

Automatic differentiation

JJV

Oct 7, 2022

Optimization

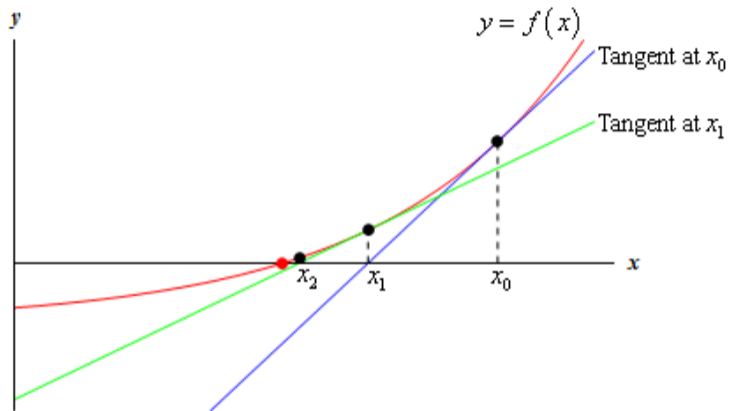
Find “the best” parameters to reach a goal

Usually, minimize a differentiable **loss function**

Find the zeroes of another function (its derivative)

How to find the zeroes of a function?

Newton's method: find x such that $f(x) = 0$



$f : \mathbf{R} \rightarrow \mathbf{R}$ differentiable, $\exists x, f'(x) \neq 0$

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$

Quadratic convergence
 $\exists C > 0, |x_{t+1} - \ell| \leq C|x_t - \ell|^2$

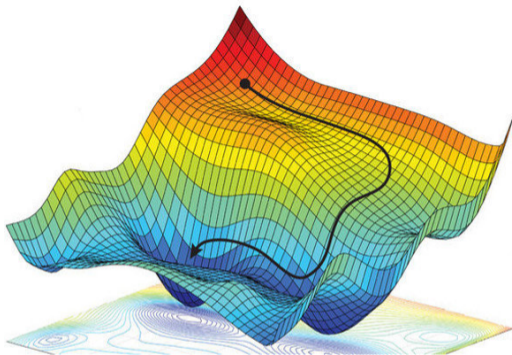
In higher dimension

Let $g : \mathbf{R}^n \rightarrow \mathbf{R}$ twice differentiable, with $n \gg 1$

What is the size of $g'(\mathbf{x})$? Usually noted $\frac{\partial g}{\partial \mathbf{x}}$ or $\nabla_{\mathbf{x}} g$.

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \underbrace{g''(\mathbf{x}_t)^{-1}}_{\in \mathbf{R}^{n \times n}, O(n^3)} \underbrace{g'(\mathbf{x}_t)}_{\in \mathbf{R}^n}$$

Gradient descent



$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t)$$

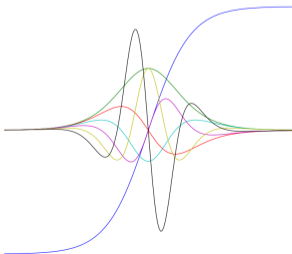
```
1 for each epoch:  
2   for  $\mathbf{x}$ ,  $y$  in dataset:  
3     compute gradients  $\frac{\partial \mathcal{L}}{\partial \theta}(f_{\theta}(\mathbf{x}), y)$  # also noted  $\nabla_{\theta} \mathcal{L}$   
4      $\theta \leftarrow \theta - \gamma \nabla_{\theta} \mathcal{L}$  #  $\gamma$  is the learning rate
```

SGD in PyTorch

```
1 import torch
2
3 # define model, n_epochs, trainloader
4 criterion = torch.nn.CrossEntropyLoss()
5 optimizer = torch.optim.SGD(model.parameters(), lr=0.001)
6
7 for _ in range(n_epochs):
8     for batch_inputs, batch_labels in trainloader:
9         outputs = model(batch_inputs)
10        loss = criterion(outputs, batch_labels)
11
12        optimizer.zero_grad()
13        loss.backward()
14        optimizer.step()
```

How to compute gradients automatically?

```
>>> from autograd import elementwise_grad as egrad # vectorize over inputs
>>> import matplotlib.pyplot as plt
>>> x = np.linspace(-7, 7, 200)
>>> plt.plot(x, tanh(x),
...          x, egrad(tanh)(x), # first derivative
...          x, egrad(egrad(tanh))(x), # second derivative
...          x, egrad(egrad(egrad(tanh)))(x), # third derivative
...          x, egrad(egrad(egrad(egrad(tanh))))(x), # fourth derivative
...          x, egrad(egrad(egrad(egrad(egrad(tanh)))))(x), # fifth derivative
...          x, egrad(egrad(egrad(egrad(egrad(egrad(tanh)))))))(x) # sixth derivative)
>>> plt.show()
```



Existing methods and their limitations

Numerical differentiation

$$\frac{f(x+h) - f(x)}{h}$$

Round-off errors

Symbolic differentiation

Have to keep symbolic expressions at each step of the process

Automatic differentiation

Chain rule

$$(f \circ g)' = g' \cdot (f' \circ g)$$

Generalized chain rule

$$\frac{df_1}{dx} = \frac{df_1}{df_2} \frac{df_2}{df_3} \cdots \frac{df_n}{dx}$$

Reminder: Jacobians

Consider differentiable $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$, its Jacobian contains its first-order partial derivatives $(J_f)_{ij} = \frac{\partial f_i}{\partial x_j}$:

$$J_f = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla^T f_1 \\ \vdots \\ \nabla^T f_m \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Example: linear layer

W is a $m \times n$ matrix

If $f(\mathbf{x}) = W\mathbf{x}$

$f_i = W_i\mathbf{x} = \sum_j W_{ij}x_j$ where W_i is i th row of W

$(J_f)_{ij} = \frac{\partial f_i}{\partial x_j} = W_{ij}$ so $J_f = W$

Generalized multivariate chain rule

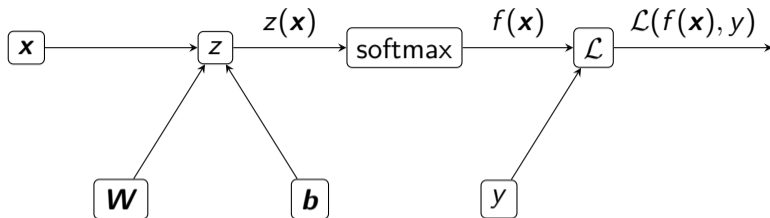
Consider differentiable $f : \mathbf{R}^m \rightarrow \mathbf{R}^k$, $g : \mathbf{R}^n \rightarrow \mathbf{R}^m$ and $\mathbf{a} \in \mathbf{R}^n$.

$$D_{\mathbf{a}}(f \circ g) = D_{g(\mathbf{a})}f \circ D_{\mathbf{a}}g$$

So the Jacobians verify:

$$J_{f \circ g} = (J_f \circ g)J_g$$

Computation graph



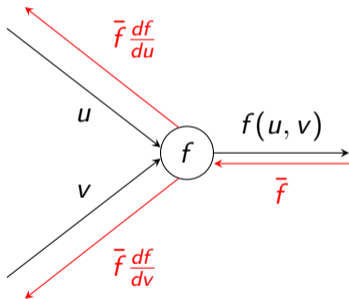
$$\begin{aligned} \frac{d\mathcal{L}}{db_1} &= \frac{d\mathcal{L}}{df} \frac{df}{db_1} = \frac{d\mathcal{L}}{df} \left(\frac{df}{dz} \frac{dz}{db_1} \right) \quad (\text{forward}) \\ &= \frac{d\mathcal{L}}{dz} \frac{dz}{db_1} = \left(\frac{d\mathcal{L}}{df} \frac{df}{dz} \right) \frac{dz}{db_1} \quad (\text{backward}) \end{aligned}$$

Properly written: $J_{\mathcal{L} \circ \text{softmax} \circ z} = (J_{\mathcal{L}} \circ f)(J_{\text{softmax}} \circ z)J_z$

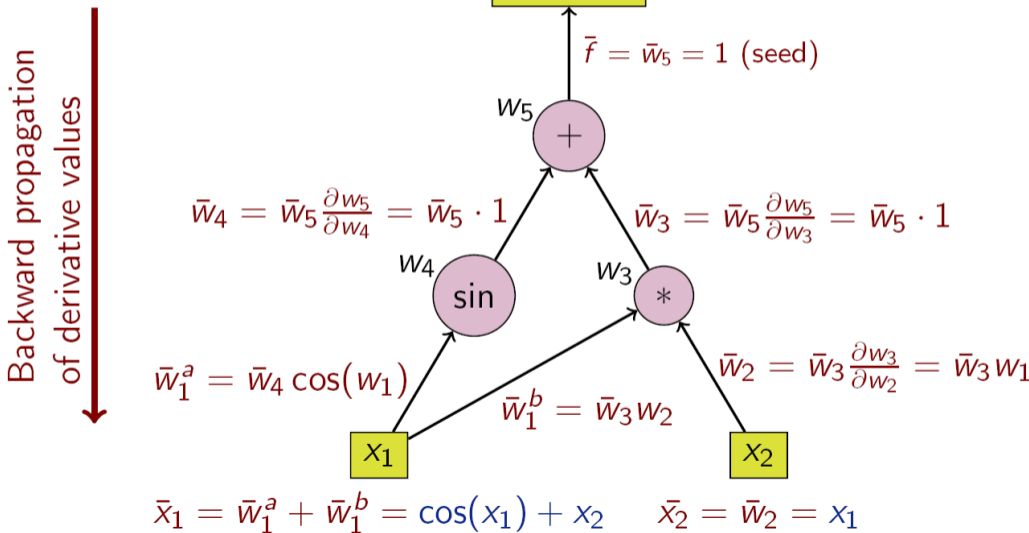
Given that $\mathbf{x} \in \mathbf{R}^d$, $z(\mathbf{x})$, $f(\mathbf{x}) \in \mathbf{R}^{d_2}$, $\mathcal{L}(f(\mathbf{x}), y) \in \mathbf{R}$,
which order is better?

Reverse accumulation (in \mathbf{R})

Let us note the adjoint $\bar{f} \triangleq \frac{d\mathcal{L}}{df}$.

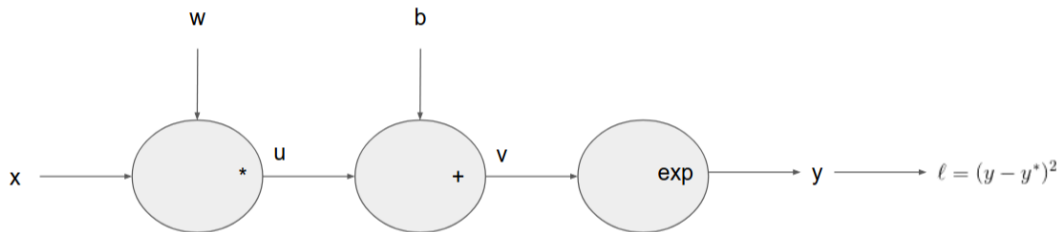


A complete example (Wikipedia)



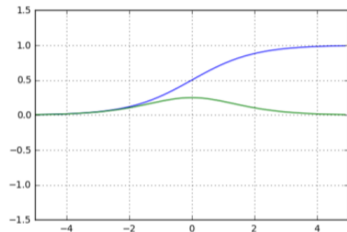
Please compute gradients

FORWARD PASS



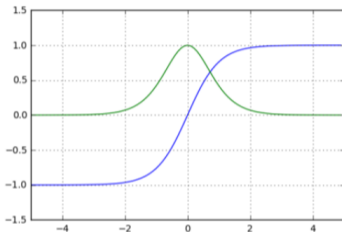
$$\begin{aligned}u &= wx \\v &= u + b \\y &= e^v \\ \ell &= (y - y^*)^2\end{aligned}$$

Examples of link functions (Ollion & Grisel)



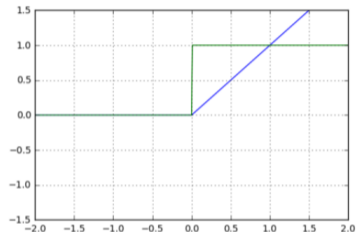
$$\text{sigm}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{sigm}'(x) = \text{sigm}(x)(1 - \text{sigm}(x))$$



$$\text{tanh}(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\text{tanh}'(x) = 1 - \text{tanh}(x)^2$$



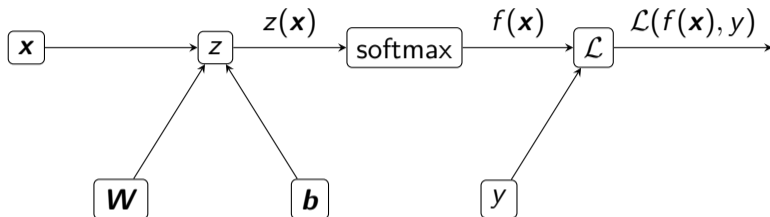
$$\text{relu}(x) = \max(0, x)$$

$$\text{relu}'(x) = 1_{x>0}$$

More interesting link functions

	Binary	Multiclass
f	softplus : $x \mapsto \log(1 + \exp(x))$	logsumexp ⁺ : $\mathbf{x} \mapsto \log(1 + \sum_c \exp(x_c))$
f'	sigmoid : $x \mapsto 1/(1 + \exp(-x))$	softmax : $\mathbf{x} \mapsto \exp(\mathbf{x}) / \sum_c \exp(x_c)$
f''	$x \mapsto \text{sigmoid}(x)(1 - \text{sigmoid}(x))$? $\mathbf{x} \mapsto \text{diag}(s(\mathbf{x})) - s(\mathbf{x})s(\mathbf{x})^T$

Computation graph: classifier



Here we define \mathcal{L} as cross-entropy:

$$\mathcal{L}(f(\mathbf{x}), y) = - \sum_{c=1}^K \mathbf{1}_{y=c} \log f(\mathbf{x})_c = - \log f(\mathbf{x})_y$$

Compute $\frac{d\mathcal{L}}{dz_c}$