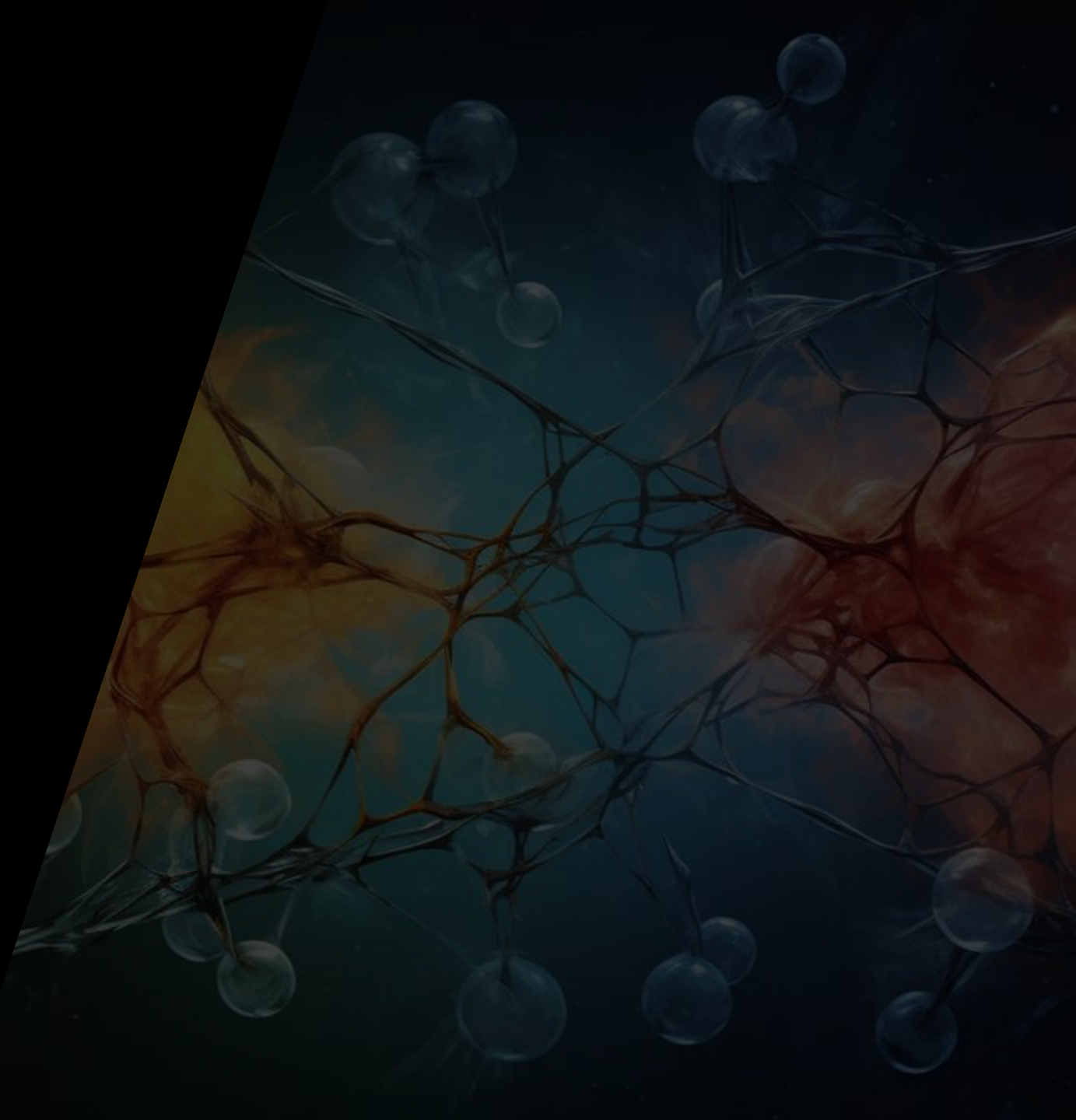


V - Optimisation with ERM

Contents:

- *Empirical risk minimization for parameter learning*
- *Gradient descent*
- *Stochastic gradient descent*
- *Backpropagation for computing gradients in neural networks*

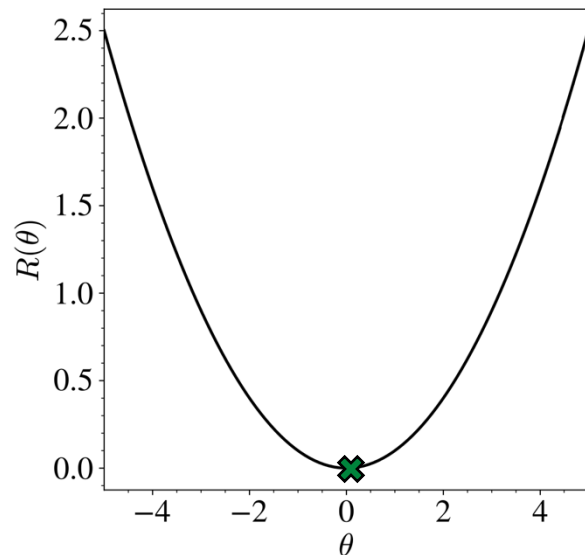


- The aim of training in supervised learning is to compute

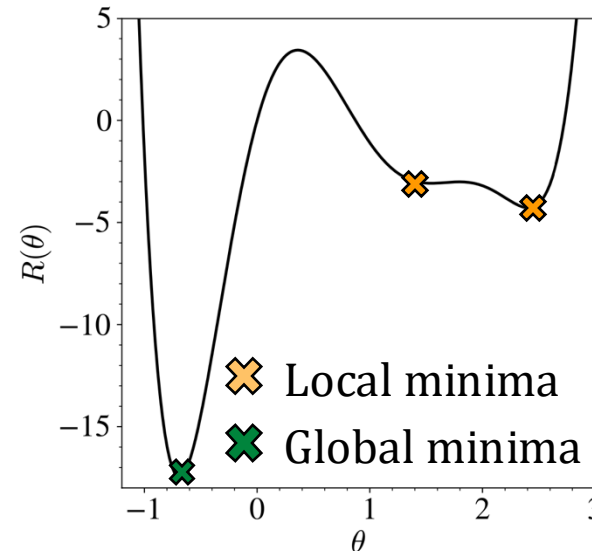
$$\hat{\theta} = \operatorname{argmin}_{\theta} R(\theta) = \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(\phi^{(i)}), y^{(i)}).$$

- In linear regression, we could directly optimise the empirical risk in closed-form, but **in general it is not possible**.
- How to perform the **Empirical Risk Minimisation** (ERM) in general?

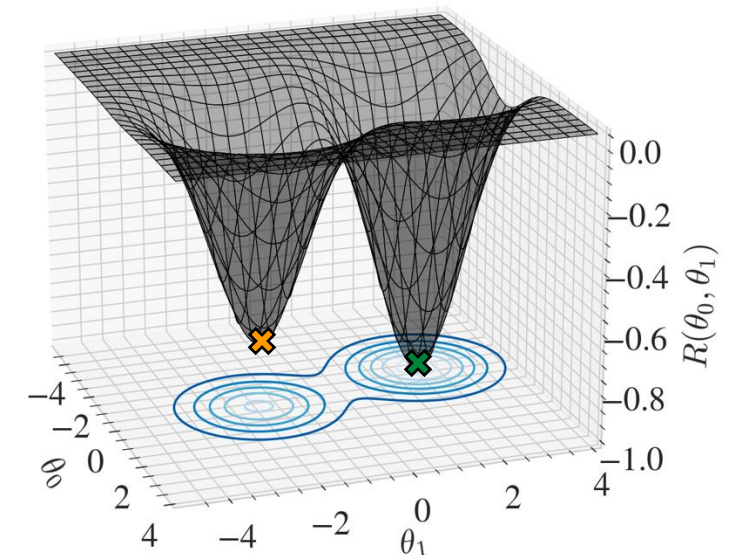
1D convex



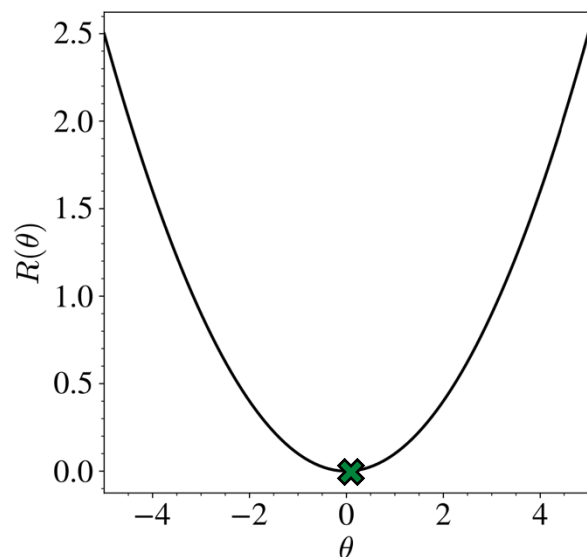
1D non-convex



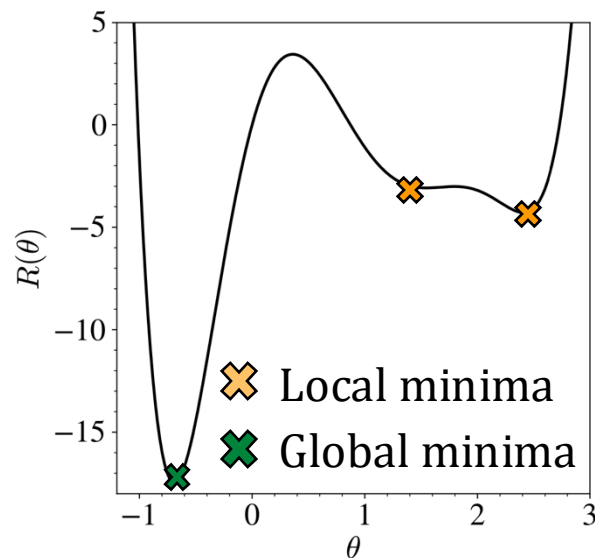
2D non-convex



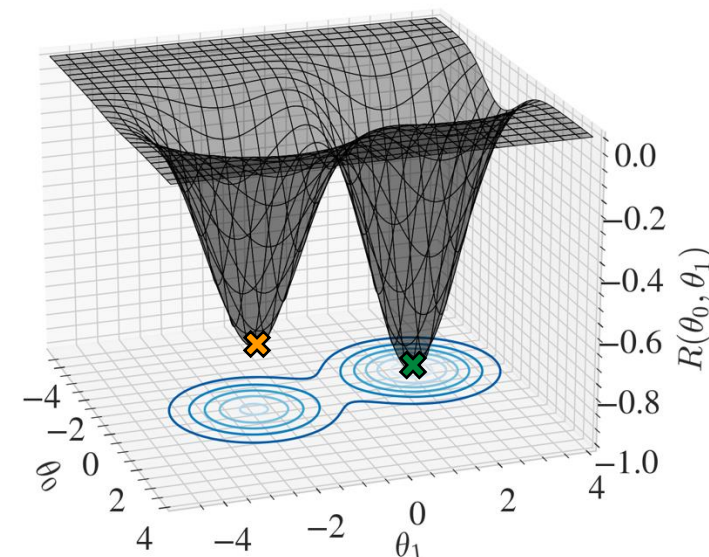
1D convex



1D non-convex



2D non-convex



- A naïve way of minimizing such functions could be to uniformly pave the parameter space and choose the one with the smallest value as being the minimum: this is **global optimization (grid search optimization)**
- This is where the **curse of dimensionality** kicks in (again). In general, we optimize models over many parameters $p \gg 1$ (imagine the pixels of an image) and all the points are far away from each other in high dimensions
- Sampling uniformly $[0, 1]^{10}$ with a step of 0.01 requires 10^{20} evaluations (think of GPT-3 and its 175 billion parameters!)

- A solution: local, directed search to navigate in the landscape

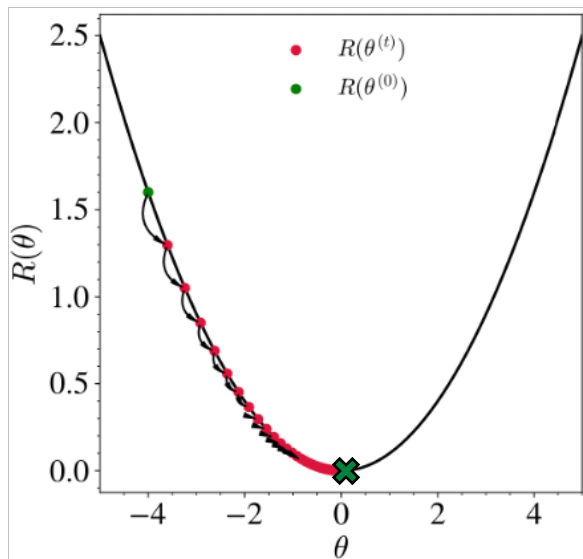
↪ Numerical optimisation by **gradient descent**

Note: the superscript does not have to do with the training example here but with the time step (θ is a parameter).

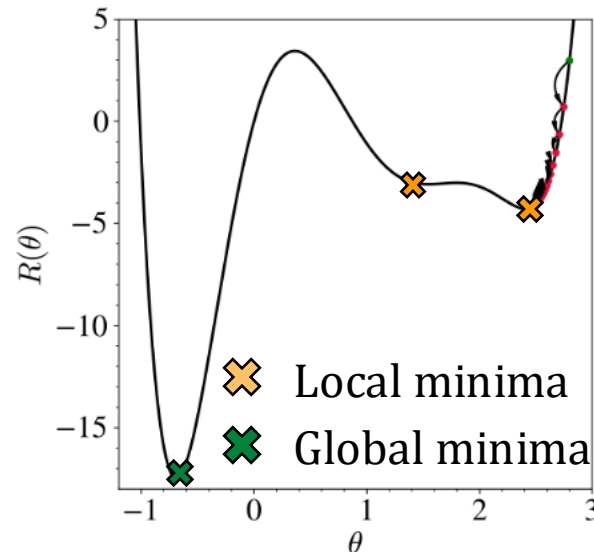
Algorithm: Gradient descent

1. Initialise θ_0 randomly
2. Compute $\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} R(\theta^{(t)})$
3. Repeat 2 until a stopping criterion is satisfied

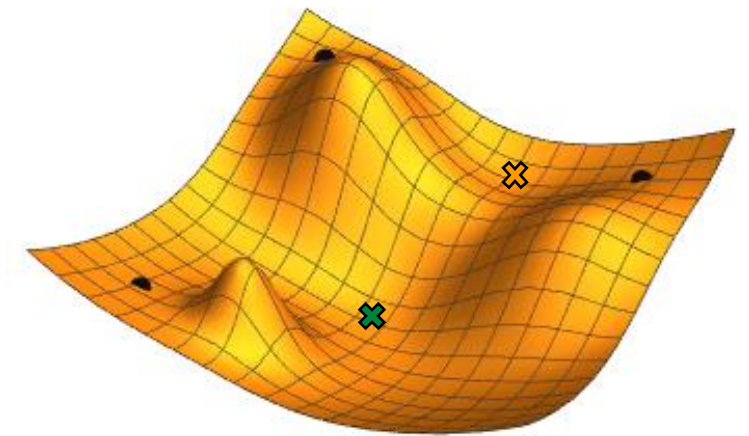
1D convex



1D non-convex



2D non-convex



Solutions in complex non-convex models depend on θ_0 : it is an **hyperparameter**

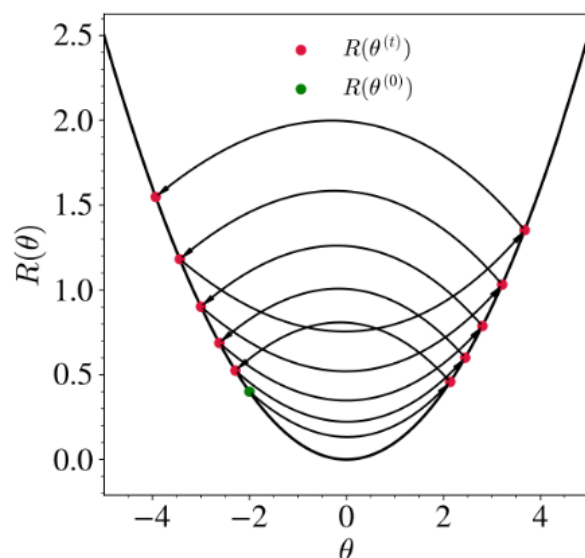
- Example of stopping criterion: $\|\theta^{(t+1)} - \theta^{(t)}\|_2^2 \leq \epsilon$, or fixed number of steps

- One **hyperparameter**: the learning rate η
- A value that is too large can lead to divergence while, when too small, the computational cost explodes (+ stuck in small asperities)

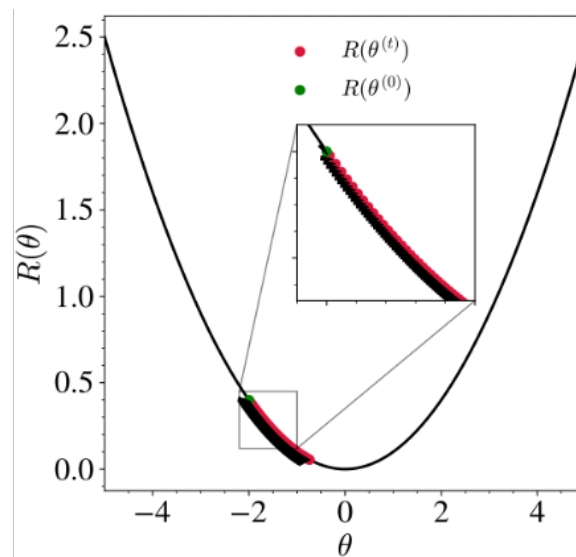


Hyperparameter: parameter that is **not learned** during the optimisation.
Ex: depth of tree in DTs, # of trees in RFs, learning rates, etc.

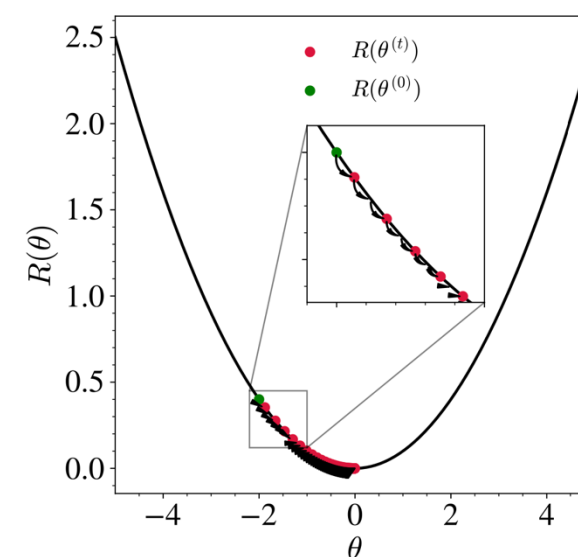
η too large



η too small



appropriate η



- Usually, we use **grid search** to find the hyperparameter performing best on a third dataset: the **validation set (or use cross-validation)**

- In linear regression, the empirical risk is

$$R(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left(\theta_0 + \theta_1 x_1^{(i)} - y^{(i)} \right)^2 \quad \boldsymbol{\theta} = [\theta_0, \theta_1]^T$$

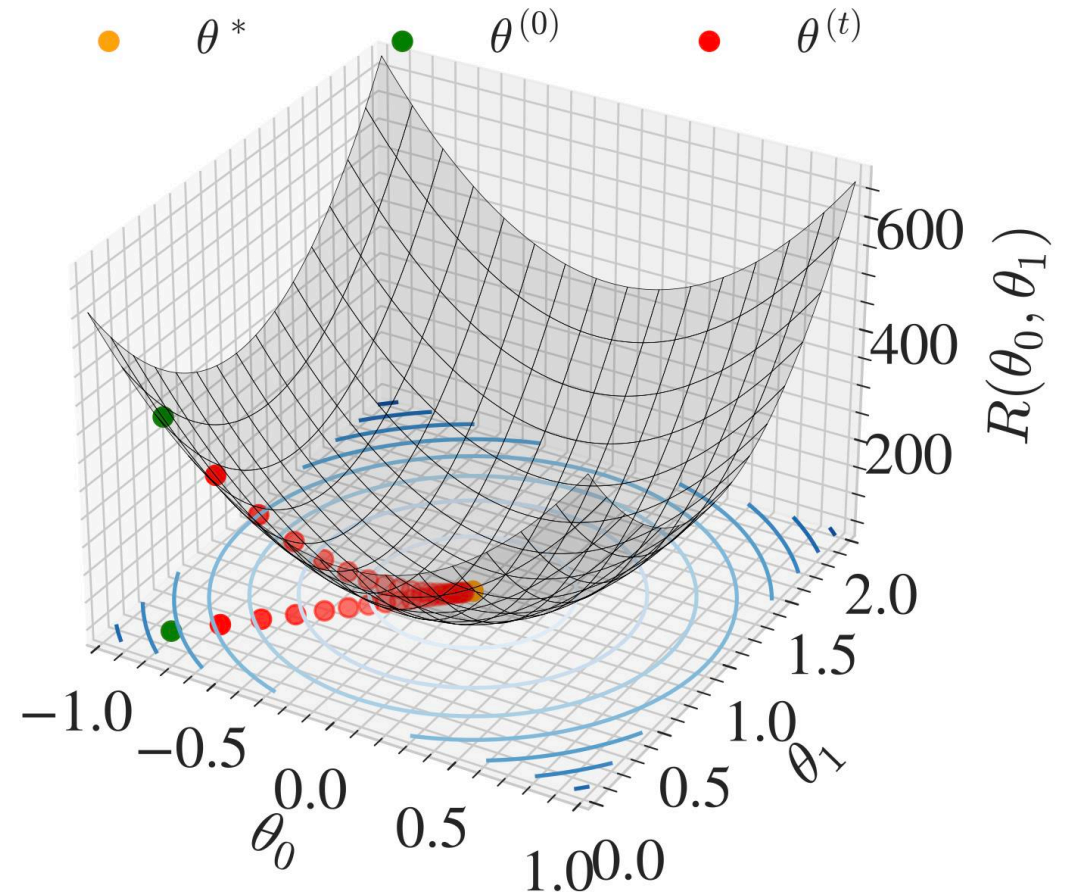
- Let's apply the gradient descent algorithm starting from $\boldsymbol{\theta}^{(0)}$ random
- Then, we need to compute the gradient $\nabla_{\boldsymbol{\theta}} R(\boldsymbol{\theta})$

$$\frac{\partial R(\boldsymbol{\theta})}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n \left(\theta_0 + \theta_1 x_1^{(i)} - y^{(i)} \right),$$
$$\frac{\partial R(\boldsymbol{\theta})}{\partial \theta_1} = \frac{2}{n} \sum_{i=1}^n x_1^{(i)} \left(\theta_0 + \theta_1 x_1^{(i)} - y^{(i)} \right)$$

- Therefore, the update is

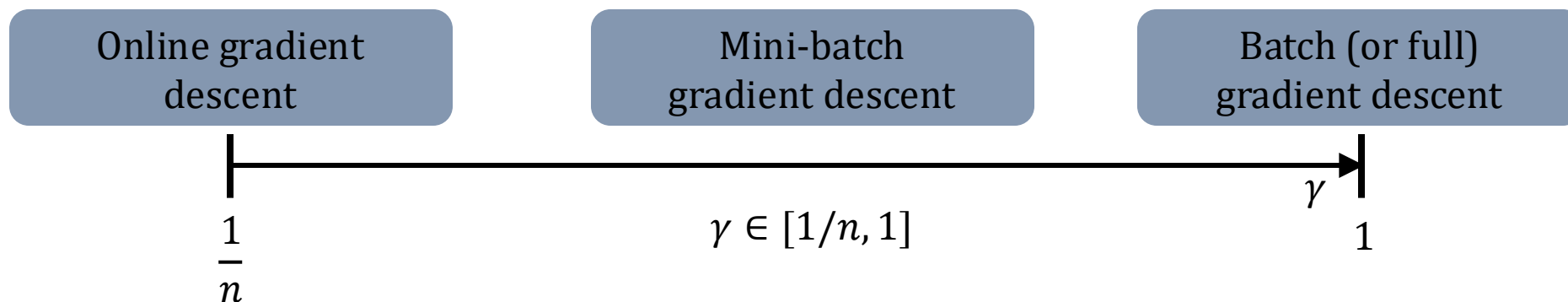
$$\theta_j^{(t+1)} = \theta_j^{(t)} - \frac{2\eta}{n} \sum_{i=1}^n x_j^{(i)} (\theta_0 + \theta_1 x_1^{(i)} - y^{(i)})$$

where $\forall i \in [1, \dots, n] \ x_0^{(i)} = 1$



Note that I **standardized** the features. This is sometimes required when optimising some models. In GD, it allows faster convergence.

- Problem of gradient descent: we **need the entire dataset** to compute $\nabla_{\theta} R(\theta^{(t)})$
- Solution: what about using only a fraction γ of the dataset chosen randomly?



Algorithm: Stochastic gradient descent

1. Initialise θ_0 randomly
2. For $e \in [1, \dots, E]$
 - 2.1. Shuffle the dataset
 - 2.2. For $i \in [1, \dots, \lfloor \gamma n \rfloor]$
 - 2.2.1. Compute $\theta^{(i+1)} = \theta^{(i)} - \eta \widehat{\nabla}_{\theta} R(\theta^{(i)})$

$$\text{where } \widehat{\nabla} R(\theta) = \sum_{j=i\gamma n}^{i\gamma n + \gamma n} \nabla R_j(\theta)$$

- e is called an **epoch** and a set of γn training examples is called a **mini-batch**
- Usually, **SGD often converges faster than full-batch** GD
- It may however oscillate around a true minimum
- Under some technical assumptions, **SGD provides an almost sure convergence** to a local (resp. global) in non-convex (resp. convex) landscapes

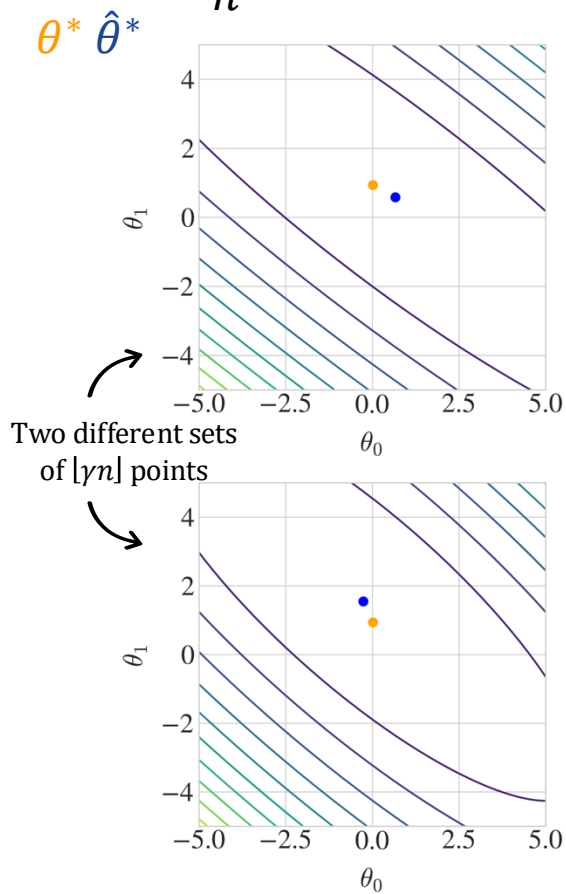
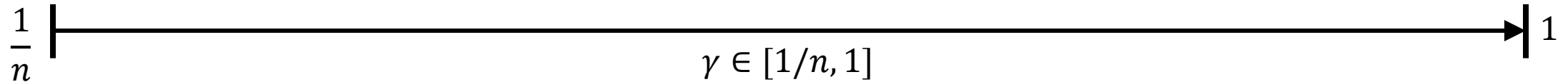
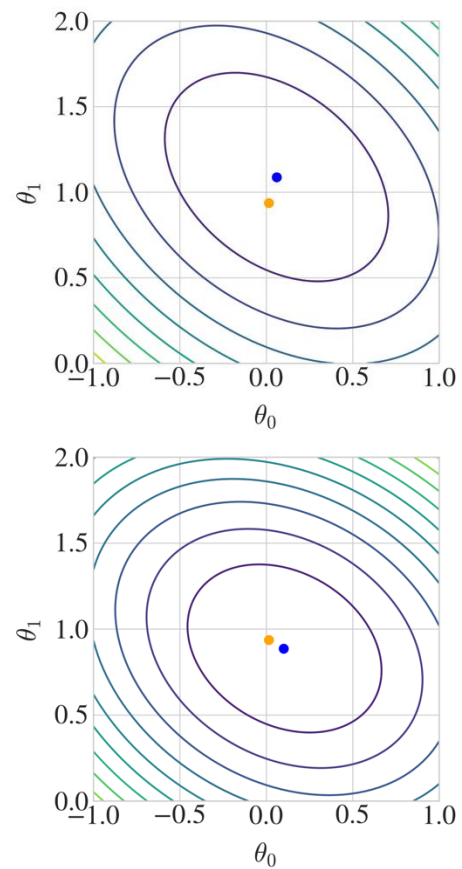
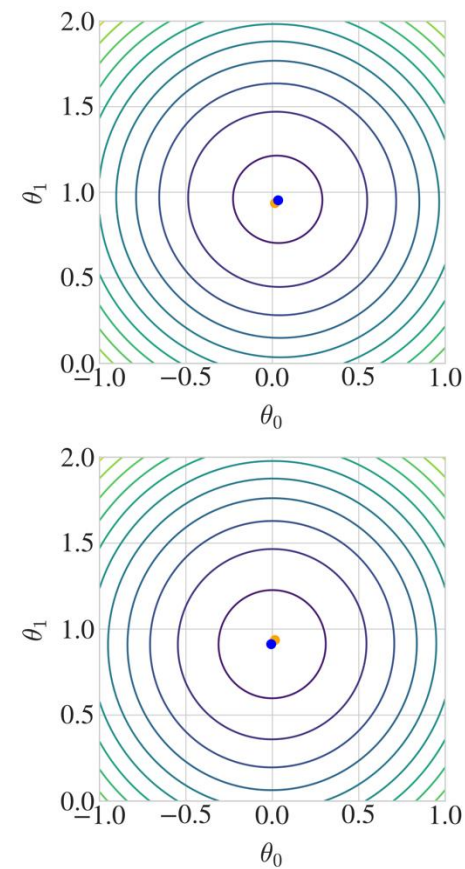
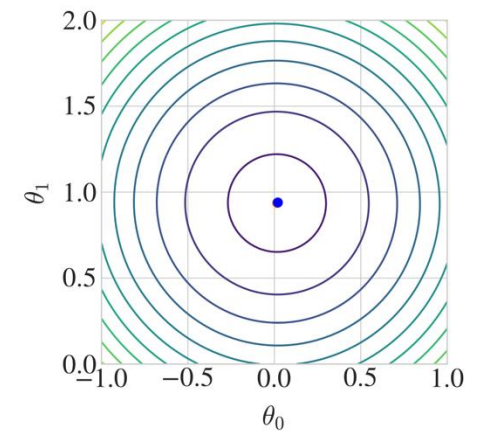
Linear models

Principles of learning

Trees and ensembling

Neural networks

Risk optimization

Online gradient
descentMini-batch
gradient descentBatch (or full)
gradient descent $\gamma = 2/n$  $\gamma = 0.1$  $\gamma = 0.5$  $\gamma = 1$

SGD on our linear regression

85

Linear models

Principles of learning

Trees and ensembling

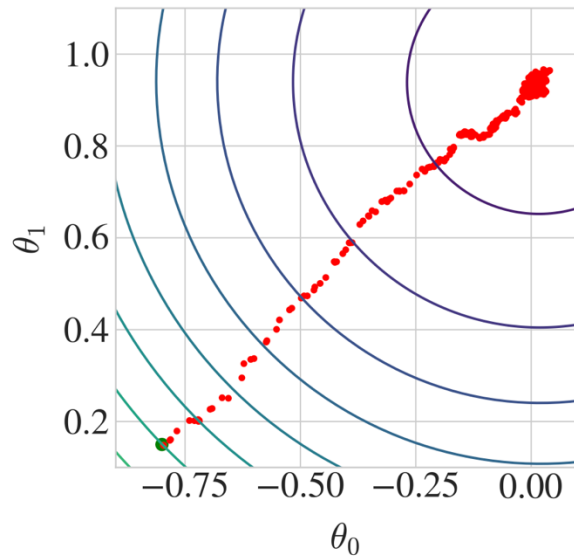
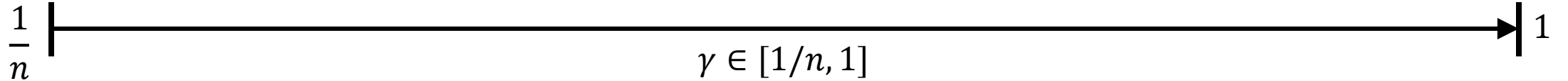
Neural networks

Risk optimization

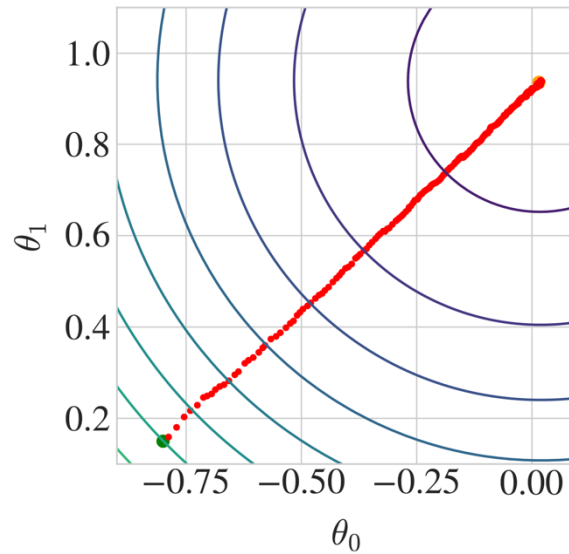
Online gradient
descent

Mini-batch
gradient descent

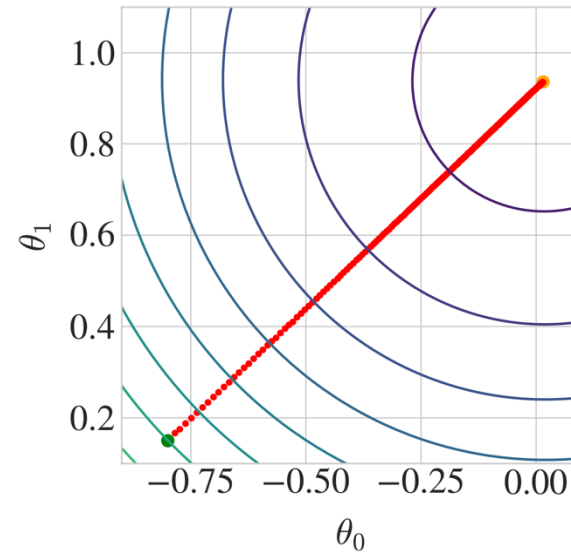
Batch (or full)
gradient descent



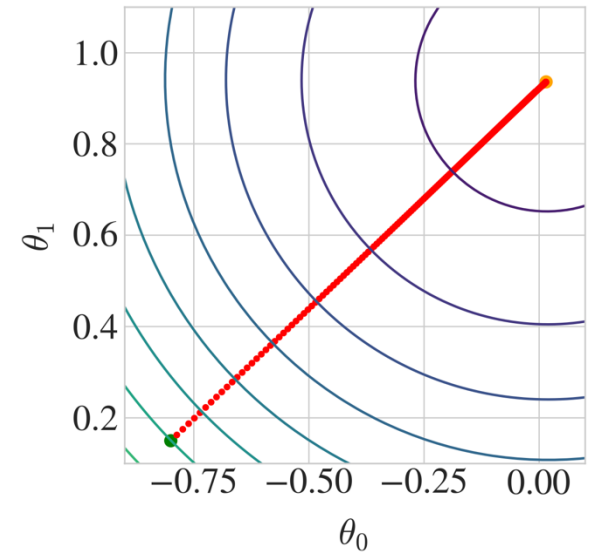
$\gamma = 2/n$



$\gamma = 0.1$



$\gamma = 0.5$



$\gamma = 1$

- The problem with neural networks is they are compositions of non-linear functions

$$y_j = s^{[L+1]} \left(\mathbf{W}^{[L+1]} s^{[L]} \left(\mathbf{W}^{[L]} \dots s^{[1]} \left(\mathbf{W}^{[1]} \mathbf{x} \right) \right) \right)$$

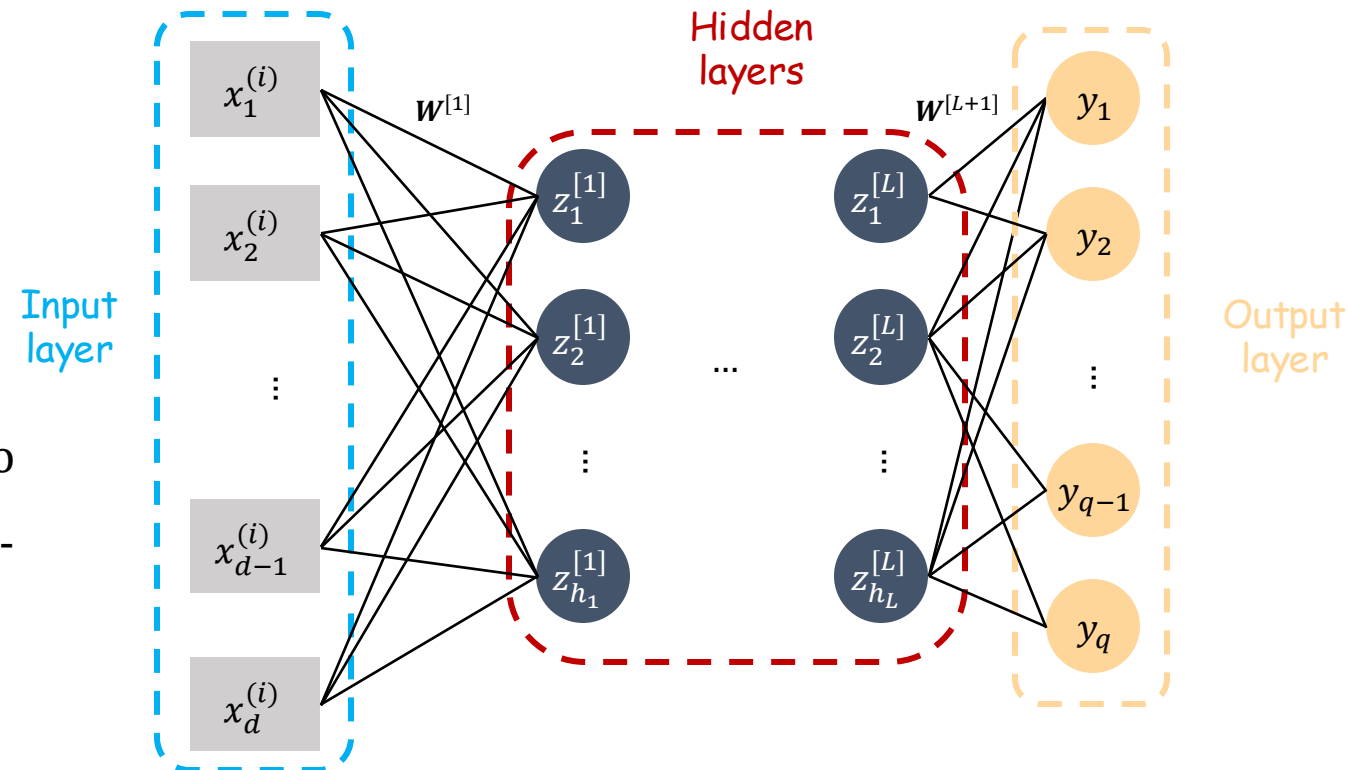
$$\mathbf{W}^{[l]} = [\mathbf{w}_1^{[l]}, \mathbf{w}_2^{[l]}, \dots, \mathbf{w}_{h_l}^{[l]}]^T \in \mathbb{R}^{h_l \times (h_{l-1} + 1)}$$

$$\mathbf{w}_j^{[l]} = [b_j^{[l]}, w_{1j}^{[l]}, w_{2j}^{[l]}, \dots, w_{h_{l-1}j}^{[l]}]$$

- We however need to optimize the cost function to obtain the “best” values of \mathbf{W} producing the closer-to-optimal target values

- How to compute $\frac{\partial \ell(\mathbf{W}^{[0]}, \dots, \mathbf{W}^{[L+1]})}{\partial w_{ij}^{[l]}}$?

- The backpropagation of errors:** an application of the **chain rule**!



- Take the example of a fully-connected network with $d = 3$, one hidden layer, and two output neurons:

$$\hat{\mathbf{y}} = s\left(\mathbf{W}^{[2]}s\left(\mathbf{W}^{[1]}\mathbf{x}\right)\right) \quad \text{with } \mathbf{x} = [1, x_1, x_2, \dots, x_d]^T$$

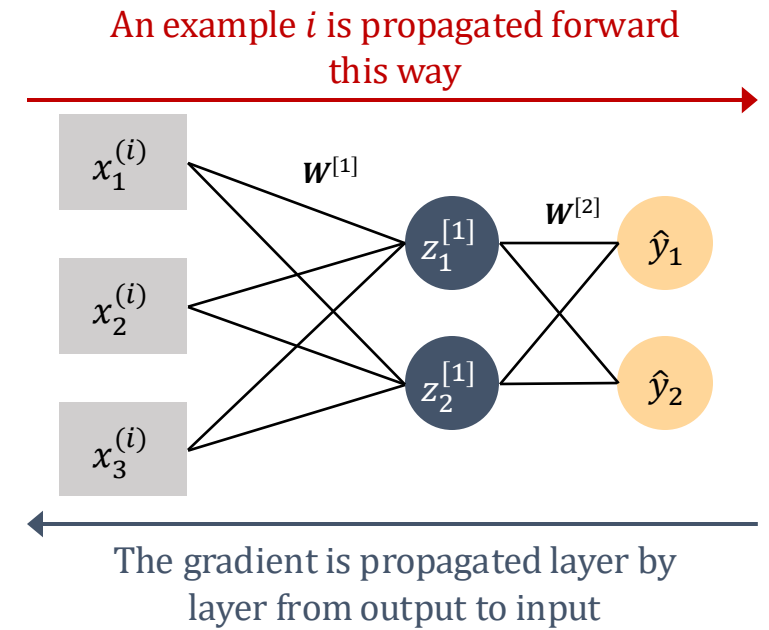
$$\mathbf{W}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix} \quad \mathbf{W}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

- Consider the squared error loss function for each training example i

$$\ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

- From the equation of $\hat{\mathbf{y}}$, we see the contribution of $\mathbf{W}^{[2]}$ is “closer” to the output than $\mathbf{W}^{[1]}$

- Let's compute $\frac{\partial \ell(\mathbf{W}^{[1]}, \mathbf{W}^{[2]})}{\partial w_{11}^{[2]}}$



- Using the chain rule

$$\frac{\partial \ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]})}{\partial w_{11}^{[2]}} = \frac{\partial \ell}{\partial u_1^{[2]}} \times \frac{\partial u_1^{[2]}}{\partial w_{11}^{[2]}}$$

where $u_1^{[2]}$ is the pre-activation of the unit 1 in layer 2 (here output layer)

$$u_1^{[2]} = w_{11}^{[2]} z_1^{[1]} + w_{21}^{[2]} z_2^{[1]} + b_1^{[2]}$$

- The second term is then easy to compute as $\frac{\partial u_1^{[2]}}{\partial w_{11}^{[2]}} = z_1^{[1]}$
- For the first term, use the chain rule again leads to

$$\delta_1^{[2]} = \frac{\partial \ell}{\partial u_1^{[2]}} = \frac{\partial \ell}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial u_1^{[2]}}$$

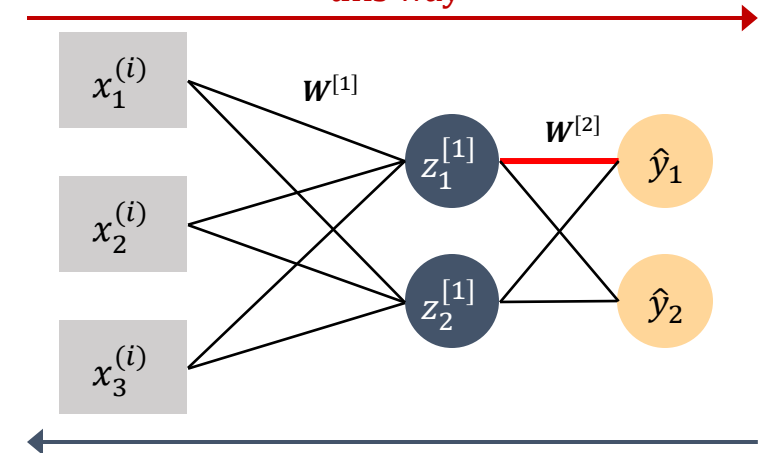
with $\hat{y}_1 = s(u_1^{[2]}) = s(w_{11}^{[2]} z_1^{[1]} + w_{21}^{[2]} z_2^{[1]} + b_1^{[2]})$

Finally,

$$\frac{\partial \ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]})}{\partial w_{11}^{[2]}} = -(y_1 - \hat{y}_1) s'(u_1^{[2]}) z_1^{[1]}$$

$$\ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{w}^{[2]} s(\mathbf{w}^{[1]} \mathbf{x}))$$

$$\mathbf{w}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

$$\mathbf{w}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

- You can proceed the same for all the weights linked to the output layer
- What about parameters in the hidden layers?
- Let's compute

$$\frac{\partial \ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]})}{\partial w_{11}^{[1]}} = \frac{\partial \ell}{\partial u_1^{[1]}} \times \frac{\partial u_1^{[1]}}{\partial w_{11}^{[1]}}$$

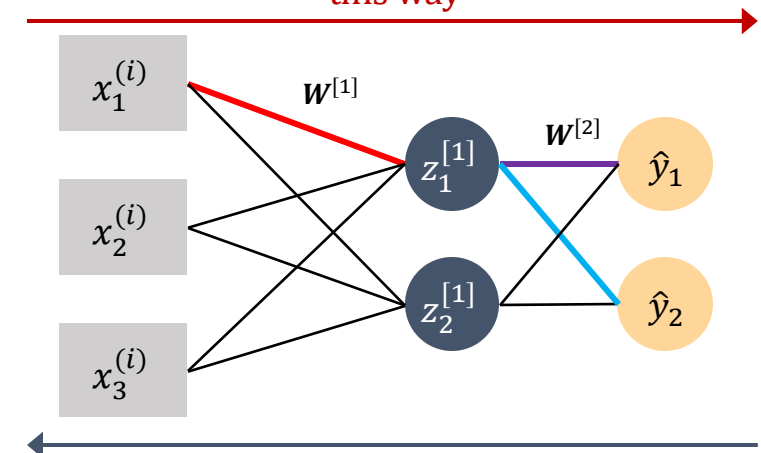
- The second term is still easy to compute
- For the first term, we have $u_1^{[1]} = w_{11}^{[1]}x_1 + w_{21}^{[1]}x_2 + b_1^{[1]}$ and observe that there are two paths to reach the weight $w_{11}^{[1]}$:

$$\delta_1^{[1]} = \frac{\partial \ell}{\partial u_1^{[2]}} \frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} + \frac{\partial \ell}{\partial u_2^{[2]}} \frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$$

Already computed in the previous layer (hence the name "Backpropagation")

$$\ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{w}^{[2]} s(\mathbf{w}^{[1]} \mathbf{x}))$$

$$\mathbf{w}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

$$\mathbf{w}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$

$$\delta_1^{[1]} = \frac{\partial \ell}{\partial u_1^{[2]}} \frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} + \frac{\partial \ell}{\partial u_2^{[2]}} \frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$$

Already computed in the previous layer (hence the name "Backpropagation")

- To compute $\frac{\partial u_1^{[2]}}{\partial u_1^{[1]}}$ and $\frac{\partial u_2^{[2]}}{\partial u_1^{[1]}}$, we can use the same chain rule again giving

$$\frac{\partial u_1^{[2]}}{\partial u_1^{[1]}} = \frac{\partial u_1^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial u_1^{[1]}} = w_{11}^{[2]} s'(u_1^{[1]})$$

$$\frac{\partial u_2^{[2]}}{\partial u_1^{[1]}} = \frac{\partial u_2^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial u_1^{[1]}} = w_{21}^{[2]} s'(u_1^{[1]})$$

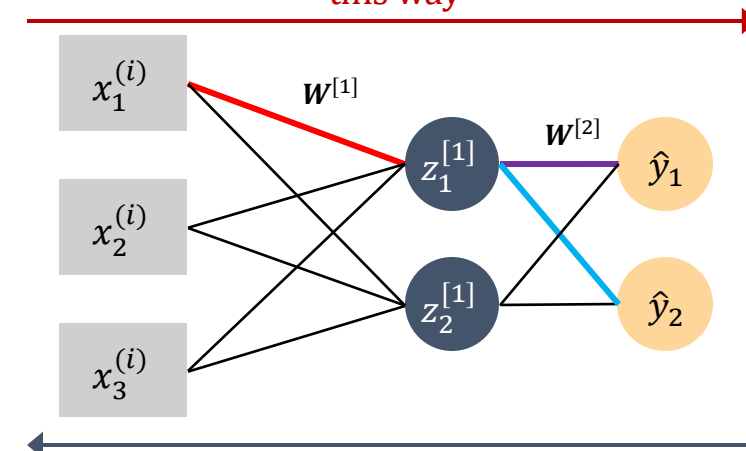
- Hence,

$$\frac{\partial \ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]})}{\partial w_{11}^{[1]}} = \left(\delta_1^{[2]} w_{11}^{[2]} + \delta_2^{[2]} w_{21}^{[2]} \right) s'(u_1^{[1]}) x_1$$

- Where everything is known!

$$\ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{w}^{[2]} s(\mathbf{w}^{[1]} \mathbf{x}))$$

$$\mathbf{w}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

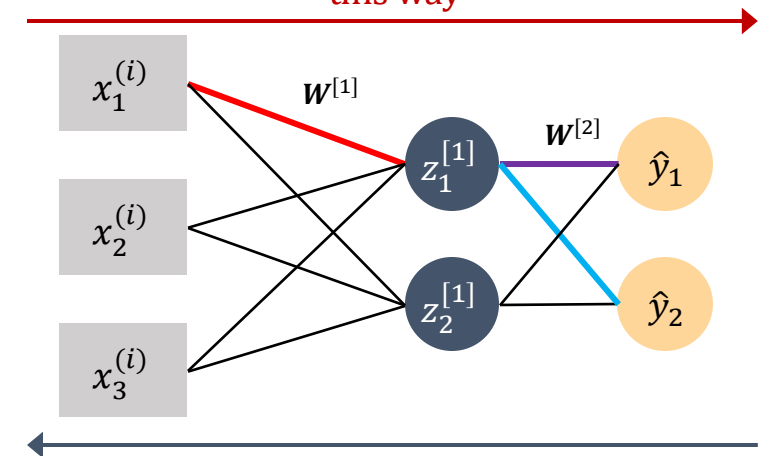
$$\mathbf{w}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$



1. Writing those equations **propagating the errors from output to input recursively** for both weights and biases leads to the **backpropagation algorithm**.
2. Together with the SGD algorithm they allow to train efficiently networks with possibly many parameters.
3. Modern programming tools implement an automatic version of this algorithm called **automatic differentiation** generating the code for backward equations from forward propagation declaration.

$$\ell(\mathbf{w}^{[1]}, \mathbf{w}^{[2]}) = \frac{1}{2} \sum_{j=1}^2 (y_j - \hat{y}_j)^2$$

An example i is propagated forward this way



The gradient is propagated layer by layer from output to input

$$\hat{\mathbf{y}} = s(\mathbf{w}^{[2]} s(\mathbf{w}^{[1]} \mathbf{x}))$$

$$\mathbf{w}^{[1]} = \begin{bmatrix} b_1^{[1]} & w_{11}^{[1]} & w_{12}^{[1]} & w_{13}^{[1]} \\ b_2^{[1]} & w_{21}^{[1]} & w_{22}^{[1]} & w_{23}^{[1]} \end{bmatrix}$$

$$\mathbf{w}^{[2]} = \begin{bmatrix} b_1^{[2]} & w_{11}^{[2]} & w_{12}^{[2]} \\ b_2^{[2]} & w_{21}^{[2]} & w_{22}^{[2]} \end{bmatrix}$$