

Mini-batch optimization

$$X = \left[\begin{array}{c|c|c} \boxed{\begin{array}{c} \vec{x}^{(1)} \\ \vec{x}^{(2)} \end{array}} & \boxed{\dots} & \boxed{\vec{x}^{(n_s)}} \\ \text{mini-batch 1} & \text{mini-batch 2} & \end{array} \right]$$

Feed the data one mini-batch at a time.

One full pass through the data is called an "epoch".

For i in range(n -epochs):

For j in range(n -batch):

$$n_batch = \frac{n_s}{\text{batch-size}}$$

Feed-forward

Back-prop

update w & b

If we assume that the time of an epoch is mostly dominated by the n_s and has very little dependence on the # batches, we can roughly take the time ^{of an} epoch as the unit of time.

Assignment: Check this!

Plot the time of an epoch vs n_s for diff batch-sizes.

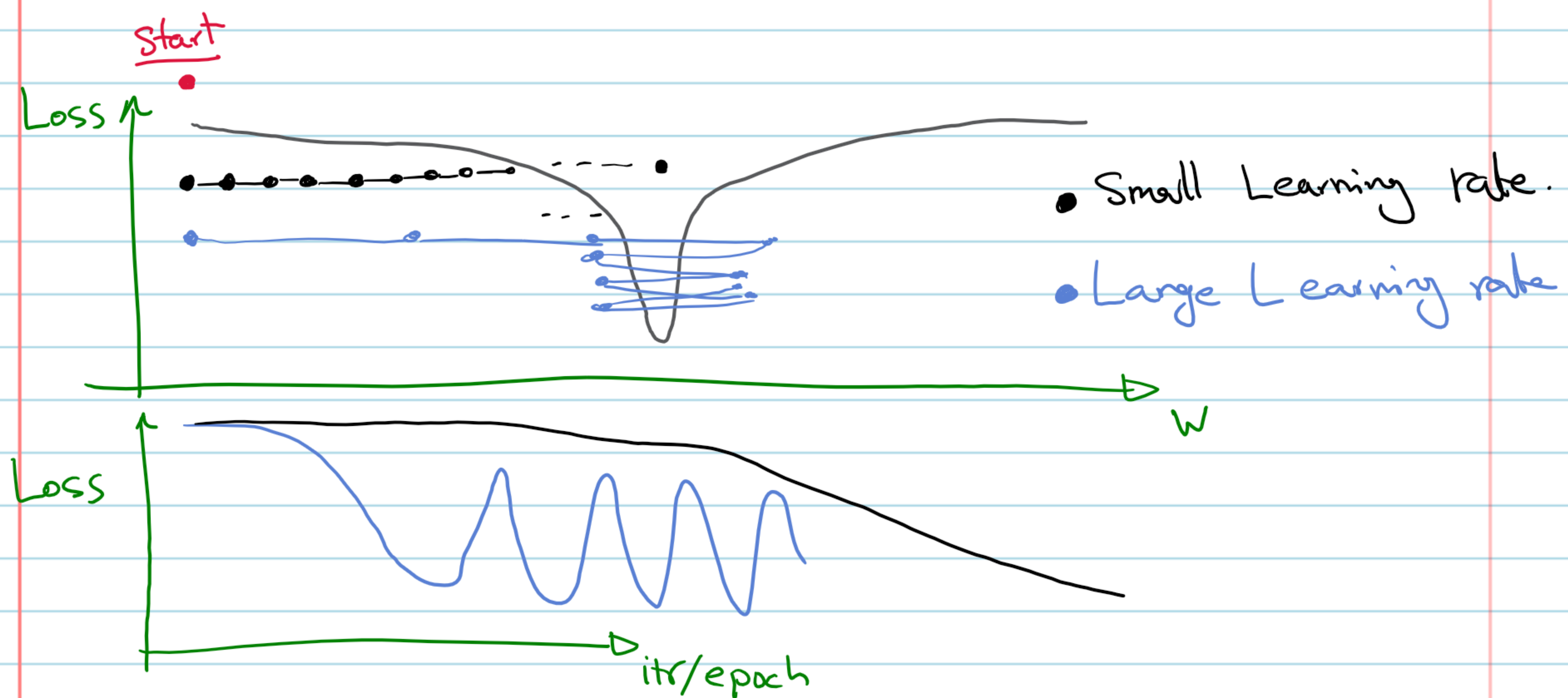
The advantage is that

- * First, we use the data more efficiently, i.e. do more GD steps in each epoch.

- * Second, the mini-batch GD adds stochasticity to the optimization process that can sometimes help escape local minimum.

Adaptive Learning Rates

Consider the following optimization problem:



With a small lr , it takes too many itr/epochs to converge to the minimum. Specially for the flat part of the loss, ($\frac{dL}{dw} \sim 0$) the training would be slow and the progress would only start after too many itr/epochs.

A large lr would be helpful to get out of the flat part of the loss function but then it cannot converge.

Ideally, we would need to use an Adaptive lr that depending on where it is, it would use a small or large value for the lr .

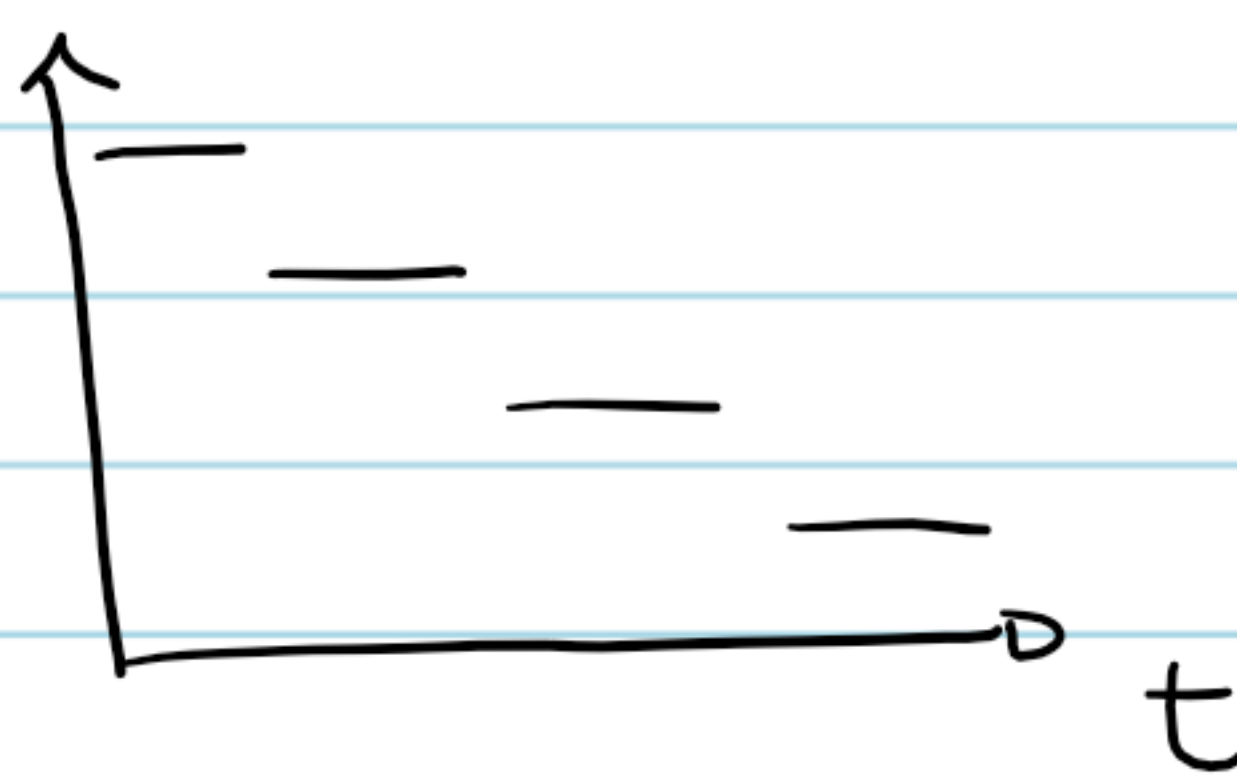
Learning Rate decay:

One solution is to decrease the lr.

1) $\alpha = \frac{k}{\sqrt{t}} \alpha_0$ k some const. (hyper param)
 t time \rightarrow itr/epoch

2) $\alpha = e^{-kt} \alpha_0$ \rightarrow Exponentially decaying

3) α : Stair case



But, this does not fully resolve the problem:

* This still gives the same lr for all the params in the problem and they can be on diff. scales.

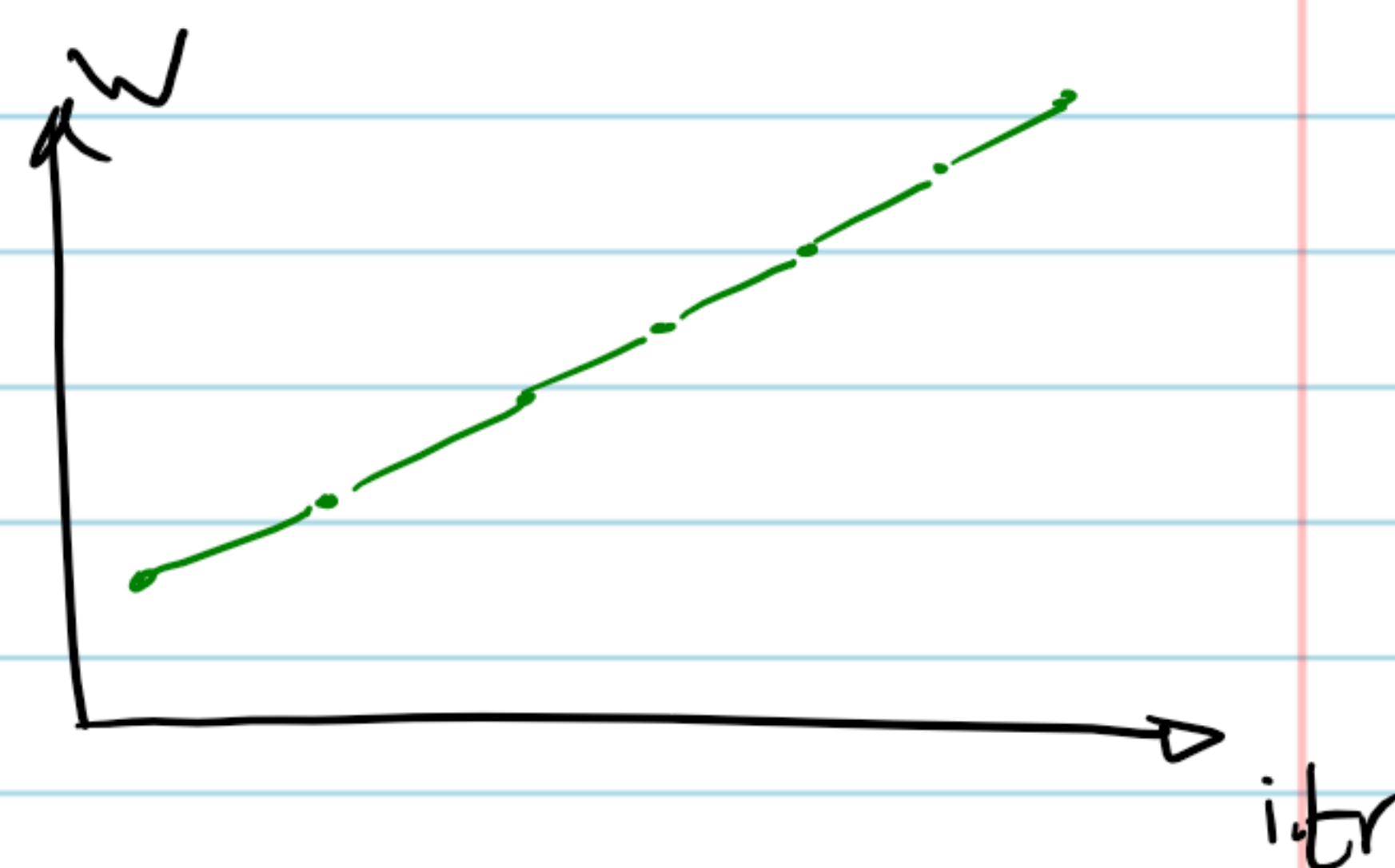
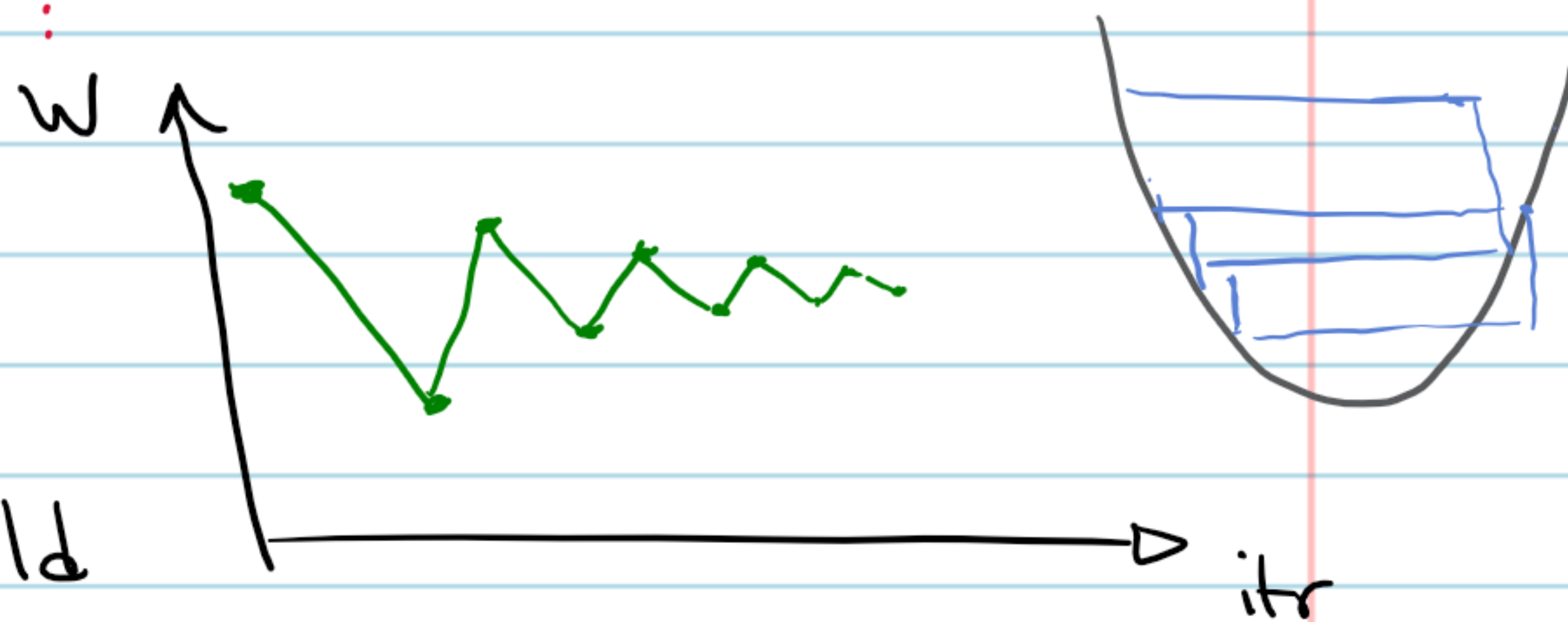
* Whether the lr should decrease or not depends on the loss and there could be situations where this does not help.

Using the history for adaptive lr:

Consider this situation.

From the oscillations in W , one could guess that the lr should decrease.

Also if, we can guess that a larger lr maybe helpful.



How can we implement our intuition?

We need to look at the history, but probably the more recent history. To do this, we apply a weighted average to the history (of changes, i.e. dW) and take that for the update:

$$W = W - \eta dW \rightarrow W = W - \eta \overline{dW}$$

Exponential-weighted average is a good choice, i.e. $\frac{w_{t-j}}{w_{t-j+1}} = \beta, 0 < \beta < 1$

At each step $t+1$, we calculate dW_{t+1} (from the loss func).

We also keep the $\{dW_t\}$ list which has the full history of dW , then calculate

$$\overline{dW}_{t+1} = \sum_{t'=0}^t w_{t'} dW_{t'} + w_{t+1} dW_{t+1}$$

$$W = W - \eta \overline{dW}$$

For the ^{exp} weights, we really need to specify the ballance between the history and the current value. This is b/c

\overline{dW}_t has all the history with the right ratio.

$$\overline{dW}_{t+1} = \beta \overline{dW}_t + (1-\beta) dW_{t+1}$$

With the recursion, we get

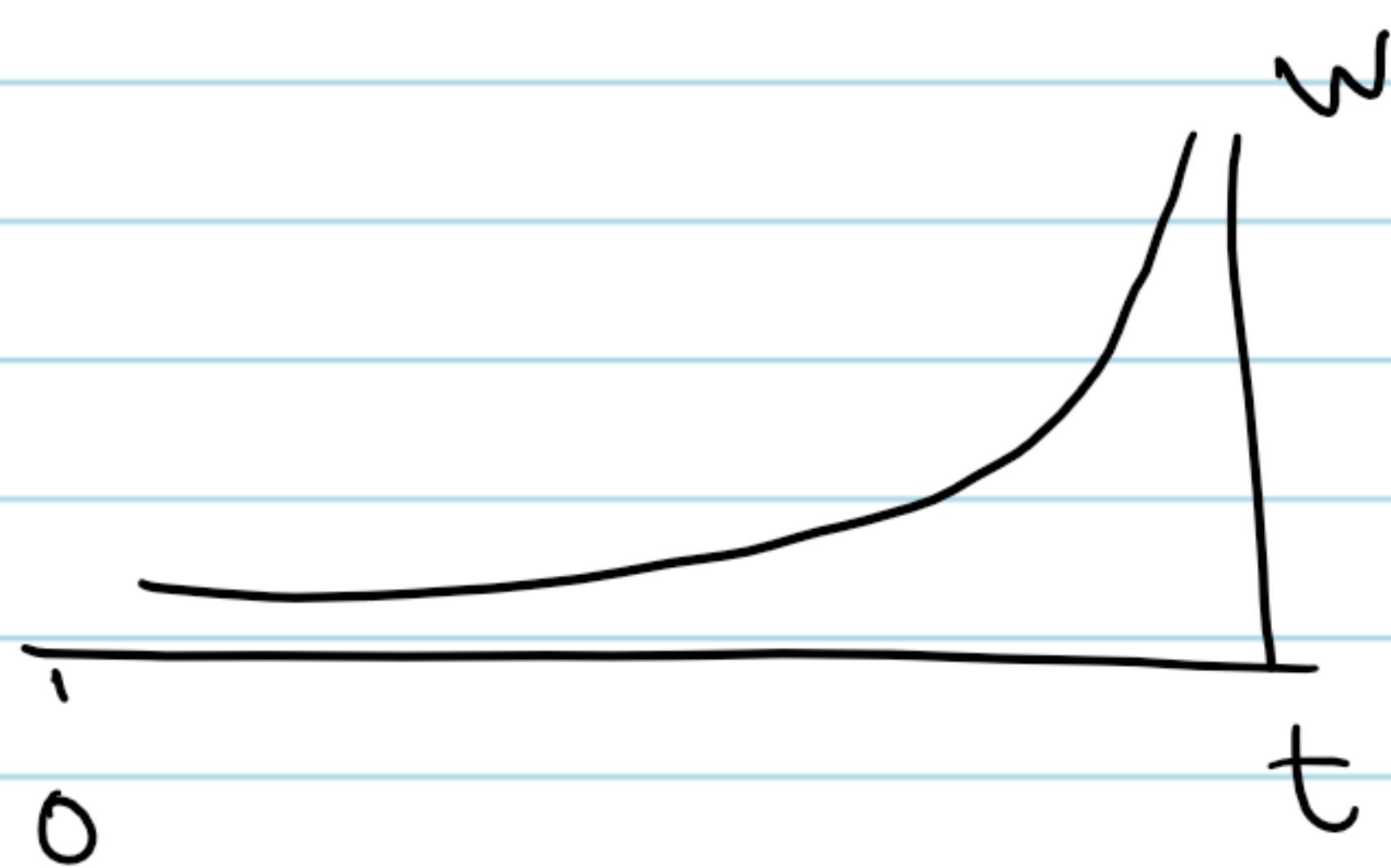
$$\overline{dW}_{t+1} = \beta^2 \overline{dW}_{t-1} + \beta(1-\beta) dW_t + (1-\beta) dW_{t+1}$$

\vdots

$$= \underbrace{dW_0}_0 + (1-\beta) \beta^t dW_1 + \beta^{t-1} (1-\beta) dW_2 + \dots \beta^{t-k} (1-\beta) dW_{k+1} + \dots (1-\beta) dW_{t+1}$$

If we plot these weight, they would make an exp function

$$W_t \propto \beta^{-t}$$



So, this would mean for each step

Feed Forward

Backprop $\rightarrow dW_{t+1}$

$$\overline{dW} = \beta \overline{dW} + (1-\beta) dW_{t+1}$$

$$W \leftarrow W + \eta \overline{dW}$$

β would be a hyperparameter that determines how much we would want to weigh the history.

$\beta = 0 \rightarrow$ No history

$(\beta \rightarrow 1) \rightarrow$ slower decay

remember $\frac{w_{t-1}}{w_t} = \beta$.



With this, we can pick a large β & if it starts fluctuating, the weighted average would be suppress the large values of dW .

\rightarrow See the example in Jupyter NB.

Example:

$dW: \{ 0.5, -0.45, 0.42, -0.4 \} , dW_{t+1} = 0.38$

$\beta = 0.9$ $(0.5) \overset{?}{\rightarrow} \frac{\overline{dW}}{0.05}$

$$(0.05)(0.9) + (-0.45)(0.1) \rightarrow 0$$

$$(0.45)(0.9) + (0.42)(0.1) \rightarrow 0.64$$

\vdots

\rightarrow Needs some work.

RMS Prop:

Alternatively, one could change the lr based on the variance of the dW .

A small var, indicates that the lr is good, on the other hand, a large var, indicates signals that the lr may be large & we may need to decrease it.

It means

$$W \rightarrow W - \eta \frac{dW}{\sqrt{\text{Var}}} \quad \overline{\text{Var}}_{t+1} = \beta \overline{\text{Var}}_t + (1-\beta) \text{Var}_{t+1} \downarrow dW^2$$

We add ϵ to avoid divergence:

$$W - \eta \frac{dW}{\sqrt{\text{Var} + \epsilon}}$$

Similar to the momentum, we use exp-weighted averaging to calculate the variance.

Adam

This combines the rms prop & momentum:

In step $(t+1)$

$$\overline{dW} = (\beta_1 \overline{dW} + (1-\beta_1) dW_{t+1}) / (1-\beta_1^t)$$

This is to compensate for the bias in the first iterations.

$$\overline{Var} = (\beta_2 \overline{Var} + (1-\beta_2) dW_{t+1}^{**2}) / (1-\beta_2^t)$$

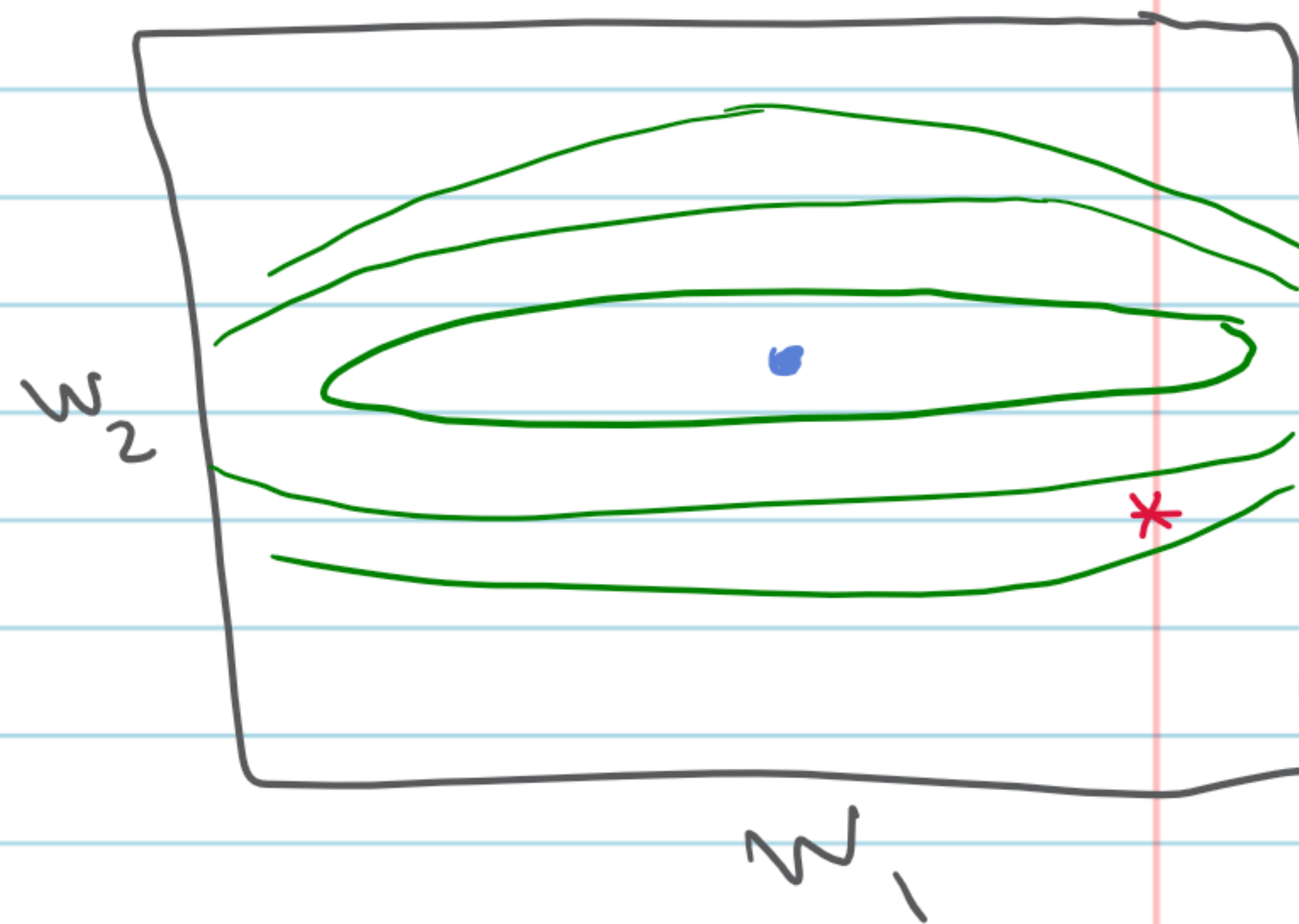
$$W := \eta \frac{\overline{dW}}{\sqrt{\overline{Var} + \epsilon}}$$

The hyperparams are $\{\eta, \beta_1, \beta_2, \epsilon\}$.

\downarrow \downarrow \downarrow
0.9 0.99 10^{-8}

Batch normalization

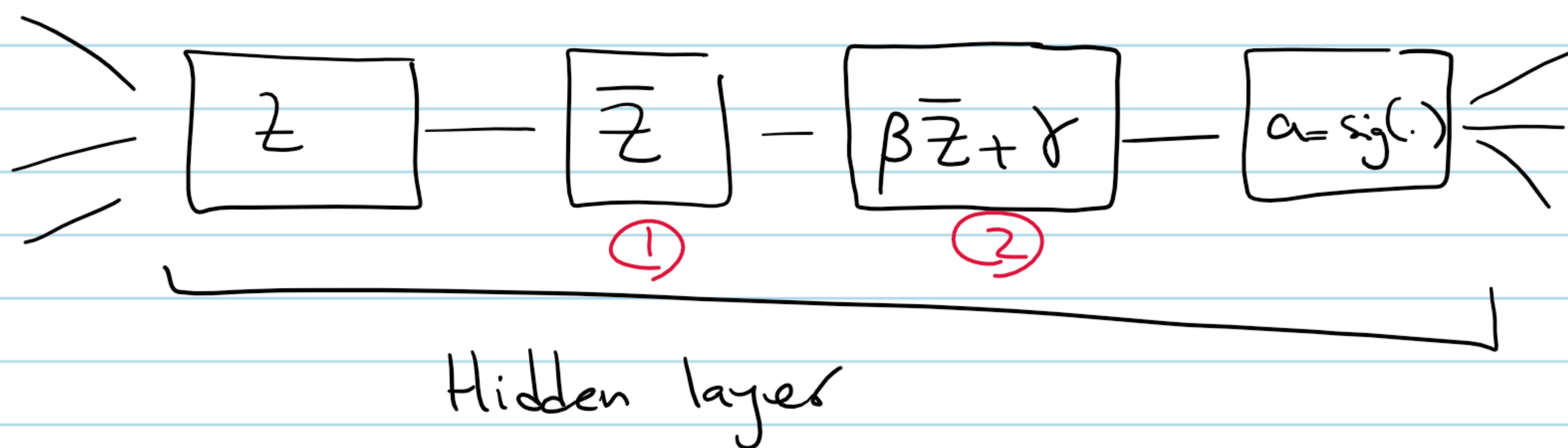
Consider the following situation. This would be hard b/c a large step with w_1 would be small for w_2 .



This is why we usually normalize / standardize our inputs.

But this still could happen for a NN since the hidden layers are not constrained.

But this can be done:



① Standardize z

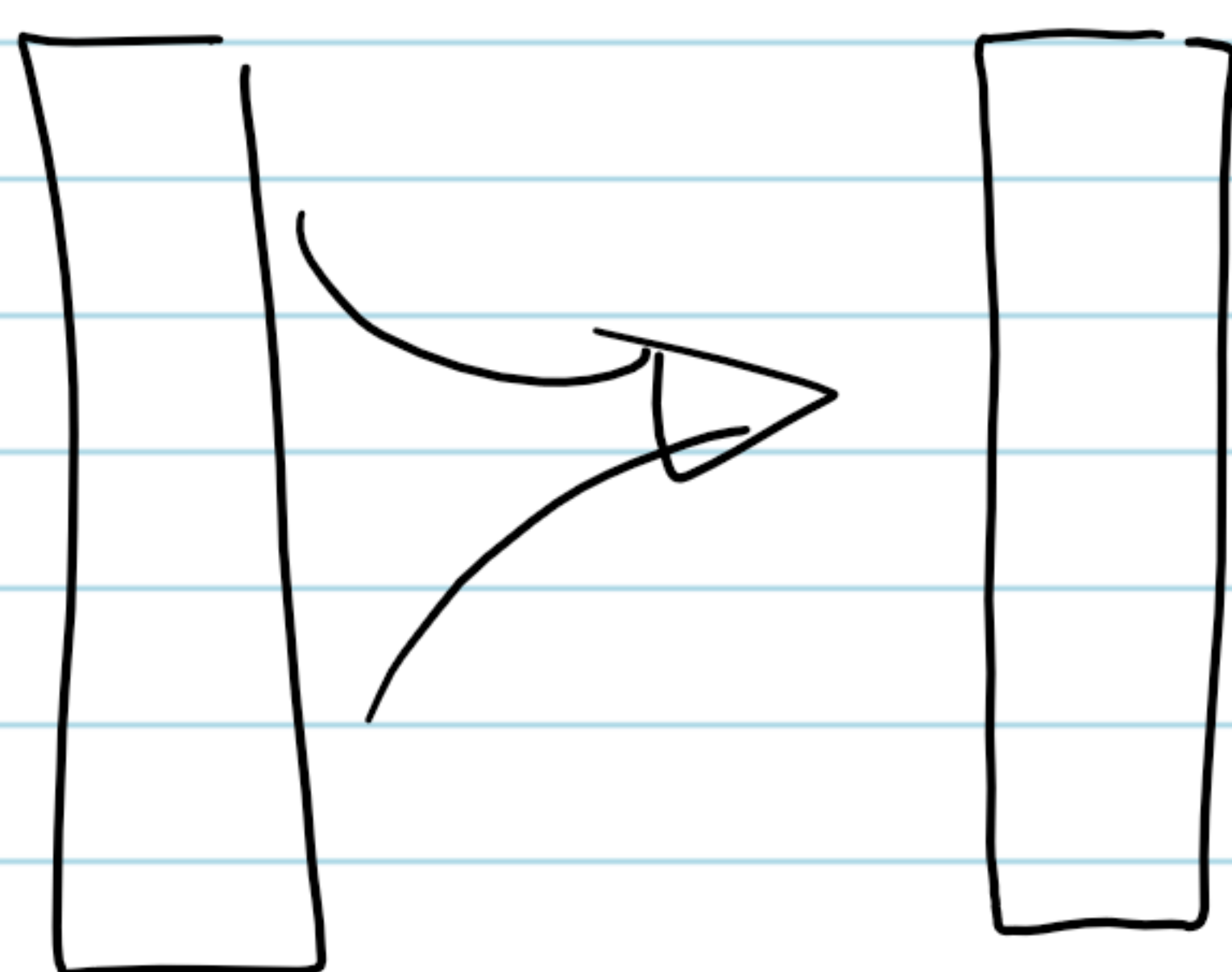
② Rescale with β & Shift by γ

↳ Params: $\{w_1, \beta, \gamma\}$ → For optimization.

Advantages

* In principle should be easier/faster to optimize/train.

* Quasi-indep. training of hidden layers



w/o the input to the next layer
could drastically change in each iteration.

But w normalization, it would have
the same mean & var.