# Lecture 9: Backpropagation

## Optimization for data sciences

### Rémy Sun

*remy.sun@inria.fr*

# Course organization
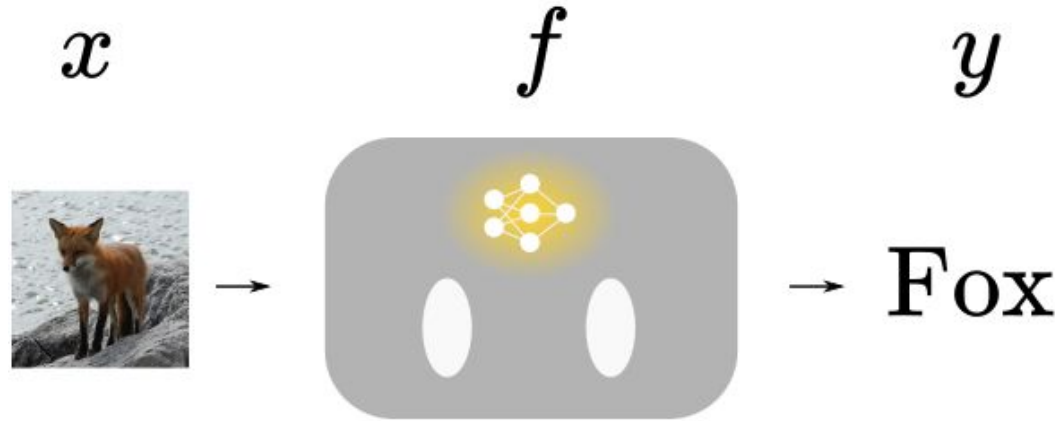
# Course organization

- Introduction to optimization
  - A few problems of interest
  - Quick mathematical refresher

- Convex problems (Following Stephen Boyd)
  - Convex sets
  - Convex functions
  - Convex problems
  - Simplex algorithm for Linear Programming

- Duality (for convex problems)
  - Lagrangian  and dual function
  - Dual problem
  - Qualification constraints
  - KKT conditions

- Newton's Descent and Barrier methods for convex case
  - Descent for the unconstrained problems
  - Equality constrained problems
  - Interior point methods
  - Lab session!

- What about the real (neural) world?
  - Problem statement
  - Let's try to solve it!
  - Gradient descent with(out) convexity
  - Gradient descent variants

- **Backpropagation**
  - **Chain rule derivation**
  - **Dynamic programming**
  - **Backpropagation**
  - **Lab session  2!**

- Reports on lab sessions
  - Labs on jupyter notebooks
    - Not every session
  - Explain the code done in the session
  - Summarize what is done in the practical

- Written Exam
  - Theoretical questions
  - We will do exercises in class
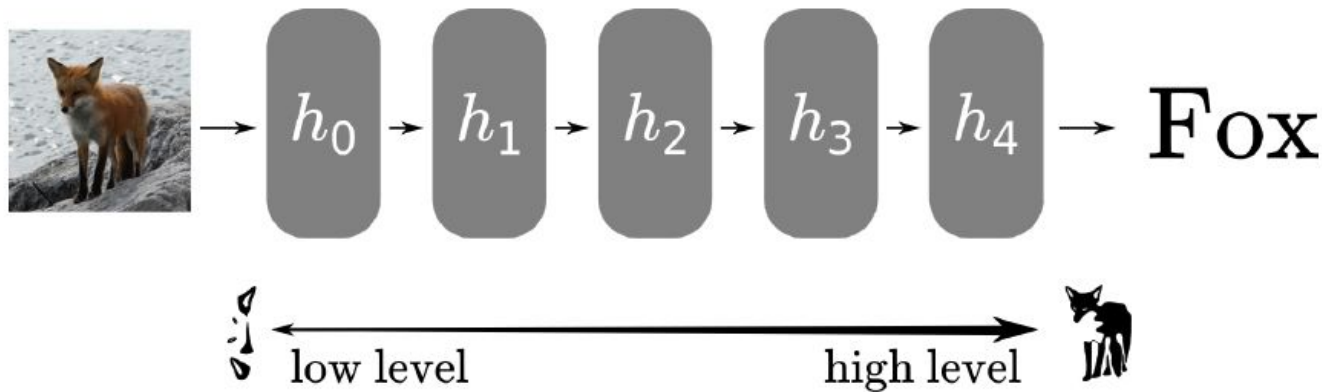
Refresher on convexity!

- Find (robot) f that classifies images well
  - Often based on neural networks
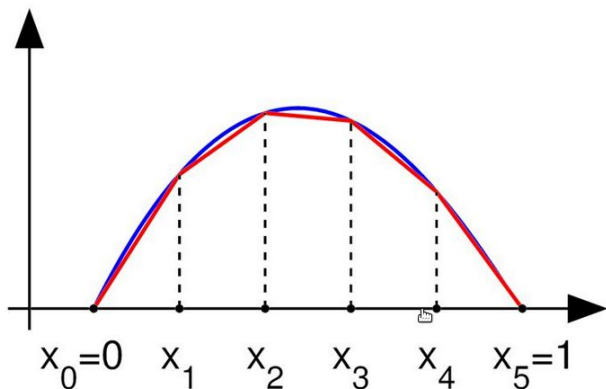
$$\forall(x, y) \in \mathcal{D}, f(x) = y$$

- Problem: we do not know $\mathcal{D}$ !
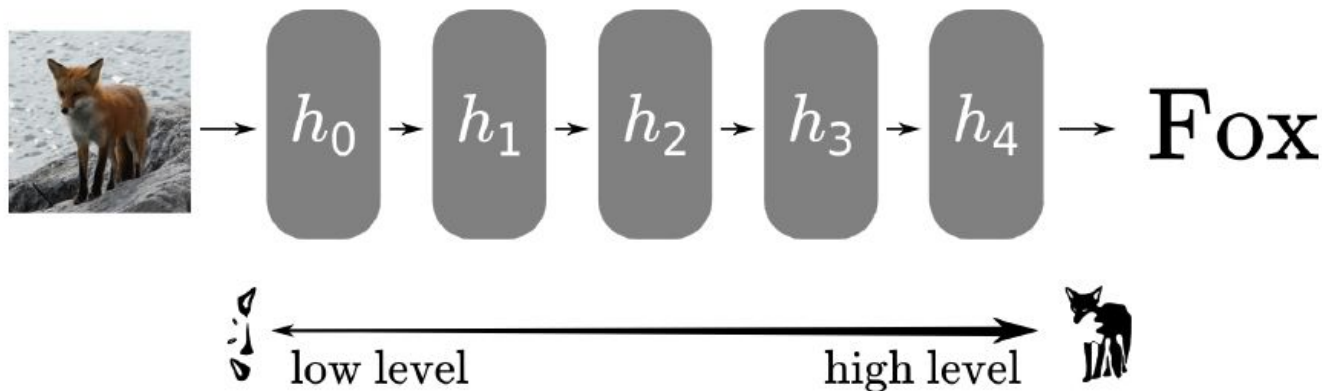  - Solved problem otherwise…
  - Evaluating the risk requires this distribution

- Solution: Use a dataset D of (x,y) sampled from $\mathcal{D}$
  - ***Empirical Risk Minimization***
  - If the (x,y) are i.i.d drawn from $\mathcal{D}$ can be expressed as a mean over the dataset

$$min_\theta \hat{\mathcal{R}}_\theta = \frac{1}{N} \sum_{i=0,\ldots,N-1} l(f_\theta(x_i), y_i)$$

low level         high level

- Neural networks are sequences of simple functions

$$f_\theta = h_\theta^0 \circ h_\theta^1 \circ \cdots \circ h_\theta^{L-1}$$

$h_0 \rightarrow h_1 \rightarrow h_2 \rightarrow h_3 \rightarrow h_4 \rightarrow$ Fox

low level      high level

$x_0 = 0 \quad x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 = 1$

- Highly expressive
  - Can fit many types of distributions

$$\begin{aligned}
\text{minimize} \quad & f_0(x) \\
\text{subject to} \quad & f_i(x) \leq 0, \quad i = 1, \ldots, m \\
& h_i(x) = 0, \quad i = 1, \ldots, p
\end{aligned}$$

- Cost function is given by neural network and loss
- No constraints!
  - Easy unconstrained problem!

- **Not convex, for a number of reasons.**

UNIVERSITÉ
CÔTE D'AZUR

- What do we need?
  - Step size
    - Fixed step size
  - Gradient
    - Could get pretty expensive too…
    - Obtained through backpropagation!
  - No Hessian!

- Under reasonable L-smooth assumption
  - Gradient descent converges!
    - To something
    - With a minuscule fixed step size

- General deep learning heuristics
  - Most local minimums are similarly good/bad
  - Big networks have very few really bad mins

# 1. Some context

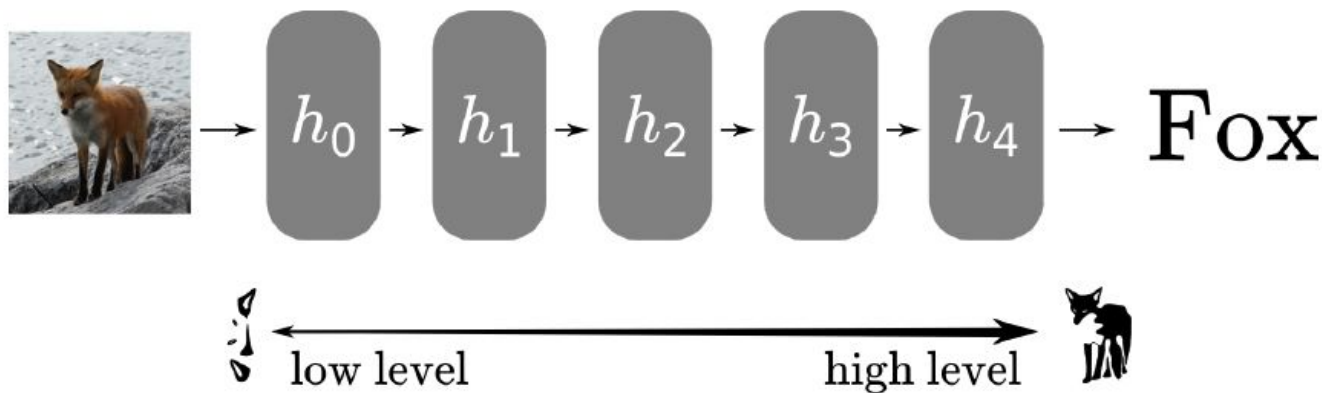$$\theta^{t+1} := \theta^t - \eta \nabla_\theta \hat{\mathcal{R}_\theta}(B)$$

- Requires finding the risk gradient wrt parameters

$$\nabla_\theta \hat{\mathcal{R}_\theta}(B) = \frac{1}{\#B} \sum_{k=0,\ldots,B-1} \nabla_\theta l(f_\theta(x_k), y_k)$$

- Boils down to computing gradients for one sample

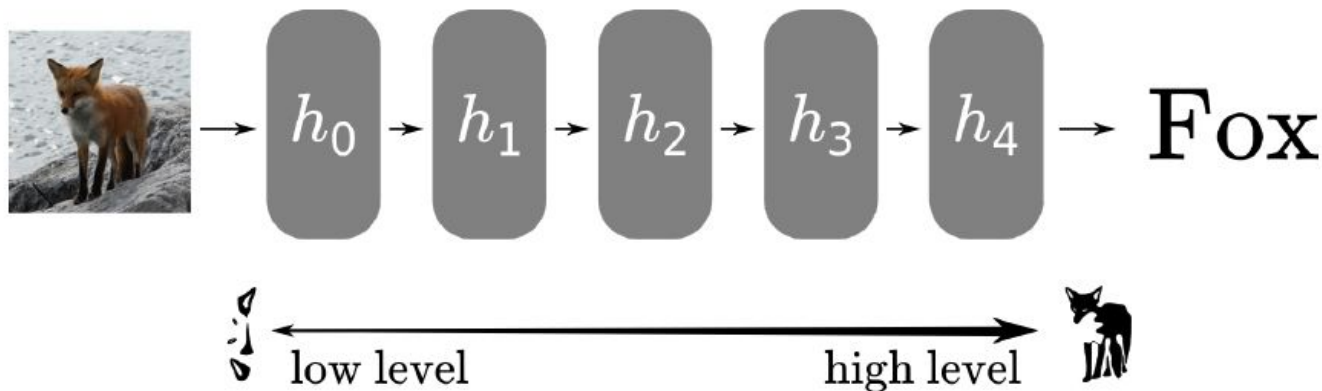$$\nabla_\theta l(f_\theta(x), y)$$

low level       high level

- Neural networks embed complex functions
  - Reason we use them

UNIVERSITÉ
CÔTE D'AZUR



- Neural networks embed complex functions
  - Reason we use them

- **How do we get the gradients?!**

# Gradients for neural networks

- Not an easy problem originally

- First neural networks do not use gradients
  - "Hard" value perceptron
  - Rosenblatt's algorithm

- 1970: Modern Backpropagation (reverse mode)
  - Computes gradients
  - Uses Chain rule derivation
  - With bottom-up dynamic programming

# 2. Chain rule

low level      high level

- The complete network is complicated

$$f_\theta = h_\theta^0 \circ h_\theta^1 \circ \cdots \circ h_\theta^{L-1}$$

- But individual derivatives are simple!

$$f_\theta = h_\theta^0 \circ h_\theta^1 \circ \cdots \circ h_\theta^{L-1}$$

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

- Derivatives can be broken down into simpler parts
  - As a product!

- Allows us to use the very simple functions
  - E.g. $\tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \, x_j + b_i^h$ has trivial partial derivatives

- Known for a very long time (17th century)

- Chain rule derivation powers backpropagation
  - Leverages simple component functions
  - Easy to break down

- Works as long as we know how to derivate layers

- Backpropagation is not just the chain rule
  - Exploding costs
  - Lots of redundant computations

- Chain rule derivation powers backpropagation
  - Leverages simple component functions
  - Easy to break down

- Works as long as we know how to derivate layers

- **Backpropagation is not just the chain rule**
  - Exploding costs
  - Lots of redundant computations
  - Backprop is **efficient** derivation

# 3. Dynamic Programming

# Dynamic programming

- Efficient computations of complex problem
    - Bellman (1950)
    - Break down into simpler (re-occuring) problems
    - And remember your solutions

- Avoid redundant computations
    - By using additional memory to store results
    - By properly structuring computations

● $F(0)=0$, $F(1)=1$, $F(N+2)=F(N)+F(N+1)$

Lots and lots and lots of computations…

● Memoization: Store sub-problem solutions as needed

- Tabulation: Start from basic problems
  - Tabulate results as we go along

- One single forward pass this way!

- Memoization vs. Tabulation depends on the problem

- Efficient computations of complex problem
  - Bellman (1950)

- Avoid redundant computations
  - By using additional memory to store results
  - By properly structuring computations

- Top-down (Memoization): Start from final problem

- Bottom-up (Tabulation): Start from basic problem

# 4. Backpropagation

- Modern Backpropagation: Linnainmaa (1970)
  - In his Master's thesis!
  - Some precursor efforts before
  - Theorized by Rosenblatt for perceptrons

- Efficient differentiation in a computational graph
  - Reverse mode autodiff
  - Not for neural network per se
  - Applied to neural networks later on

- Networks are complex but made of simple parts!
  - Simple gradients of component functions
  - Chain-rule allows decomposition into simple gradients

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial a}\frac{\partial a}{\partial w}$$

  - Remember the computed values

- Need to store intermediate activations "a" to evaluate partial derivatives

$$\frac{\partial a}{\partial w}$$

- Only one pass (in backward)!

Entrée
Couche cachée
Couche de sortie

- Simple 1 hidden layer MLP
  - 2 inputs
  - 2 outputs
  - 4 hidden activations

- Classification problem
  - Outputs probabilities
  - Cross-entropy loss

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \, x_j + b_i^h \\[2mm] h_i = \tanh(\tilde{h}_i) \\[2mm] \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \, h_j + b_i^y \\[2mm] \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



Couche cachée
Couche de sortie
Entrée

$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \, x_j + b_i^h \\[2mm] h_i = \tanh(\tilde{h}_i) \\[2mm] \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \, h_j + b_i^y \\[2mm] \hat{y}_i = \mathrm{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

UNIVERSITÉ
CÔTE D'AZUR

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



$$\frac{\partial \mathcal{L}}{\partial \tilde{h}} = \qquad \frac{\partial \mathcal{L}}{\partial h} = \qquad \frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \tilde{y}} \qquad \frac{\partial \mathcal{L}}{\partial \hat{y}}$$

$$\frac{\partial \mathcal{L}}{\partial W_h} = \qquad \frac{\partial \mathcal{L}}{\partial W_y} = \frac{\partial \mathcal{L}}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial W_y}$$

Couche cachée
Couche de sortie

Entrée



$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \; x_j + b_i^h \\[2mm] h_i = \tanh(\tilde{h}_i) \\[2mm] \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \; h_j + b_i^y \\[2mm] \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum\limits_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

40

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



Couche
cachée

Couche
de sortie

Entrée

$$\frac{\partial \mathcal{L}}{\partial \tilde{h}} = \qquad \frac{\partial \mathcal{L}}{\partial h} = \frac{\partial \mathcal{L}}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial h} \qquad \frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \tilde{y}} \qquad \frac{\partial \mathcal{L}}{\partial \hat{y}}$$

$$\frac{\partial \mathcal{L}}{\partial W_h} = \qquad \frac{\partial \mathcal{L}}{\partial W_y} = \frac{\partial \mathcal{L}}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial W_y}$$

$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \ x_j + b_i^h \\ h_i = \tanh(\tilde{h}_i) \\ \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \ h_j + b_i^y \\ \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

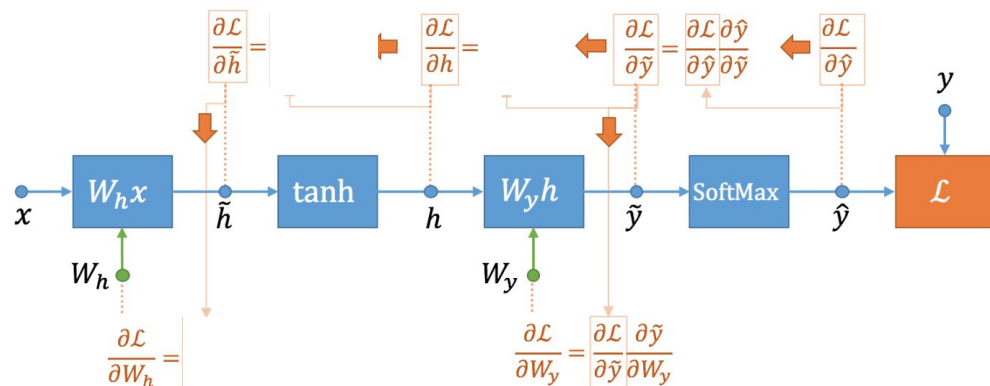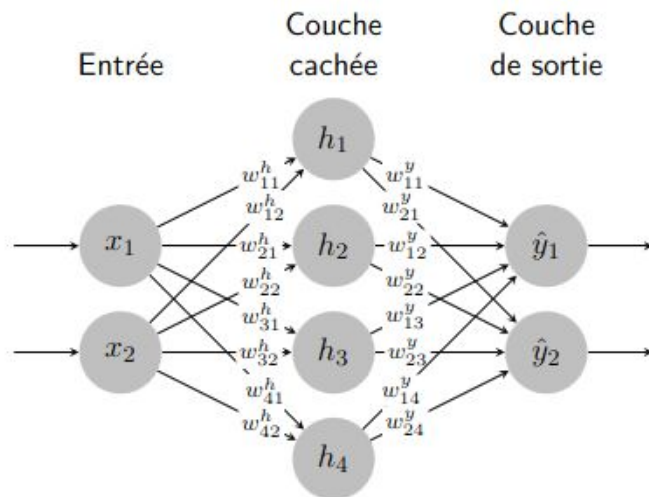$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



$$\frac{\partial \mathcal{L}}{\partial \tilde{h}} = \frac{\partial \mathcal{L}}{\partial h}\frac{\partial h}{\partial \tilde{h}} \qquad \frac{\partial \mathcal{L}}{\partial h} = \frac{\partial \mathcal{L}}{\partial \tilde{y}}\frac{\partial \tilde{y}}{\partial h} \qquad \frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}}\frac{\partial \hat{y}}{\partial \tilde{y}} \qquad \frac{\partial \mathcal{L}}{\partial \hat{y}}$$

$$\frac{\partial \mathcal{L}}{\partial W_h} = \qquad\qquad \frac{\partial \mathcal{L}}{\partial W_y} = \frac{\partial \mathcal{L}}{\partial \tilde{y}}\frac{\partial \tilde{y}}{\partial W_y}$$

Entrée — Couche cachée — Couche de sortie

$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \; x_j + b_i^h \\[2mm] h_i = \tanh(\tilde{h}_i) \\[2mm] \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \; h_j + b_i^y \\[2mm] \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$
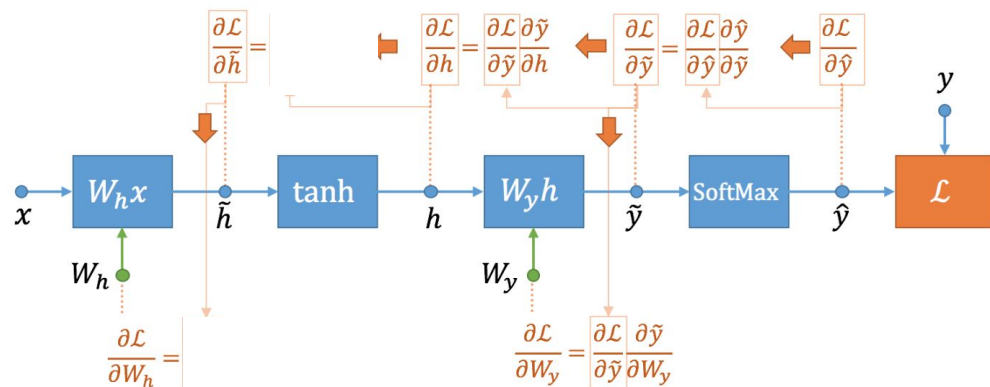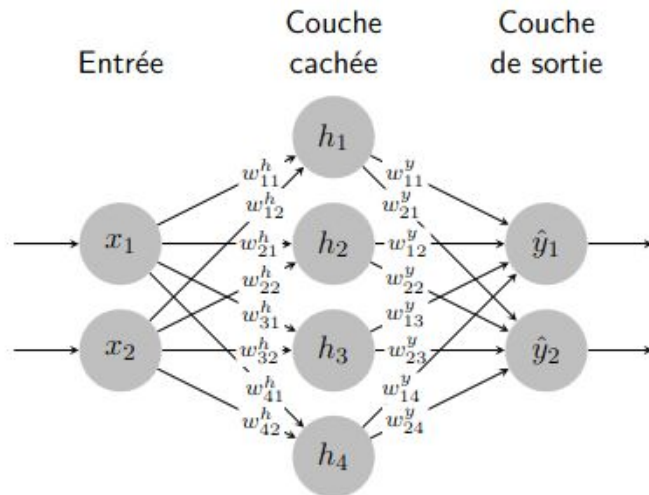
$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



$$\frac{\partial \mathcal{L}}{\partial \tilde{h}} = \frac{\partial \mathcal{L}}{\partial h} \frac{\partial h}{\partial \tilde{h}} \qquad \frac{\partial \mathcal{L}}{\partial h} = \frac{\partial \mathcal{L}}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial h} \qquad \frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \tilde{y}} \qquad \frac{\partial \mathcal{L}}{\partial \hat{y}}$$

$$\frac{\partial \mathcal{L}}{\partial W_h} = \frac{\partial \mathcal{L}}{\partial \tilde{h}} \frac{\partial \tilde{h}}{\partial W_h} \qquad \frac{\partial \mathcal{L}}{\partial W_y} = \frac{\partial \mathcal{L}}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial W_y}$$

Couche cachée — Couche de sortie

Entrée

$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \, x_j + b_i^h \\ h_i = \tanh(\tilde{h}_i) \\ \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \, h_j + b_i^y \\ \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

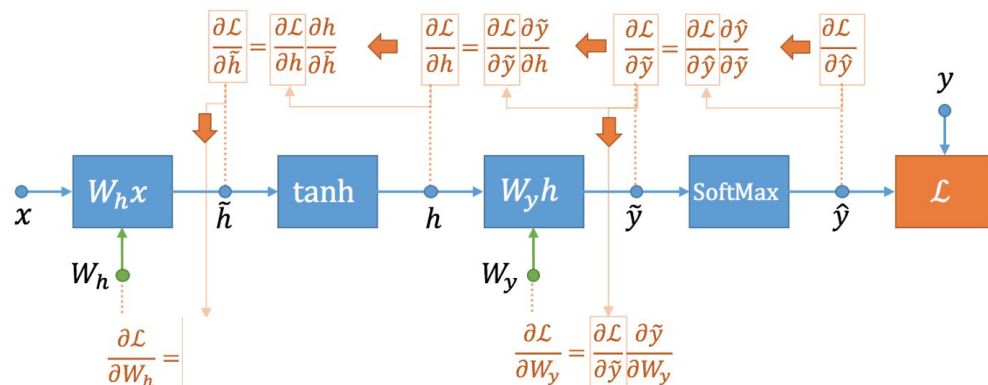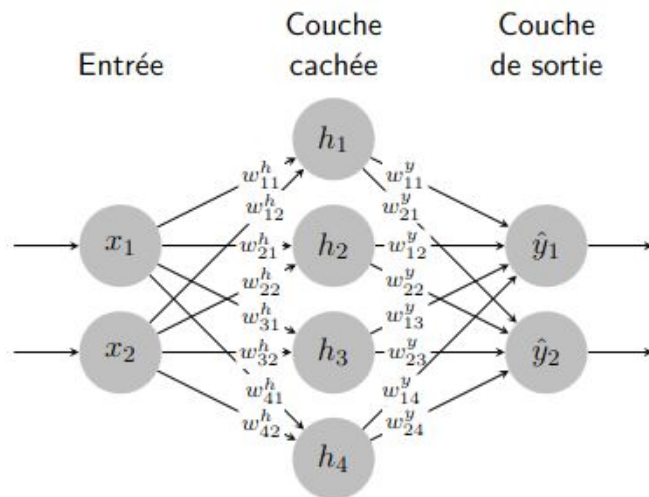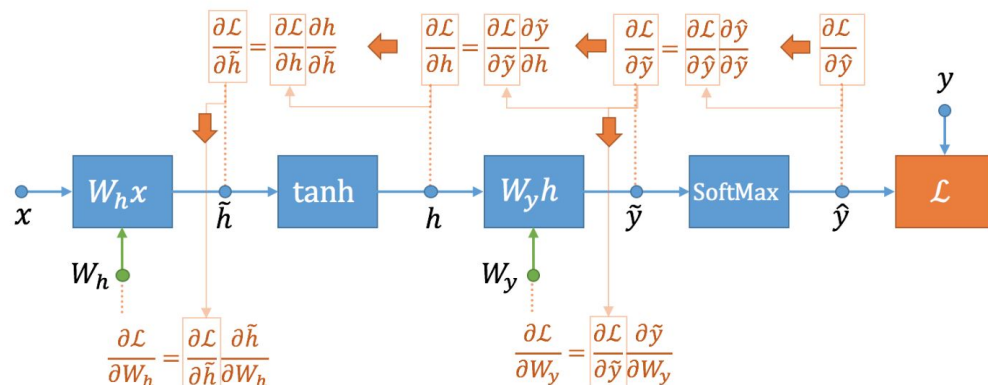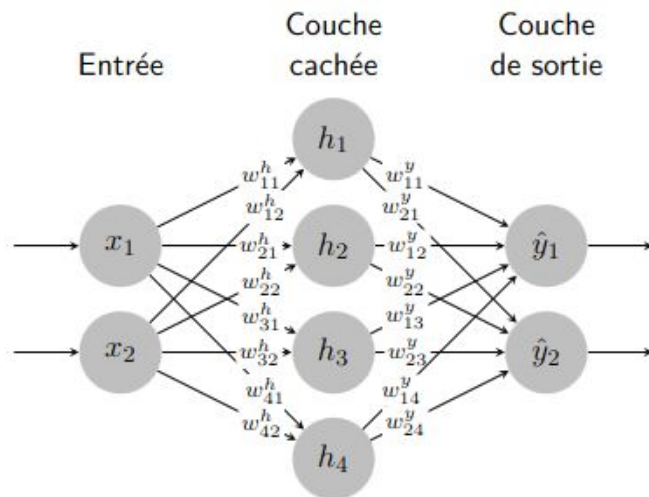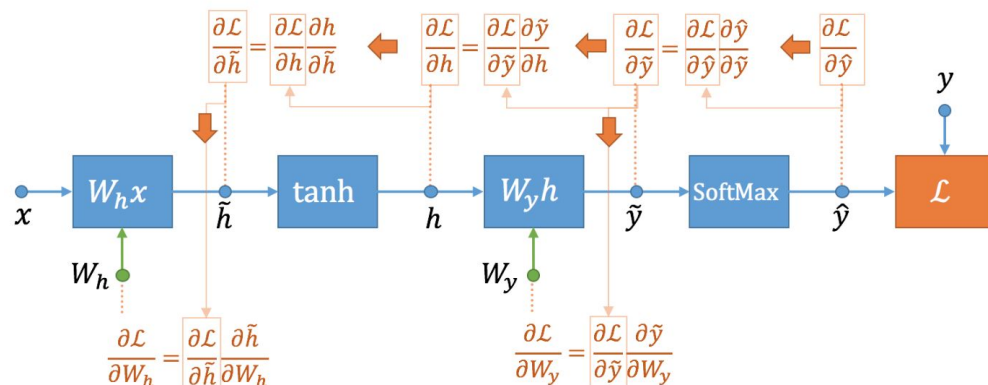$$\begin{cases} \delta_i^y = \dfrac{\partial \ell}{\partial \tilde{y}_i} = \\ \dfrac{\partial \ell}{\partial W_{i,j}^y} = \\ \dfrac{\partial \ell}{\partial b_i^y} = \end{cases}$$

43

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$

$$\frac{\partial \mathcal{L}}{\partial \tilde{h}} = \frac{\partial \mathcal{L}}{\partial h}\frac{\partial h}{\partial \tilde{h}} \qquad \frac{\partial \mathcal{L}}{\partial h} = \frac{\partial \mathcal{L}}{\partial \tilde{y}}\frac{\partial \tilde{y}}{\partial h} \qquad \frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}}\frac{\partial \hat{y}}{\partial \tilde{y}} \qquad \frac{\partial \mathcal{L}}{\partial \hat{y}}$$

$$x \quad\boxed{W_h x}\quad \tilde{h} \quad\boxed{\text{tanh}}\quad h \quad\boxed{W_y h}\quad \tilde{y} \quad\boxed{\text{SoftMax}}\quad \hat{y} \quad\boxed{\mathcal{L}} \quad y$$

$$W_h \qquad\qquad W_y$$

$$\frac{\partial \mathcal{L}}{\partial W_h} = \frac{\partial \mathcal{L}}{\partial \tilde{h}}\frac{\partial \tilde{h}}{\partial W_h} \qquad\qquad \frac{\partial \mathcal{L}}{\partial W_y} = \frac{\partial \mathcal{L}}{\partial \tilde{y}}\frac{\partial \tilde{y}}{\partial W_y}$$

Couche cachée    Couche de sortie

Entrée

$x_1$ $x_2$ $h_1$ $h_2$ $h_3$ $h_4$ $\hat{y}_1$ $\hat{y}_2$

$w_{11}^h$ $w_{12}^h$ $w_{21}^h$ $w_{22}^h$ $w_{31}^h$ $w_{32}^h$ $w_{41}^h$ $w_{42}^h$

$w_{11}^y$ $w_{21}^y$ $w_{12}^y$ $w_{22}^y$ $w_{13}^y$ $w_{23}^y$ $w_{14}^y$ $w_{24}^y$

$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h\, x_j + b_i^h \\[2mm] h_i = \tanh(\tilde{h}_i) \\[2mm] \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y\, h_j + b_i^y \\[2mm] \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum\limits_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$
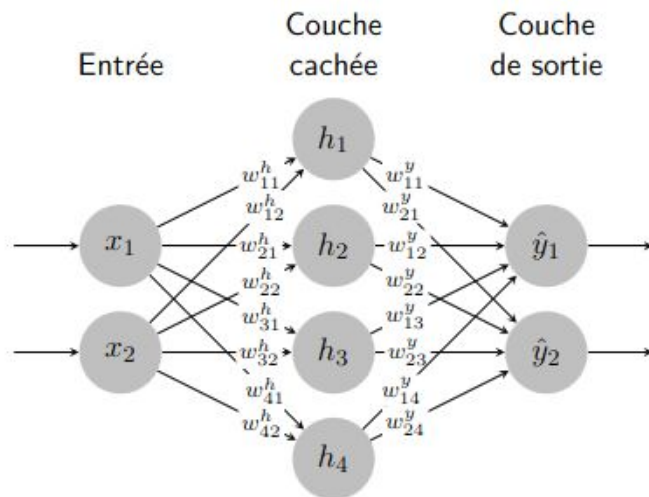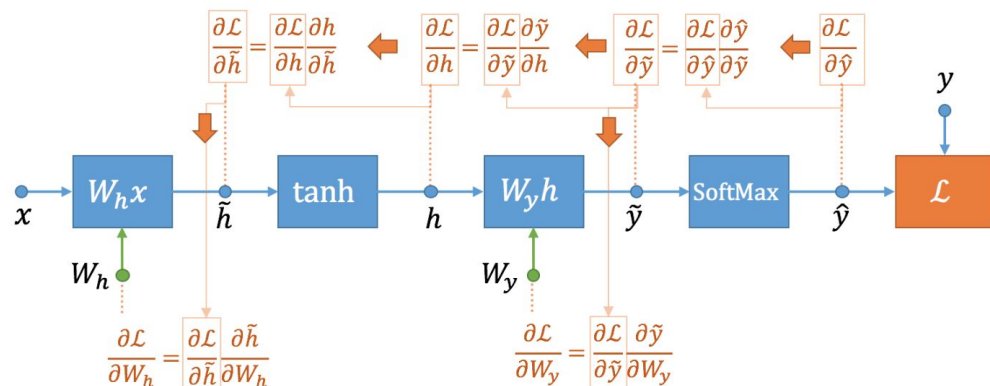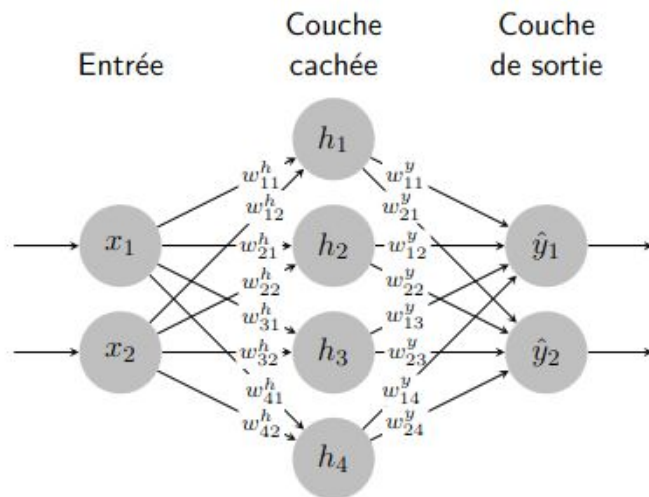
$$\begin{cases} \delta_i^y = \dfrac{\partial \ell}{\partial \tilde{y}_i} = \hat{y}_i - y_i \\[2mm] \dfrac{\partial \ell}{\partial W_{i,j}^y} = \delta_i^y\, h_j \\[2mm] \dfrac{\partial \ell}{\partial b_i^y} = \delta_i^y \end{cases}$$

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



Couche
cachée

Couche
de sortie

Entrée



$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \, x_j + b_i^h \\ h_i = \tanh(\tilde{h}_i) \\ \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \, h_j + b_i^y \\ \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

$$\begin{cases} \delta_i^y = \dfrac{\partial \ell}{\partial \tilde{y}_i} = \hat{y}_i - y_i \\ \dfrac{\partial \ell}{\partial W_{i,j}^y} = \delta_i^y \, h_j \\ 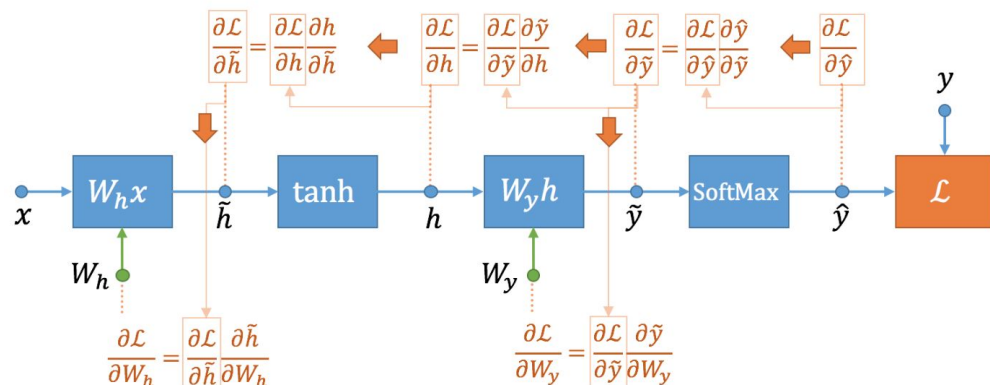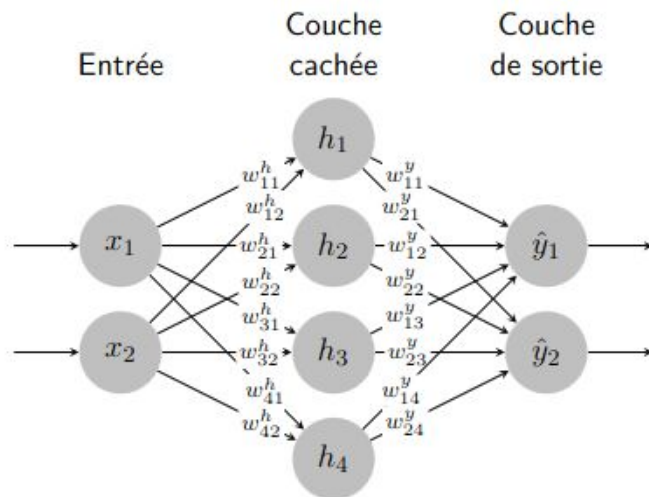\dfrac{\partial \ell}{\partial b_i^y} = \delta_i^y \\ \delta_i^h = \dfrac{\partial \ell}{\partial \tilde{h}_i} = \\ \dfrac{\partial \ell}{\partial W_{i,j}^h} = \\ \dfrac{\partial \ell}{\partial b_i^h} = \end{cases}$$
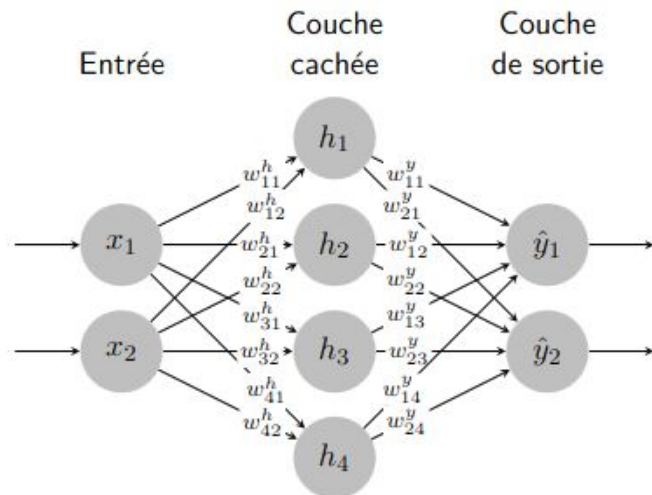
UNIVERSITÉ CÔTE D'AZUR

$$l_{CE}(\hat{y}, y) = - \sum_{i=0}^{\#Classes-1} y_i \log(\hat{y}_i)$$



$$\frac{\partial \mathcal{L}}{\partial \tilde{h}} = \frac{\partial \mathcal{L}}{\partial h}\frac{\partial h}{\partial \tilde{h}} \quad\quad \frac{\partial \mathcal{L}}{\partial h} = \frac{\partial \mathcal{L}}{\partial \tilde{y}}\frac{\partial \tilde{y}}{\partial h} \quad\quad \frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}}\frac{\partial \hat{y}}{\partial \tilde{y}} \quad\quad \frac{\partial \mathcal{L}}{\partial \hat{y}}$$

$$\frac{\partial \mathcal{L}}{\partial W_h} = \frac{\partial \mathcal{L}}{\partial \tilde{h}}\frac{\partial \tilde{h}}{\partial W_h} \quad\quad \frac{\partial \mathcal{L}}{\partial W_y} = \frac{\partial \mathcal{L}}{\partial \tilde{y}}\frac{\partial \tilde{y}}{\partial W_y}$$

Couche cachée — Couche de sortie
Entrée

$$\begin{cases} \tilde{h}_i = \sum_{j=1}^{n_x} W_{i,j}^h \, x_j + b_i^h \\[2mm] h_i = \tanh(\tilde{h}_i) \\[2mm] \tilde{y}_i = \sum_{j=1}^{n_h} W_{i,j}^y \, h_j + b_i^y \\[2mm] \hat{y}_i = \text{SoftMax}(\tilde{y}_i) = \dfrac{e^{\tilde{y}_i}}{\sum_{j=1}^{n_y} e^{\tilde{y}_j}} \end{cases}$$

$$\begin{cases} \delta_i^y = \dfrac{\partial \ell}{\partial \tilde{y}_i} = \hat{y}_i - y_i \\[2mm] \dfrac{\partial \ell}{\partial W_{i,j}^y} = \delta_i^y \, h_j \\[2mm] \dfrac{\partial \ell}{\partial b_i^y} = \delta_i^y \\[2mm] \delta_i^h = \dfrac{\partial \ell}{\partial \tilde{h}_i} = (1 - h_i^2)\sum_{j=1}^{n_y} \delta_j^y W_{j,i}^y \\[2mm] \dfrac{\partial \ell}{\partial W_{i,j}^h} = \delta_i^h \, x_j \\[2mm] \dfrac{\partial \ell}{\partial b_i^h} = \delta_i^h \end{cases}$$

- Core problem: Find gradient updates

- Gradient can be computed efficiently with backpropagation
  - Chain rule starting from the "end"
  - Re-use computed gradients (Bottom-up DP)
  - Keep forward activations for gradients
  - Simple layers mean simple gradient blocks

$$\begin{cases} \tilde{\mathbf{h}} = \mathbf{x}\mathbf{W}^{h\top} + \mathbf{b}^h \\ \mathbf{h} = \tanh(\tilde{\mathbf{h}}) \\ \tilde{\mathbf{y}} = \mathbf{h}\mathbf{W}^{y\top} + \mathbf{b}^y \\ \hat{\mathbf{y}} = \mathrm{SoftMax}(\tilde{\mathbf{y}}) \end{cases} \qquad \begin{cases} \nabla_{\tilde{\mathbf{y}}} = \hat{\mathbf{y}} - \mathbf{y} \\ \nabla_{\mathbf{W}^y} = \nabla_{\tilde{\mathbf{y}}}^{\top}\mathbf{h} \\ \nabla_{\mathbf{b}^y} = \nabla_{\tilde{\mathbf{y}}}^{\top} \\ \nabla_{\tilde{\mathbf{h}}} = (\nabla_{\tilde{\mathbf{y}}}\mathbf{W}^y) \odot (1 - \mathbf{h}^2) \\ \nabla_{\mathbf{W}^h} = \nabla_{\tilde{\mathbf{h}}}^{\top}\mathbf{x} \\ \nabla_{\mathbf{b}^h} = \nabla_{\tilde{\mathbf{h}}}^{\top} \end{cases}$$

- Lab2 on Moodle
  - Implement this by hand with basic torch!
  - Careful with batch dimension!

```python
def backward(params, outputs, Y):
    bsize = Y.shape[0]
    grads = {}

    Y_tilde_grad = outputs["yhat"] - Y
    h_tilde_grad = torch.mm(Y_tilde_grad, params['Wy']
                            ) * (1 - torch.pow(outputs['h'], 2))

    grads["Wy"] = torch.mm(Y_tilde_grad.T, outputs["h"])
    grads["Wh"] = torch.mm(h_tilde_grad.T, outputs['X'])
    grads["by"] = Y_tilde_grad.sum(dim=0,keepdim=True).T
    grads["bh"] = h_tilde_grad.sum(0, keepdim=True).T

    grads['Wy'] /= bsize
    grads['by'] /= bsize
    grads['Wh'] /= bsize
    grads['bh'] /= bsize

    return grads
```

- Torch.tensor object
  - Np.array like
  - Tracked on a computational graph
  - .grad variable to track gradients
  - .backward to backpropagate gradients through the graph
  - Activate .autograd!

# Lab 2: Autograd backward

```python
def backward(params, outputs, Y):
    bsize = Y.shape[0]
    grads = {}

    Y_tilde_grad = outputs["yhat"] - Y
    h_tilde_grad = torch.mm(Y_tilde_grad, params['Wy']
                            ) * (1 - torch.pow(outputs['h'], 2))

    grads["Wy"] = torch.mm(Y_tilde_grad.T, outputs["h"])
    grads["Wh"] = torch.mm(h_tilde_grad.T, outputs['X'])
    grads["by"] = Y_tilde_grad.sum(dim=0,keepdim=True).T
    grads["bh"] = h_tilde_grad.sum(0, keepdim=True).T

    grads['Wy'] /= bsize
    grads['by'] /= bsize
    grads['Wh'] /= bsize
    grads['bh'] /= bsize

    return grads
```

```python
params['Wh'] = torch.randn(nh, nx) * 0.3
params['Wh'].requires_grad = True
params['bh'] = torch.zeros(nh, 1, requires_grad=True)
params['Wy'] = torch.randn(ny, nh) * 0.3
params['Wy'].requires_grad = True
params['by'] = torch.zeros(ny, 1, requires_grad=True)
```

```python
with torch.no_grad():
    params['Wy'] -= eta * params['Wy'].grad
    params['Wh'] -= eta * params['Wh'].grad
    params['by'] -= eta * params['by'].grad
    params['bh'] -= eta * params['bh'].grad

    params['Wy'].grad.zero_()
    params['Wh'].grad.zero_()
    params['by'].grad.zero_()
    params['bh'].grad.zero_()
```