

Machine Learning

Neural Networks and Deep Learning: Training

Gabriele Russo Russo Francesco Lo Presti

Laurea Magistrale in Ingegneria Informatica – A.Y. 2023/24



TOR VERGATA
UNIVERSITÀ DEGLI STUDI DI ROMA

Dipartimento di Ingegneria Civile e
Ingegneria Informatica

Training Neural Networks

- ▶ In principle, training an ANN is not different than training the other parametric models you studied so far
- ▶ The goal is determining the model parameters θ that optimize a cost function $J(\theta)$
 - ▶ e.g., via stochastic gradient descent
- ▶ In this case, θ includes:
 - ▶ the connection weights, $\mathbf{W}^{(\ell)}, \forall 1 \leq \ell \leq L$
 - ▶ the bias vectors, $\mathbf{b}^{(\ell)}, \forall 1 \leq \ell \leq L$

Training Neural Networks: Cost function

Which is the cost/objective function to optimize?

- ▶ Let's consider a regression task and the squared error as the loss function

$$\mathcal{L} = \frac{1}{2}(t - y)^2 \quad (1)$$

where t is the prediction target and y is the NN prediction.

- ▶ The cost $J(\boldsymbol{\theta})$ may comprise a regularization term

$$J(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{1}{2}(t^{(i)} - y^{(i)})^2 + \lambda \Omega(\boldsymbol{\theta}) \quad (2)$$

- ▶ e.g., assuming ℓ_2 regularization:

$$J(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{1}{2}(t^{(i)} - y^{(i)})^2 + \lambda \|\boldsymbol{\theta}\|_2 \quad (3)$$

Training Neural Networks: SGD

Cannot we just keep using SGD?

Short answer: Yes

NN training via SGD (a naive algorithm!)

```
1 for epoch  $\leftarrow 1, \dots, N_{epochs}$  do  
    /* Loop over training samples */  
2     for  $i = 1, \dots, N$  do  
3          $\hat{g} \leftarrow \nabla_{\theta} J^{(i)}(\theta)$       /* Compute the gradient */  
4          $\theta \leftarrow \theta - \alpha \hat{g}$   
5     end  
6 end
```

Training Neural Networks: SGD (2)

Cannot we just keep using SGD?

- ▶ A few more challenges involved...
- ▶ Main issue: nonlinearity of NNs cause loss functions to become **nonconvex**
 - ▶ We must use gradient-based methods to drive the cost function to a low value
 - ▶ No global convergence guarantees!
- ▶ How to (efficiently) compute the gradient?
 - ▶ For single-layer NNs (e.g., Perceptron), the loss is just a function of the weights: easy
 - ▶ For multi-layer NNs, loss is a composite function of the weights of the previous layers
 - ▶ Deriving an analytical expression is easy, but evaluation can be computationally expensive

Backpropagation

- ▶ **Backpropagation** (or, **backprop**) is an algorithm to compute the gradient of a function, including the gradient of the NN cost function $\nabla_{\theta} J(\theta)$
 - ▶ Rumelhart, Hinton, Williams (1986). *Learning representations by back-propagating errors*, Nature, 323 (6088): 533–536.
- ▶ The term “backpropagation” is often used loosely to refer to the entire training algorithm – including how the gradient is used, such as by SGD – but this is not strictly correct

Backpropagation: Overview

The algorithm consists of 2 phases:

- ▶ **Forward phase:** A training instance is fed into the NN as input, resulting in a cascade of computations across the layers. The final output is used to compute the cost function.
- ▶ **Backward phase:** The gradient of the cost function with respect to the parameters is computed, traversing the NN in the opposite direction (starting from the output layer).

To efficiently compute the gradient, backprop leverages **dynamic programming** and the **chain rule** of calculus.

Forward Propagation

- ▶ **Forward propagation** (or **forward pass**): computation of intermediate variables and outputs in order, from the input layer to the output layer
- ▶ We saw forward propagation in action in the [xor](#) notebook

For a NN with a single hidden layer, forward prop. computes:

- ▶ $z = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
- ▶ $\mathbf{h} = \phi(z)$
- ▶ $\mathbf{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$

Recall: Chain Rule of Calculus

- ▶ Let $f : \mathbb{R} \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}$
- ▶ The derivative of the composite function $f(g(x))$, according to the chain rule, is computed as:

$$\frac{d}{dx}f(g(x)) = f'(g(x))g'(x) \quad (4)$$

- ▶ Intuitively: increasing x by some infinitesimal quantity h_1 “causes” g to change by the infinitesimal $h_2 = g'(x)h_1$. This in turn causes f to change by $f'(g(x))h_2 = f'(g(x))g'(x)h_1$.
- ▶ Equivalently, letting $y = g(x)$ and $z = f(y)$:

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} \quad (5)$$

Recall: Chain Rule of Calculus (2)

- ▶ Consider a multivariate function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and $g_i : \mathbb{R} \rightarrow \mathbb{R}$

$$\frac{\partial f(g_1(x), \dots, g_n(x))}{\partial x} = \sum_{k=1}^n \frac{\partial f(g_1(x), \dots, g_n(x))}{\partial g_k(x)} \cdot \frac{\partial g_k(x)}{\partial x} \quad (6)$$

- ▶ Intuitively: changing x slightly causes each g_k to change slightly. Each of these effects causes a slight change to f . For infinitesimal changes, these effects combine additively.
- ▶ Alternatively, $g : \mathbb{R} \rightarrow \mathbb{R}^n$ and

$$\frac{\partial f(g(x))}{\partial x} = \sum_{k=1}^n \frac{\partial f(g(x))}{\partial g_k(x)} \cdot \frac{\partial g_k(x)}{\partial x} \quad (7)$$

Recall: Chain Rule of Calculus (3)

- ▶ Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and $g : \mathbb{R}^m \rightarrow \mathbb{R}^n$
- ▶ In this case, we consider the partial derivatives w.r.t. each x_i

$$\frac{\partial f(g(\mathbf{x}))}{\partial x_i} = \sum_{k=1}^n \frac{\partial f(g(\mathbf{x}))}{\partial g_k(\mathbf{x})} \cdot \frac{\partial g_k(\mathbf{x})}{\partial x_i} \quad (8)$$

- ▶ In vector notation:

$$\nabla_{\mathbf{x}} f(g(\mathbf{x})) = \left(\frac{\partial \mathbf{g}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{g}} f(\mathbf{g}(\mathbf{x})) \quad (9)$$

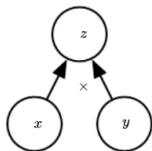
where $\left(\frac{\partial \mathbf{g}}{\partial \mathbf{x}} \right)$ is the Jacobian matrix of $g(\mathbf{x})$.

- ▶ Same rule exploited if we have multi-dimensional tensors (just think that you can flatten the tensor into a vector).

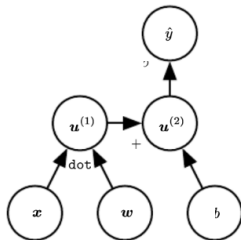
Computational Graph

- ▶ It is useful to describe NNs using a more precise formalism, namely **computational graphs**
- ▶ Each **node** indicates a **variable**
 - ▶ A scalar, a vector, a matrix, or a tensor
- ▶ An **operation** is a function of 1+ variables
 - ▶ We assume operations to return a single output variable (which may be a vector though!)
- ▶ An **edge** (x, y) indicates that y is computed by applying an operation to x
 - ▶ We may annotate y with the operation

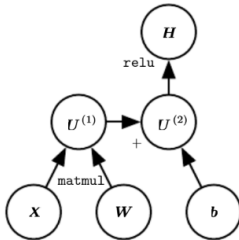
Examples



(a)

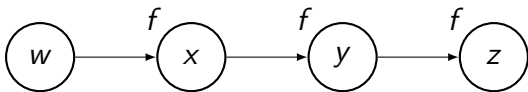


(b)



(c)

Example: Chain Rule



$$z = f(y) = f(f(x)) = f(f(f(w)))$$

$$\begin{aligned}\frac{\partial z}{\partial w} &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial w} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} = \\ &= f'(y) f'(x) f'(w) = \\ &= f'(f(\textcolor{red}{f}(\textcolor{red}{w}))) f'(\textcolor{red}{f}(\textcolor{red}{w})) f'(w)\end{aligned}$$

Applying the rule, we often need to evaluate subexpressions more than once. We either store them or re-compute them

Backpropagation with Scalar Variables

- ▶ Consider a computational graph with nodes $u^{(1)}, u^{(2)}, \dots, u^{(n)}$, in topological order, all representing **scalar** variables
- ▶ The first $k < n$ variables $u^{(1)}, \dots, u^{(k)}$ represent input variables
- ▶ $u^{(n)}$ is the output variable
- ▶ We want to compute the gradient w.r.t. input

$$\frac{\partial u^{(n)}}{\partial u^{(i)}}, \forall i \in \{1, \dots, k\}$$

- ▶ Note: when training a NN, we usually are not interested in computing the gradient w.r.t. input, but to parameters

Backpropagation with Scalar Variables (2)

- ▶ Each node $u^{(i)}$ is associated with an operation $f^{(i)}$

$$u^{(i)} = f^{(i)}(\mathcal{A}^{(i)}), \quad \mathcal{A}^{(i)} = \{u^{(j)} \mid j \in \text{Parent}(u^{(i)})\}$$

Forward pass (with scalar variables only)

```
1 for  $i = 1, \dots, k$  do  
2   |    $u^{(i)} \leftarrow x_i$                                 /* Input values */  
3 end  
4 for  $i = k + 1, \dots, n$  do  
5   |    $u^{(i)} \leftarrow f^{(i)}(\mathcal{A}^{(i)})$   
6 end  
7 return  $u^{(n)}$ 
```


Backpropagation with Scalar Variables (3)

- ▶ We use a lookup table `grad` to store the computed value of $\frac{\partial u^{(n)}}{\partial u^{(i)}}$ for every $u^{(i)}$
- ▶ For convenience, let $children(j) = \{i | j \in Parent(u^{(i)})\}$

Backward propagation (with scalar variables only)

- 1 Run forward pass (previous slide) to evaluate variables
- 2 Initialize empty lookup table `grad`
- 3 `grad[u(n)] ← 1`
- 4 **for** `j=n-1,n-2,...,1` **do**
 - 5 `/* Apply chain rule` `*/`
`grad[u(j)] ← $\sum_{i \in children(j)} grad[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$`
- 6 **end**
- 7 **return** `{grad[u(i)], i = 1, ..., k}`

Backpropagation for NNs

- ▶ We now focus on the specific case of backpropagation for multi-layer neural networks
- ▶ To keep things simple, we assume that a single training instance x and the associated label t are used to compute the cost $J(\theta)$
- ▶ We aim to compute $\nabla_{\theta} J(\theta)$

Forward Pass

Forward pass

```
1  $\mathbf{h}^{(0)} \leftarrow \mathbf{x}$ 
2 for  $k=1,\dots,L$  do
3    $\mathbf{a}^{(k)} \leftarrow \mathbf{W}^{(k)} \mathbf{h}^{(k-1)} + \mathbf{b}^{(k)}$ 
4    $\mathbf{h}^{(k)} \leftarrow f^{(k)}(\mathbf{a}^{(k)})$ 
5 end
6  $\mathbf{y} \leftarrow \mathbf{h}^{(L)}$ 
7 return  $J = L(\mathbf{y}, \mathbf{t}) + \lambda \Omega(\boldsymbol{\theta})$ 
```

Backward Pass

Backpropagation for multi-layer NNs

```
1  $\mathbf{g} \leftarrow \nabla_{\mathbf{y}} J = \nabla_{\mathbf{y}} L(\mathbf{y}, \mathbf{t})$ 
2 for  $k=L, \dots, 2, 1$  do
3    $\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f^{(k)'}(\mathbf{a}^{(k)})$ 
4    $\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\boldsymbol{\theta})$ 
5    $\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)T} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta})$ 
   /* Propagate the gradients */
6    $\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)T} \mathbf{g}$ 
7 end
8 return  $\nabla_{\mathbf{W}^{(k)}} J, \nabla_{\mathbf{b}^{(k)}} J, \forall k = 1, \dots, L$ 
```

Backward Pass (2)

$$1 \quad g \leftarrow \nabla_y J = \nabla_y \mathcal{L}(y, t)$$

- ▶ g is initialized as the gradient of the cost function with respect to the output vector y
- ▶ The cost function may include a regularization term, which does not depend on the output
- ▶ We need just differentiate the loss function
- ▶ This step can be performed analytically and depends on the loss function in use
 - ▶ e.g., if $\mathcal{L} = \frac{1}{2} \|y - t\|_2^2$, we get $\nabla_y \mathcal{L} = (y - t)$

Backward Pass (3)

```
2 for k=L,..,2,1 do
3    $\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot \mathbf{f}^{(k)'}(\mathbf{a}^{(k)})$ 
4   ...
```

- ▶ We iterate **backward** from the L -th layer to the first hidden layer
- ▶ At the beginning of each iteration, \mathbf{g} stores the gradient of the cost with respect to the k -th layer
- ▶ e.g., at first iteration, \mathbf{g} is the gradient w.r.t. the output layer

Backward Pass (4)

```
2 for  $k=L,..,2,1$  do
3   |  $\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f^{(k)'}(\mathbf{a}^{(k)})$ 
```

- ▶ We compute the gradient w.r.t. the pre-activation value
- ▶ Being $f^{(k)}$ the activation function of the k -th layer:

$$\mathbf{h}^{(k)} = f^{(k)}(\mathbf{a}^{(k)})$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(k)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(k)}} \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{a}^{(k)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(k)}} \odot f^{(k)'}(\mathbf{a}^{(k)}) = \mathbf{g} \odot f^{(k)'}(\mathbf{a}^{(k)})$$

- ▶ Note: \odot denotes the element-wise product

Remark: Softmax

```
2 for k=L,...,2,1 do
3   |  $\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot \mathbf{f}^{(k)'}(\mathbf{a}^{(k)})$ 
```

- ▶ The simple expression of line 3 above is not correct if $f^{(k)}$ is the softmax
- ▶ In that case, the output of the layer depends on **all** the pre-activation variables. Therefore, the Jacobian must be used

$$\mathbf{g} \leftarrow \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(k)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{h}^{(k)}} \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{a}^{(k)}} = \left(\frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{a}^{(k)}} \right)^T \cdot \mathbf{g}$$

Backward Pass (5)

$$4 \quad \nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\boldsymbol{\theta})$$

$$5 \quad \nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)T} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta})$$

- ▶ Recall: \mathbf{g} is the gradient w.r.t. $\mathbf{a}^{(k)} = \mathbf{W}^{(k)} \mathbf{h}^{(k-1)} + \mathbf{b}^{(k)}$
- ▶ We now take into account the regularization term as well:

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{W}^{(k)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(k)}} + \lambda \frac{\partial \Omega(\boldsymbol{\theta})}{\partial \mathbf{W}^{(k)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(k)}} \frac{\partial \mathbf{a}^{(k)}}{\partial \mathbf{W}^{(k)}} + \nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) = \\ &= \mathbf{g} \mathbf{h}^{(k-1)T} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) \end{aligned}$$

- ▶ The same reasoning holds for $\mathbf{b}^{(k)}$

Backward Pass (6)

$$6 \quad g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)T} g$$

- ▶ We use g to **back-propagate** the gradient
- ▶ At the next iteration g must hold the gradient w.r.t. $h^{(k-1)}$

$$a^{(k)} = W^{(k)} h^{(k-1)} + b^{(k)}$$

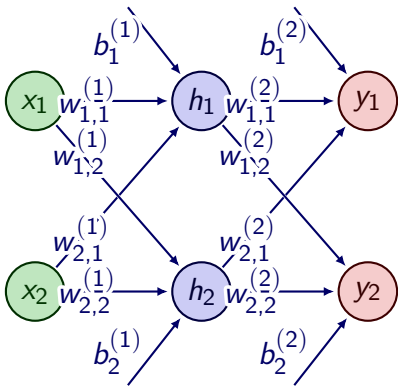
$$\frac{\partial \mathcal{L}}{\partial h^{(k-1)}} = \frac{\partial \mathcal{L}}{\partial a^{(k)}} \frac{\partial a^{(k)}}{\partial h^{(k-1)}} = W^{(k)T} g$$

Derivative of Activation Functions: Examples

	$f(x)$	$f'(x)$
Logistic (sigmoid)	$\frac{1}{1+e^{-x}}$	$f(x)(1-f(x))$
ReLU	$\max\{0, x\}$	$\begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$
Hyp. tan	$\frac{e^x - e^{-x}}{e^x + e^{-x}}$	$1 - f(x)^2$

Exercise

Given the following NN specification and the input vector $\mathbf{x} = \{0.5, 0.1\}$, apply backprop algorithm



$$\mathbf{W}^{(1)} = \begin{bmatrix} 0.15 & 0.25 \\ 0.2 & 0.3 \end{bmatrix}$$

$$\mathbf{W}^{(2)} = \begin{bmatrix} 0.4 & 0.5 \\ 0.5 & 0.55 \end{bmatrix}$$

$$\mathbf{b}^{(1)} = [0.35, 0.35]$$

$$\mathbf{b}^{(2)} = [0.6, 0.6]$$

$$\mathbf{t} = [0.01, 0.99]$$

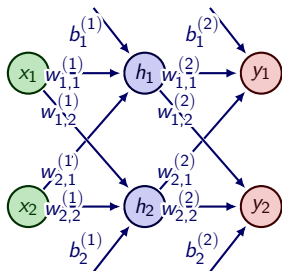
$$\mathcal{L} = \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|_2^2$$

All the units use the **logistic activation**.

 `ex_backprop.ipynb` (solution)

Exercise: Solution (partial)

Viene usata sigmoide: $\frac{1}{1 + e^{(-z)}}$

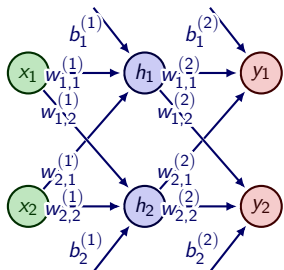


$$\begin{aligned} \mathbf{a}^{(1)} &= \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} = \\ &= \begin{bmatrix} 0.15 & 0.25 \\ 0.2 & 0.3 \end{bmatrix} \cdot \begin{bmatrix} 0.5 \\ 0.1 \end{bmatrix} + \begin{bmatrix} 0.35 \\ 0.35 \end{bmatrix} = \\ &= \begin{bmatrix} 0.095 \\ 0.155 \end{bmatrix} + \begin{bmatrix} 0.35 \\ 0.35 \end{bmatrix} = \begin{bmatrix} 0.445 \\ 0.505 \end{bmatrix} \end{aligned}$$

$$\mathbf{h}^{(1)} = \sigma(\mathbf{a}^{(1)}) = \begin{bmatrix} 0.6094 \\ 0.6236 \end{bmatrix}$$

Note: numerical results may be inaccurate due to approximations...check the notebook for correct results!

Exercise: Solution (partial)

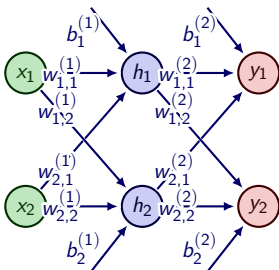


$$\begin{aligned} \mathbf{a}^{(2)} &= \mathbf{W}^{(2)} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} = \\ &= \begin{bmatrix} 0.4 & 0.5 \\ 0.5 & 0.55 \end{bmatrix} \cdot \begin{bmatrix} 0.6094 \\ 0.6236 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.6 \end{bmatrix} = \\ &= \begin{bmatrix} 0.5244 \\ 0.6478 \end{bmatrix} + \begin{bmatrix} 0.6 \\ 0.6 \end{bmatrix} = \begin{bmatrix} 1.1244 \\ 1.2478 \end{bmatrix} \end{aligned}$$

$$\mathbf{y} = \mathbf{h}^{(2)} = \sigma(\mathbf{a}^{(2)}) = \begin{bmatrix} 0.7548 \\ 0.7770 \end{bmatrix}$$

$$\mathcal{L} = \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|_2^2 = \frac{1}{2} \left\| \begin{bmatrix} 0.7448 \\ -0.213 \end{bmatrix} \right\|_2^2 = 0.3$$

Exercise: Solution (partial)



$$\mathbf{g} \leftarrow \nabla_{\mathbf{y}} \mathcal{L} = \mathbf{y} - \mathbf{t} = \begin{bmatrix} 0.7448 \\ -0.213 \end{bmatrix}$$

$k=2$

$$\begin{aligned} \mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(2)}} \mathcal{L} &= \mathbf{g} \odot \sigma'(\mathbf{a}^{(2)}) = \\ &= \begin{bmatrix} 0.7448 \\ -0.213 \end{bmatrix} \odot \begin{bmatrix} \sigma(1.1244)(1 - \sigma(1.1244)) \\ \sigma(1.2478)(1 - \sigma(1.2478)) \end{bmatrix} = \\ &= \begin{bmatrix} 0.7448 \\ -0.213 \end{bmatrix} \odot \begin{bmatrix} 0.1851 \\ 0.1733 \end{bmatrix} = \begin{bmatrix} 0.1378 \\ -0.0369 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{w}^{(2)}} \mathcal{L} &= \mathbf{g} \mathbf{h}^{(1)T} = \begin{bmatrix} 0.1378 \\ -0.0369 \end{bmatrix} \cdot \begin{bmatrix} 0.6094 & 0.6236 \end{bmatrix} = \\ &= \begin{bmatrix} 0.0840 & 0.0860 \\ -0.0225 & -0.0230 \end{bmatrix} \end{aligned}$$

Computational Demand

- ▶ Let's consider a generic computational graph
- ▶ Assume each operation to have roughly the same cost (actually, an operation may be a matrix multiplication!)
 - ▶ We can reason in terms of operations
- ▶ A naive application of the chain rule would have an exponential runtime because of repeated evaluation of subexpressions
- ▶ With n nodes, the forward pass of backprop executes $\mathcal{O}(n)$ operations
- ▶ Backward pass adds an operation for every edge in the original graph: $\mathcal{O}(n^2)$

Computational Demand for NNs

- ▶ For the specific case of feedforward fully-connected NNs, we can be more precise
- ▶ Consider a NN with L hidden layers, each with m neurons
- ▶ We drop the assumption of atomic operations (i.e., a matrix multiplication is not $\mathcal{O}(1)$)
- ▶ **Forward pass**: an add-and-multiply operation for every weight: $\mathcal{O}(Lm^2)$
- ▶ **Backward pass**: more operations to perform, but still: $\mathcal{O}(Lm^2)$
- ▶ Note: this is just one step in the whole training process!

Generalizing Backpropagation: Input Batches

- ▶ So far, we considered the gradient of the cost function computed on a **single input instance x**
- ▶ In practice, training is more efficient if **minibatches** of input are considered, because computational parallelism can be exploited (e.g., vectorization)
- ▶ The same algorithm can be applied: just replace x with a matrix X , where each row is a different training instance
- ▶ Operations involving gradients do not change: just imagine to “flatten” matrices (or tensors) into a vector

Training with Minibatch SGD

- ▶ Recall Minibatch SGD from the first part of the course

NN training via Minibatch SGD

```
1 while stopping criterion not met do  
2   |   Sample minibatch  $\mathcal{B}$  randomly from the training set  
3   |    $\mathbf{g} \leftarrow \frac{1}{|\mathcal{B}|} \nabla_{\boldsymbol{\theta}} J^{\mathcal{B}}(\boldsymbol{\theta})$       /* Compute the gradient */  
4   |    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \mathbf{g}$   
5 end
```

Minibatch Size

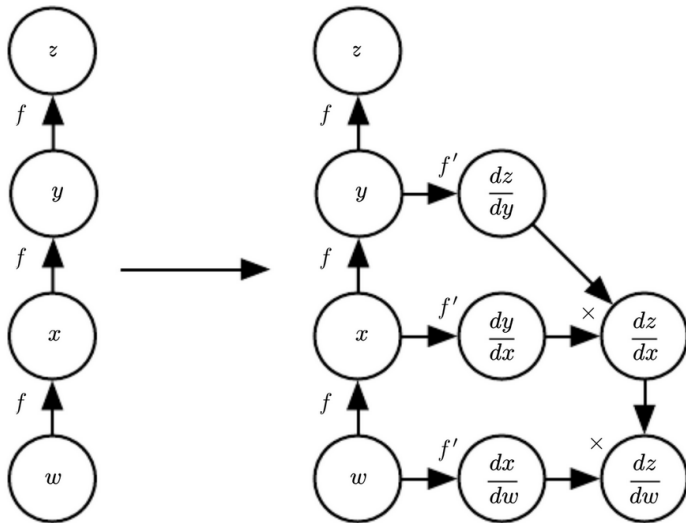
How to choose the size of the minibatch?

- ▶ Larger batches generally make training faster, but smaller batches favor convergence
- ▶ Batch size is usually a power of 2 (the reason should be obvious...)
- ▶ Typical values are, e.g., 8, 16, 32, 64
- ▶ The choice depends on the hardware in use
- ▶ **Tip:** keep it fixed while optimizing other hyperparameters, and optimize it at the end (we will come back to this issue later!)

Symbol-to-symbol Differentiation

- ▶ Computational graphs (as well as algebraic expressions) are defined over **symbols** (i.e., variables that do not have a specific value)
- ▶ An approach to backpropagation implementation: given a set of numerical values for the inputs in the graph, return the gradient at those input values
 - ▶ Called **symbol-to-number differentiation**
 - ▶ Implemented in Torch and Caffe
- ▶ Alternative approach (**symbol-to-symbol**): extend the original graph with new nodes that provide a symbolic description of the desired derivatives
 - ▶ Implemented in Theano and TensorFlow

Example: Symbol-to-symbol Differentiation



General Backprop and Custom Operations

- ▶ Backpropagation can be implemented in a way that enable easy integration of new operations on vectors, matrices and tensors
- ▶ The algorithm does not require a specification of which operations are performed, nor how to compute their derivatives
- ▶ It only needs a suitable `op` interface, providing a `backprop()` method to backpropagate the gradient
- ▶ Deep learning libraries provide built-in implementations of most used operations, but also support custom operations
- ▶ More details: Goodfellow–Bengio–Courville Sec. 6.5.6

Example with scikit-learn

- ▶ scikit-learn provides `MLPClassifier` and `MLPRegressor`, based on multi-layer feedforward NNs
- ▶ Very easy to use, but not recommended to work with NNs
 - ▶ No GPU support
 - ▶ Less flexibility compared to other frameworks (e.g., TensorFlow)
- ▶ We see a simple classification example on the Iris data set

 `mlpclassifier.ipynb`

Exercise

Spend some time on
<http://playground.tensorflow.org>

3/10/2023

Introduction to Keras and TensorFlow

TensorFlow

- ▶ Library for ML focused on DNNs
- ▶ Developed by [Google Brain](#) team
 - ▶ Many scientists and developers involved, including Geoffrey Hinton (2018 Turing Award winner)
- ▶ First released in 2015; major update (TensorFlow 2) in 2019
- ▶ Built-in support for GPU execution, as well as distributed execution



TensorFlow

Installing TensorFlow

- ▶ Main version available for Linux; Windows; macOS (no GPU)
 - ▶ tensorflow.js for browsers
 - ▶ TensorFlow Lite for embedded/mobile devices
- ▶ A command should be enough:

```
pip install tensorflow
```

- ▶ If you use Google Colab, it is already installed

Keras

Keras semplifica la definizione di modelli ed architettura rete, ma non offre training per la rete, per questo ha bisogno di un backend (TensorFlow)

- ▶ Keras is an open-source library, first released in 2015, that provides a **Python API** for ANNs
- ▶ The goal is easing the definition of deep models
- ▶ Keras requires a **backend** for actual execution
 - ▶ TensorFlow, PyTorch, JAX, ...
- ▶ Since 2017, TensorFlow has integrated its own implementation of Keras (`tf.keras`)
 - ▶ we will use this Keras implementation
 - ▶ no need to install it explicitly

Quindi installando TensorFlow, avremo incluso anche una versione di Keras (`tf.Keras`)

TensorFlow APIs

Abbiamo tre diverse interfacce per definire modelli:

Le prime due vengano da Keras, e useremo quasi esclusivamente loro.

La Core Api è di più basso livello, e la uso se devo introdurre qualcosa di particolare, quindi in generale sarà poco usata.

- ▶ As part of `tf.keras`:
 - ▶ **Sequential API**: easiest to use; OK for most apps
 - ▶ **Functional API**: more flexibility
- ▶ **Core API**: only for advanced stuff
- ▶ We will mostly use the Sequential API

Basics

► Let's start with the basics of TensorFlow

- Tensors è un array multidimensionale (tipo numpy). Se voglio inserire dei numeri ho prefisso "constant". "shape" = size, dtype = tipo
- Variables
- Graphs
- Automatic differentiation se metto una variabile in una funzione, tensorflow calcola automaticamente la derivata. Ciò porta alla creazione di un graph.
- Modules

Modules: se spengo server, non voglio perdere dati, allora salvo il modulo, cioè tutte le variabili, pesi etc... che posso ricaricare quando voglio.




tf_basics.ipynb

abbiamo anche `tf.nn.softmax`, applicabile ai livelli di uscita.

Il tensore è IMMUTABILE, ovvero, ad ogni operazione viene creato un nuovo tensore, quindi ciò che posso fare è aggiornare la variabile associata al tensore.
Oppure usare una "Variable", tensore mutabile.

con `@tf.function` si ottimizza la funzione, creando un grafo delle operazioni che verrà richiamato le volte successive per velocizzare il codice, usando tale grafo. Se passo da interi a virgola mobili, il grafo viene ricompilato. Queste operazioni le vediamo poco, lavoriamo a livello più alto.

- ▶ 2 core data structures: layers and models (no neuroni, è scatola nera in cui metto input e produce output)
- ▶ **Layer**: simple input/output transformation
 - ▶ encapsulates a state (weights) and some computation
 - ▶ may represent a layer of a DNN, but also a data preprocessing step
- ▶ **Model**: directed acyclic graph (DAG) of layers
 - ▶ e.g., **Sequential** model: a linear stack of layers
 - ▶ **fit()** and **predict()** methods (similar to scikit-learn)
grado diretto aciclico di livelli, quindi grafo dei livelli.

 `tf_keras.ipynb`

"Dense", indica livello fully connected. (vedi foto)

tali metodi sono quelli offerti dagli stimatori in scikit learn, addestramento (con loop di addestramento visto prima) e predizione

NNs for Regression: Hints

7/11/2023

come impostare hyperparametri modello, vediamo REGOLE GENERALI, o LINEE GUIDA.

Se lavoro in 3 dimensioni -> 3 neuroni. Se parlo di automobili, metto un neurone per ogni feature (cavalli, posti, anno...). Il feature engineering va fatto anche con reti neurali, anche se abbiamo visto che alcune cose "le fanno da sole". Se posso applicarne qualcuna, perchè sono facili ad es, meglio farlo.

Hyperparameter	Possible/typical values
# input neurons	1 per input feature
# hidden layers	?? (1-5 for many tasks)
# hidden units per layer	?? (10-100 for most tasks)
# output units	1 per predicted dimension
Hidden activation	ReLU (quantomeno per iniziare uso lei)
Output activation	None; ReLU for positive outputs; tanh for bounded outputs
Loss	MSE or MAE MAE = errore medio assoluto MSE = pesa di più se sbaglio di tanto

Quanti livelli nascosti? da 1 a 5 sono normalmente sufficienti, dipende dalla complessità.

Ciò si lega anche alle unità nel livello nascosto. Se ho 1000 feature, non sembra sensato mettere solo 10 unità per livello. Spesso si va a tentativi, detto anche "hyperparameter optimization".

NNs for **Classification**: Hints

Hyperparameter	Binary	Multiclass
Input and hidden layers	<i>Same as regression</i>	
# output units	1	1 per class
Output activation	Sigmoid	Softmax
Loss	Cross entropy	Cross entropy

Sigmoid per Binary, perchè output compreso tra 0 e 1, ok per probabilità.

Se multiclasse, usiamo la softmax, prende i valori pre attivazione e produce vettore in uscita tra 0 e 1.

Remark: Cross-Entropy Loss

- ▶ With 2 classes, in Keras you have BinaryCrossentropy

*t è 0 o 1, solo un termine sarà diverso da 0. Log(y) è la probabilità, se y=1 loss è 0 (ln(1)=0).
Altrimenti più mi allontano
più il logaritmo diventa piccolo,
mettiamo il meno davanti e allora cresce.*

$$\mathcal{L} = -(t \log y + (1 - t) \log (1 - y))$$

- ▶ With $C > 2$ classes and, hence, $y = (y_1, \dots, y_C)$, the loss is:

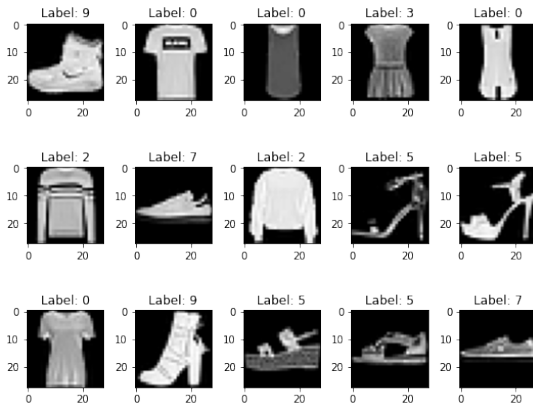
solo un termine sarà diverso da 0

$$\mathcal{L} = - \sum_{c=0}^C t_c \log y_c$$


- ▶ CategoricalCrossentropy if targets are given as one-hot vectors
- ▶ SparseCategoricalCrossentropy if targets are integers between 0 and $(C - 1)$

Example: Fashion MNIST

- ▶ We consider the **Fashion MNIST** dataset
- ▶ 60,000 grayscale images of 28×28 pixels each to classify, with 10 classes



Tuning Hyperparameters

- ▶ Flexibility of NNs is also one of their drawbacks: many **hyperparameters** to tweak!
 - ▶ Number of layers, units
 - ▶ Optimization algorithm (and its hyperparameters!)
 - ▶ ...
- ▶ Exhaustive explorations (i.e., grid search) usually not admissible
- ▶ You can perform a randomized search using scikit-learn facilities
 -  `tf_fashion_hyper.ipynb`
- ▶ Specialized libraries for hyperparameter tuning, e.g., **Hyperopt**, **Talos**, **Spearmint**, ...

References

- ▶ Goodfellow et al.: 6.5, 7.1, 7.4, 7.8, 7.11, 7.12, 8.1–8.5, 8.7.1
- ▶ D2L: 5.3–5.6, 8.5, 12
- ▶ Hands-on ML: Chapters 10-11