

# Deep Learning: Training Algorithms

**Ozan Özdenizci**

Institute of Theoretical Computer Science

[ozan.ozdenizci@igi.tugraz.at](mailto:ozan.ozdenizci@igi.tugraz.at)

Deep Learning VO - WS 23/24

Lecture 5 - November 6th, 2023

# Recap: Stochastic Gradient Descent (SGD)

- The **gradient** points in the direction of the steepest increase of the error.

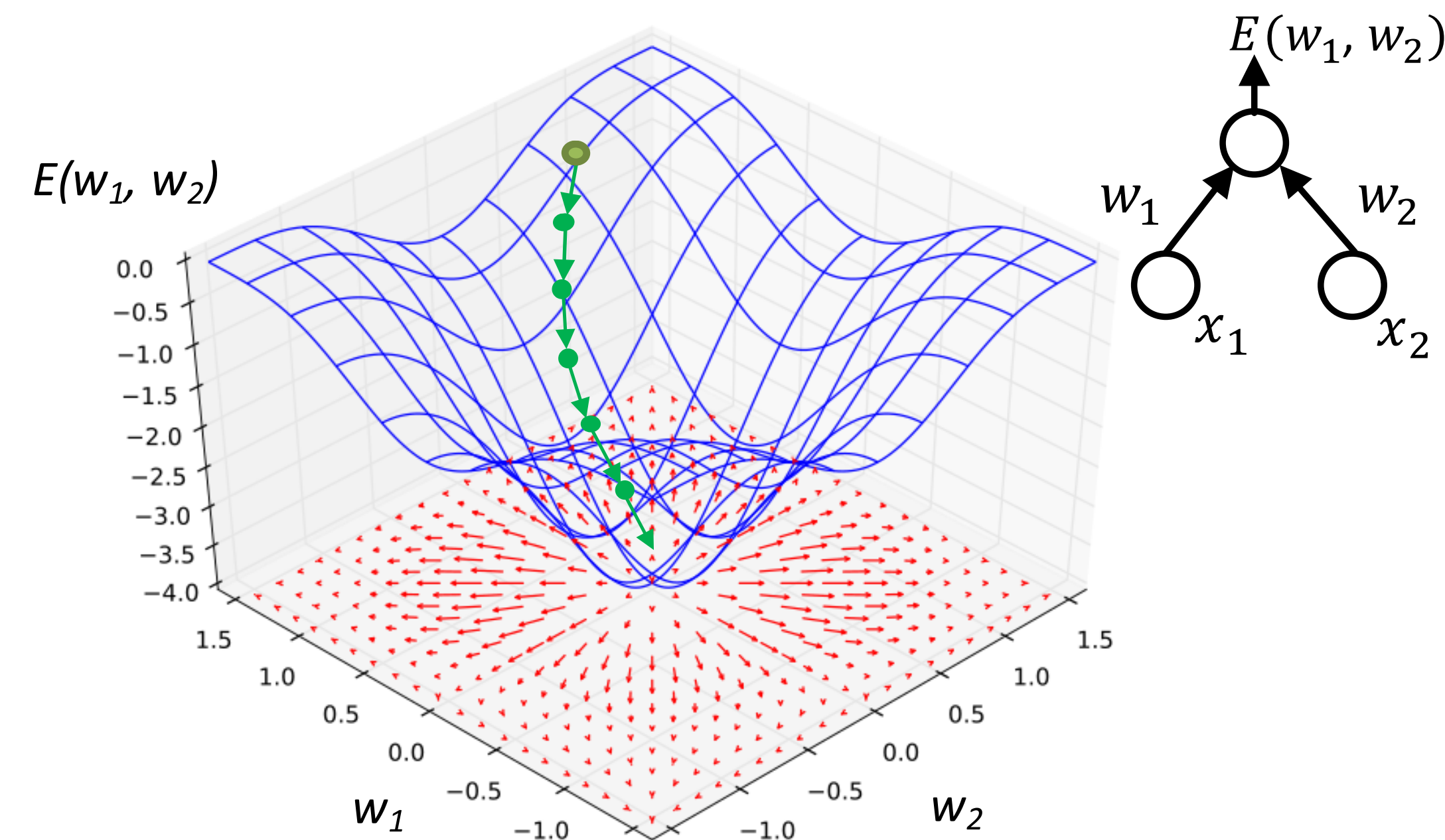
$$\nabla_{\mathbf{w}} E = \begin{pmatrix} \frac{\partial E}{\partial w_1} \\ \vdots \\ \frac{\partial E}{\partial w_D} \end{pmatrix}$$

- We slightly change parameters to reduce error:

$$\mathbf{w}_{new} = \mathbf{w}_{old} - \epsilon \nabla_{\mathbf{w}} E(\mathbf{w}_{old}) \quad \epsilon > 0 \text{ is the learning rate.}$$

- Stochastic Gradient Descent:**

- Divide the training set into **mini-batches**.
- Approximate the true gradient by the gradient over the mini-batch.



```
1 Choose  $\mathbf{w}$  randomly
2 Let B be the set of minibatches
2  $i \leftarrow 0$ 
3 REPEAT
4   FOR  $b$  in B
4      $\mathbf{w} \leftarrow \mathbf{w} - \epsilon(i) \nabla_{\mathbf{w}} E(\mathbf{w}, b)$ 
5      $i \leftarrow i+1$ 
6 UNTIL Stopping Criterion
```

# Today

---

- ❑ Basic Training Algorithms
- ❑ Common Pitfalls & Training Considerations
- ❑ Searching Hyperparameters

# (1) SGD with decreasing learning rate

Due to the batch-induced noise in stochastic GD, gradient does not become 0 at minimum.

- One therefore reduces the learning rate  $\epsilon_k$  for each update  $k$ .
- Sufficient condition to guarantee convergence of SGD:

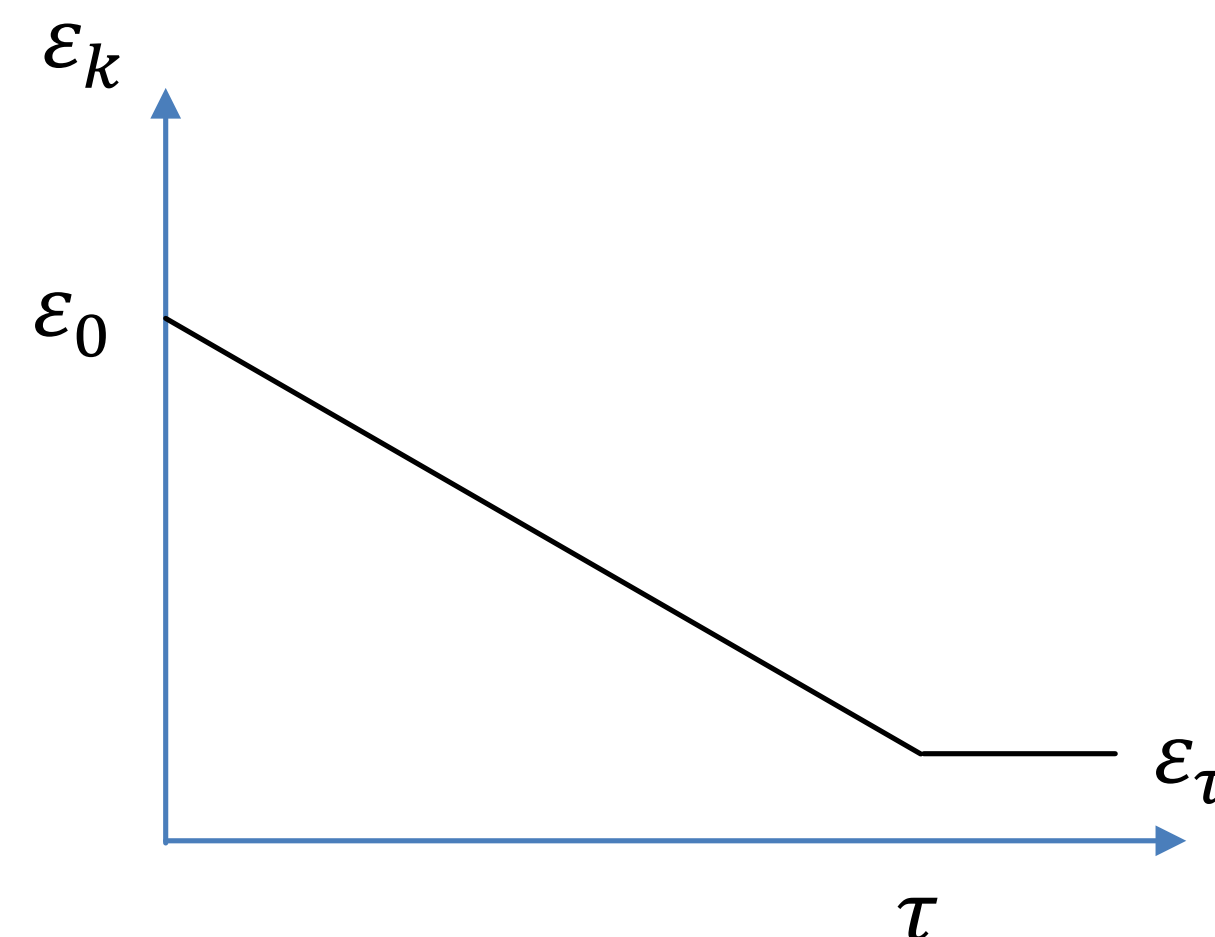
$$\sum_{k=0}^{\infty} \epsilon_k = \infty \quad \text{and} \quad \sum_{k=0}^{\infty} \epsilon_k^2 < \infty$$

**A commonly used approach (in practice):**

- Start with initial rate  $\epsilon_0$
- Decrease linearly until update  $\tau$

$$\epsilon_k = \left(1 - \frac{k}{\tau}\right) \epsilon_0 + \frac{k}{\tau} \epsilon_\tau$$

- Then keep constant  $\epsilon_\tau$



**Choosing the meta-parameters:**

$\tau$  : A few passes through the training examples (e.g., 10 epochs)

$\epsilon_\tau$  : about 1% of  $\epsilon_0$

$\epsilon_0$  : Monitor first several iterations and choose learning rate higher than the best performing, but not so high so that it does not cause severe instability.

## (2) Momentum

Use a running average of the gradient for the update:

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} L(\theta)$$

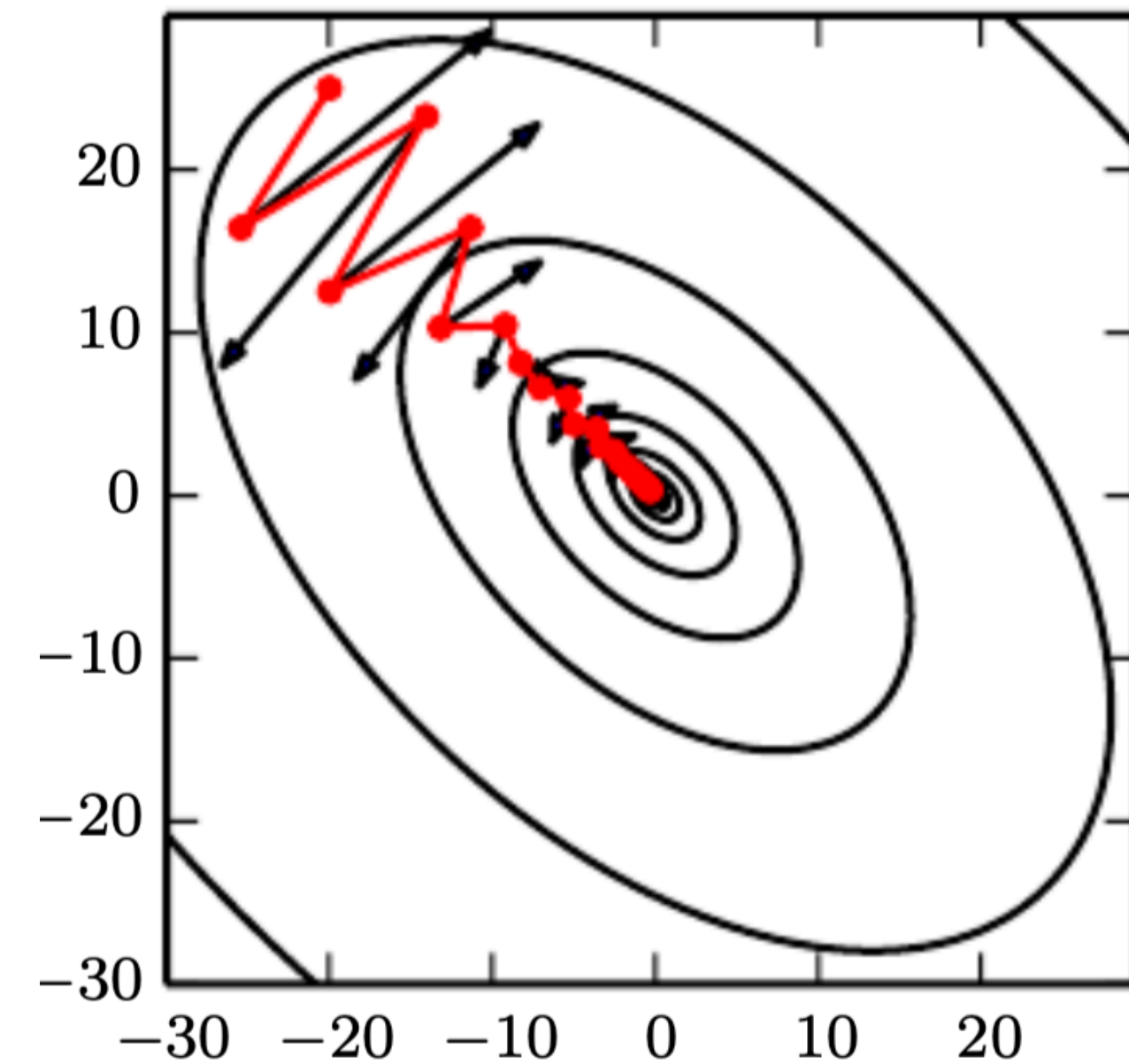
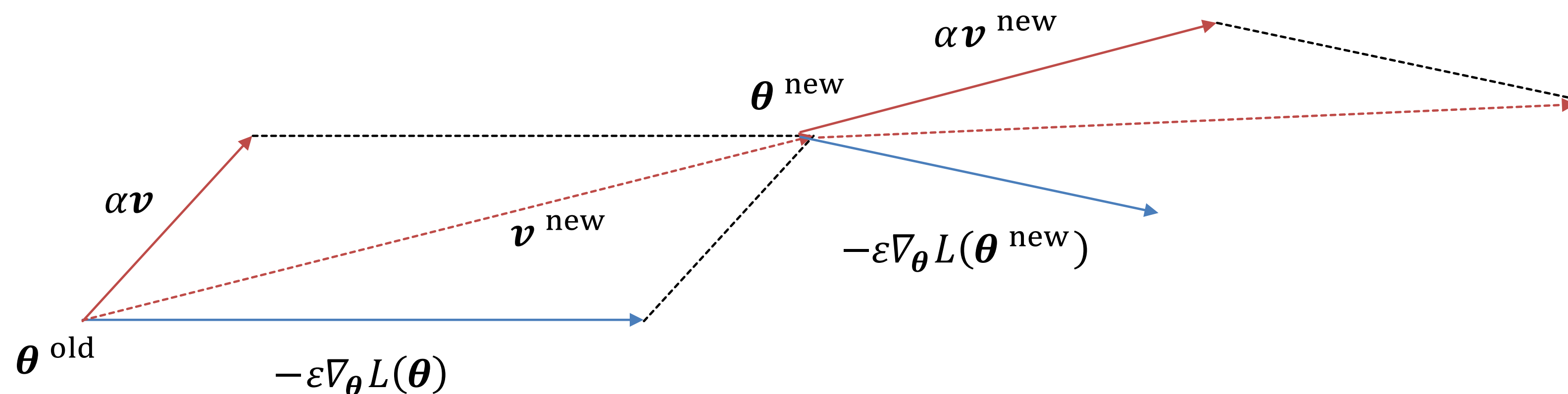
$\mathbf{v}$  : velocity

$$\theta \leftarrow \theta + \mathbf{v}$$

$L$  : objective function (loss)

**Basic idea:** Give the parameter vector a **mass** that moves on the error landscape.

- Reduces gradient noise of minibatches
- Reinforces consistent gradient directions



**Choosing the meta-parameters:**

$\alpha \in [0,1)$  : How quickly previous contributions decay (common choices: 0.5, 0.9, 0.99).

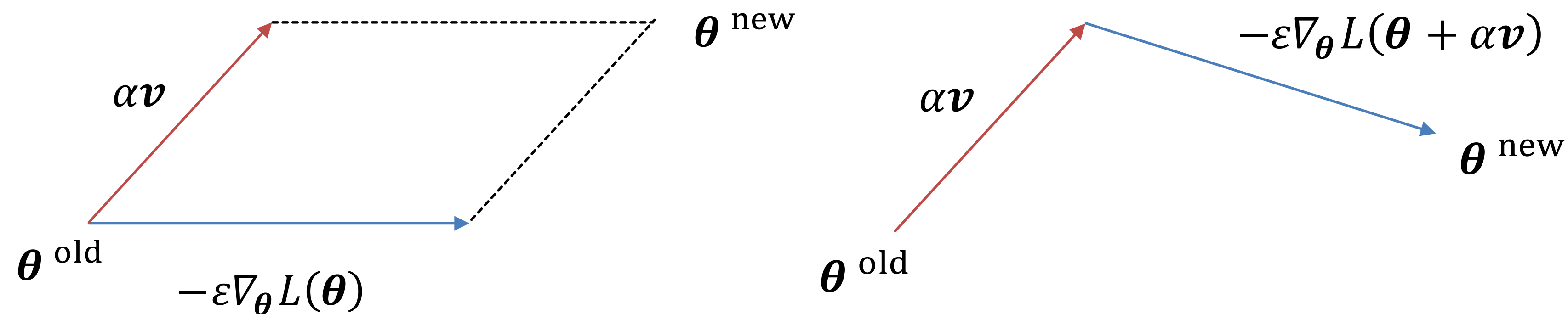
Often adapted over time, starting with small value and then increased.



### (3) Nesterov Momentum

Evaluate at the location where parameters would be if you continued with velocity (no gradient force).

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} L(\theta + \alpha \mathbf{v})$$
$$\theta \leftarrow \theta + \mathbf{v}$$



- The Nesterov-Momentum variant has stronger theoretical convergence properties on convex functions and slightly better results in practice.

# (4) Algorithms with Adaptive Learning Rates

---

The error is often highly sensitive to some directions in parameter space and insensitive to others.

- Hence, adapting a *global learning rate* is less helpful.
- Instead, one often tries to find **individual learning rates** for each parameter.

## Algorithms:

- ▶ (4.1) AdaGrad
- ▶ (4.2) RMSProp
- ▶ (4.3) Adam

# (4.1) AdaGrad

Scale learning rates by dividing them by the square root of the sum of all historical squared gradient values.

---

**Algorithm 8.4** The AdaGrad algorithm

---

**Require:** Global learning rate  $\epsilon$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$ .

    Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---



# (4.1) AdaGrad

Scale learning rates by dividing them by the square root of the sum of all historical squared gradient values.

## Algorithm 8.4 The AdaGrad algorithm

**Require:** Global learning rate  $\epsilon$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$   
corresponding targets  $\mathbf{y}^{(i)}$ .

Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$ .

Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root are  
element-wise)

Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

- This is quite useful for convex optimization.
- *Limitation:* For non-convex optimization, early gradients (accumulated at the beginning of training) are often not relevant later.

## AdaGrad - A brief summary:

$$\mathbf{g} = \begin{pmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_D} \end{pmatrix} \rightarrow \text{Gradient computed over one mini batch.}$$

$$\mathbf{r} = \begin{pmatrix} \sum_j \left( \frac{\partial L^{(j)}}{\partial w_1} \right)^2 \\ \vdots \\ \sum_j \left( \frac{\partial L^{(j)}}{\partial w_D} \right)^2 \end{pmatrix} \rightarrow \text{Squared gradients accumulated across all minibatches so far}$$

$$\Delta w_k = -\frac{\epsilon}{\delta + \sqrt{\sum_j \left( \frac{\partial L^{(j)}}{\partial w_k} \right)^2}} \cdot \frac{\partial L}{\partial w_k}$$

individual update for a single parameter  $w_k$

## (4.2) RMSProp

Like AdaGrad, but using a running average to compute the square root of the sum of historical squared gradients.

---

**Algorithm 8.5** The RMSProp algorithm

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers

Initialize accumulation variables  $\mathbf{r} = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$ .

    Compute parameter update:  $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\delta + \mathbf{r}}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---

## (4.2) RMSProp

Like AdaGrad, but using a running average to compute the square root of the sum of historical squared gradients.

---

### Algorithm 8.5 The RMSProp algorithm

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers

Initialize accumulation variables  $\mathbf{r} = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$ .

    Compute parameter update:  $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\delta + \mathbf{r}}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---

### RMSProp

- Does not shrink the learning rate according to the entire history of the squared gradients.
- Instead, uses a decaying average to discard history from the extreme past.



## (4.2) RMSProp with Nesterov Momentum

RMSProp can also be combined with the idea of momentum.

---

**Algorithm 8.6** RMSProp algorithm with Nesterov momentum

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$ , momentum coefficient  $\alpha$

**Require:** Initial parameter  $\theta$ , initial velocity  $v$

Initialize accumulation variable  $r = 0$

**while** stopping criterion not met **do**

Sample a minibatch of  $m$  examples from the training set  $\{x^{(1)}, \dots, x^{(m)}\}$  with corresponding targets  $y^{(i)}$ .

Compute interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$ .

Compute gradient:  $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$ .

Accumulate gradient:  $r \leftarrow \rho r + (1 - \rho) g \odot g$ .

Compute velocity update:  $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$ . ( $\frac{1}{\sqrt{r}}$  applied element-wise)

Apply update:  $\theta \leftarrow \theta + v$ .

**end while**

---

## (4.3) Adam

Similar to RMSProp with momentum, but with corrected moments.

---

### Algorithm 8.7 The Adam algorithm

---

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ .  
(Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant  $\delta$  used for numerical stabilization (Suggested default:  $10^{-8}$ )

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$

Initialize time step  $t = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

    Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

    Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

    Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

    Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

    Compute update:  $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$  (operations applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$

**end while**

---

### Adam

- Initially, the moments are biased towards the **initialization**. This is corrected in Adam.
- Compared to SGD, Adam requires much less hyper-parameter tuning. Default hyper-parameters often work well (in practice).

# Today

---

- ☒ Basic Training Algorithms
- ☐ Common Pitfalls & Training Considerations
- ☐ Searching Hyperparameters



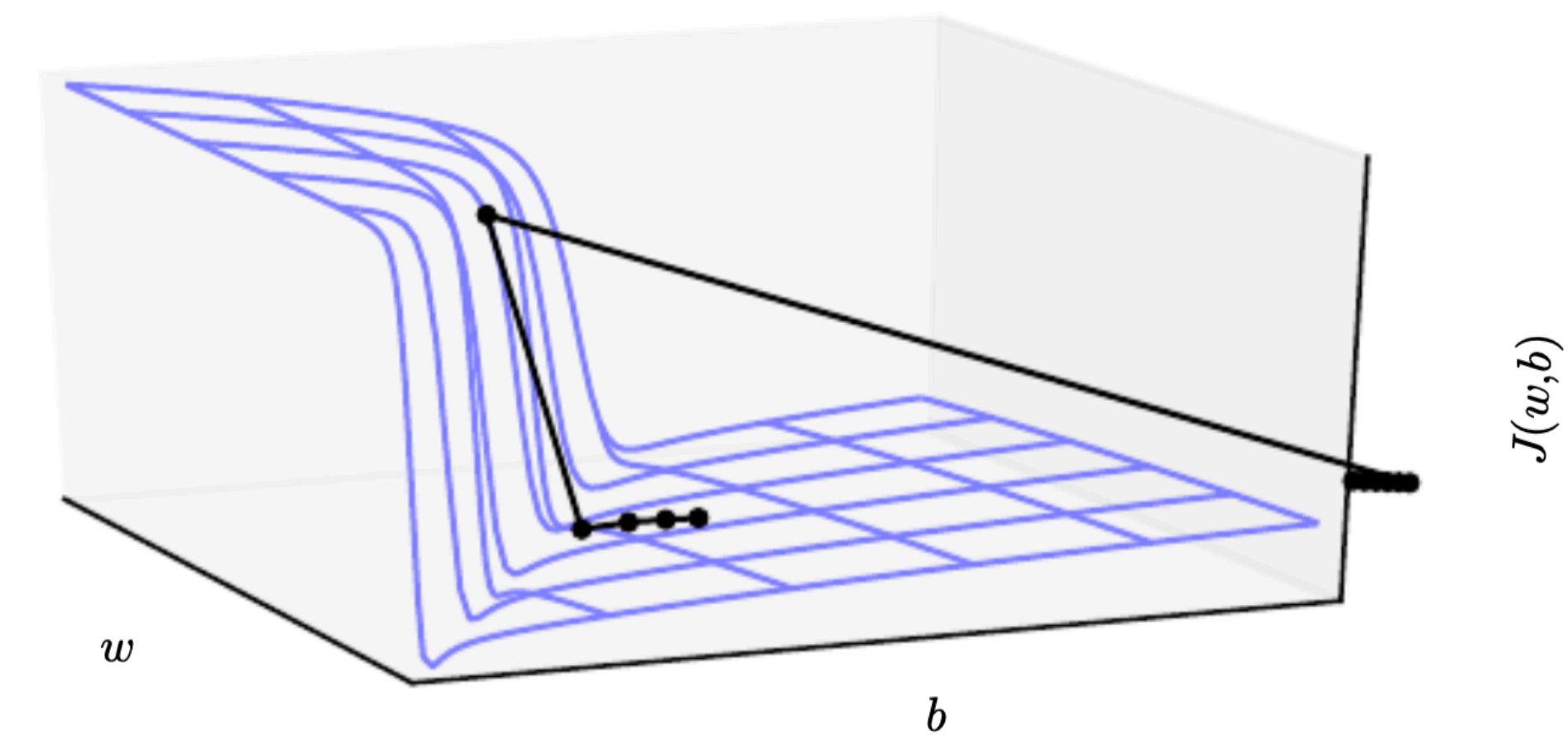
# Gradient Clipping

Deep networks often have extremely steep regions in the error surface, resembling to cliffs.

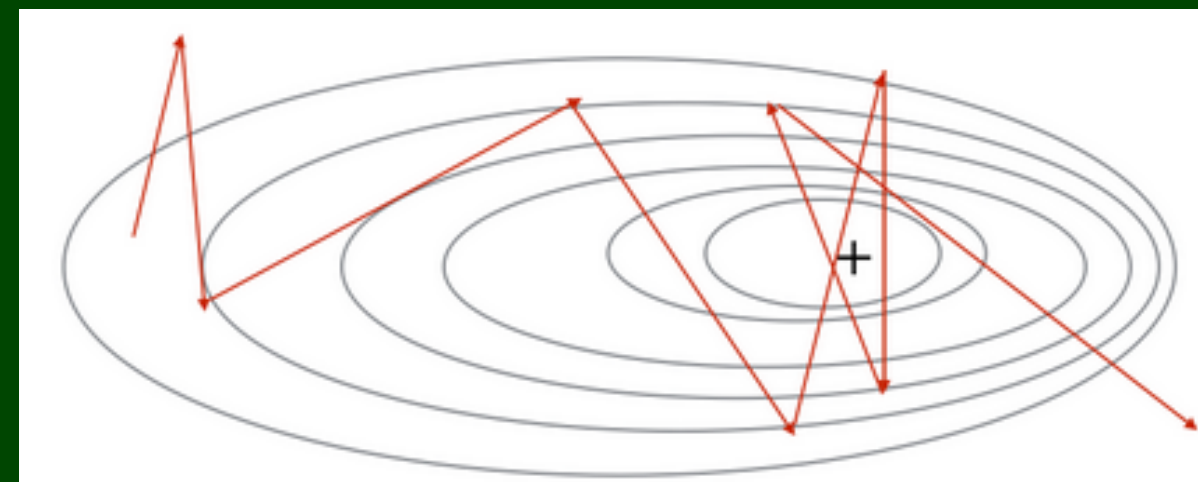
There, the gradient can become very large, which leads to instabilities and to large jumps in parameter space.

**Gradient clipping:** Clip large gradients by a threshold.

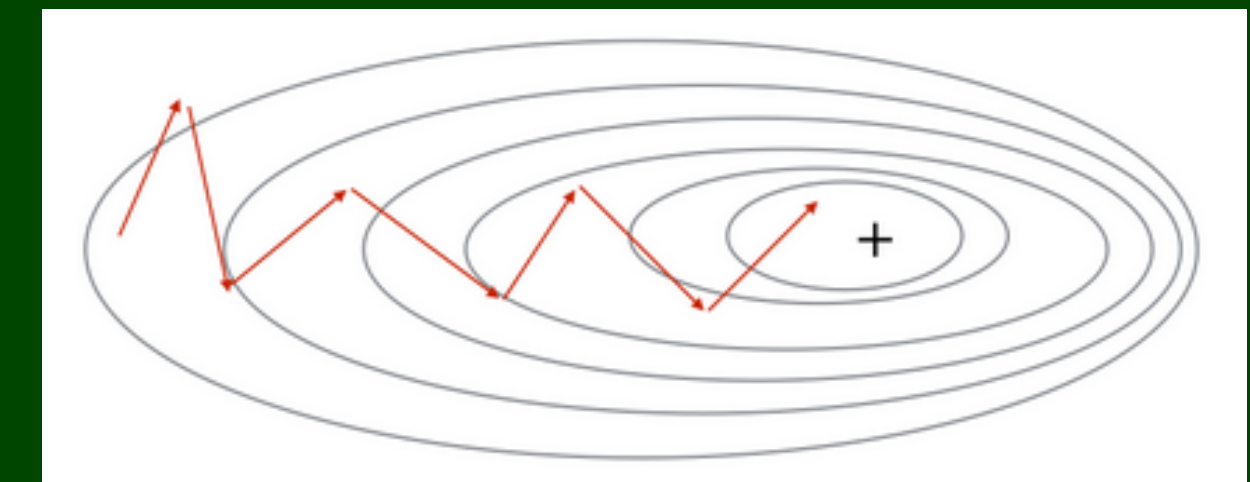
$$\text{if } ||\mathbf{g}|| > \nu \text{ then } \mathbf{g} \leftarrow \frac{\nu \mathbf{g}}{||\mathbf{g}||}$$



without gradient clipping



with gradient clipping

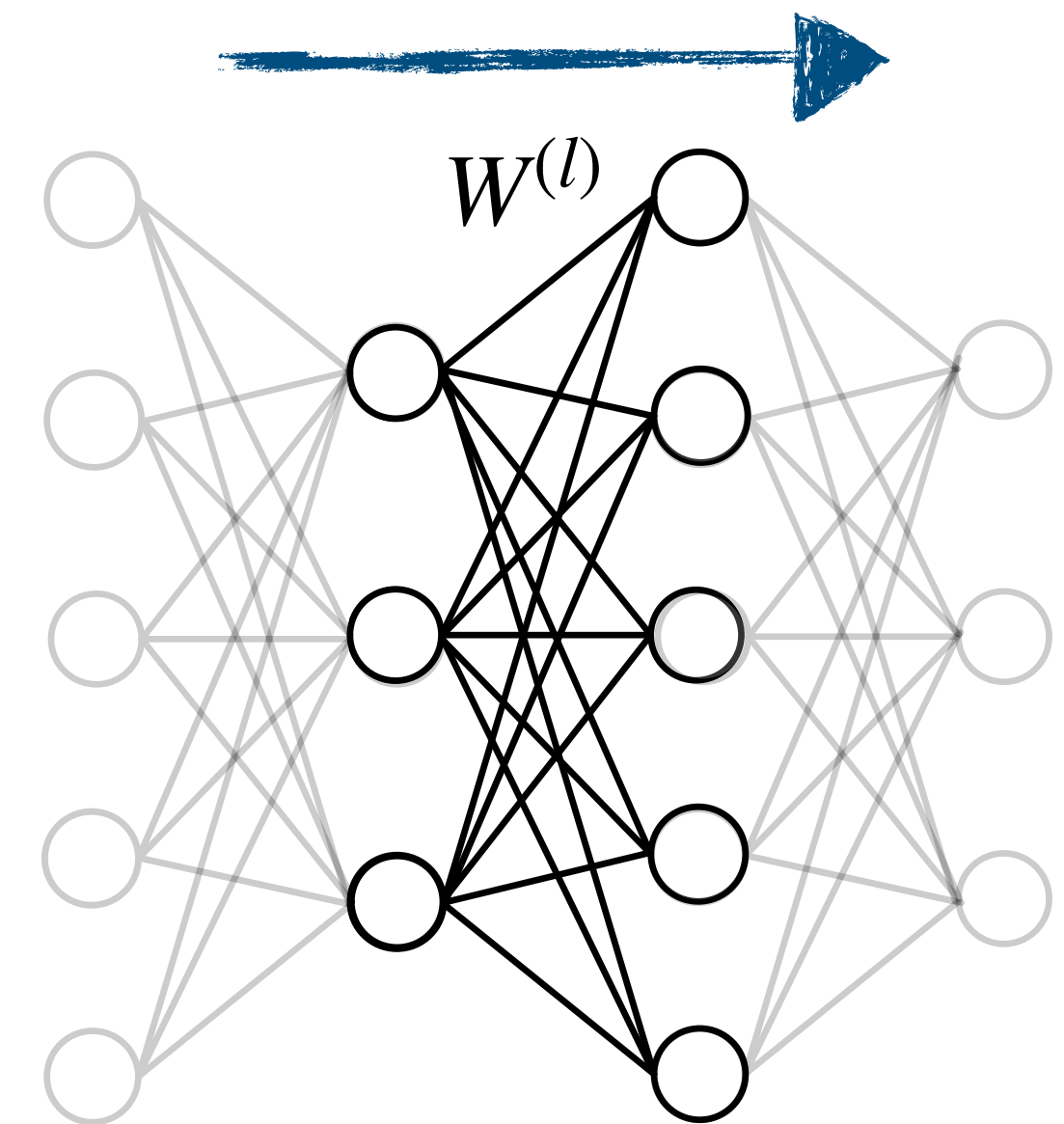


# Parameter Initialization

Initial parameters can determine whether learning converges at all, or how quickly it converges.

Important considerations for network weights:

- **Breaking the symmetry:**
  - Units should compute different functions.
  - Initialize weights randomly (e.g., using a Gaussian or uniform distribution).
- **Sizes of initial weights:**
  - Large
    - Stronger symmetry breaking and avoids losing the signal.
  - Too large
    - Exploding values during forward/backward pass, or saturated activation functions (bad).
  - Small
    - Regularization: We want the weights to be small, so start with small ones.



# Parameter Initialization - Heuristics

Consider the weights of a fully connected layer with  $m$  inputs and  $n$  outputs.

Heuristic 1:

$$w_{ij} \sim \text{Uniform} \left( -\frac{1}{\sqrt{m}}, \frac{1}{\sqrt{m}} \right)$$

Heuristic 2:

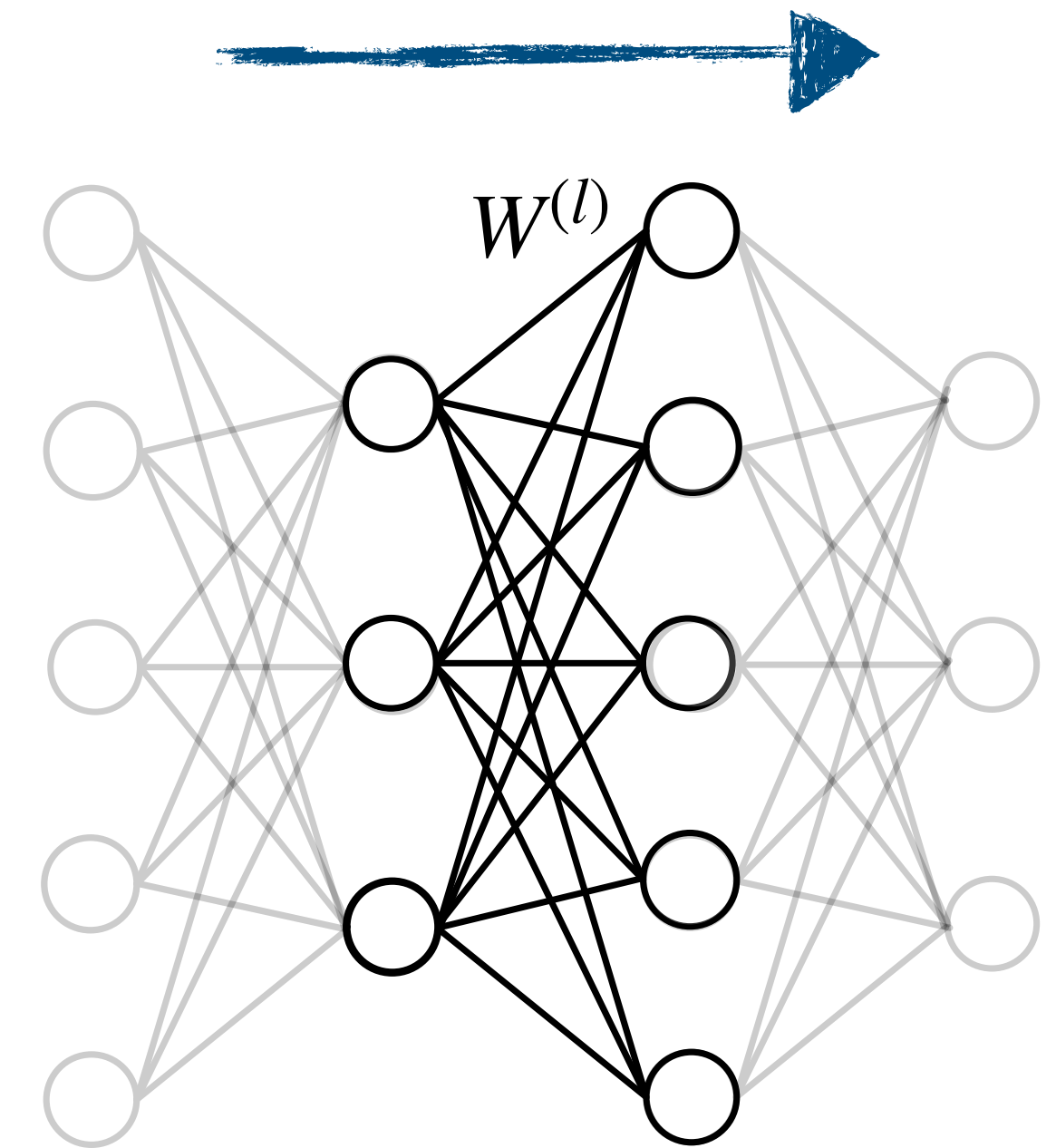
$$w_{ij} \sim \text{Uniform} \left( -\frac{6}{\sqrt{m+n}}, \frac{6}{\sqrt{m+n}} \right)$$

[Glorot & Bengio, AISTATS 2010]

**Potential problem for large layers:** Each individual weight becomes very small.

Alternative: **Sparse initialization** - each unit has initially exactly  $k$  non-zero weights.

**Empirical:** Look at distribution of outputs and gradients on a single minibatch.



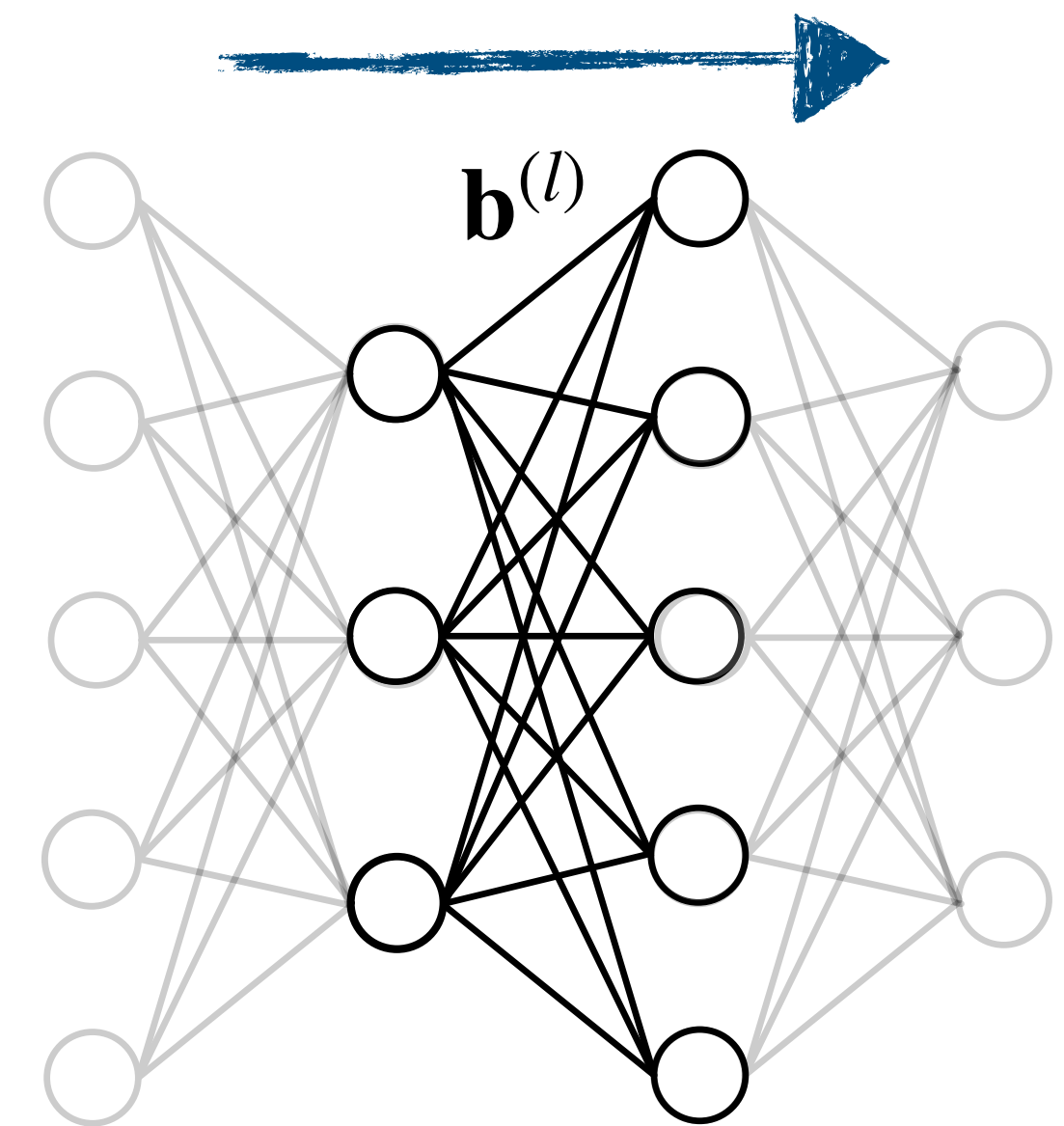
$W^{(l)}$  is an  $n \times m$  matrix  
 $m$  : inputs  
 $n$  : outputs

# Parameter Initialization

Initial parameters can determine whether learning converges at all, or how quickly it converges.

## Important considerations for network biases:

- Can all be set to the same heuristically chosen constant.
- For the output, one can choose biases to obtain right target means.
- Avoid saturation
  - e.g., ReLUs make activations mostly positive (small positive bias, e.g., 0.01)



## Further techniques:

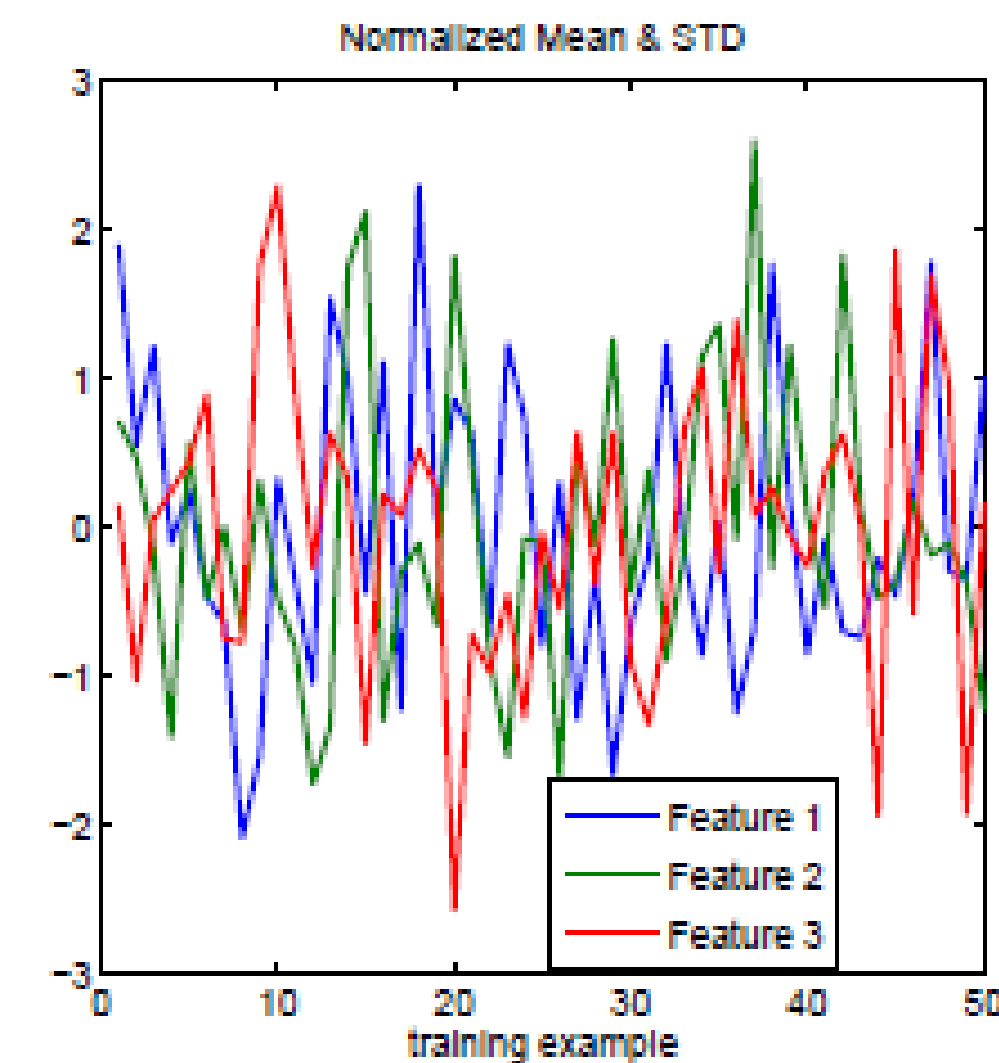
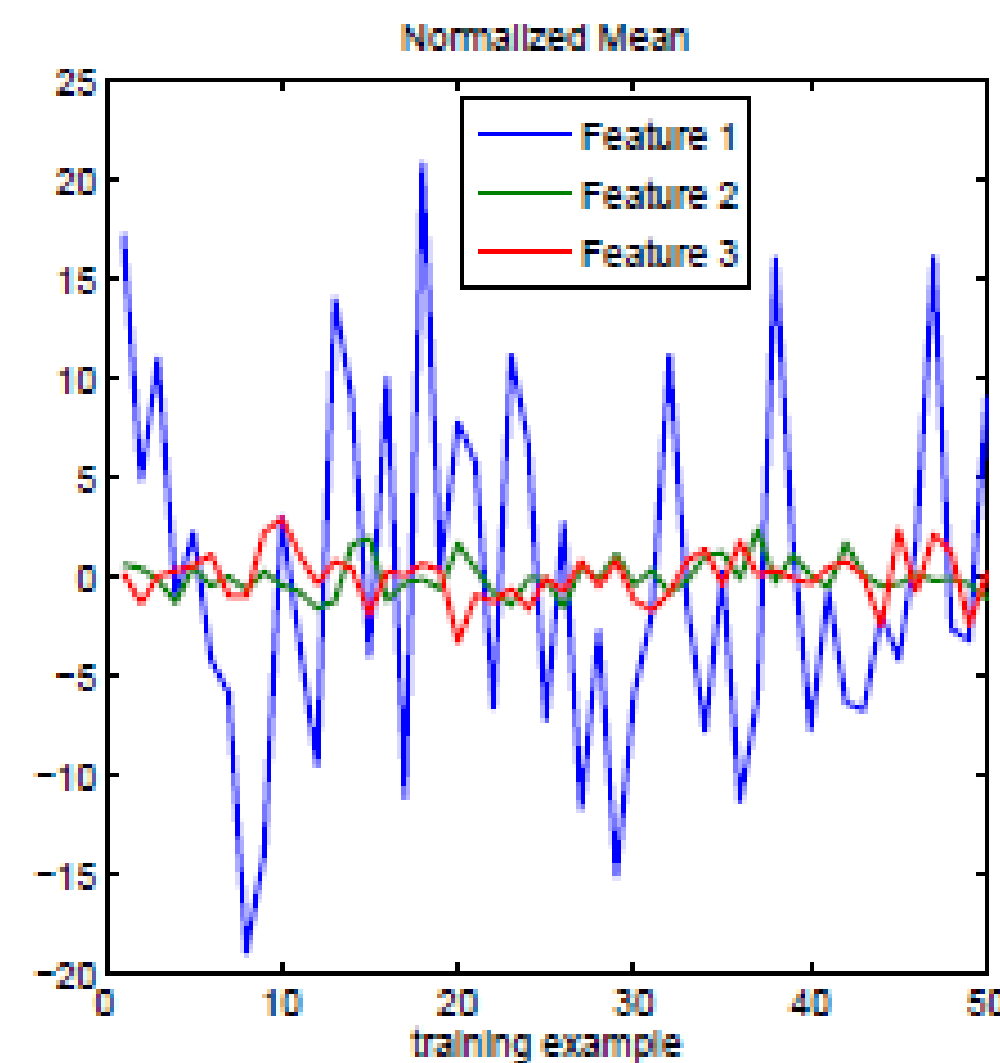
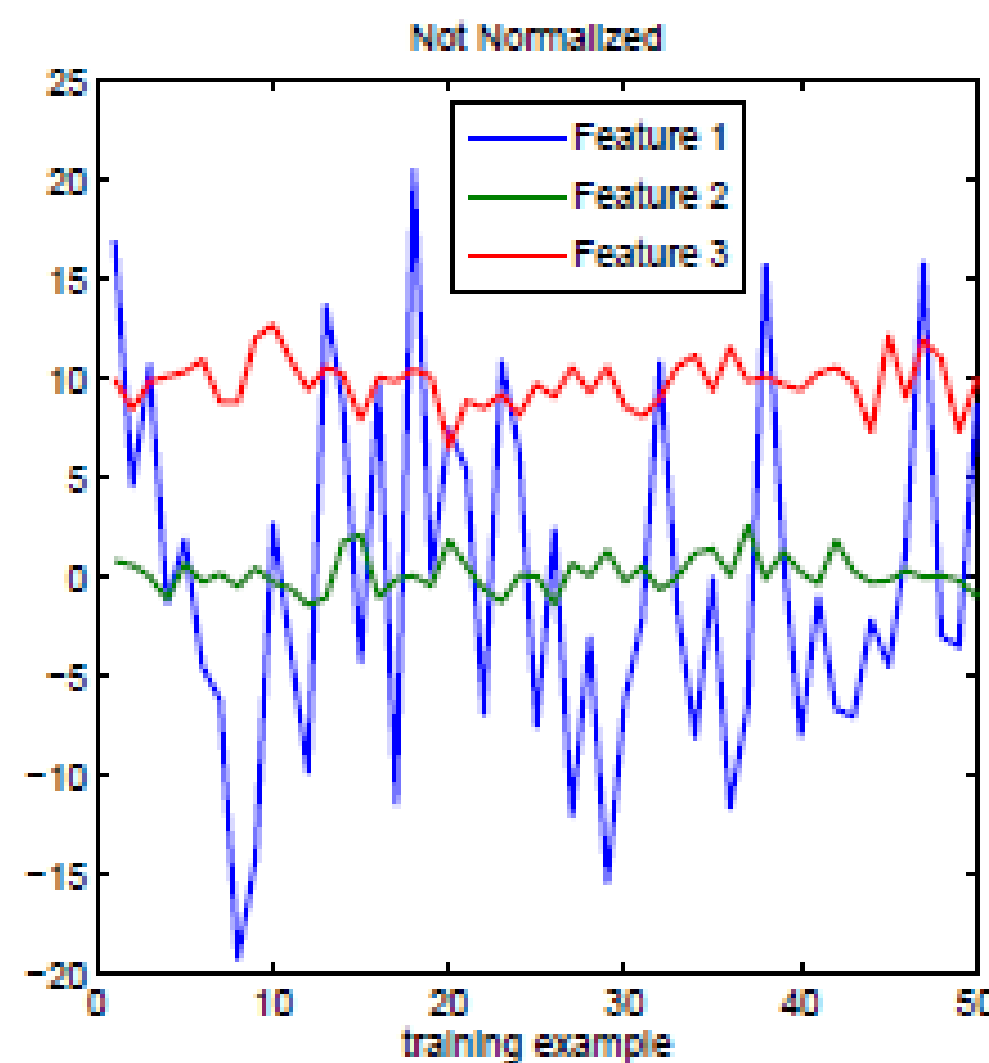
- Unsupervised pre-training
- Pre-training on similar tasks (e.g., for image-related tasks one can use networks pre-trained on the ImageNet dataset, which will already capture many aspects of natural images)

# Activity Normalization Methods

## Input normalization

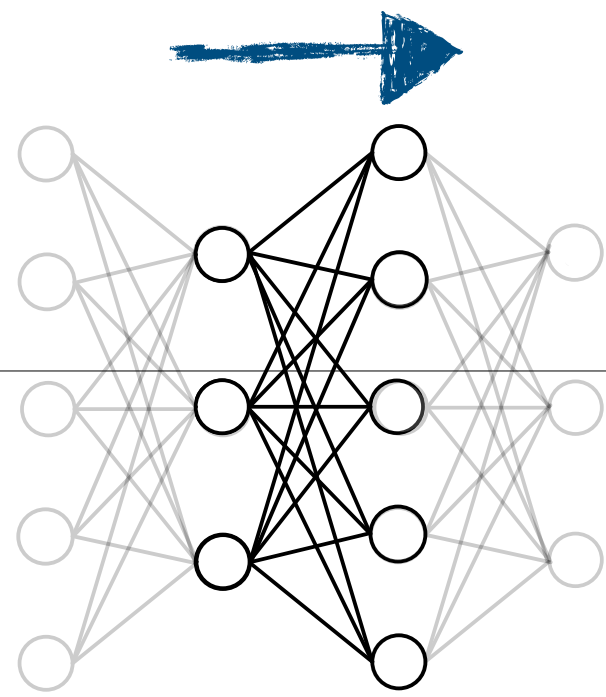
- We usually normalize each input feature to have zero-mean and unit-variance.
- This can significantly speed up learning.
- Mean and std. statistics are computed over the training set, and then applied to training, validation and test sets.

$$\tilde{x}_i^{(n)} = \frac{x_i^{(n)} - \text{mean}(\{x_i^{(1)}, \dots, x_i^{(N)}\})}{\text{std}(\{x_i^{(1)}, \dots, x_i^{(N)}\})}$$





# Batch Normalization



We can interpret the output of a layer as the input to the next layer.

In Batch-Normalization, one does this normalization on the output of each layer for the current batch.

$$\tilde{z}_i^{(n)} = \frac{z_i^{(n)} - \text{mean}(\{z_i^{(1)}, \dots, z_i^{(N')}\})}{\sqrt{\text{var}(\{z_i^{(1)}, \dots, z_i^{(N')}\}) + \epsilon}}$$

$z_i^{(1)}, \dots, z_i^{(N')}$  : output of unit  $i$  in the batch

$\epsilon$  : small constant to avoid division by zero

- This can be seen as an intermediate layer that linearly transforms the layers' output.
- During gradient descent, the error is backpropagated through this transformation.
- **When testing a single example**, one can use an average of the means and standard deviations in a number of previous batches.



# Batch Normalization

$$\tilde{z}_i^{(n)} = \frac{z_i^{(n)} - \text{mean}(\{z_i^{(1)}, \dots, z_i^{(N')}\})}{\sqrt{\text{var}(\{z_i^{(1)}, \dots, z_i^{(N')}\}) + \epsilon}}$$

$z_i^{(1)}, \dots, z_i^{(N')}$  : output of unit  $i$  in the batch

$\epsilon$  : small constant to avoid division by zero

**Normalizing the output of units can reduce the expressive power of the network.**

- Using trained (learnable) parameters  $\gamma_i, \beta_i$ , we often replace  $\tilde{z}_i^{(n)}$  as:

$$\gamma_i \tilde{z}_i^{(n)} + \beta_i$$

- The resulting network is equivalent to the original one (without batch normalization).
- However, it has better learning dynamics. The mean and std. dev. of a unit's output is now determined by two parameters independent of other network parameters.
- Batch-Normalization is not applicable to small batch sizes and recurrent neural networks.

# Layer Normalization

---

$$\tilde{z}_i^{(n)} = \frac{z_i^{(n)} - \text{mean}(\{z_1^{(n)}, \dots, z_m^{(n)}\})}{\sqrt{\text{var}(\{z_1^{(n)}, \dots, z_m^{(n)}\}) + \epsilon}}$$

**Here, we normalize over all outputs of a layer.**

$z_1^{(n)}, \dots, z_m^{(n)}$  : outputs of all units of one layer for training example  $n$

- This is not just a reparametrization. Every example will have its individual transform.
- Can be computed online, also with batch-size 1.
- Works also easily for recurrent neural networks.

# Today

---

- ☒ Basic Training Algorithms
- ☒ Common Pitfalls & Training Considerations
- ☐ Searching Hyperparameters

# Hyperparameter search

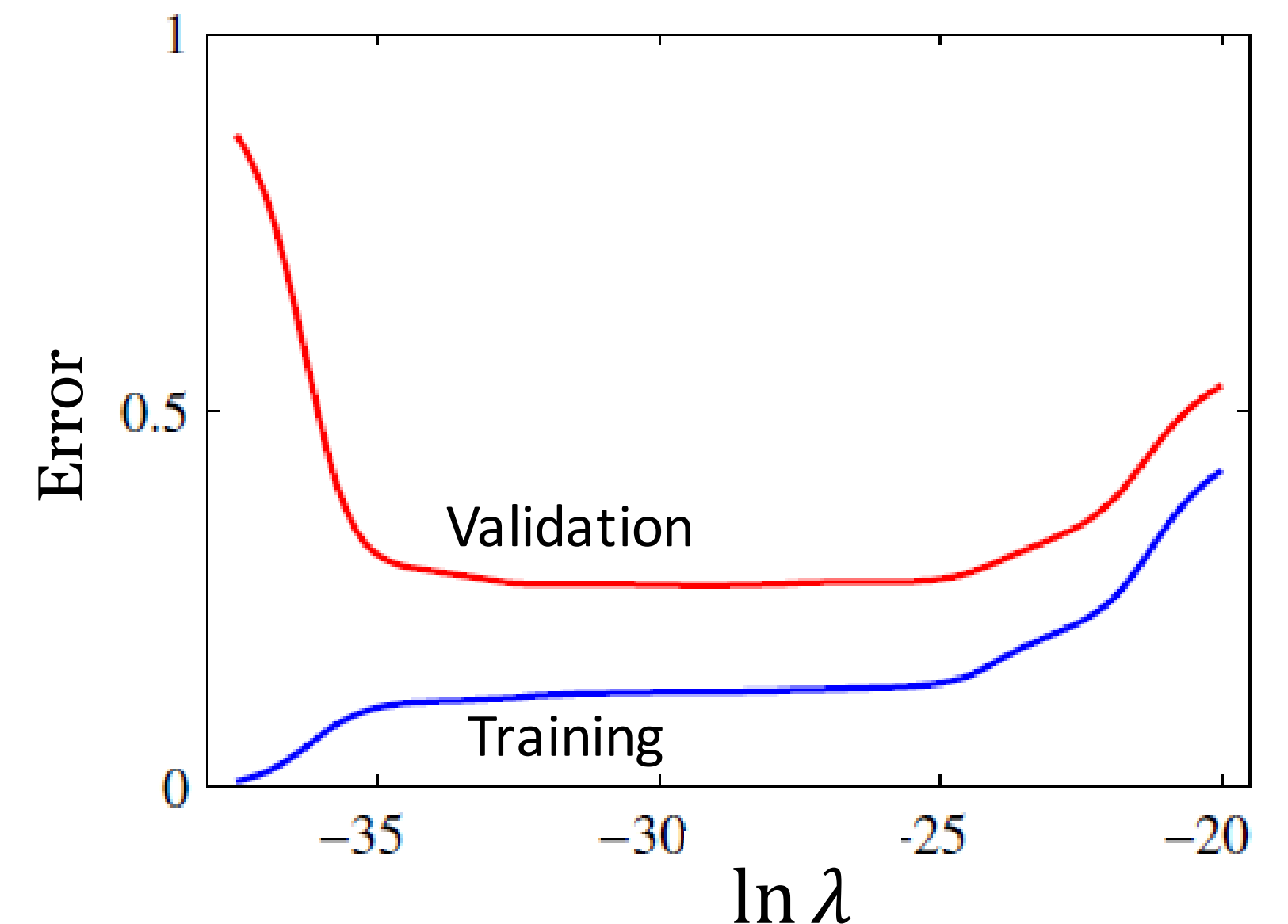
For each application, a set of hyperparameters (h-parameters) has to be defined, such as:

- Network architecture,
- Learning rate, number of iterations, batch size,
- Learning algorithm and its parameters (e.g. momentum),
- Regularization hyperparameters (*we will see next lecture!*)

## Validation

- Performances of different h-parameter settings are compared on a **validation set**.
- When the best setting was determined, one can train the model with this set included in the training data.
- *Example:* Tuning the weight decay regularization parameter:

$$E(w) = \sum_{n=1}^N \{\hat{y}(x^{(n)}) - t^{(n)}\}^2 + \lambda \|w\|^2$$



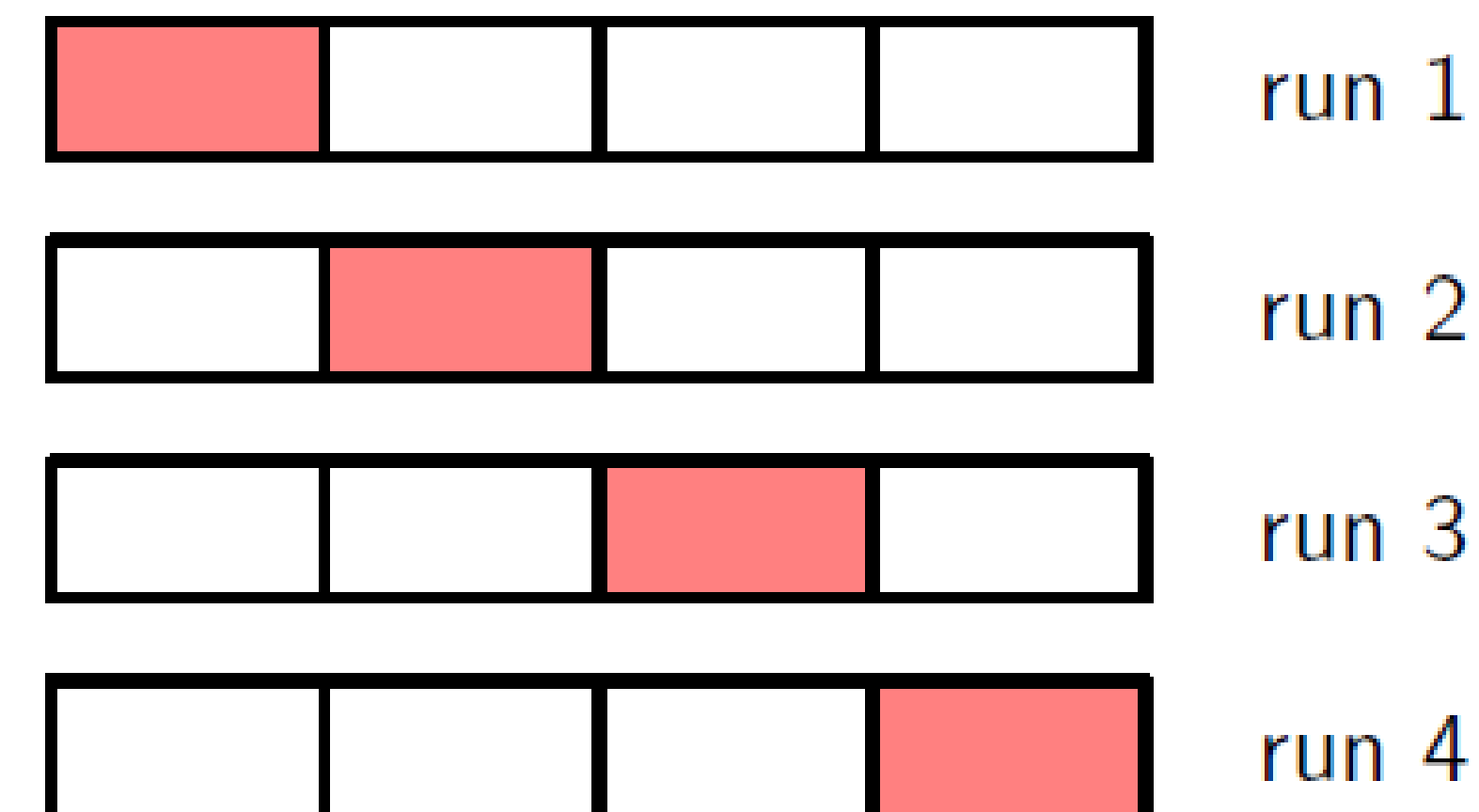
# Cross-validation

## Problem with small data sets:

- When validation set is small, we get a bad estimate of generalization error.
- When training set is small, the model is not well-trained.
- For small data sets, one can use *m-fold cross validation*.

## m-fold Cross Validation:

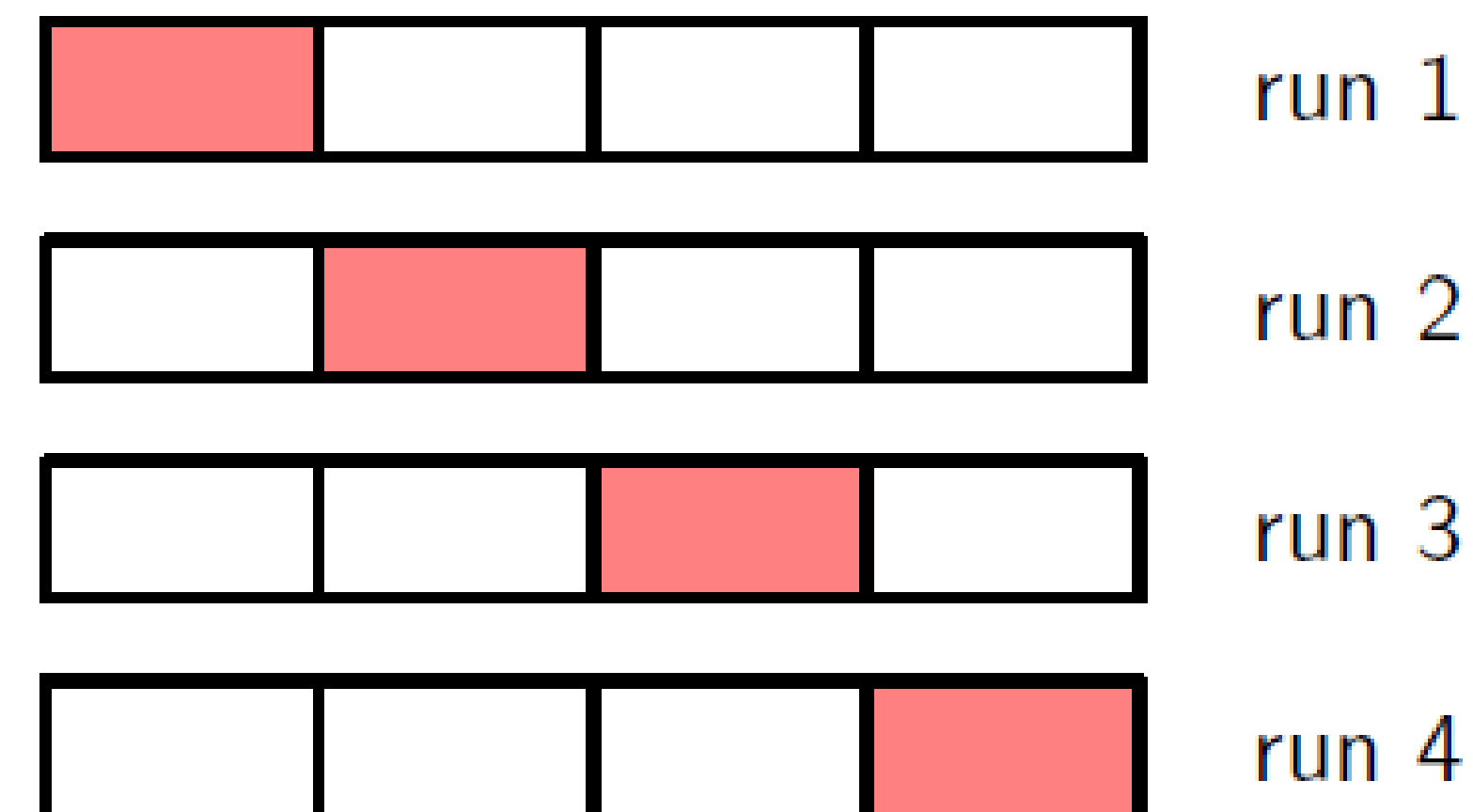
- Divide data  $\mathcal{D}$  into (roughly) equally sized subsets  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_m$
- For  $i = 1$  TO  $m$ 
  - Training set is given by  $\mathcal{D}_1 \cup \dots \cup \mathcal{D}_{i-1} \cup \mathcal{D}_{i+1} \cup \dots \cup \mathcal{D}_m$
  - Validation set is given by  $\mathcal{D}_i$
- One obtains  $m$  validation errors. Compute mean and variance of those.



# Cross-validation

## Example application: Determine regularization parameter $\lambda$

- Make one CV run for each of a set of reasonable  $\lambda$  values
- Take the  $\lambda$  with the best average validation error
- Retrain the model with this  $\lambda$  and the whole data (excluding test set)
- Get an estimate generalization error using the test set.



**Leave-one-out CV:** Special case where  $m$  is equal to the size of the data set.

**Stratified CV:** Distribute the data on the folds such that targets (labels) have approximately the same distribution.



# Options for Hyperparameter search

---

- ▶ **Grid-search** (test many h-parameter settings)
- ▶ **Random search** (improves on Grid-search)
- ▶ **Neural Architecture Search (NAS)** (automating architecture search)
- ▶ **Manual tuning** (often used)

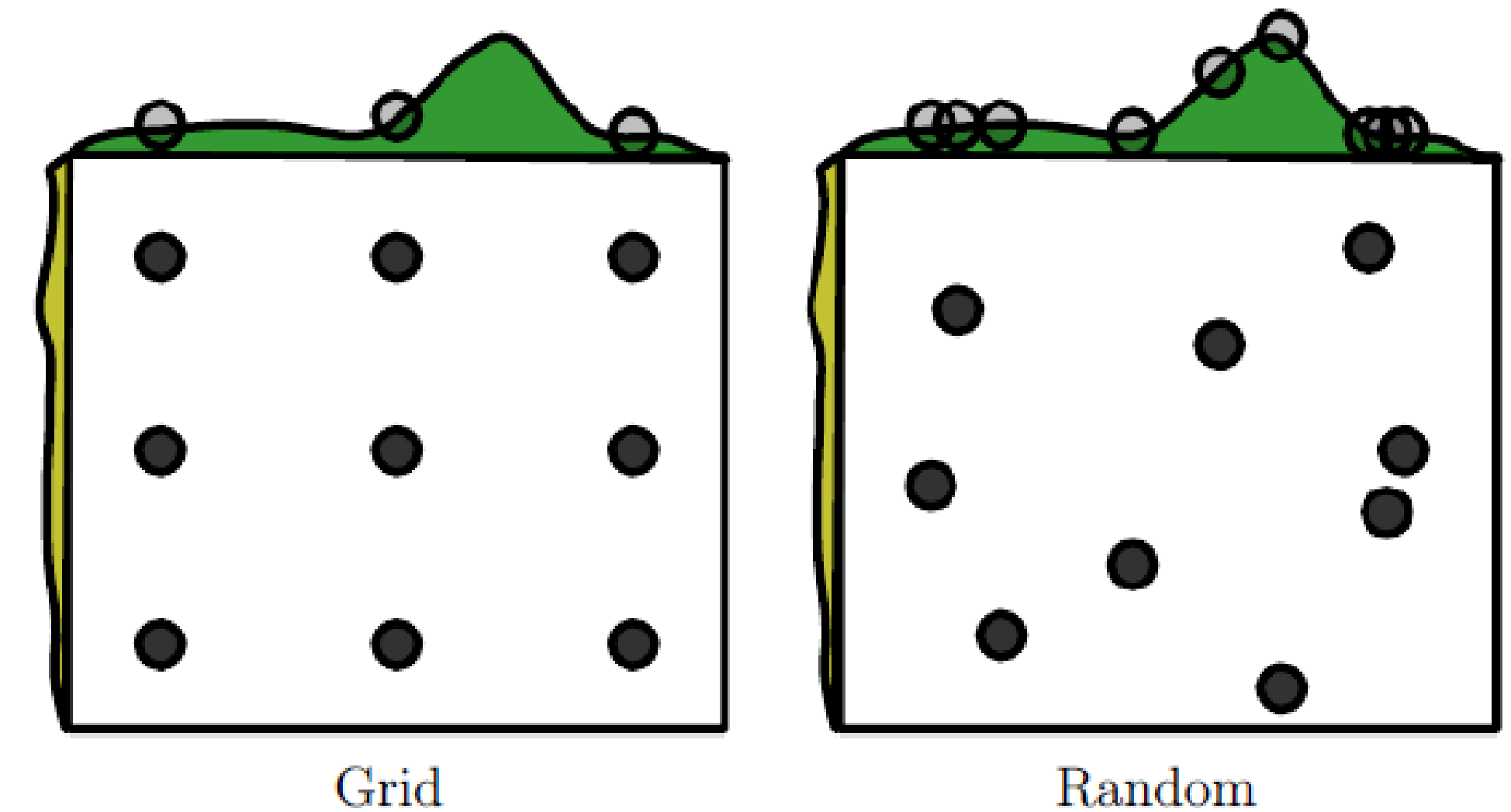
# Grid-Search & Random Search

## Grid-search

- Define a grid on h-parameters values, test each point.

## Random search

- Define interval and distribution for each h-parameter.
- Draw each h-parameter independently.
- Train model with h-parameter setting.
- Advantages: Tests more values for each individual h-parameter.
- Often used distributions: Uniform or uniform on log-scale.



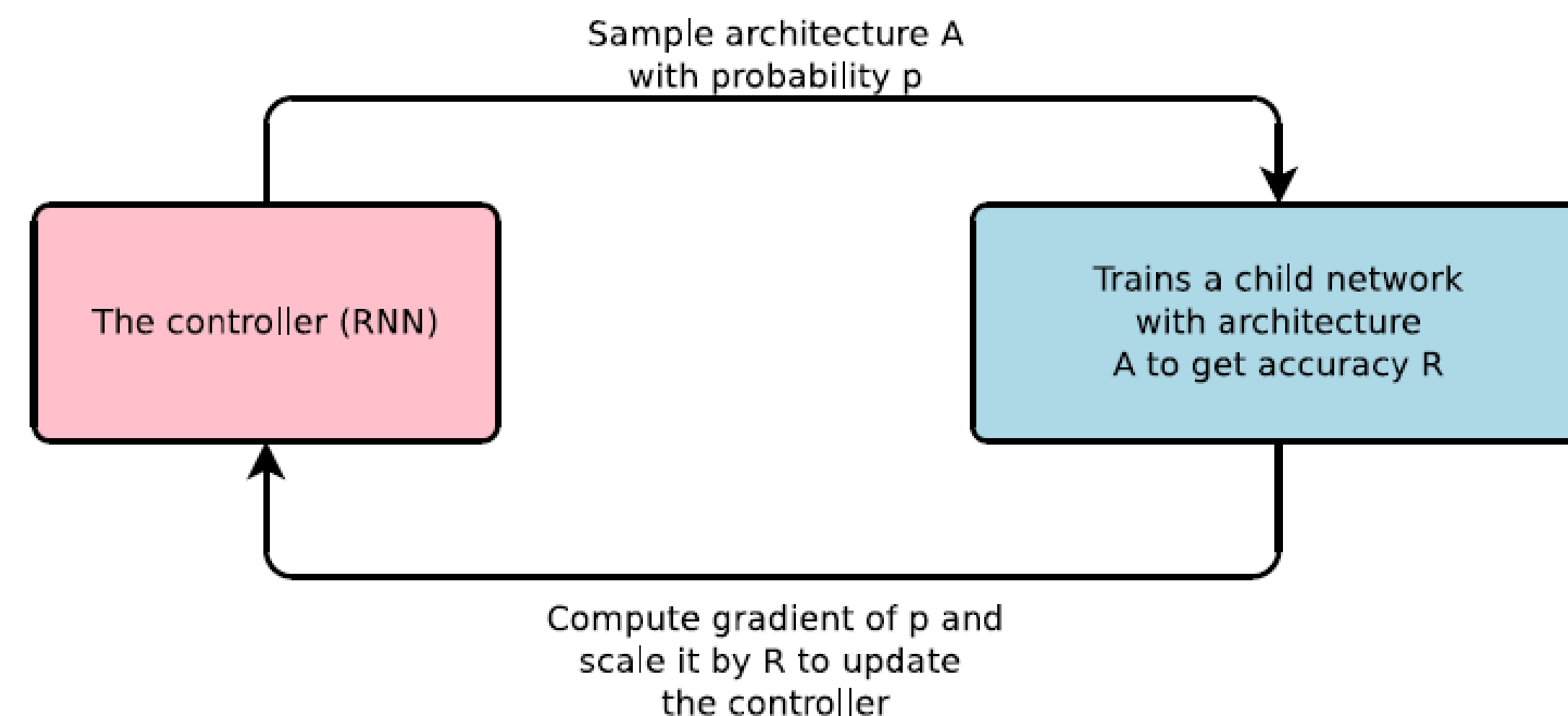
If the data set is small, one can still use m-fold cross validation for each point.

# Neural Architecture Search (NAS)

- A technique for automating the design of neural networks.
- Resulting networks are on par or outperform hand-designed architectures.

## General idea:

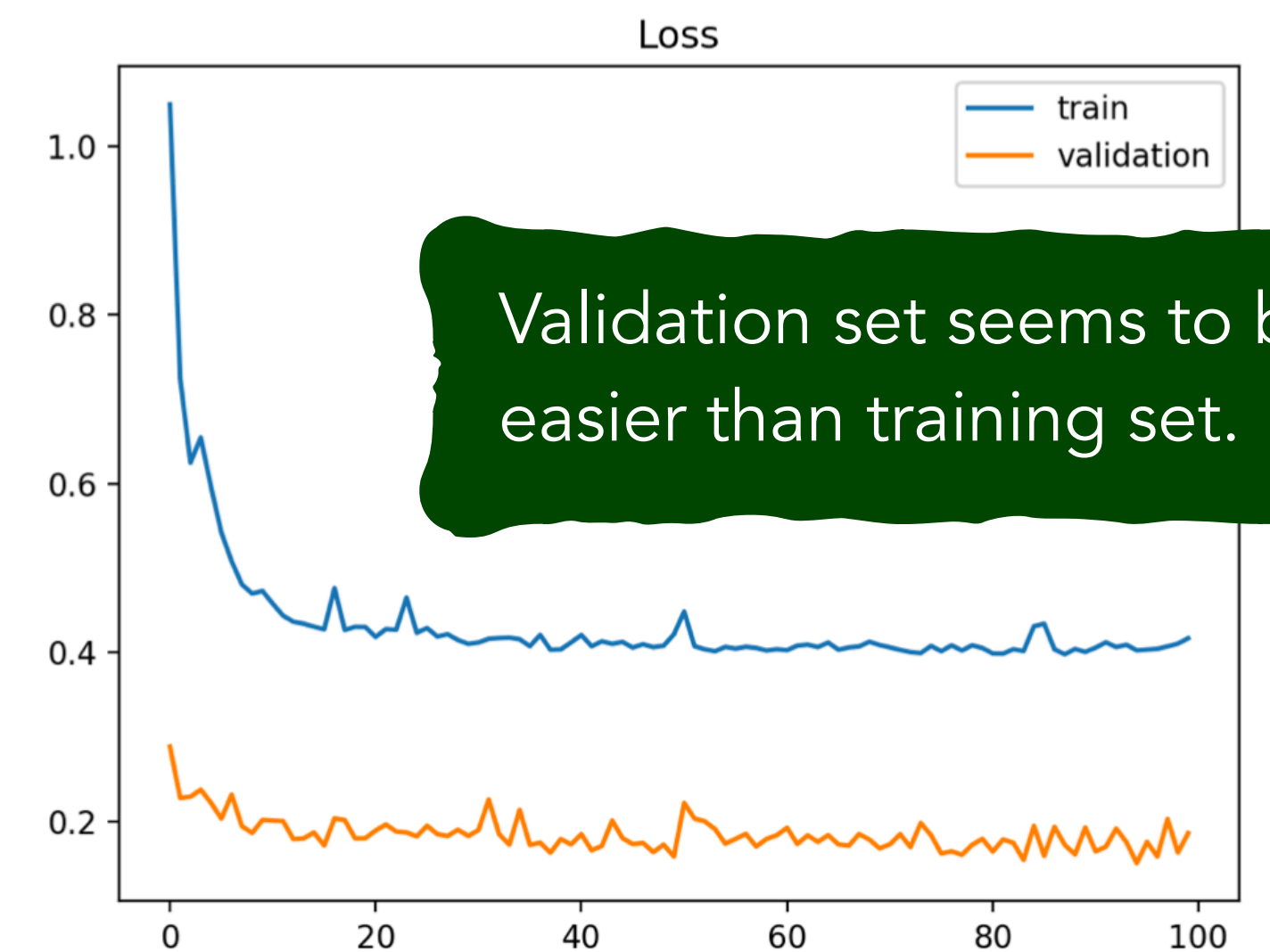
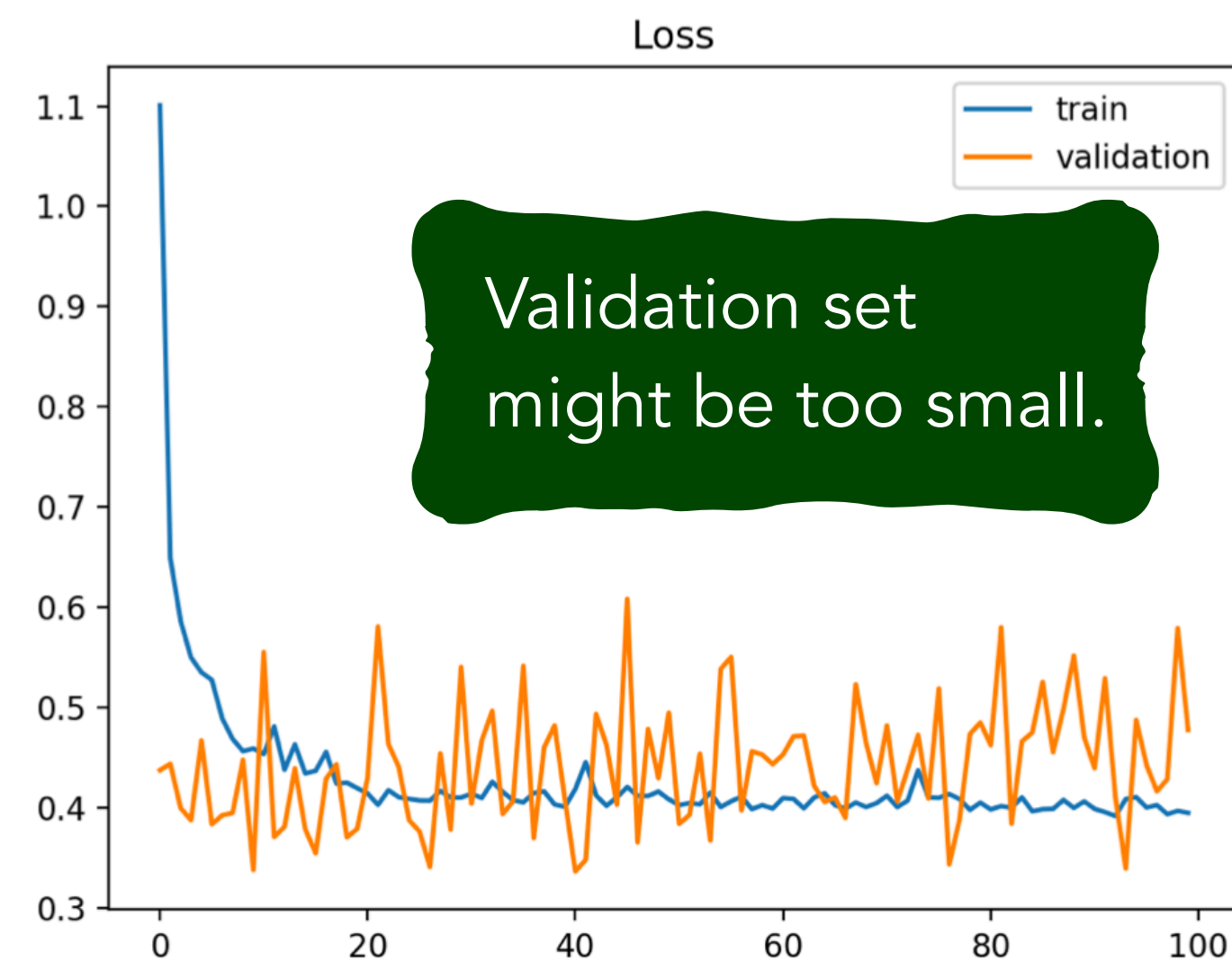
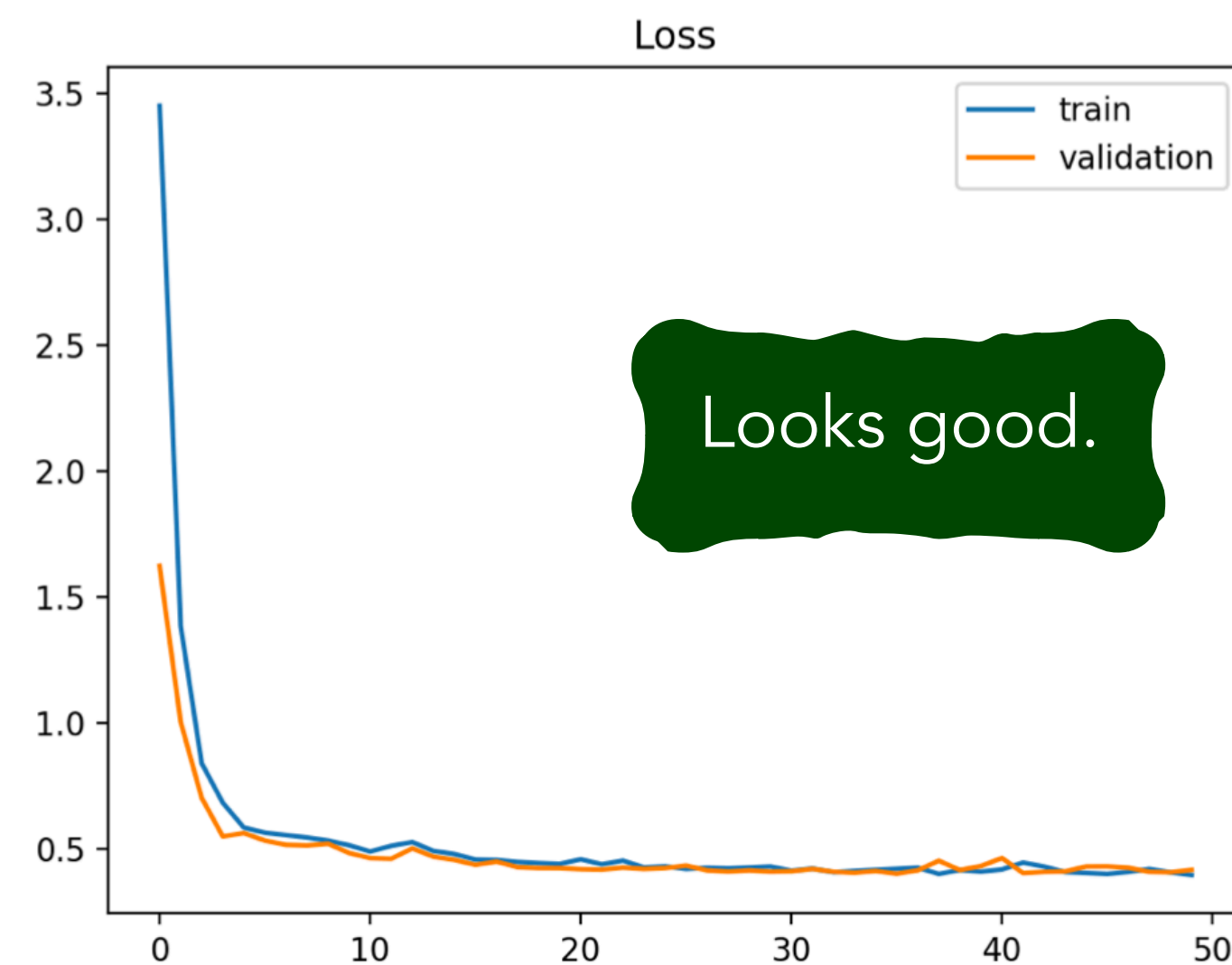
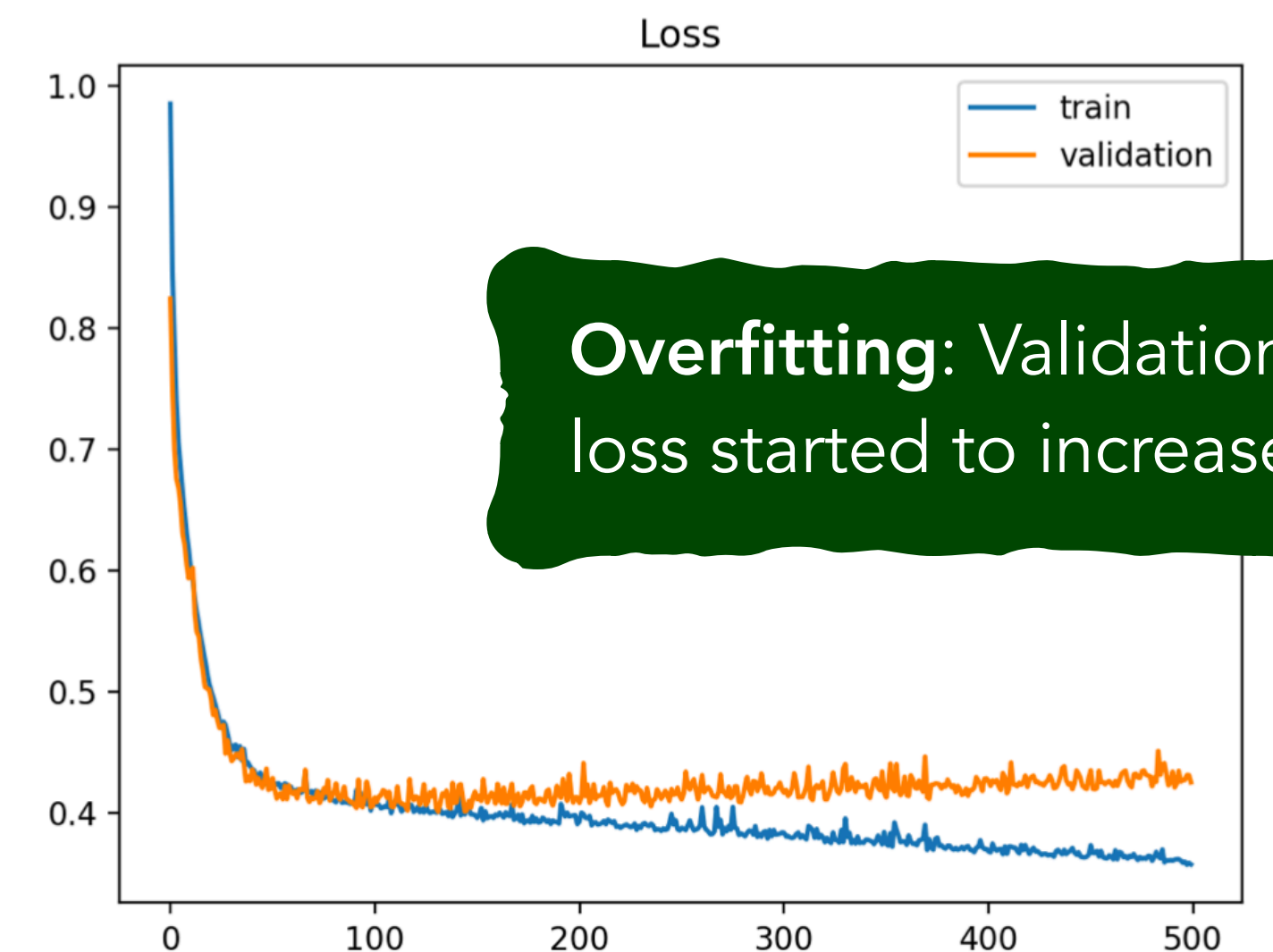
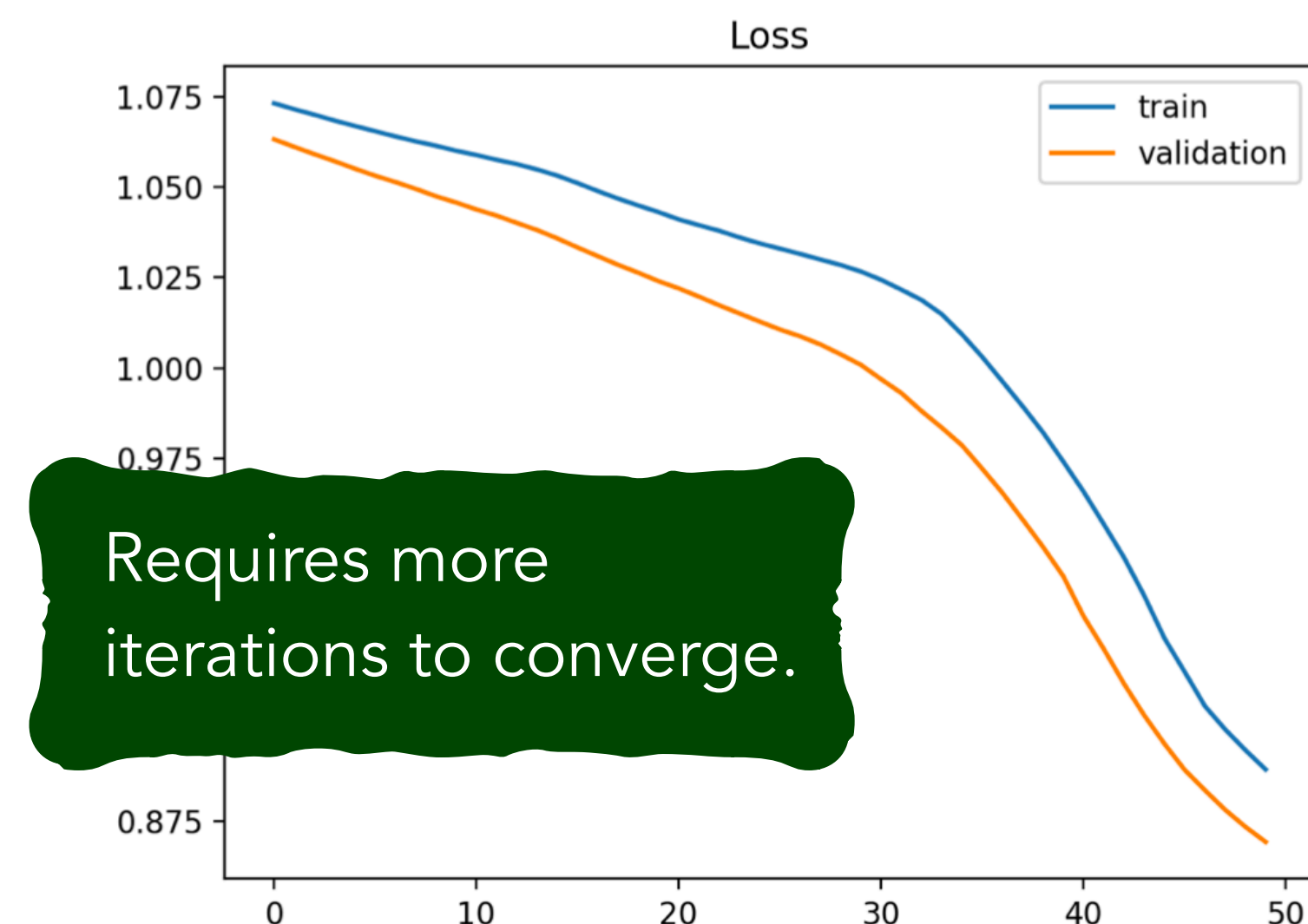
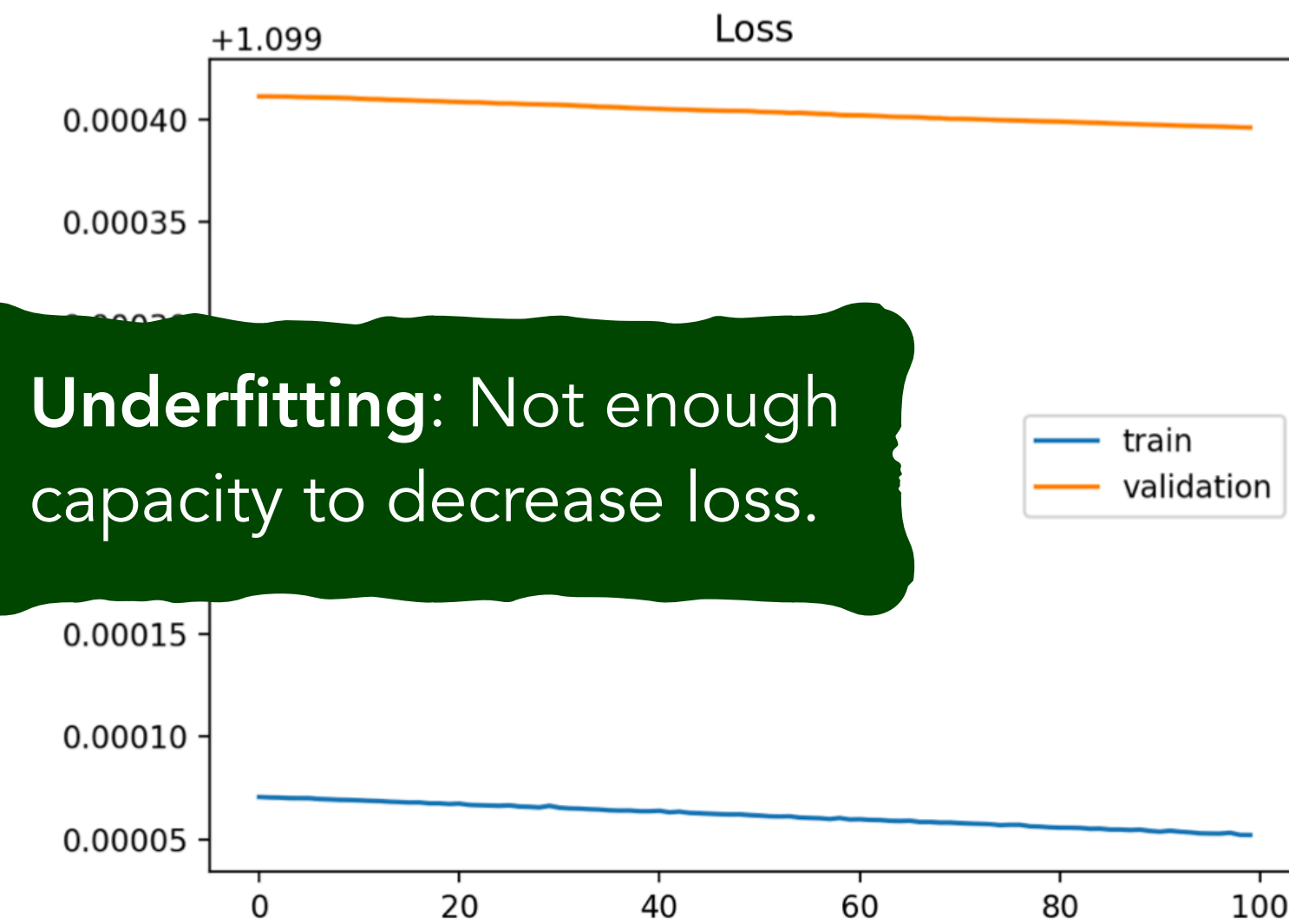
- Use a recurrent neural network (RNN, the controller) that samples network architectures.
- Performance is evaluated, better performance leads to larger reward
- The RNN is optimized with reinforcement learning to maximize the expected reward
- Disadvantage: Very compute-hungry



Zoph, B., & Le, Q. V. (2016). "Neural architecture search with reinforcement learning." *arXiv preprint arXiv:1611.01578*.

# Monitoring the Training Process

slide figures from  
notes by A. Geiger



# Today

---

- ☑ Basic Training Algorithms
- ☑ Common Pitfalls & Training Considerations
- ☑ Searching Hyperparameters

## Questions?