term.

$$\text{Weighted Average} = x_1 w_1 + x_2 w_2 + \dots + x_n w_n$$

Example: If $w_1 = .6$, that means **60%** of the average comes from x_1 .

Note that, since we're talking about **proportions**, they need to **sum to 1**: it wouldn't make sense to have more than 100% of the average.

At each time step, we're adding one new gradient: the **present** one.

We'll simplify our average to those two terms: the **present** gradient, versus all the **past** gradients.

These weights are **separate** from the weights inside our neural network.

They do, however, represent the same type of concept: the NN weights scale the **input**, while these weights scale the **gradients**.

- We represent the importance (**weight**) of our **past** gradients using the variable γ .
- We want the two terms to add to 1: so, the importance of **current** gradient is $1 - \gamma$.

$$\underbrace{\text{Average}}_{\bar{A}_t} = \underbrace{\text{Old gradients}}_{\gamma G_{t-1}} + \underbrace{\text{New gradient}}_{(1-\gamma)g_t} \quad (\text{B.1})$$

Now, we have **control** over how much the present or past gradient matters: we just have to adjust γ .

B.1.1.3 Running Average

We still have some work to do: first, we haven't made it clear how we're incorporating our old gradients: we lumped them into one term.

Let's try building up from $t = 1$. We'll assume our previous gradients are 0, for simplicity.

$$A_0 = g_0 = 0 \quad (\text{B.2})$$

Our first step will average this with our **first** gradient:

$$A_1 = \gamma g_0 + (1 - \gamma) g_1 \quad (\text{B.3})$$

Simplifying to:

$$A_1 = (1 - \gamma) g_1 \quad (\text{B.4})$$

What about our second step?

$$A_2 = \overbrace{\gamma G_{t-1}}^{\text{Old gradients}} + (1 - \gamma) g_2 \quad (\text{B.5})$$

We *could* just plug in g_1 . But, A_1 contains the information about our first gradient g_1 , **and** the gradient before it, g_0 .

$$A_2 = \overbrace{\gamma A_1}^{\text{Contains } g_1, g_0} + (1 - \gamma) g_2 \quad (\text{B.6})$$

We can repeat this process:

$$A_3 = \overbrace{\gamma A_2}^{\text{Contains } g_2, g_1, g_0} + \overbrace{(1 - \gamma) g_3}^{\text{New gradient}} \quad (\text{B.7})$$

And so, we've created a general way to **average** as our program **runs** through different gradients.

$$A_t = \gamma A_{t-1} + (1 - \gamma) g_t \quad (\text{B.8})$$

To allow more flexibility, we'll allow γ to **vary in time**, as γ_t .

$$A_t = \gamma_t A_{t-1} + (1 - \gamma_t) a_t \quad (\text{B.9})$$

- We call this a **running average**.

Definition 4

A **running average** is a way to average past data with present data **smoothly**.

Our **initial** value for the average is typically zero:

$$A_0 = 0$$

Then, we begin introducing **new** data points.

- You use the parameter γ_t to indicate how much you want to prioritize **past data**.
- Thus, $1 - \gamma_t$ indicates the value of **new data**.

$$\underbrace{A_t}_{\text{Average}} = \underbrace{\gamma_t A_{t-1}}_{\text{Old gradients}} + \underbrace{(1 - \gamma_t) g_t}_{\text{New gradient}}$$

- Note that instead of γ , we write γ_t : this "discount factor" can vary with time.

Clarification 5

This is technically only one kind of **running average**: here, we use an "**exponential moving average**".

There are different ways to average past data points, with different **weighting** schemes.

- For example, you could do a "**simple moving average**", where you average equally over the last n data points.

B.1.1.4 Running Averages: The Distant Past

So, how does this "running average" approach affect our different data points, further in the past? Let's find out.

For simplicity, let's assume $\gamma_t = \gamma$: it's a **constant**.

$$A_t = \gamma A_{t-1} + (1 - \gamma) g_t \tag{B.10}$$

We can expand A_{t-1} to see further in the past:

$$A_t = \gamma \left(\gamma A_{t-2} + (1-\gamma)g_{t-1} \right) + (1-\gamma)g_t \quad (\text{B.11})$$

And even further:

This is starting to get messy: don't worry if it's hard to read.

$$A_t = \gamma \left(\gamma \left(\gamma A_{t-3} + (1-\gamma)g_{t-2} \right) + (1-\gamma)g_{t-1} \right) + (1-\gamma)g_t \quad (\text{B.12})$$

Let's rewrite this.

$$\gamma^3 A_{t-3} + \gamma^2 (1-\gamma)g_{t-2} + \gamma (1-\gamma)g_{t-1} + (1-\gamma)g_t \quad (\text{B.13})$$

We see a "stacking" effect for γ :

- We only partly include our newest data point: we scale it by $1-\gamma$, to make room for the past.

$$(1-\gamma)g_t \quad (\text{B.14})$$

- But if your gradient is 1 time unit in the past, we apply γ **once**, "forgetting" some more of that gradient.

$$\gamma (1-\gamma)g_{t-1} \quad (\text{B.15})$$

- But if your gradient is 2 units in the past, we apply γ **twice**: we've "forgotten" some of it twice.

$$\gamma^2 (1-\gamma)g_{t-2} \quad (\text{B.16})$$

Each time we do an average, we scale down our older data points by γ . So, the further in the past, the less effect they have.

- This is exactly the kind of design we wanted!

Concept 6

A **running average** tends to pay less attention to data further in the **past**.

- In general, if you are k time units in the past, we apply a factor of γ^k .

Because $\gamma < 1$, this **exponentially** decays to 0.

If we want to fully expand A_t , it's easiest to use a sum:

$$A_T = \sum_{t=0}^T \gamma^{(T-t)} \cdot (1 - \gamma) \cdot g_t$$

You can compare this formulation against what we computed above.

B.1.1.5 Momentum

Applying a running average to our gradients gives us **momentum**.

- This analogy to **physics** represents how our point "moves" through the **weight space**, to optimize J .
- The gradient gives us a "direction" of motion. So, our **momentum** represents the direction we were "already moving": the **previous** averaged gradient.

We use M to represent the "averaged gradient" that we use to move. Our initial momentum is 0:

$$M_0 = 0 \tag{B.17}$$

And we want to average our current gradient with past gradients:

$$M_t = \gamma M_{t-1} + (1 - \gamma) g_t \tag{B.18}$$

What is our **past** gradient g_t ? Well, we want to use W to modify J : $\frac{\partial J}{\partial W}$, or $\nabla_W J$.

And we're moving through the **weight space**, so our input to the gradient is the previous set of weights, W_{t-1} .

$$g_t = \nabla_W J(W_{t-1}) \tag{B.19}$$

And finally, M_t , our "averaged gradient", determines how we move.

Gradient descent adjusts our weights: we go from one list of weights, to another. That's why we say we're moving through the "weight space".

Because our hypothesis is determined by our weights, this is also a "hypothesis space".

$$W_t = W_{t-1} - \eta M_t \quad (\text{B.20})$$

Definition 7

Momentum is a technique for gradient descent where we do a **running average** between our current gradient, and our older gradients.

This approach reduces **oscillation**, and thus aims to improve the speed of convergence for our models.

- Our initial "momentum" (averaged gradient) is **0**.
- The amount we value new data is given by γ_t .
- Old data is scaled by $1 - \gamma_t$.

$$M_0 = 0$$

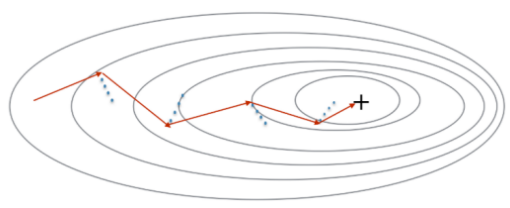
$$M_t = \underbrace{\gamma_t M_{t-1}}_{\text{Old gradients}} + \underbrace{(1 - \gamma_t) \nabla_W J(W_{t-1})}_{\text{New gradient}}$$

We use this momentum to take our step:

$$W_t = W_{t-1} - \eta M_t$$

This approach puts more emphasis on newer data points, and less on older ones.

Example: We can see this "dampened oscillation" in action below:



The blue lines indicate the more "severe" path that we would've taken with normal gradient descent. the orange line is the line we take with momentum.

Notice that generally, the orange path is closer to correct! We still oscillate somewhat, but much less than we would have otherwise.

B.1.2 Adadelta

B.1.2.1 Per-weight step sizes

Earlier, we discussed creating **separate** step sizes for different weights in our network:

- Different weights could have different **sensitivities**: one weight could have a much larger/smaller impact on J , based on **structure**.
- We already have the challenge of figuring out what **step sizes** cause **slow convergence/divergence**: it would be even harder to find a step size that fit **all** of our weights, at the same time.

So, instead, we decided that each weight gets its **own step size**.

- But, we never elaborated on **how** we compute that (adjustable) step size.

Notation 8

Our base step size is indicated as η .

We'll call the **modified** step size for weight j , η_j .

B.1.2.2 Scaling

Our goal now, is to come up with a **systematic** way to **assign step sizes**. This allows us to **adjust**, rather than run our whole gradient descent with a "bad" step size.

Why do we adjust our step size? To avoid slow convergence, oscillation, or divergence.

So, we have less of a problem if we choose η poorly.

- We might expect **slow convergence** if the derivative is too **small**: we carefully take small steps, but we aren't having much of an impact on J .
 - Across this "flatter" region of the surface, in one direction, we might expect it to be safer to move further.
- We might expect **divergence** or **oscillation** if the derivative is too **large**.
 - We might end up "missing" possible solutions/local minima by **overshooting** them.
 - So, "steeper" regions might be riskier.

Concept 9

Our goal is to **scale** our step size, so that it **adapts** to the situation:

- We want to **shrink** our step for **large** gradients
- We want to **increase** our step for **small** gradients

And we're interested in the **magnitude** of these gradients.

This sort of behavior is easily captured by including a factor of $1/\|g_t\|$. However, this has a smoothness problem: so, we'll use $1/\|g_t\|^2$ instead.

Though, we need to keep it separate for each of our data points.

Notation 10

The gradient for **weight j** at **time t** is given as

$$g_{t,j} = \nabla_w J(W_{t-1})_j$$

Note the double-subscript.

- By isolating weight j, we have a constant, not a vector.

We can now write our gradient update rule:

$$W_{t,j} = W_{t-1,j} - \eta_j \cdot g_{t,j} \quad (\text{B.21})$$

And we're currently using the step size

$$\eta_j = \frac{\eta}{g_{t,j}^2} \quad (\text{B.22})$$

Don't save this equation! It isn't our final formula.

B.1.2.3 Averaging

But, we don't necessarily know how "steep" our region **generally** is, based on the current gradient g_t . g_t only gives us **one point** in space.

It would be helpful to include information from the **past**: we'll be re-using the **weighted average**, once again.

$$G_t = \gamma G_{t-1} + (1 - \gamma) g_t^2 \quad (\text{Maybe?})$$

Once again, it's helpful that the weighted average gradually "forgets" older information: we care less about gradients which are "further" from the present.

- Note that we're averaging the **squared** magnitude.

Technically this is still **incorrect**: we need the j notation.

$$G_{t,j} = \gamma G_{t-1,j} + (1 - \gamma) g_{t,j}^2 \quad (\text{Fixed!})$$

It's incorrect because g_t is a vector: we can't square it directly, we have to square its magnitude.

B.1.2.4 Division by zero

There's a problem with our weight adjustment:

$$\eta_j = \frac{\eta}{G_{t,j}} \quad (\text{B.23})$$

What happens if the denominator is near zero? It'll explode to a huge number! And at zero, it's undefined.

To solve this, we'll add a very small constant, ϵ .

$$\eta_j = \frac{\eta}{G_{t,j} + \epsilon} \quad (\text{B.24})$$

We won't prescribe any particular choice of ϵ here.

Now, our scaling factor will never be bigger than $1/\epsilon$.

B.1.2.5 Square root

One last concern, to wrap up: currently, we're diving by the **squared** gradient. This is actually somewhat overkill.

Remember that our goal is to do the following operation:

$$W_{t,j} = W_{t-1,j} - \underbrace{\eta_j \cdot g_{t,j}}_{\text{Update}} \quad (\text{B.25})$$

With our current formula, our update has

- a factor of $g_{t,j}$ in the **numerator**
- from η_j , a factor proportional to $g_{t,j}^2$ in the **denominator**

This is "scaling" our gradient by more than we want to. So, we'll take the **square root** of the denominator.

$$\eta_j = \frac{\eta}{\sqrt{G_{t,j} + \epsilon}} \quad (\text{B.26})$$

Our goal is to make the scales of different axes more similar, not to neglect dimensions with high gradient (high effect on loss)

This is the completed form of our **adadelta step size rule**!

Definition 11

Adadelta is a technique for **adaptive step size**, which:

- **Decreases** step size in dimensions with a history of **high-magnitude** gradients
- **Increases** step size in dimensions with a history of **low-magnitude** gradients

Suppose our gradient for weight W_j at time t is represented by

$$g_{t,j} = \nabla_w J(W_{t-1})_j$$

This is accomplished by **scaling** the step size η to create η_j :

$$\eta_j = \frac{\eta}{\sqrt{G_{t,j} + \epsilon}}$$

Where $G_{t,j}$ is a "**running average**" of the previous gradients for weight W_j .

$$\begin{aligned} G_{0,j} &= 0 \\ G_{t,j} &= \gamma G_{t-1,j} + (1 - \gamma) g_{t,j}^2 \end{aligned}$$

So, our completed gradient descent rule takes the form:

$$W_{t,j} = W_{t-1,j} - \frac{\eta}{\sqrt{G_{t,j} + \epsilon}} \cdot g_{t,j}$$

B.1.2.6 Sparse Data

One major advantage of adadelta is its use for **sparse datasets**, where many variables only show up in a small percentage of the data.

- If a variable is much less frequent, then the weighted average G_t will be much smaller.
- So, when those data points **do** appear, the step size is much larger.

This allows our model to learn more from variables that don't show up as frequently.

Concept 12

Adadelta often works well with **sparse data**: datasets where many variables rarely show up as non-zero.

The step sizes for these variables become much larger, so our model can learn more from less-common information.

However, this can run the risk of paying attention to variables that are sparse, but **not** especially meaningful.

- It's important to choose your variables carefully, so the model doesn't "learn" from noise.

B.1.3 Adagrad

Originally, adadelta derives from a **simpler** method, known as **adagrad**, or "adaptive gradient".

- The main difference with this approach is, rather than find the **weighted average** G_t , we simply **sum** the previous gradients.

The main problem with this approach is that, over time, G_t becomes too **large**. So, our step size becomes too small, and our algorithm slows down.

By averaging, we don't run into this problem of G_t "accumulating": older data is gradually **forgotten**.

B.1.4 Adam

Momentum and Adadelata both bring some benefits:

- **Momentum** averages our current gradient with **previous gradients**, to reduce **oscillation** and make a more direct path to the solution.
- **Adadelata** modifies our **step sizes**: it takes smaller steps in directions of high gradient (reduce overshooting) and takes bigger steps in directions of low gradient (converge faster).

There's nothing structurally incompatible between them. So, why not incorporate both?

- This combination is called **Adam**: it has become the most popular way to handle step sizes in neural networks.

Concept 13

Adam integrates the techniques of both **momentum** and **adadelata**.

B.1.4.1 Momentum and Adadelata

We used two different **running averages**, both using γ for their "**discount factor**": how much we discount the effect of older data.

- We'll replace γ with B_1 (momentum) and B_2 (adadelata).

It's "discounting", or reducing the effect of older data, because $\gamma < 1$.

Notation 14

In the adam algorithm, B_1 and B_2 are **discount factors** replacing γ from momentum and adadelata.

We use the same notation for gradients:

$$g_{t,j} = \nabla_w J(W_{t-1})_j \quad (\text{B.27})$$

First, we'll keep track of our **averaged gradient**, $m_{t,j}$. This will be the **direction** we plan to move.

$$m_{t,j} = B_1 m_{t-1,j} + (1 - B_1) g_{t,j} \quad (\text{B.28})$$

Then, we'll keep track of the **average squared gradient**, $v_{t,j}$. This will tell us whether to scale up/down our **step size** on each variable.

$$v_{t,j} = B_2 v_{t-1,j} + (1 - B_2) g_{t,j}^2 \quad (\text{B.29})$$

We can combine these into our equation, the way you use momentum and adadelta:

$$W_{t,j} = W_{t-1,j} - \frac{\eta}{\sqrt{v_{t,j}} + \epsilon} \cdot m_{t,j} \quad (\text{B.30})$$

We're not quite done yet, though.

B.1.4.2 Normalization

There's one issue with our running average, that we should address.

Suppose we have a sequence of numbers:

$$[1, 1, 1] \quad (\text{B.31})$$

What's the running average of this sequence of numbers, with $\gamma = .1$? You'd expect it to be 1 for all elements, right?

- We start with $a_0 = 0$.

$$\begin{aligned} a_1 &= .1 * 0 + .9 * 1 = .9 \\ a_2 &= .1 * .9 + .9 * 1 = .99 \\ a_3 &= .1 * .99 + .9 * 1 = .999 \end{aligned} \quad (\text{B.32})$$

It turns out not to be true! Why is that?

Because of our initial term, a_0 : it has nothing to do with the data. Because it'll always be 0, it **deflates** the value of all the remaining averages.

How much does it deflate the average? Well, let's consider the general equation:

$$A_T = \sum_{t=0}^T \gamma^{(T-t)} (1 - \gamma)^t g_t \quad (\text{B.33})$$

What "fraction" of the total is missing, thanks to $A_0 = 0$?

- Well, all of our weighted terms $\gamma^{(T-t)} (1 - \gamma)^t$ should add up to 1: each one represents the "percent/fraction" of the average which comes from that g_t term.

A_0 matches $t = 0$, so we find:

$$A_T = \overbrace{A_0 \gamma^T}^{A_0=0} + \sum_{t=1}^T \gamma^{(T-t)} (1 - \gamma)^t g_t \quad (\text{B.34})$$

So, out of our total 1, or 100%, we're missing γ^T .

- Our new total is $1 - \gamma^T$.
- In order to correct for the "zeroed out" part of our average, we multiply by

$$\frac{\text{Desired total}}{\text{Real total}} = \frac{1}{1 - \gamma^T} \quad (\text{B.35})$$

Concept 15

We've chosen $A_0 = 0$: this is **unrelated** to our real data.

As a result, A_t doesn't accurately reflect our weighted average: it's slightly **smaller**.

- A_0 is "included" as a 0, even though it's not a **real** data point.
- So, whatever **percent** of our data is represented by A_0 , is "falsely empty".

A_0 is scaled by γ^t , so that's the amount **removed** from our "true" weighted average.

The "real" weighted average, that ignores A_0 , has to un-do the deflation caused by including fake, starting data:

$$\hat{A}_t = A_t \cdot \frac{1}{1 - \gamma^t}$$

We'll make this correction for our running averages:

$$\begin{aligned} \hat{m}_{t,j} &= \frac{m_{t,j}}{1 - B_1^t} \\ \hat{v}_{t,j} &= \frac{v_{t,j}}{1 - B_2^t} \end{aligned} \quad (\text{B.36})$$

B.1.4.3 Adam

Now, we have our complete equation for adam:

Definition 16

Adam is a technique for improving **gradient descent** that integrates two other techniques. Our computed gradient is given as

$$g_{t,j} = \nabla_w J(W_{t-1})_j$$

- **Momentum** averages our previous gradients with our current one, to reduce oscillation and get a "averaged" view of data. B_1 gives the weight of old data.

$$m_{0,j} = 0 \quad m_{t,j} = B_1 m_{t-1,j} + (1 - B_1) g_{t,j}$$

- **Adadelta** estimates how "steep" our gradient has been recently, and uses it to adjust our step size, for better convergence. B_2 gives the weight of old data.

$$v_{0,j} = 0 \quad v_{t,j} = B_2 v_{t-1,j} + (1 - B_2) g_{t,j}^2$$

We include a **correction** factor for the zeroed initial conditions: $m_0 = v_0 = 0$.

$$\hat{m}_{t,j} = \frac{m_{t,j}}{1 - B_1^t} \quad \hat{v}_{t,j} = \frac{v_{t,j}}{1 - B_2^t}$$

Finally, our update rule takes the form:

$$W_{t,j} = W_{t-1,j} - \frac{\eta}{\sqrt{\hat{v}_{t,j}} + \epsilon} \cdot \hat{m}_{t,j}$$

Impressively, adam is not very sensitive to the initial conditions for ϵ , B_1 and B_2 .

To implement this for NNs, we just need to keep track of each value, for each weight. If you do it systematically, this is easier than it sounds.

It's able to "self-correct" its step sizes and gradient based on data.

B.2 Batch Normalization Details

Let's review the general premise

of batch normalization:

Definition 17

Batch Normalization is a process where we

- Standardize the pre-activation for each layer using the mean μ_i and standard deviation σ_i (for the i^{th} dimension). Select ϵ : $0 < \epsilon \ll 1$.

$$\bar{z}_{ij} = \frac{z_{ij} - \mu_i}{\sigma_i + \epsilon}$$

- Choose the new mean and standard deviation for the pre-activation using $(n \times 1)$ vectors G and B

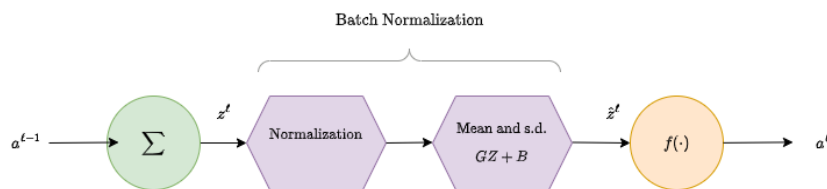
$$\hat{z}_{ik} = G_i * \bar{z}_{ij} + B_i$$

B.2.1 Applying batch normalization to backprop

Remark (Optional)

The following section mostly deals with the details of computing batch normalization derivatives.

As we showed with our figure,



Batch normalization adds two units that **interrupt** our chain of nested functions.

That means we need to figure out how to do backprop, bridging across the new gap between z and \hat{z} .

So, that only leaves a couple problems:

- "Bridging the gap" between derivatives before and after BN, with $\partial \hat{z}^\ell / \partial z^\ell$.
- Gradients for G and B : they're parameters now, too, so we need to train them.

Concept 18

Introducing batch normalization add new functions **in between** our old ones.

- Our input data has to travel through those additional layers.
- This **changes** the relationship between our current value, and the output loss.

So, in order to do backprop correctly, we have to figure out the **derivatives** of those functions.

B.2.1.1 Bridging the gap

We want to connect the start and end of batch normalization:

$$\frac{\partial \hat{Z}}{\partial Z} \quad (\text{B.37})$$

As usual, with the chain rule, we'll connect them by considering any values/function **between** them.

In this case, Z is normalized (\bar{Z}), and **then** we apply G and B (\hat{Z}). We're missing the "normalized" step.

$$\frac{\partial \hat{Z}}{\partial Z} = \frac{\partial \hat{Z}}{\partial \bar{Z}} \cdot \frac{\partial \bar{Z}}{\partial Z} \quad (\text{B.38})$$

Let's compare any two of these : \hat{Z} and \bar{Z} , for example.

- These are both $(n \times k)$ matrices.
- That means that each has two dimensions of variables. If we were to take the derivative between them, we would need $2 * 2$ axes: a **4-tensor**.

That sounds terrible. Instead, we'll compute these derivatives **element-wise**.

$$\frac{\partial \hat{Z}_{ab}}{\partial Z_{ef}} = \frac{\partial \hat{Z}_{ab}}{\partial \bar{Z}_{cd}} \cdot \frac{\partial \bar{Z}_{cd}}{\partial Z_{ef}} \quad (\text{B.39})$$

~~~~~

**B.2.1.2 Indexing**

First, let's simplify these indices: only some pairs of elements matter.

$$\hat{Z}_{ik} = G_i * \bar{Z}_{ik} + B_i \quad (\text{B.40})$$

It seems that  $\hat{Z}_{ik}$  is only affected by the element with the same indices:  $\bar{Z}_{ik}$ .

### Concept 19

$\hat{Z}_{ik}$  is only a function of the same index element  $\bar{Z}_{ik}$ .

- Any other elements from  $\bar{Z}$  has no effect.

$$(a \neq c \text{ or } b \neq d) \implies \frac{\partial \hat{Z}_{ab}}{\partial \bar{Z}_{cd}} = 0$$

So, we write our remaining derivatives as  $\partial \hat{Z}_{ik} / \partial \bar{Z}_{ik}$ .

What about the other derivative?

$$\bar{Z}_{ik} = \frac{Z_{ik} - \mu_i}{\sigma_i + \epsilon}$$

$\mu_i$  and  $\sigma_i$  include various different data points  $Z_{ik}$ , but only the  $i^{\text{th}}$  dimension.  $Z_{ik}$  requires the exact same indices.

### Concept 20

$\bar{Z}_{ik}$  is only a function of elements in the same dimension  $i$ ,  $Z_{ij}$ .

$$c \neq e \implies \frac{\partial \bar{Z}_{cd}}{\partial Z_{ef}} = 0$$

Our remaining derivatives take the form  $\partial \bar{Z}_{ik} / \partial Z_{ij}$ .

If we boil this down, we get all of our non-zero derivatives:

$$\frac{\partial \hat{Z}_{ik}}{\partial Z_{ij}} = \frac{\partial \hat{Z}_{ik}}{\partial \bar{Z}_{ik}} \cdot \frac{\partial \bar{Z}_{ik}}{\partial Z_{ij}} \quad (\text{B.41})$$

#### B.2.1.3 Computing $\partial \hat{Z}_{ik} / \partial \bar{Z}_{ik}$

We return to our previous equation:

$$\hat{Z}_{ik} = G_i * \bar{Z}_{ik} + B_i \implies \frac{\partial \hat{Z}_{ik}}{\partial \bar{Z}_{ik}} = G_i \quad (\text{B.42})$$

#### B.2.1.4 Computing $\partial \bar{Z}_{ik} / \partial Z_{ij}$

For the other derivative:

$$\bar{Z}_{ik} = \frac{Z_{ik} - \mu_i}{\sigma_i + \epsilon} \quad (\text{B.43})$$

This gets a bit complicated, because  $Z_{ij}$  can affect three different terms:  $Z_{ik}$ ,  $\mu_i$ , and  $\sigma_i$ .

We'll solve this by using the multi-variable chain rule.

We're linearly adding the effect due to each of these three variables, separately.

$$\frac{\partial \bar{Z}_{ik}}{\partial Z_{ij}} = \overbrace{\frac{\partial \bar{Z}_{ik}}{\partial Z_{ik}} \cdot \frac{dZ_{ik}}{dZ_{ij}}}^{\text{Z}_{ik}\text{'s effect}} + \overbrace{\frac{\partial \bar{Z}_{ik}}{\partial \mu_i} \cdot \frac{d\mu_i}{dZ_{ij}}}^{\mu_i\text{'s effect}} + \overbrace{\frac{\partial \bar{Z}_{ik}}{\partial \sigma_i} \cdot \frac{d\sigma_i}{dZ_{ij}}}^{\sigma_i\text{'s effect}} \quad (\text{B.44})$$

In each of these terms, we treat the other two variables as "constant".

### B.2.1.5 Lots of mini-derivatives

Let's compute the  $\bar{Z}$  derivatives.

$$\bar{Z}_{ik} = \frac{Z_{ik} - \mu_i}{\sigma_i + \epsilon} \quad (\text{B.45})$$

gives us

$$\frac{\partial \bar{Z}_{ik}}{\partial Z_{ik}} = \frac{1}{\sigma_i + \epsilon} \quad \frac{\partial \bar{Z}_{ik}}{\partial \mu_i} = \frac{-1}{\sigma_i + \epsilon} \quad \frac{\partial \bar{Z}_{ik}}{\partial \sigma_i} = -\left(\frac{Z_{ik} - \mu_i}{(\sigma_i + \epsilon)^2}\right) \quad (\text{B.46})$$

Now, let's compute the  $Z_{ij}$  derivatives.

$$\boxed{\frac{\partial Z_{ik}}{\partial Z_{ij}} = \delta_{jk}} = \mathbf{1}(j = k) = \begin{cases} 1 & j = k \\ 0 & j \neq k \end{cases} \quad (\text{B.47})$$

$$\mu_i = \frac{1}{K} \sum_{j=1}^K Z_{ij} \Rightarrow \boxed{\frac{\partial \mu_i}{\partial Z_{ij}} = \frac{1}{K}} \quad (\text{B.48})$$

$$\sigma_i^2 = \frac{1}{K} \sum_{j=1}^K (Z_{ij} - \mu_i)^2 \Rightarrow \boxed{\frac{\partial \sigma_i}{\partial Z_{ij}} = \frac{Z_{ij} - \mu_i}{K\sigma_i}} \quad (\text{B.49})$$

### B.2.1.6 Assembling our derivatives

Now, we plug them in.

$$\frac{\partial \bar{Z}_{ik}}{\partial Z_{ij}} = \frac{\partial \bar{Z}_{ik}}{\partial Z_{ik}} \cdot \frac{dZ_{ik}}{dZ_{ij}} + \frac{\partial \bar{Z}_{ik}}{\partial \mu_i} \cdot \frac{d\mu_i}{dZ_{ij}} + \frac{\partial \bar{Z}_{ik}}{\partial \sigma_i} \cdot \frac{d\sigma_i}{dZ_{ij}} \quad (\text{B.50})$$

First, the  $\bar{Z}$  derivatives.

$$\frac{\partial \bar{z}_{ik}}{\partial z_{ij}} = \left( \frac{1}{\sigma_i + \epsilon} \right) \cdot \frac{dz_{ik}}{dz_{ij}} - \left( \frac{1}{\sigma_i + \epsilon} \right) \cdot \frac{d\mu_i}{dz_{ij}} - \left( \frac{z_{ik} - \mu_i}{(\sigma_i + \epsilon)^2} \right) \cdot \frac{d\sigma_i}{dz_{ij}} \quad (\text{B.51})$$

And now the  $z_{ij}$  derivatives.

$$\frac{\partial \bar{z}_{ik}}{\partial z_{ij}} = \left( \frac{1}{\sigma_i + \epsilon} \right) \cdot \delta_{jk} - \left( \frac{1}{\sigma_i + \epsilon} \right) \cdot \left( \frac{1}{K} \right) - \left( \frac{z_{ik} - \mu_i}{(\sigma_i + \epsilon)^2} \right) \cdot \left( \frac{z_{ij} - \mu_i}{K\sigma_i} \right) \quad (\text{B.52})$$

### Key Equation 21

We have found the batch normalization derivatives

$$\frac{\partial \hat{z}_{ik}}{\partial \bar{z}_{ik}} = G_i$$

$$\frac{\partial \bar{z}_{ik}}{\partial z_{ij}} = \frac{1}{K(\sigma_i + \epsilon)} \left( K\delta_{jk} - 1 - \frac{(z_{ik} - \mu_i)(z_{ij} - \mu_i)}{\sigma_i(\sigma_i + \epsilon)} \right)$$

Which we multiply together to find  $\partial \hat{z}_{ik} / \partial z_{ij}$ .

Once we've computed these derivatives, we can use them to extend the chain of  $\partial L / \partial \hat{z}_{ik}$ .

We're aiming to handle both of the batch normalization functions, with  $\partial L / \partial z_{ij}$ .

- $z_{ij}$  affects every data point  $\hat{z}_{ik}$ , and every data point affects  $L$ .
- So, we'll have to consider every data point  $k$  using the multi-variable chain rule:

We're going from  $\hat{z}_{ik}$ , which is post-BN, to  $z_{ij}$ , which is pre-BN.

$$\frac{\partial L}{\partial z_{ij}} = \underbrace{\sum_{k=1}^K}_{\text{MV Chain Rule}} \overbrace{\frac{\partial L}{\partial \hat{z}_{ik}} \cdot \frac{\partial \hat{z}_{ik}}{\partial z_{ij}}}^{\text{Data point } k\text{'s effect}} \quad (\text{B.53})$$

We can use this to go further back in our layers: as far as we want, as long as we don't forget this derivative!

### B.2.1.7 Gradients for B and G

We want  $\partial L / \partial \mathbf{B}$  and  $\partial L / \partial \mathbf{G}$ . Thanks to the work we did just now, we can travel backwards through numerous layers, to reach any  $\mathbf{B}^\ell$  and  $\mathbf{G}^\ell$ .

So, we'll assume we have  $\partial L / \partial \hat{\mathbf{Z}}^\ell$ . Once again, we'll omit the layer notation.

- We'll focus on a single bias,  $B_i$ : this biases one dimension of  $\hat{z}_{ik}$ .

$$\hat{z}_{ik} = G_i * \bar{z}_{ik} + B_i \implies \frac{\partial \hat{z}_{ik}}{\partial B_i} = 1 \quad (\text{B.54})$$



- This bias is applied to every of our K data points: every data point affects the loss L. So, we'll have to sum them with the multi-variable chain rule.

This is exactly the same as how we did  $\partial L / \partial \mathbf{z}_{ij}$ .

$$\frac{\partial L}{\partial \mathbf{B}_i} = \sum_{k=1}^K \frac{\partial L}{\partial \hat{\mathbf{z}}_{ik}} \cdot \frac{\partial \hat{\mathbf{z}}_{ik}}{\partial \mathbf{B}_i} = \boxed{\sum_{k=1}^K \frac{\partial L}{\partial \hat{\mathbf{z}}_{ik}}} \quad (\text{B.55})$$

Now, we focus on a single scaling factor  $G_i$ :

$$\hat{\mathbf{z}}_{ik} = G_i * \bar{\mathbf{z}}_{ik} + \mathbf{B}_i \implies \frac{\partial \hat{\mathbf{z}}_{ik}}{\partial G_i} = \bar{\mathbf{z}}_{ik} \quad (\text{B.56})$$

And once again, we add across different data points:

$$\frac{\partial L}{\partial G_i} = \sum_{k=1}^K \frac{\partial L}{\partial \hat{\mathbf{z}}_{ik}} \cdot \frac{\partial \hat{\mathbf{z}}_{ik}}{\partial G_i} = \boxed{\sum_{k=1}^K \frac{\partial L}{\partial \hat{\mathbf{z}}_{ik}} \bar{\mathbf{z}}_{ik}} \quad (\text{B.57})$$

We're finished.

### Key Equation 22

Here, we have the gradients for the batch scaling coefficients, B and G.

$$\frac{\partial L}{\partial \mathbf{B}_i} = \sum_{k=1}^K \frac{\partial L}{\partial \hat{\mathbf{z}}_{ik}}$$

$$\frac{\partial L}{\partial G_i} = \sum_{k=1}^K \frac{\partial L}{\partial \hat{\mathbf{z}}_{ik}} \bar{\mathbf{z}}_{ik}$$