

Introduction to Machine Learning

Part II

Judith Katzy
Hamburg, September 2024

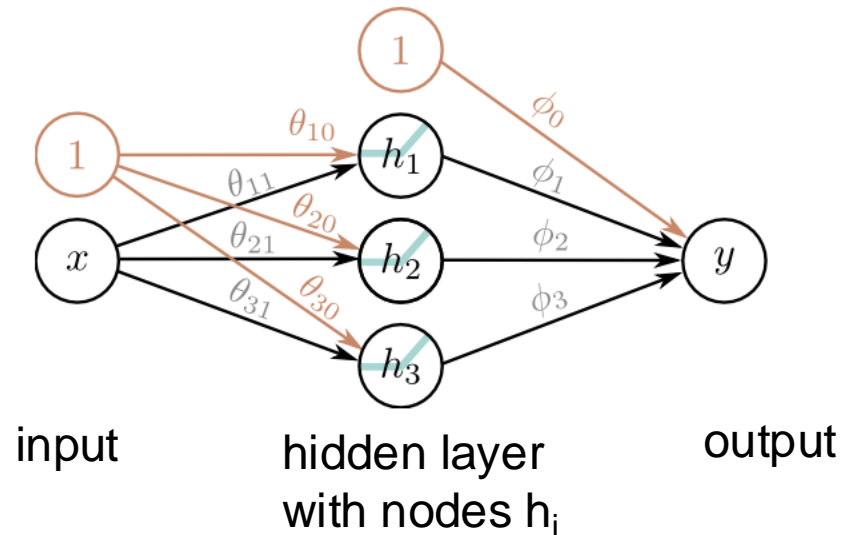


Outline

- Optimisation
- Stochastic learning
- Loss & Regularisation

Neural network family of functions

$$\begin{aligned}y &= f[x, \phi] \\ &= \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]\end{aligned}$$



$$\begin{aligned}\hat{\phi} &= \operatorname{argmin}_{\phi} [L[\phi]] \\ &= \operatorname{argmin}_{\phi} \left[\sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \right]\end{aligned}$$

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

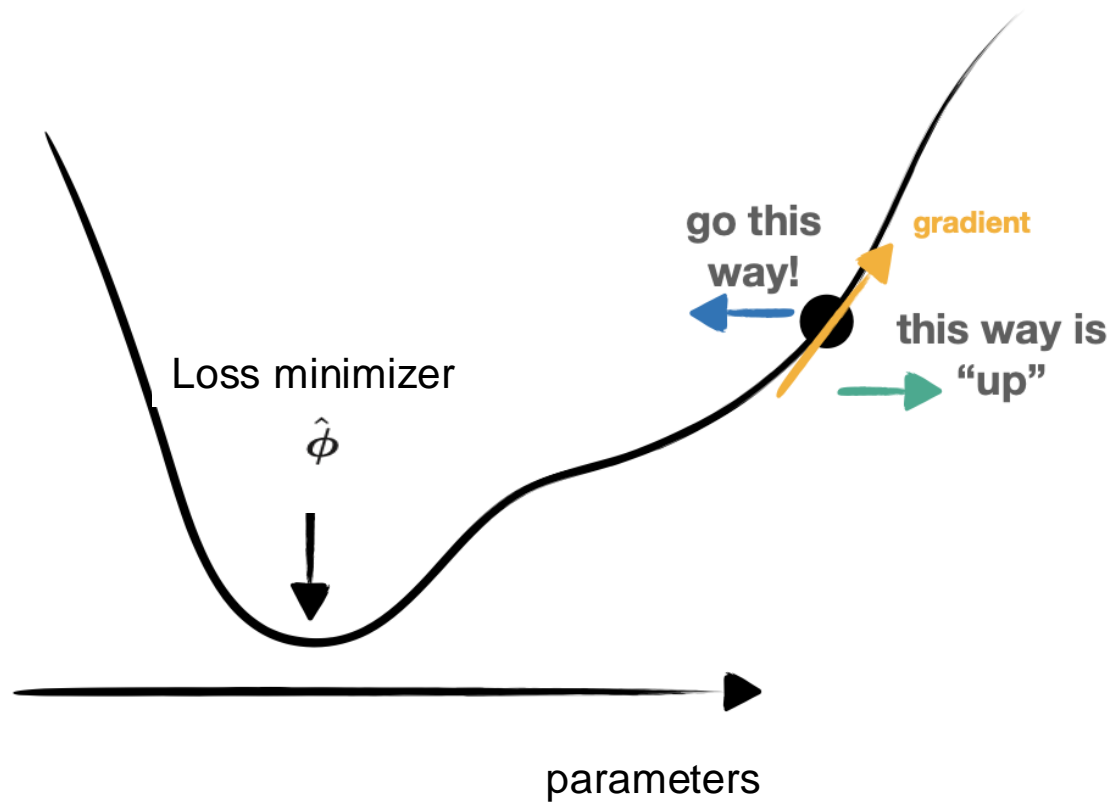
$$h_2 = a[\theta_{20} + \theta_{21}x]$$

$$h_3 = a[\theta_{30} + \theta_{31}x],$$

Iterative Optimisation

Gradient descent

A natural idea is to optimize by walking downhill

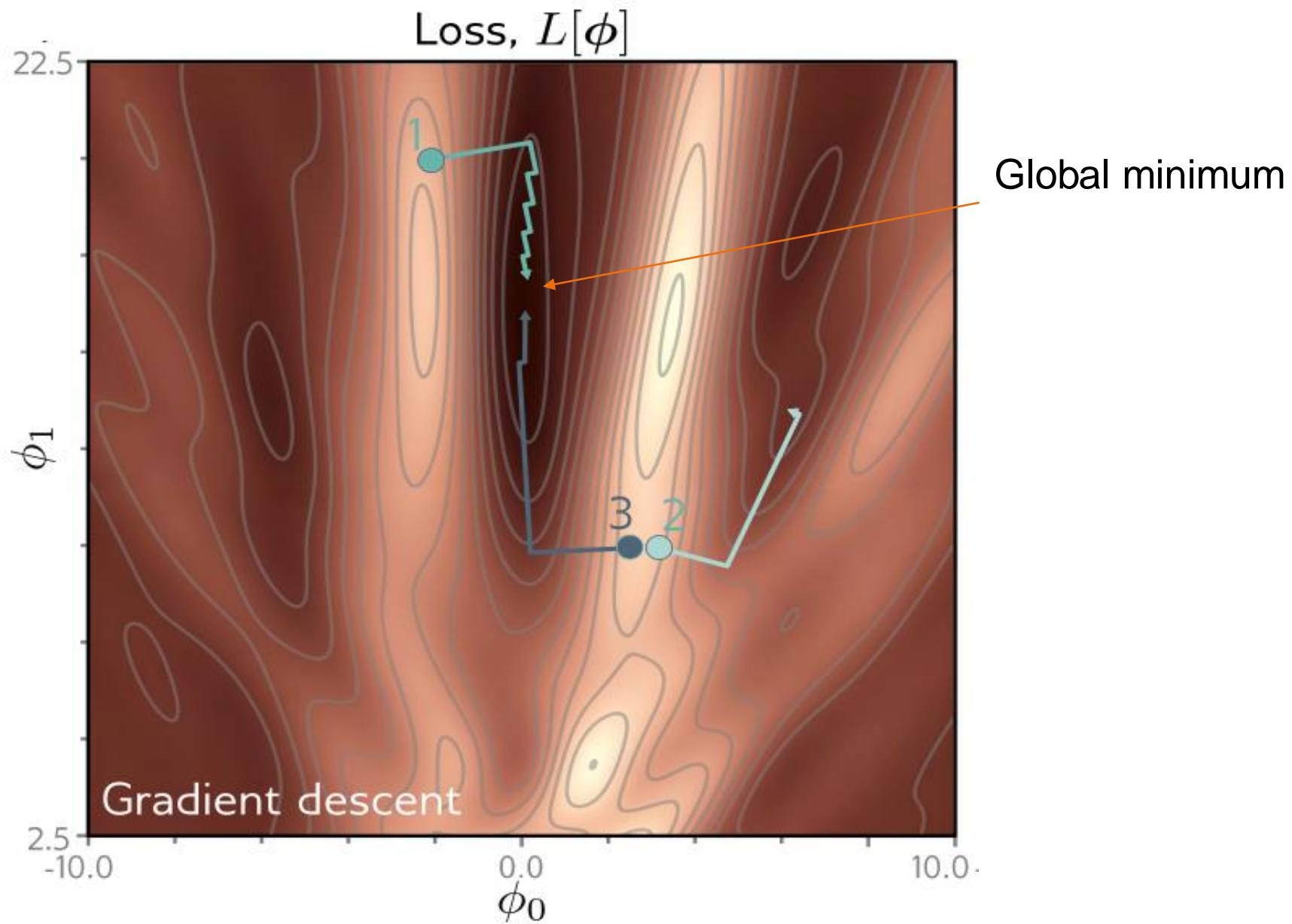


$$\phi \leftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi},$$

alpha: learning rate

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L[\phi]]$$

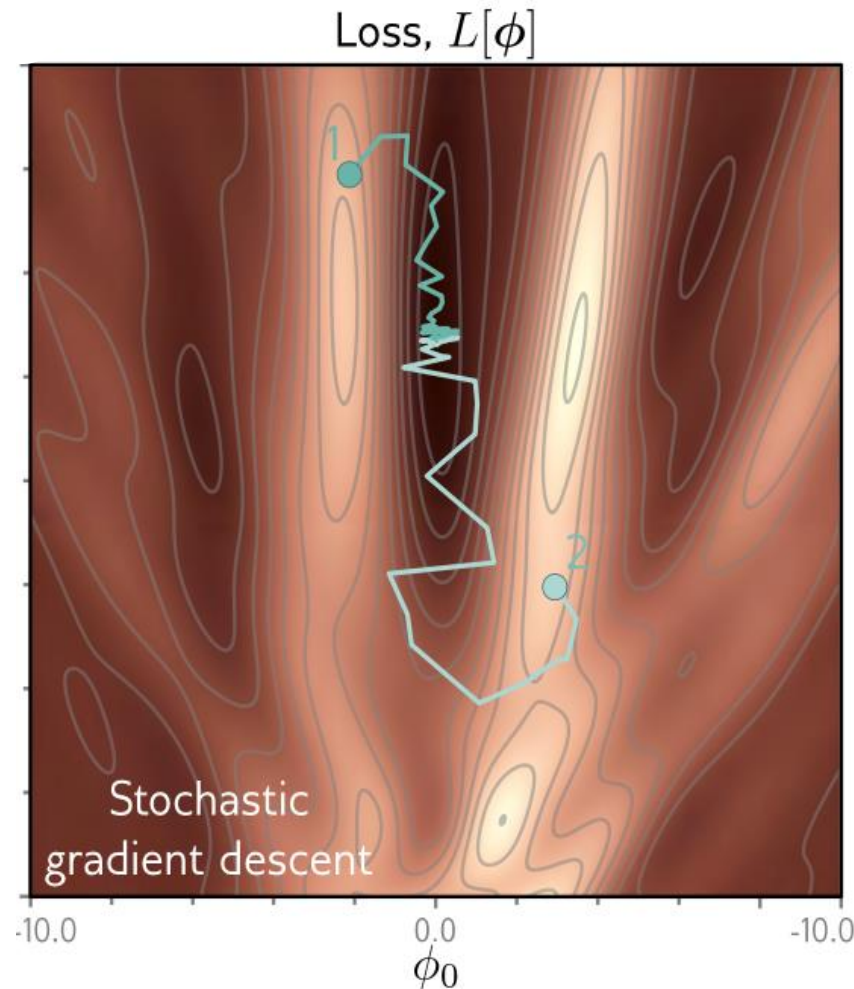
Guaranteed to work for convex functions



Stochastic gradient descent

Remember: goal is generalization not training loss

Evaluating the loss on a small “mini-batch” instead of the full data: useful noise to jump over e.g. local minima.



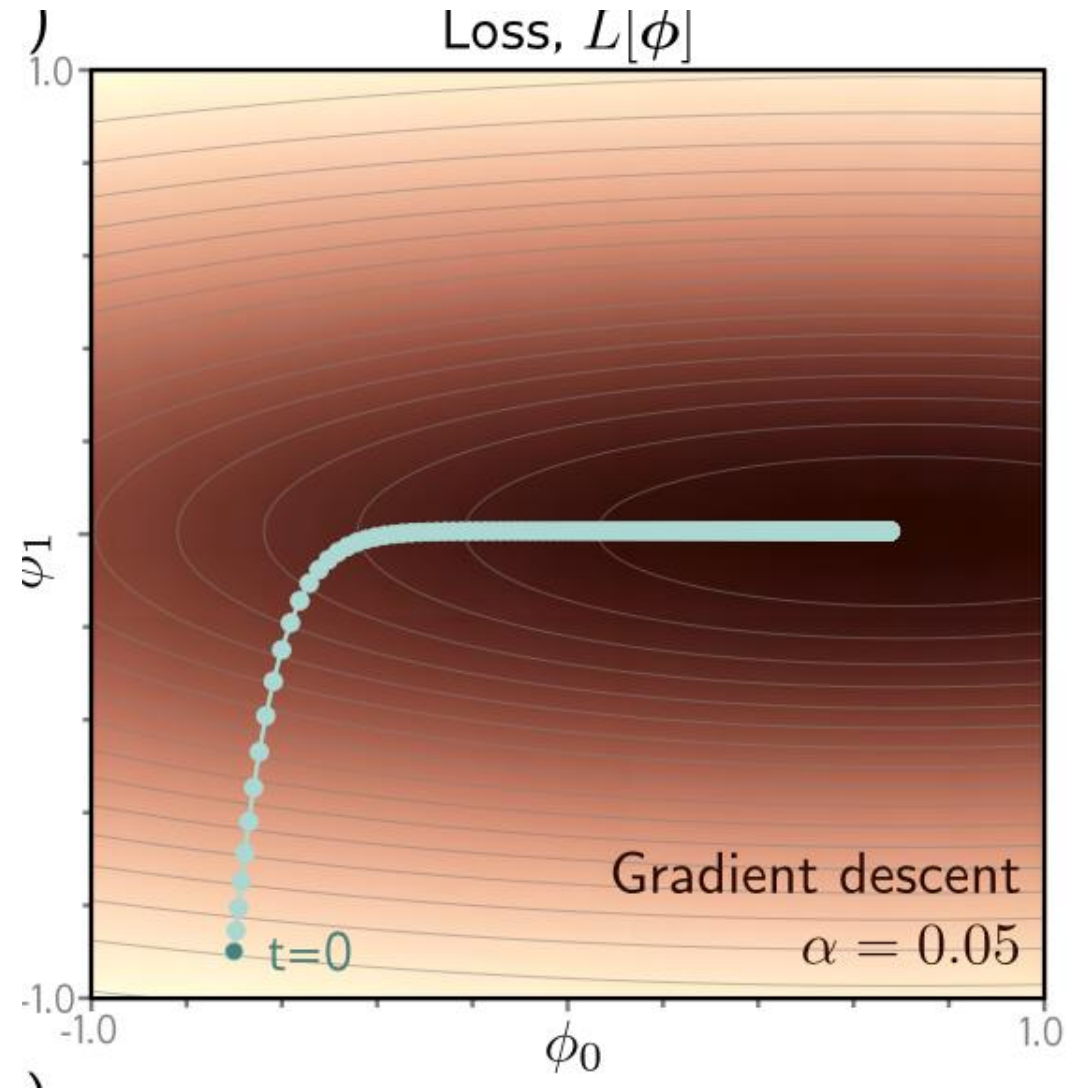
Update loss on (mini)batch

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \sum_{i \in \mathcal{B}_t} \frac{\partial l_i[\phi_t]}{\partial \phi}$$

Additional benefit:

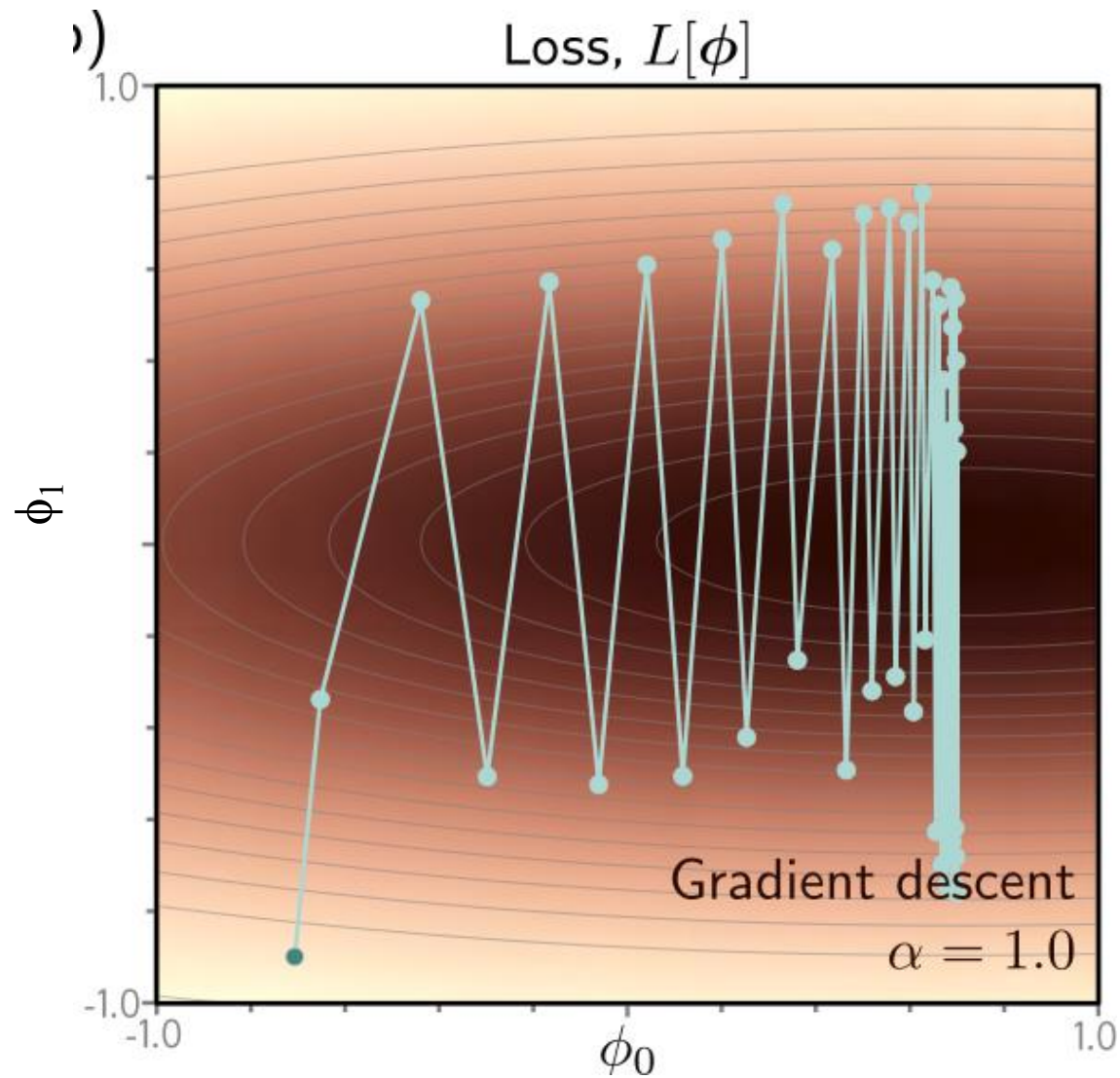
Less computationally expensive

Tuning the learning rate



$$\phi \leftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi},$$

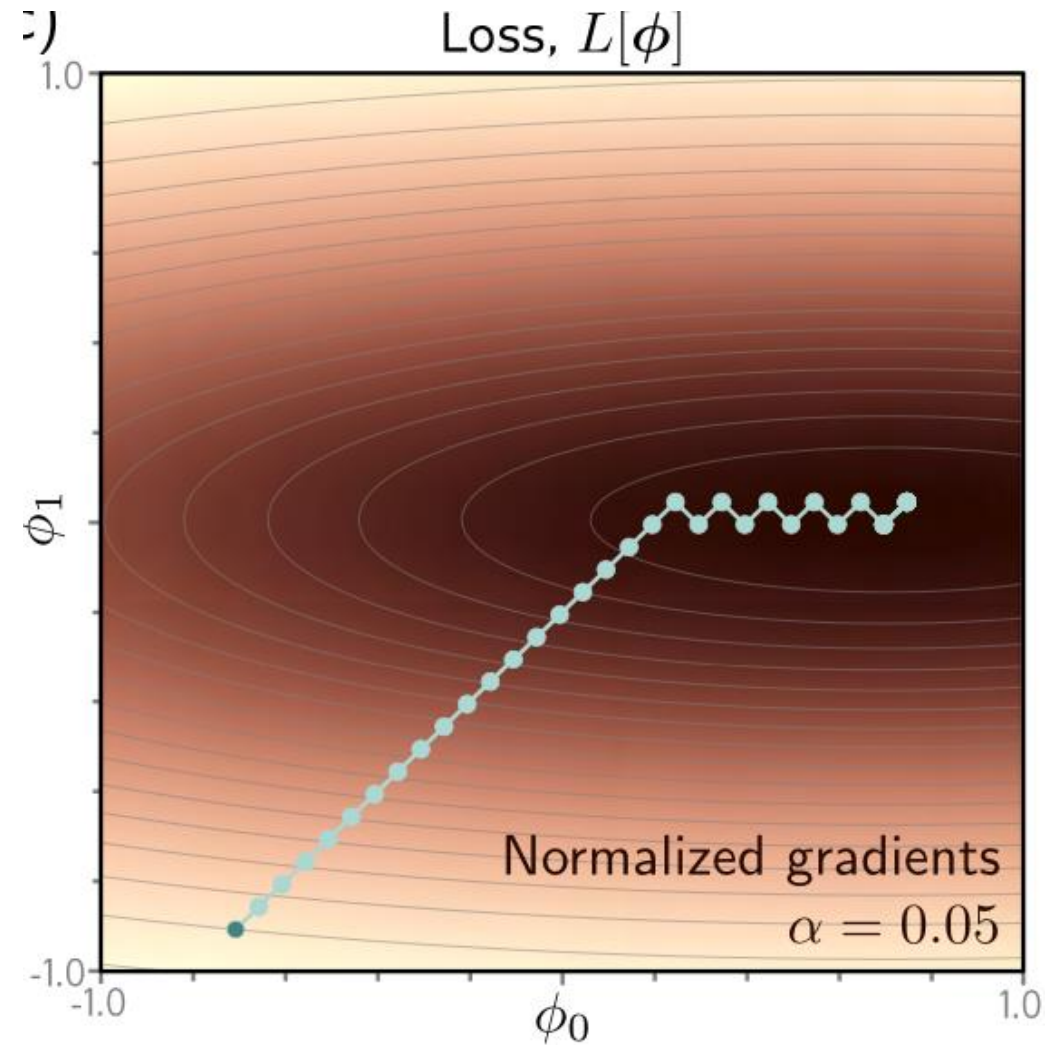
Tuning the learning rate



$$\phi \leftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi},$$

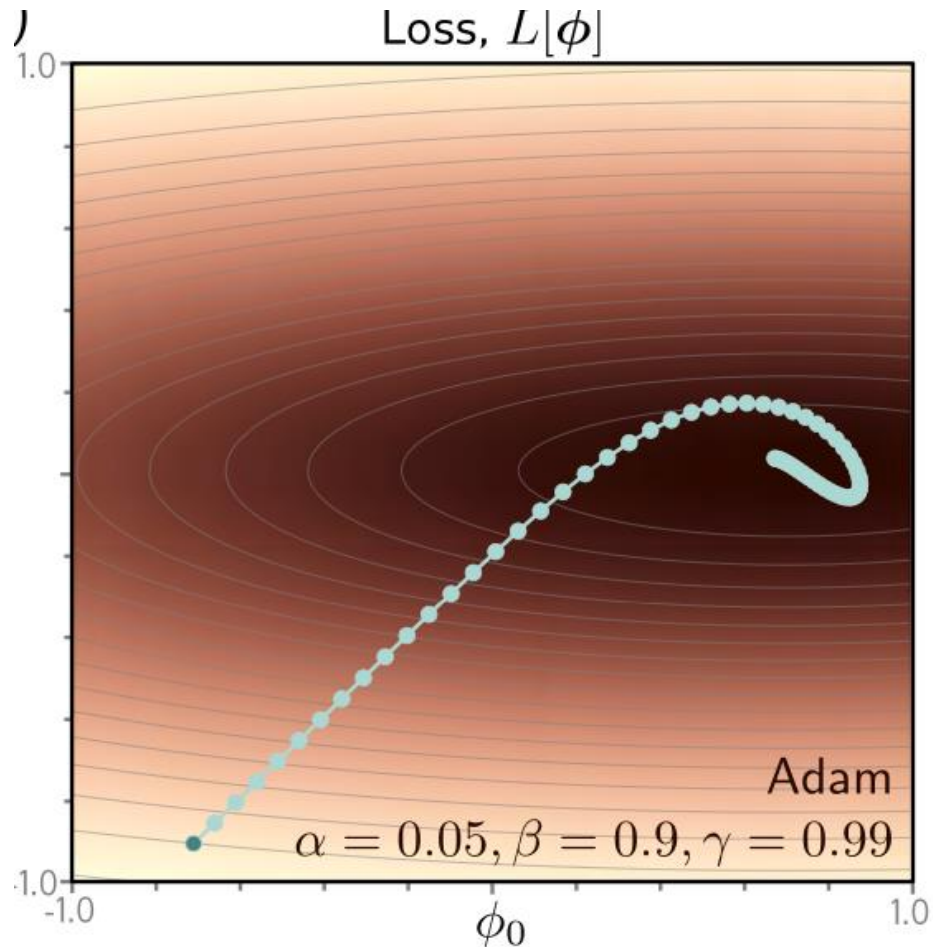
Tuning the learning rate

$$\phi \leftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi},$$



Adapting the learning rate dynamically

Many additional tricks & nuances in practical optimization algorithms to improve convergence for non-convex problems



Example of a good default: ADAM

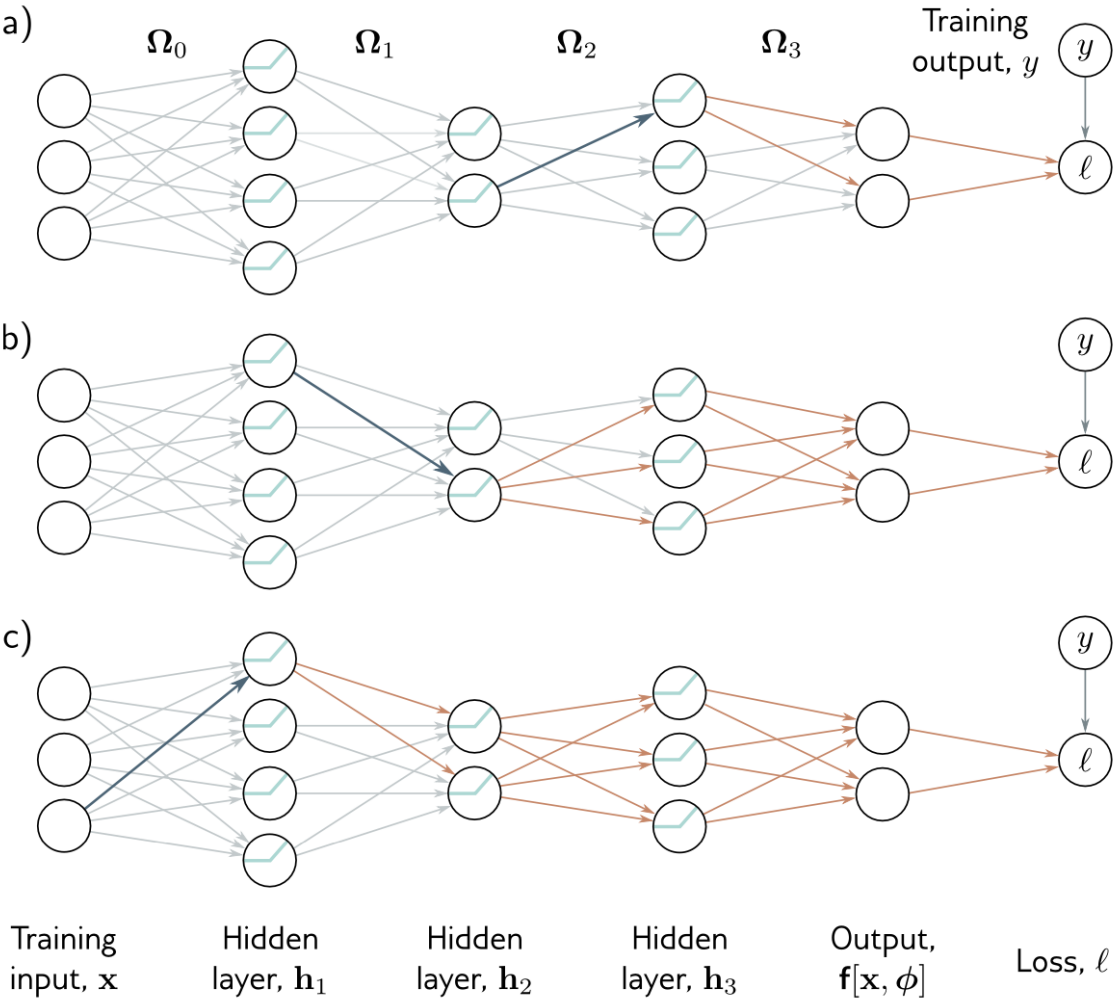
Weighted average over history

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \frac{\partial L[\phi_t]}{\partial \phi}$$

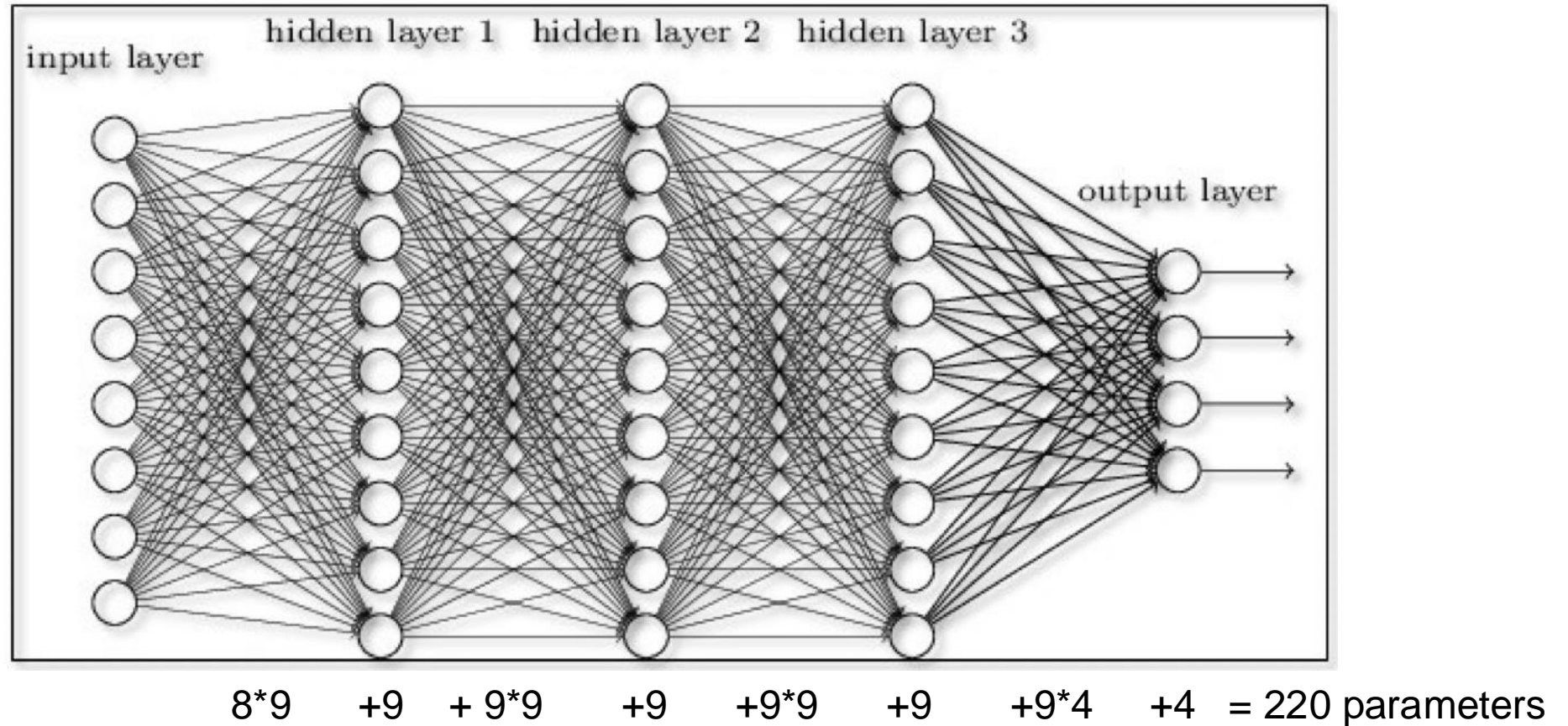
$$\mathbf{v}_{t+1} \leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma) \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1} + \epsilon}},$$

Gradient descent needs.....GRADIENTS!



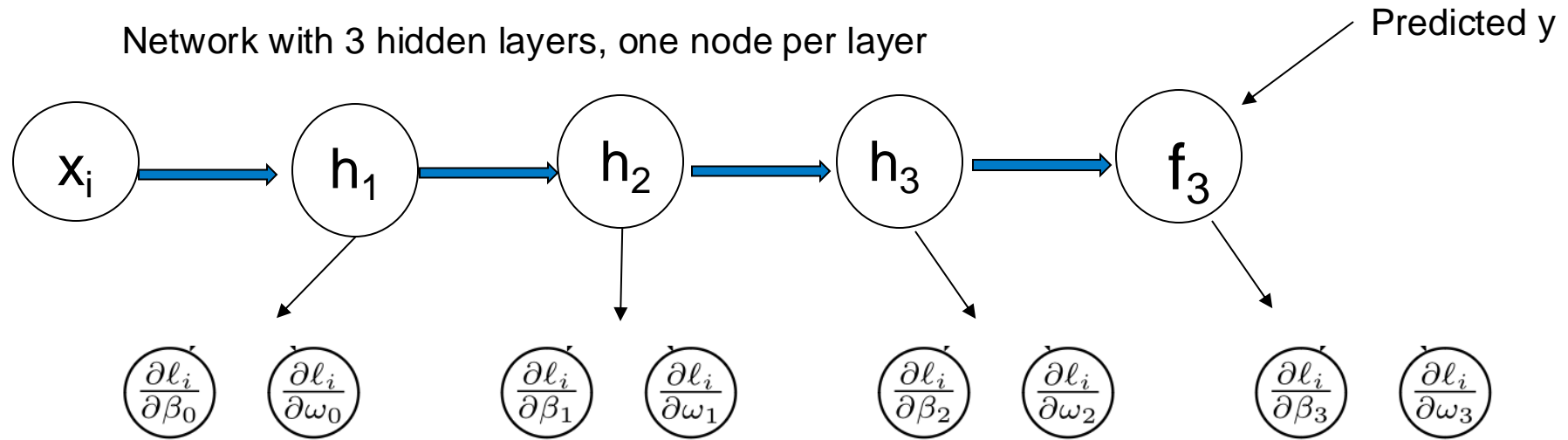
...and there are plenty of gradients



Biggest networks to date have up to $\sim 10^{12}$ parameters (e.g. ChatGPT)

Let's start simple:

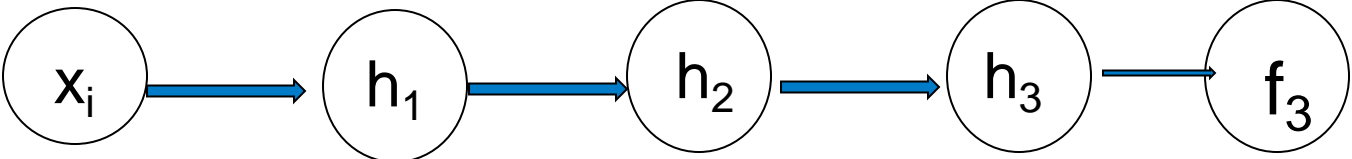
Network with 3 hidden layers, one node per layer



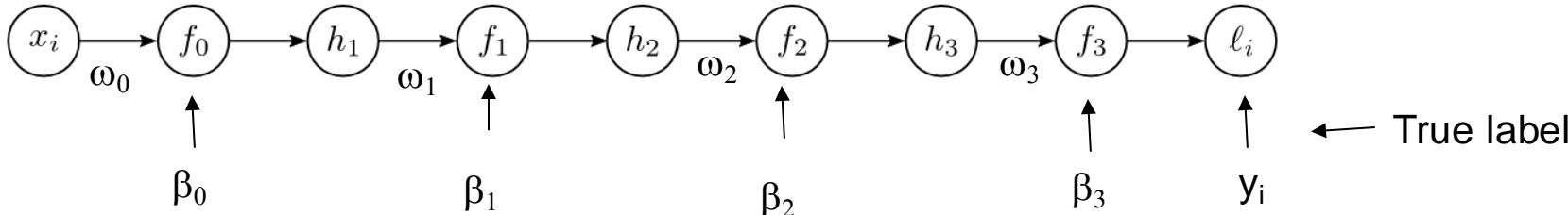
Our goal 8 gradients!

Let's start simple:

Network with 3 hidden layers, one node per layer



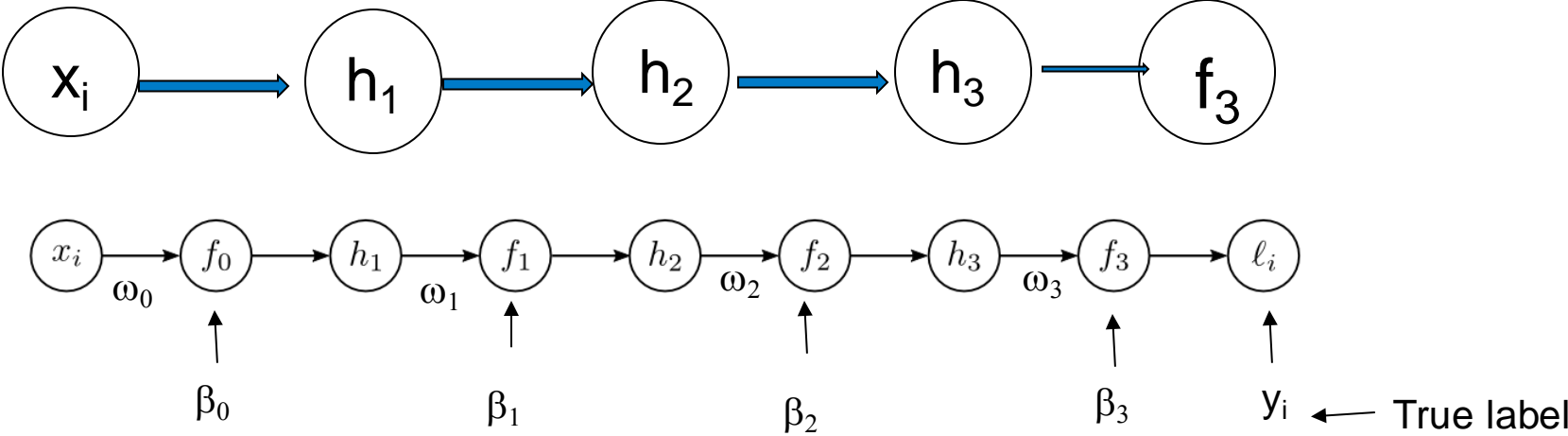
Break down in elements (algorithmic differentiation):



$$\begin{aligned}
 f_0 &= \beta_0 + \omega_0 \cdot x_i \\
 h_1 &= \mathbf{a} [f_0] \\
 f_1 &= \beta_1 + \omega_1 \cdot h_1 \\
 h_2 &= \mathbf{a} [f_1] \\
 f_2 &= \beta_2 + \omega_2 \cdot h_2 \\
 h_3 &= \mathbf{a} [f_2] \\
 f_3 &= \beta_3 + \omega_3 \cdot h_3 \\
 l_i &= (f_3 - y_i)^2.
 \end{aligned}$$

Let's start simple:

Network with 3 hidden layers, one node per layer

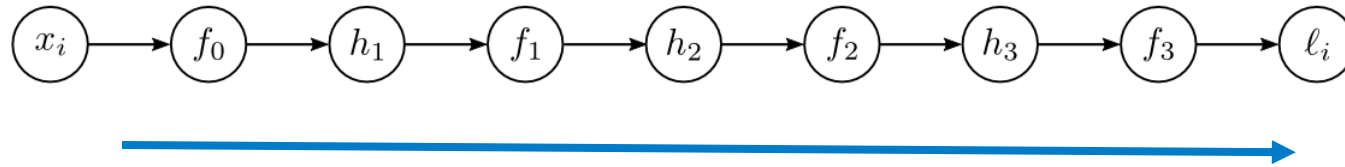


$$\begin{aligned}
 f_0 &= \beta_0 + \omega_0 \cdot x_i \\
 h_1 &= \mathbf{a} [f_0] \\
 f_1 &= \beta_1 + \omega_1 \cdot h_1 \\
 h_2 &= \mathbf{a} [f_1] \\
 f_2 &= \beta_2 + \omega_2 \cdot h_2 \\
 h_3 &= \mathbf{a} [f_2] \\
 f_3 &= \beta_3 + \omega_3 \cdot h_3 \\
 l_i &= (f_3 - y_i)^2.
 \end{aligned}$$

$$l_i (\omega_0, \beta_0, \omega_1, \beta_1, \omega_2, \beta_2, \omega_3, \beta_3) = (f_3(h_3(f_2(h_2(f_1(h_1(f_0(\omega_0, \beta_0)))))))) - y_i)^2$$

Backpropagation: collect the ingredients for gradients

Forward pass



Store these intermediate values for later use

Backpropagation: collect ingredients for gradients

Calculate & store derivatives with respect to intermediate variables

$$l_i(\omega_0, \beta_0, \omega_1, \beta_1, \omega_2, \beta_2, \omega_3, \beta_3) = (f_3(h_3(f_2(h_2(f_1(h_1(f_0(\omega_0, \beta_0)))))))) - y_i)^2$$

$$\frac{\partial l_i}{\partial f_3} = 2(f_3 - y_i).$$

$$\frac{\partial l_i}{\partial h_3} = \frac{\partial f_3}{\partial h_3} \frac{\partial l_i}{\partial f_3}.$$

$$\frac{\partial l_i}{\partial f_2} = \frac{\partial h_3}{\partial f_2} \left(\frac{\partial f_3}{\partial h_3} \frac{\partial l_i}{\partial f_3} \right)$$

$$\frac{\partial l_i}{\partial h_2} = \frac{\partial f_2}{\partial h_2} \left(\frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial l_i}{\partial f_3} \right)$$

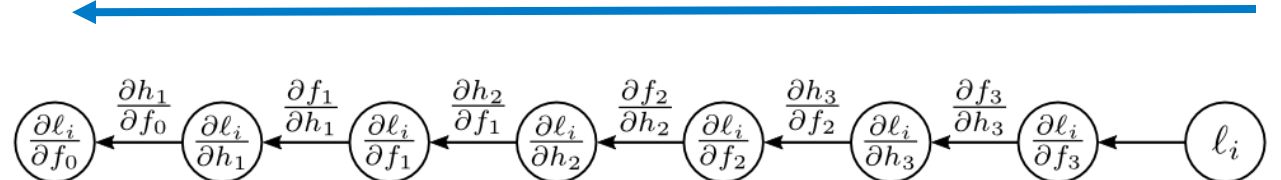
$$\frac{\partial l_i}{\partial f_1} = \frac{\partial h_2}{\partial f_1} \left(\frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial l_i}{\partial f_3} \right)$$

$$\frac{\partial l_i}{\partial h_1} = \frac{\partial f_1}{\partial h_1} \left(\frac{\partial h_2}{\partial f_1} \frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial l_i}{\partial f_3} \right)$$

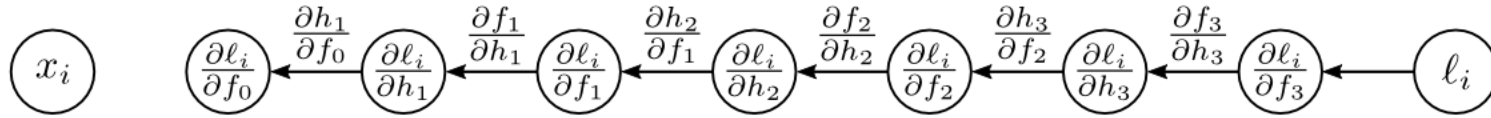
$$\frac{\partial l_i}{\partial f_0} = \frac{\partial h_1}{\partial f_0} \left(\frac{\partial f_1}{\partial h_1} \frac{\partial h_2}{\partial f_1} \frac{\partial f_2}{\partial h_2} \frac{\partial h_3}{\partial f_2} \frac{\partial f_3}{\partial h_3} \frac{\partial l_i}{\partial f_3} \right)$$

Reuse previously calculated derivative

Computationally very efficient!



Backpropagation: derivatives of simple functions

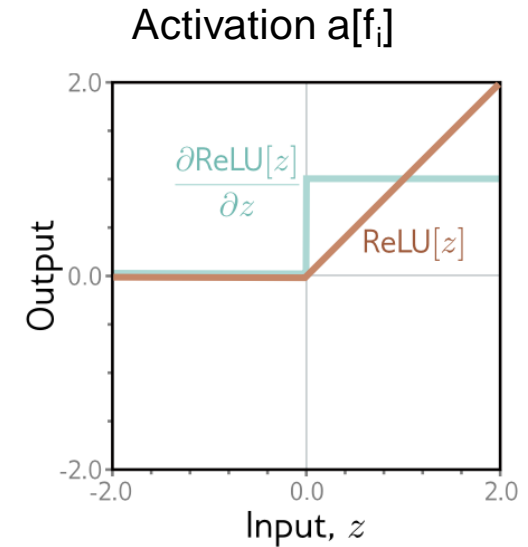


$$h_3 = a[f_2]$$

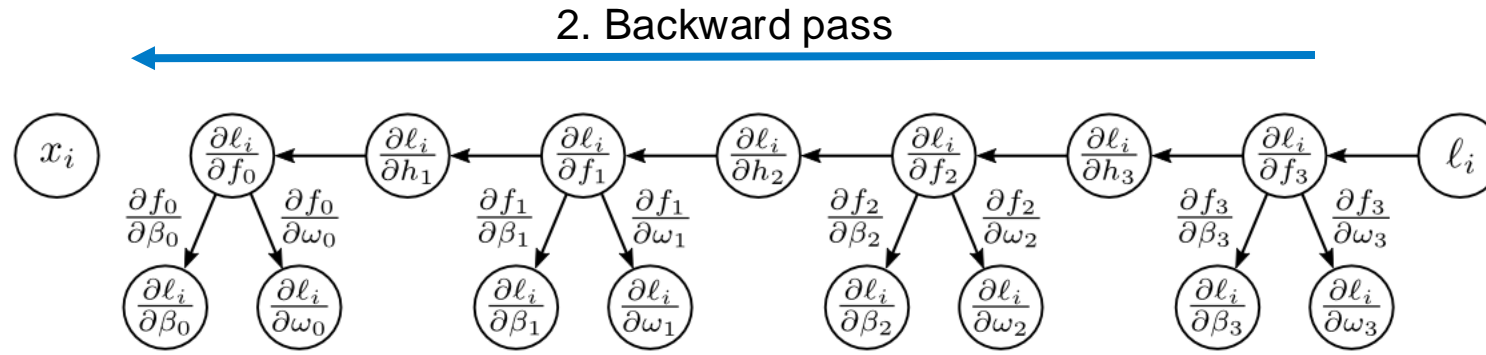
$$f_3 = \beta_3 + \omega_3 \cdot h_3$$

$$\frac{\partial h_3}{\partial f_2} = \begin{cases} 1 & \text{for } f_2 > 0 \\ 0 & \text{for } f_2 \leq 0 \end{cases}$$

$$\frac{\partial f_3}{\partial h_3} = \omega_3$$



Backpropagation: collect ingredients



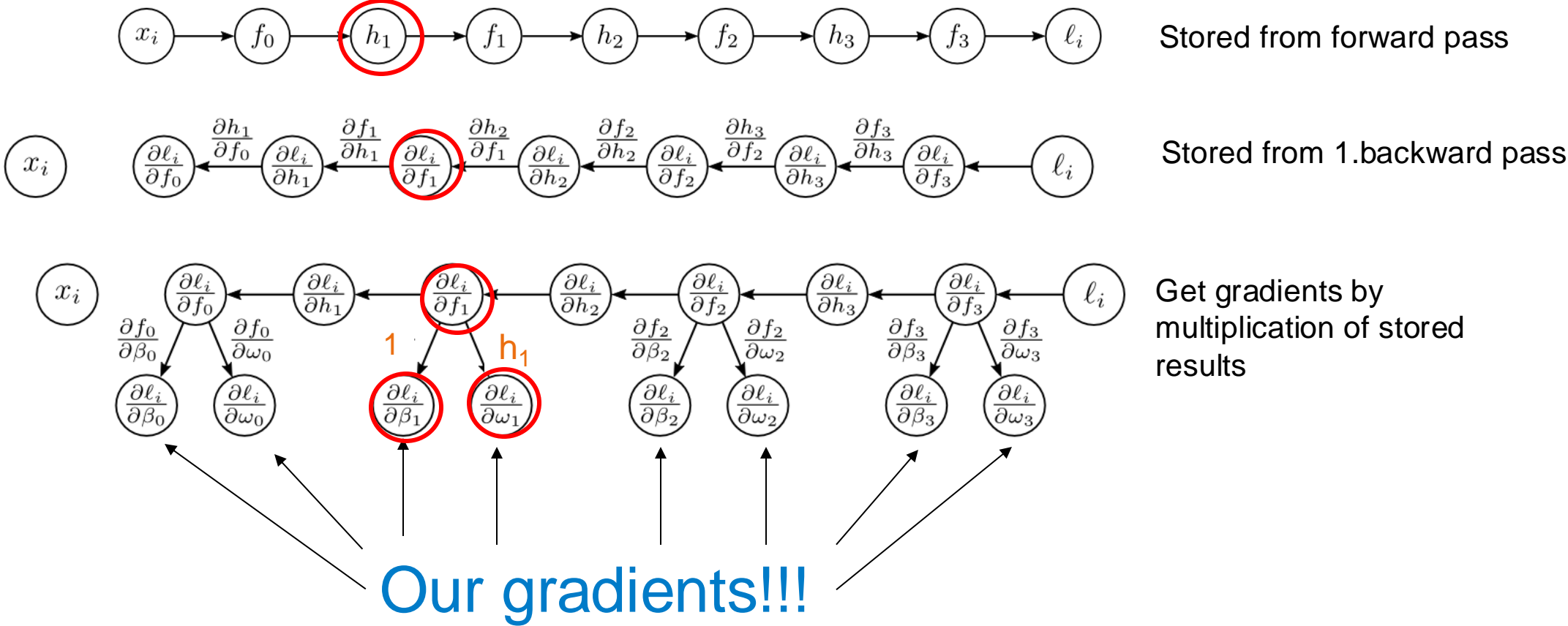
Finally gradient of loss with respect to weights and biases:

$$\frac{\partial f_k}{\partial \beta_k} = 1 \quad \text{and} \quad \frac{\partial f_k}{\partial \omega_k} = h_k.$$

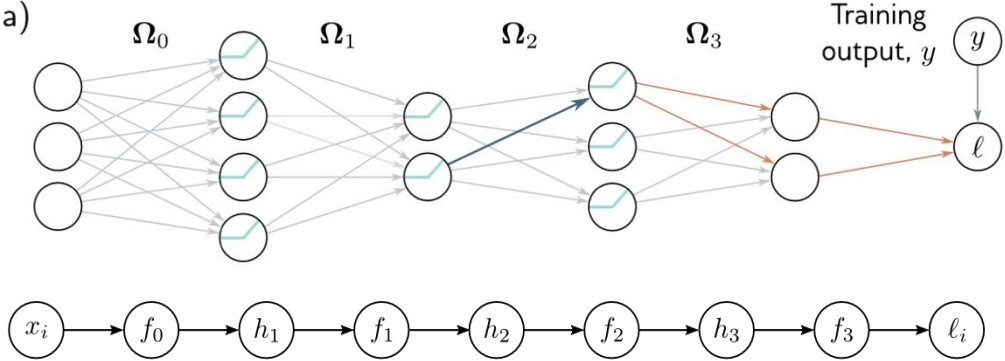
Recall: $f_k = \beta_k + \omega_k h_k$

Already stored in forward pass

Putting it all together...



Now complex network...forward pass

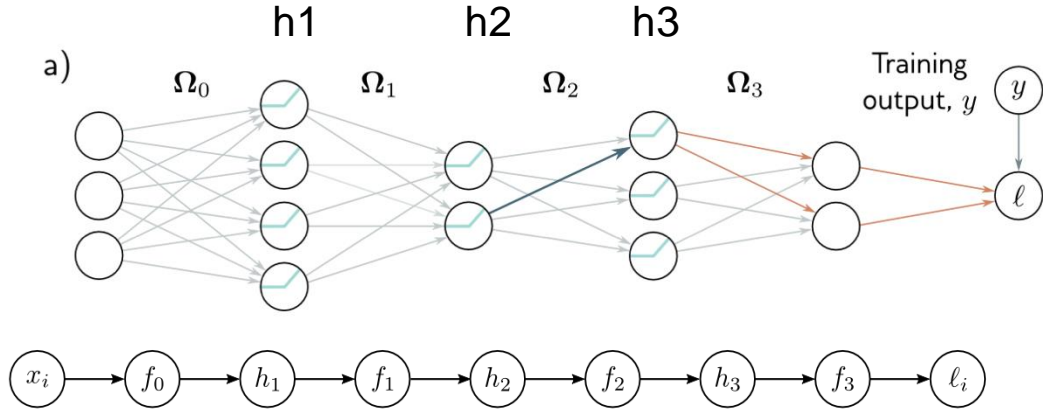


Ω_i : matrix, x_i, b_i, h_i, f_i : vector:

| | | |
|-------|----------------------------|--------------|
| x_i | | 3×1 |
| f_0 | $= \beta_0 + \Omega_0 x_i$ | 4×1 |
| h_1 | $= a[f_0]$ | 4×1 |
| f_1 | $= \beta_1 + \Omega_1 h_1$ | 2×1 |
| h_2 | $= a[f_1]$ | 2×1 |
| f_2 | $= \beta_2 + \Omega_2 h_2$ | 3×1 |
| h_3 | $= a[f_2]$ | 3×1 |
| f_3 | $= \beta_3 + \Omega_3 h_3$ | 2×1 |
| l_i | $= l[f_3, y_i]$, | 1×1 |

$$\Omega_0 = \begin{pmatrix} \omega_{11} & \omega_{12} & \omega_{13} \\ \omega_{21} & \omega_{22} & \omega_{23} \\ \omega_{31} & \omega_{32} & \omega_{33} \\ \omega_{41} & \omega_{42} & \omega_{43} \end{pmatrix} \begin{pmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{pmatrix}$$

Now complex model...backward pass



$$\begin{aligned}
 \mathbf{f}_2 &= \beta_2 + \mathbf{\Omega}_2 \mathbf{h}_2 && 3 \times 1 \\
 \mathbf{h}_3 &= \mathbf{a}[\mathbf{f}_2] && 3 \times 1 \\
 \mathbf{f}_3 &= \beta_3 + \mathbf{\Omega}_3 \mathbf{h}_3 && 2 \times 1 \\
 l_i &= l[\mathbf{f}_3, y_i], && 1 \times 1
 \end{aligned}$$

$$\frac{\partial l_i}{\partial \mathbf{f}_2} = \frac{\partial \mathbf{h}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} \frac{\partial l_i}{\partial \mathbf{f}_3}$$

Derivatives of vectors become matrices

3x3 3x2 2x1

and scalar multiplication turns into matrix multiplication

but there are still simplifications

$$\frac{\partial l_i}{\partial \mathbf{f}_2} = \frac{\partial \mathbf{h}_3}{\partial \mathbf{f}_2} \frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} \frac{\partial l_i}{\partial \mathbf{f}_3}.$$

$$\frac{\partial \mathbf{f}_3}{\partial \mathbf{h}_3} = \frac{\partial}{\partial \mathbf{h}_3} (\boldsymbol{\beta}_3 + \boldsymbol{\Omega}_3 \mathbf{h}_3) = \boldsymbol{\Omega}_3^T.$$

proof with matrix calculus

$$\frac{\partial \mathbf{h}_3}{\partial \mathbf{f}_2}$$

diagonal matrix \rightarrow replace by vector $\mathbb{I}(\mathbf{f}_2 > 0)$ and pointwise multiply
very compute efficient!

$$\frac{\partial l_i}{\partial \mathbf{f}_2} = \mathbb{I}(\mathbf{f}_2 > 0) \boldsymbol{\Omega}_3^T \frac{\partial l_i}{\partial \mathbf{f}_3}.$$

Correspondingly for all other f_i, h_i

Summary backpropagation

Neural network $f(\mathbf{x}_i, \phi)$: K hidden layers, activation function (e.g. ReLu)

Loss: $\ell_i = l[\mathbf{f}[\mathbf{x}_i, \phi], \mathbf{y}_i]$.

Forward pass: compute and store

$$\mathbf{f}_0 = \beta_0 + \mathbf{\Omega}_0 \mathbf{x}_i$$

$$\mathbf{h}_k = \mathbf{a}[\mathbf{f}_{k-1}] \quad k \in \{1, 2, \dots, K\}$$

$$\mathbf{f}_k = \beta_k + \mathbf{\Omega}_k \mathbf{h}_k. \quad k \in \{1, 2, \dots, K\}$$

Backward pass: compute and store

$$\frac{\partial \ell_i}{\partial \beta_k} = \frac{\partial \ell_i}{\partial \mathbf{f}_k} \quad k \in \{K, K-1, \dots, 1\}$$

$$\frac{\partial \ell_i}{\partial \mathbf{\Omega}_k} = \frac{\partial \ell_i}{\partial \mathbf{f}_k} \mathbf{h}_k^T \quad k \in \{K, K-1, \dots, 1\}$$

$$\frac{\partial \ell_i}{\partial \mathbf{f}_{k-1}} = \mathbb{I}[\mathbf{f}_{k-1} > 0] \odot \left(\mathbf{\Omega}_k^T \frac{\partial \ell_i}{\partial \mathbf{f}_k} \right), \quad k \in \{K, K-1, \dots, 1\}$$

Finally first layer

$$\frac{\partial \ell_i}{\partial \beta_0} = \frac{\partial \ell_i}{\partial \mathbf{f}_0}$$

$$\frac{\partial \ell_i}{\partial \mathbf{\Omega}_0} = \frac{\partial \ell_i}{\partial \mathbf{f}_0} \mathbf{x}_i^T$$

Comments

- Backpropagation is super compute efficient but not memory efficient
 - ➔ All intermediate values of the forward pass and all weight matrices are stored -> might limit size of model to be trained
- If memory allows, perform the forward and backward passes for the entire batch in parallel
 - Matrices and vectors get another index for the data event, i.e, become a multi-dimensional tensor
 - “tensor” = generalization of matrix to arbitrary dimension (vector = 1d tensor, matrix=2d tensor, etc.)
- All of this is implemented in TensorFlow and PyTorch
 - **After** this school you will just use a single line, like `loss.backward()`

Vanishing and exploding gradients

Due to the matrix multiplication

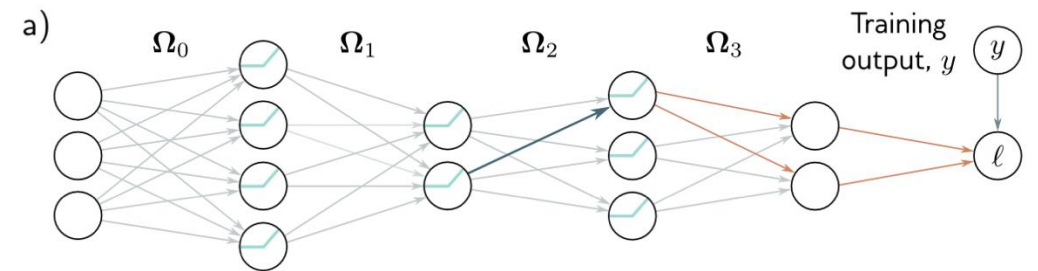
- very small weights may get to 0 and effectively remove nodes from the network “vanishing gradient”
- Large weights might get exponentiated to values beyond the precision of floating point arithmetic

Initialisation

To avoid vanishing or exploding gradients sample from normal distribution with mean 0 and variance

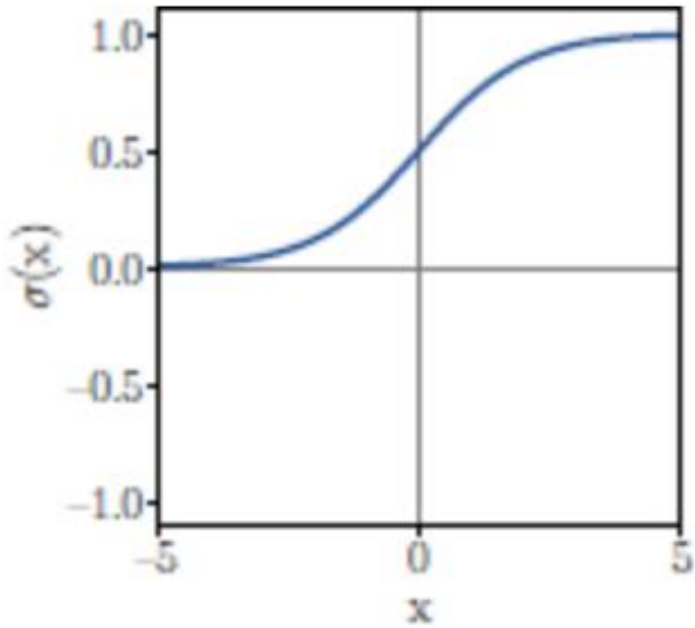
$$\sigma_{\Omega}^2 = \frac{4}{D_h + D_{h'}}$$

D_h : dimension of layer where weights are applied
 $D_{h'}$: dimension of layer to which they are fed

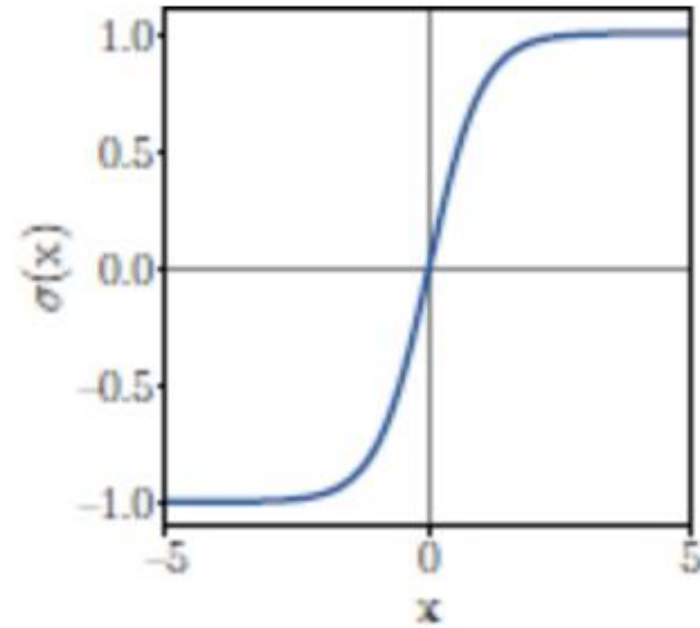


Complex network architectures may use separate simple networks to determine initialization for large complex networks

Gradients of other activation functions



Logistic sigmoid

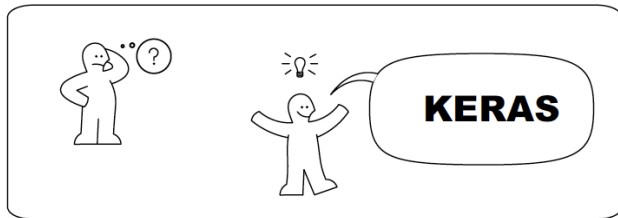


tanh

The full chain

ML frameworks like TensorFlow (with Keras API), PyTorch and JAX put a lot of the pieces together to provide a performant setup
-> See exercises with Peter

NEURALA NÄTVERK



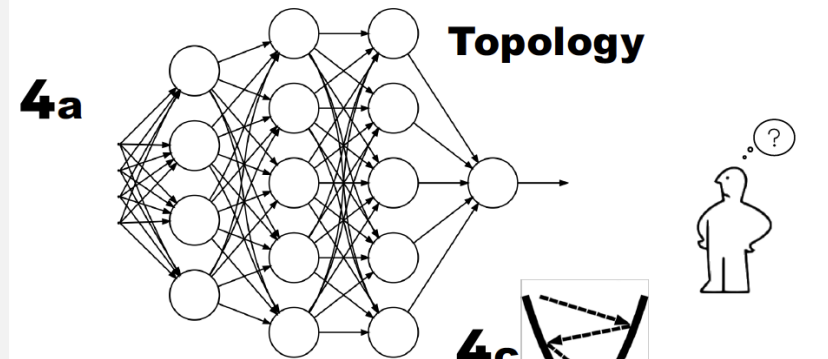
1  **Features**

2a **Train**  **Test**

2b  **Normalize**

3  **Metric**

Assembly Instruction



4b  **Loss**

5  **Fit**

6  **Overtraining**

7 **Apply!**

The network

```
def create_model(input_dim, output_dim):  
    return torch.nn.Sequential(  
        torch.nn.Linear(input_dim, 9),  
        torch.nn.ReLU(),  
        torch.nn.Linear(9, 9),  
        torch.nn.ReLU(),  
        torch.nn.Linear(9, 9),  
        torch.nn.ReLU(),  
        torch.nn.Linear(9, output_dim),  
    )
```

....improving step by step

Learn by revisiting the data often and adjusting

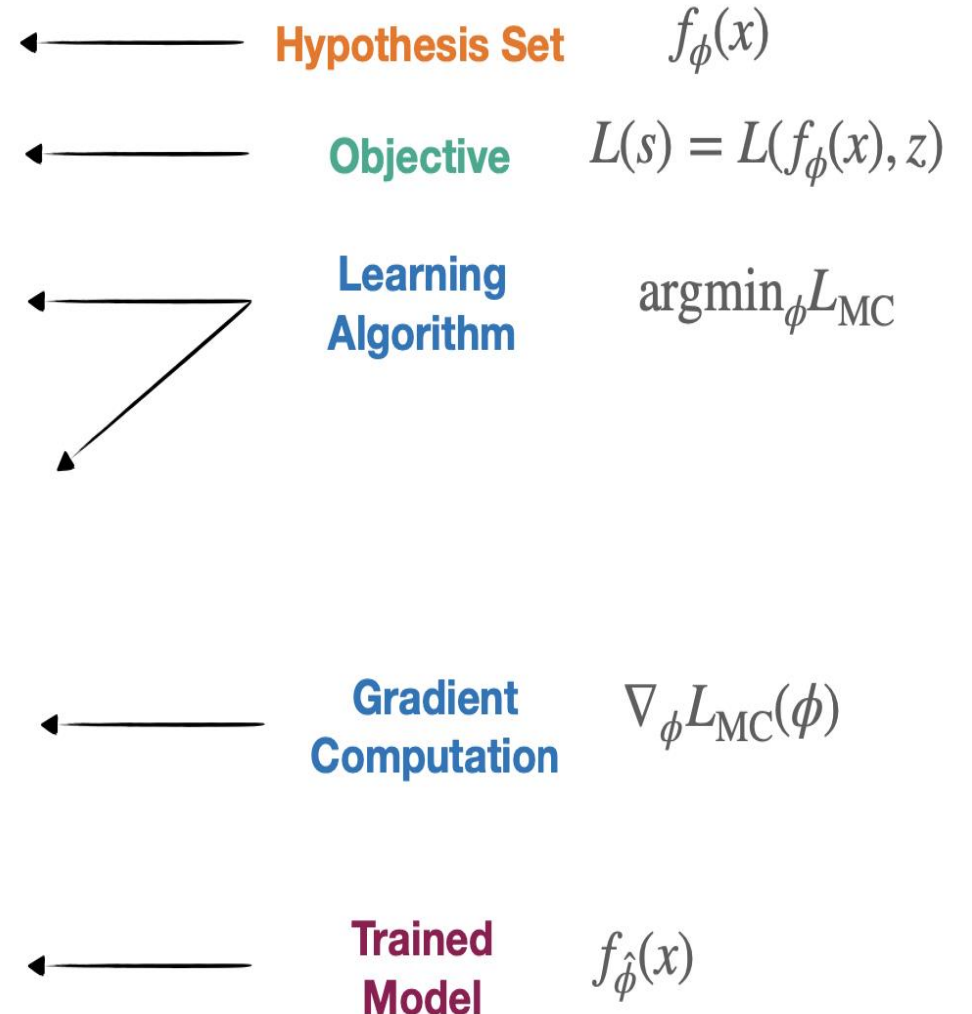
```
f = initial_guess()
for n in range(steps):
    examples ~ p(data)
    loss = evaluate(f, examples)
    adjustment = react(loss, f)
    f = new_hypo(f, adjustment)
```


A full training loop

↓ Data $s \sim p(s)$

```
def learn(samples):
    features, labels = samples
    model = MyModel()
    loss_func = torch.nn.MSELoss()
    opt = torch.optim.Adam(
        model.parameters(), lr = 1e-3
    )

    for i in range(steps):
        predictions = model(samples)
        loss = loss_func(predictions, labels)
        loss.backward()
        opt.step()
        opt.zero_grad()
    return model
```



Back-up