⑥

# Training Neural Networks (I)

**★ Mini-batch SGD**

Loop:

1. Sample a batch of data

2. Forward prop it through the graph (network), get loss.

3. Backprop to calculate the gradients
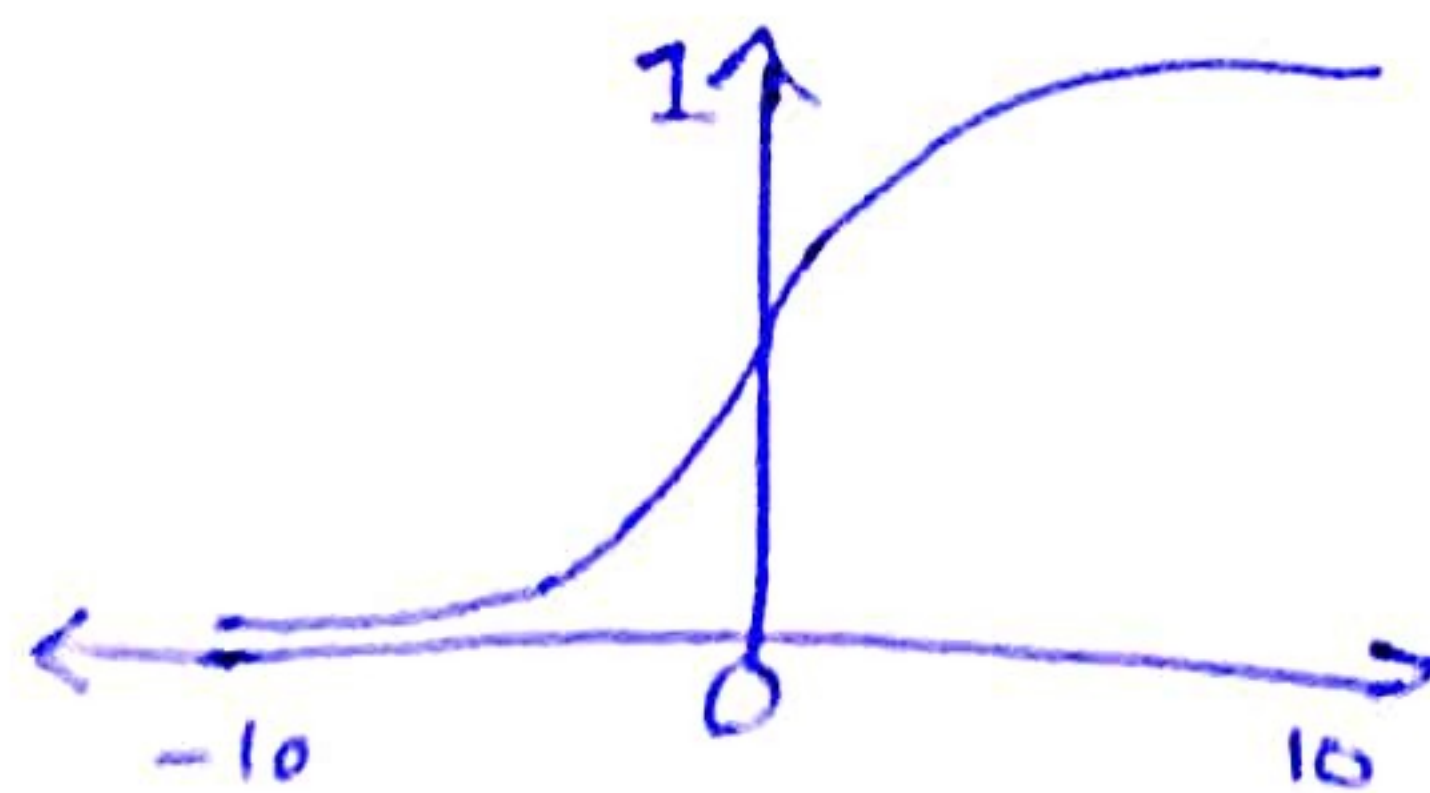
4. Update the parameters using the gradient.

### Part 1

→ Activation Functions
→ Data Preprocessing
→ Weight Initialization
→ Batch Normalization
→ Babysitting the Learning Process
→ Hyperparameter Optimization

**★ Activation Function**

**⊛ Sigmoid**

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
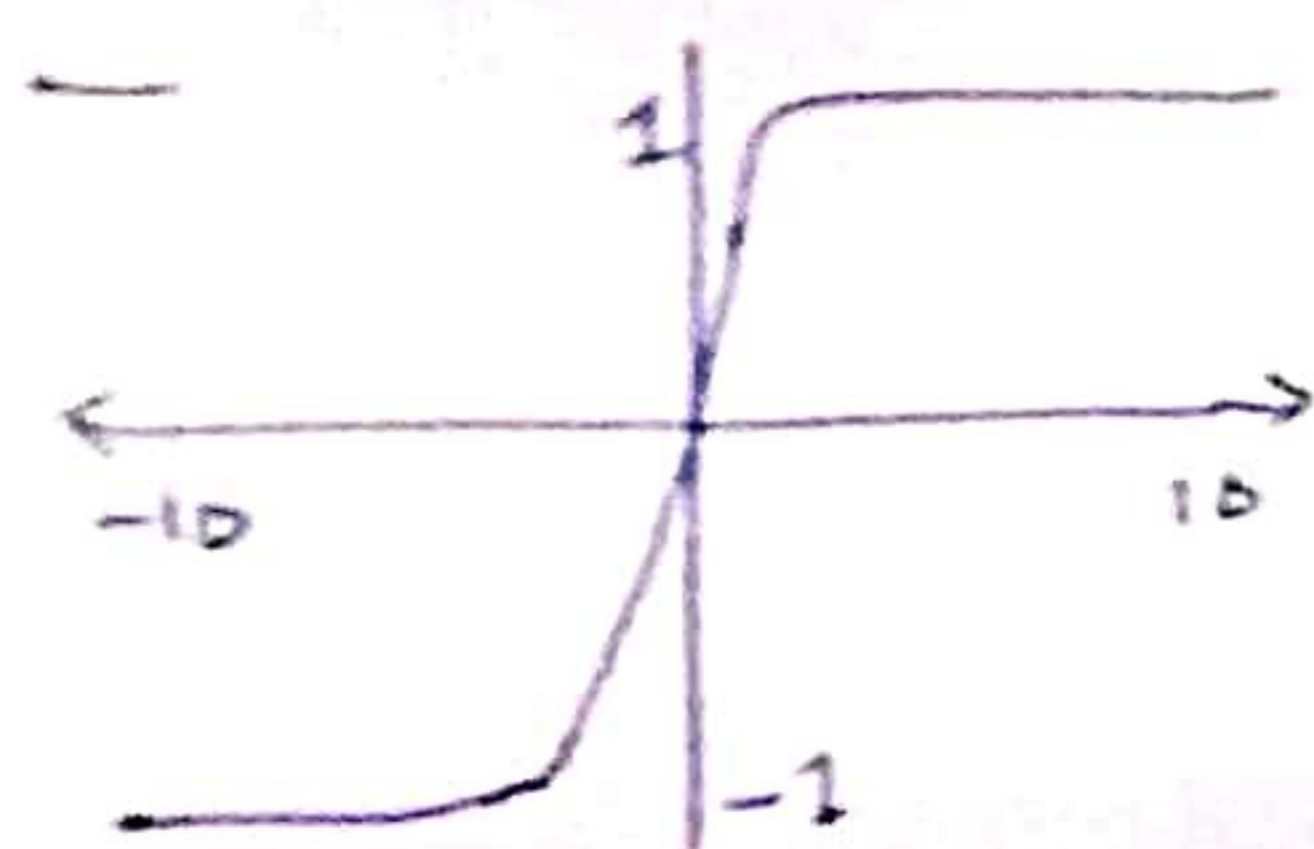


→ Squashes numbers to range [0, 1]

⇒ 3 Problems:

① Saturated neurons "kill" the gradient

② Sigmoid outputs are not zero-centered.
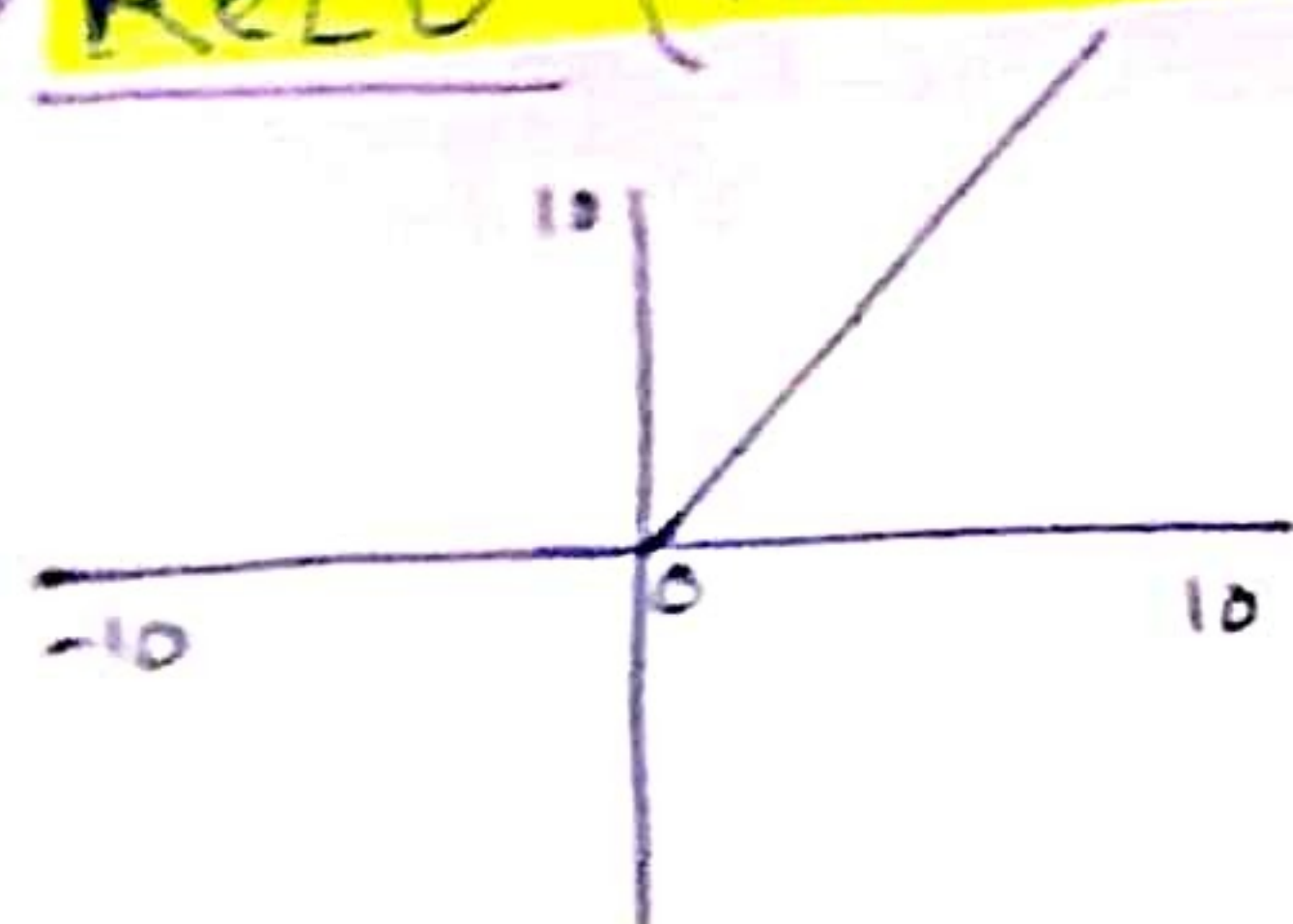
③ exp() is a bit compute expensive.

⊛ tanh(x)

→ Squashes number to range [-1,1]

→ Zero centered (nice)

→ Still kills gradients when saturated.

⊛ ReLU (Rectified Linear Unit)

→ Does not saturate (in +ve region)

→ Very computationally efficient

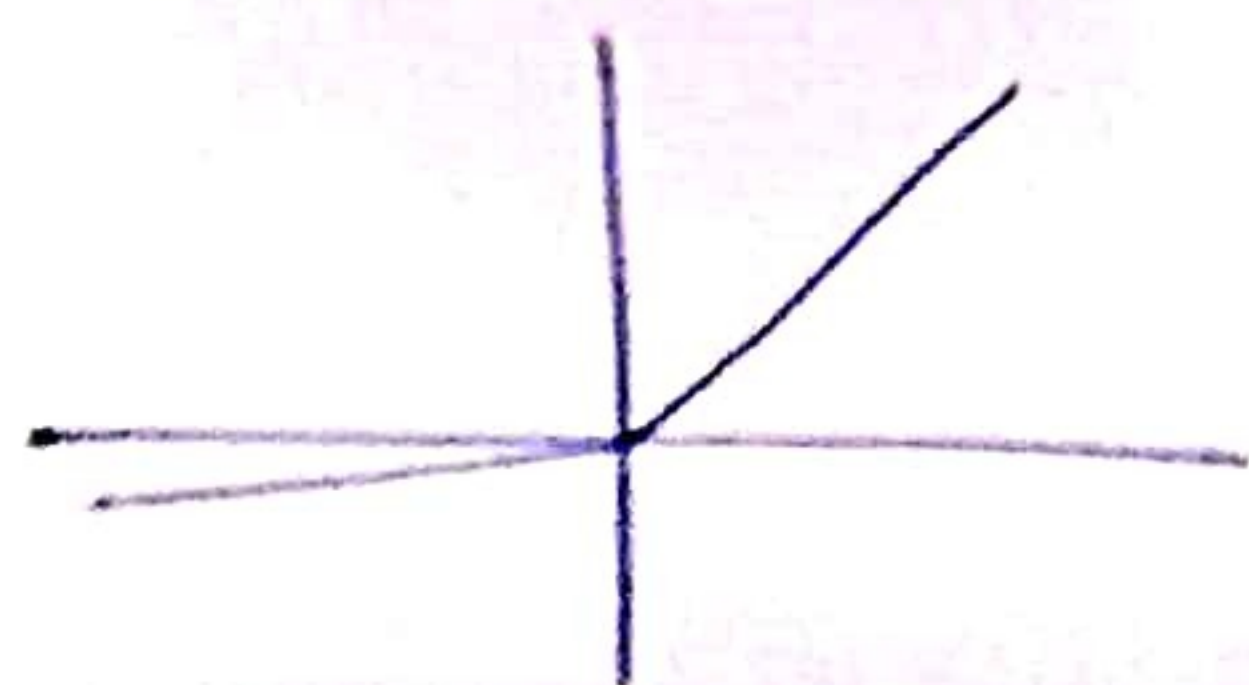→ Converges much faster than sigmoid/tanh in practice (eg 6x)

→ Actually man biologically plausible than sigmoid.

$f(x) = max(0, x)$

→ 2 Problem:

① Not zero centered output

② Gradient when x < 0 {die}

⊛ Leaky ReLU

→ Does not saturate

→ Computationally efficient
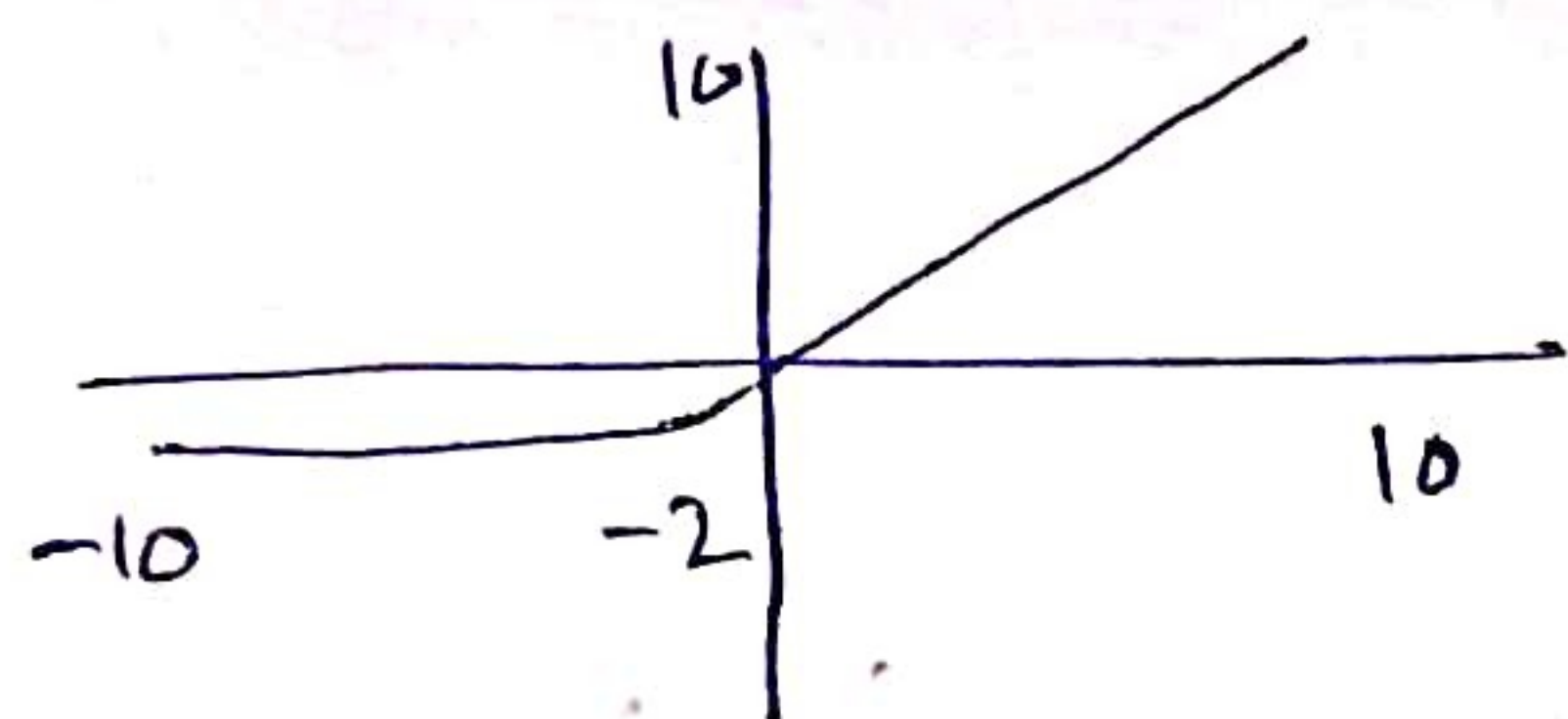
→ Converges much faster than sigmoid/tanh in practice!

→ will not die.

$f(x) = max(0.01x, x)$

⊛ ==Parametric ReLU (PReLU)==

$$f(x) = \max(\alpha x, x)$$

⊛ ==Exponential Linear Units (ELU)==



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$

→ All benefits of ReLU

→ Closer to zero mean output

→ Negative saturation regime compared with Leaky ReLU adds some robustness to noise.

→ Computation requires exp()

⊛ ==Maxout Neuron==
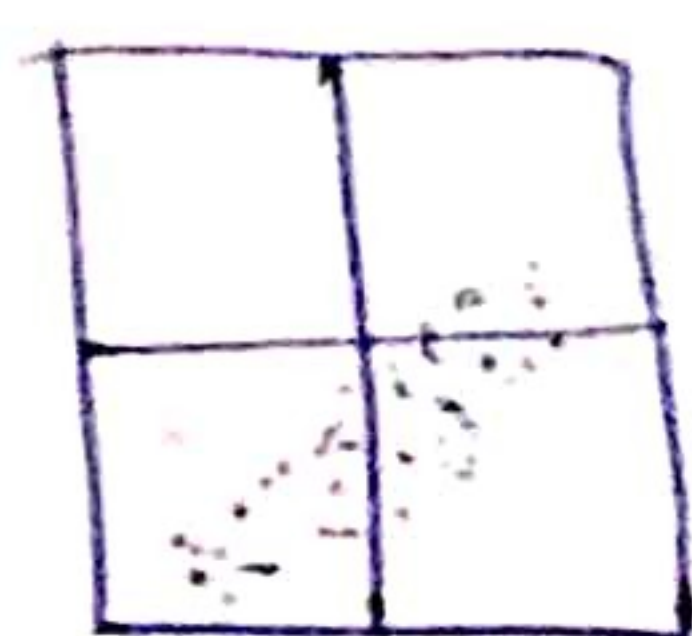
→ Generalize ReLU & Leaky ReLU

→ Linear Regime

→ Does not saturate

→ Does not Die

$$\max(\omega_1^T x + b_1, \omega_2^T x + b_2)$$
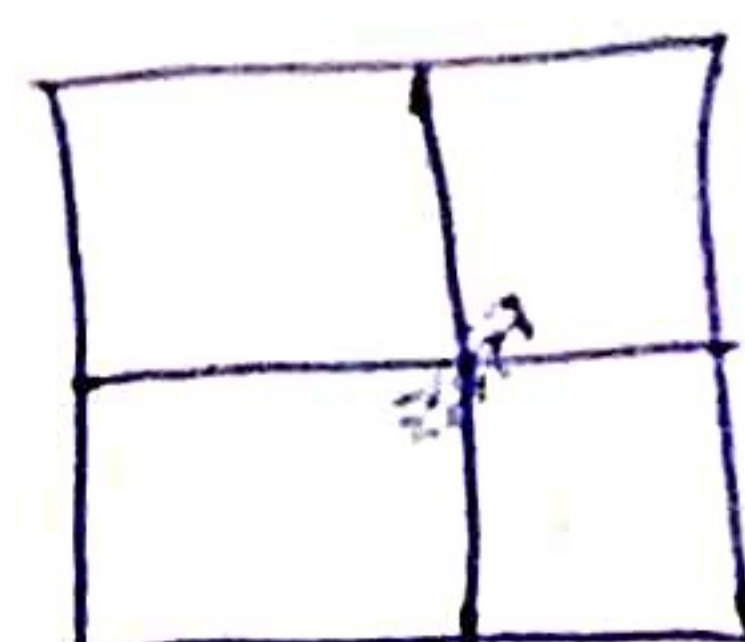
Problem: doubles the number of parameters/neuron.

# ★ Data Preprocessing



Original          Zero-Centered          Normalized
data                    data                       data

$\left\{\begin{array}{l}\rightarrow \text{Subtract the} \\ \quad \text{mean image} \\ \rightarrow \text{Subtract per-channel} \\ \quad \text{mean}\end{array}\right\}$  $\left\{\begin{array}{l}\rightarrow \text{Not Common} \\ \quad \text{for images}\end{array}\right\}$

# ★ Weight Initialization

① W = 0

⇒ All neuron will do the same thing

⇒ All neuron will always be same.

② Small random numbers for weight

( gaussian with zero mean & $10^{-2}$ standard deviation )

⇒ Works ~ Okay for small networks, but Problems with deeper networks.

③ Xavier initialization (2010)

→ Reasonable initialization.

⇒ But when using the ReLU nonlinearity it breaks.

④ He et al. (2015)

↳ Additional $\frac{1}{2}$

# ∗ Batch Normalization

⇒ "You want unit gaussian activations? Just make them so?"

⇒ Consider a batch of activations at some layer. To make each dimension unit gaussian, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{Var[x^{(k)}]}}$$

⇒ And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Note:
⇒ The network can learn

$$\gamma^{(k)} = \sqrt{Var[x^{(k)}]}$$
$$\beta^{(k)} = E[x^{(k)}]$$

to recover the identity mapping

# ∗ Babysitting the Learning Process

Step1: Pre process the data

Step2: Choose the architecture

Step3: Double check that the loss is reasonable.

ⓐ → Zero regularization
ⓑ → Some regularization

{ Loss should go up }

Step4: Start training with very small amout of data

 ↳ Turn of Regularization
 ↳ Make sure that you can overfit very small Partion of the training data.

$$\{ \text{Make the Loss} \rightarrow 0 \}$$

Step5: Start with small regularization & find learning rate that makes the loss go down.

 ↳ Loss barely changing

   ⟹ Learning rate is probably too small

Rough range for learning rate we should be Cross-validating is somewhere. $[1e^{-3} ... 1e^{-5}]$

**✱ Hyperparameter Optimization**

**⊞ Cross-validation Strategy**

 Coarse → fine Cross-validation in stages

 First stage: Only a few epochs to get rough idea of what params work
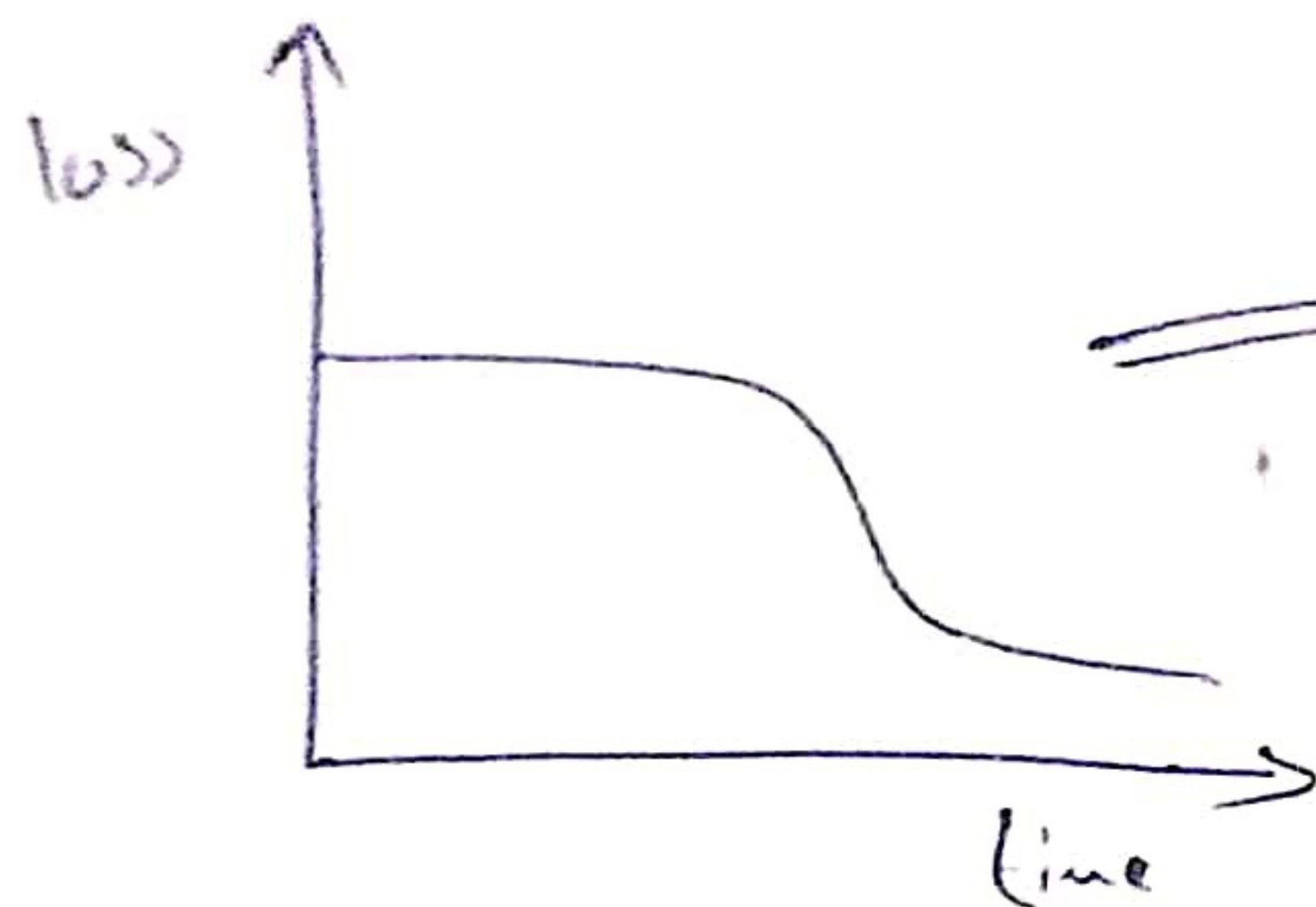
 Second Stage: longer runtime, finer search.

⟹ Random Search is better than Grid Search for hyper parameter.

⟹ Search for hyperparameter in log space when appropriate.

⇒ Hyperparameters to play with:

    ↳ Network architecture

    ↳ Learning rate, its decay schedule, update type

    ↳ regularization (L2/Dropout strength)

⇒ Monitor & Visualize the loss curve.



loss

time

⟶ Bad initialization a prime suspect