# MNIST classification in numpy

We will consider classifying MNST with a simple NN of the form

$$\hat{z}^i = W^i z^i + b^i \quad , \quad i=1,\dots, k-1$$
$$z^i = \sigma(\hat{z}^i) \quad\quad i=2,\dots, k-1$$

where $z_1 = x$, where $x$ is an input to the network, $k$ is # of layers & we let

$$h(x,\theta) := \hat{z}^k \quad\quad \leftarrow \text{ I think we want } h(x,\theta) = SM(\hat{z}^k)$$

Here, $\theta = ((W^i)_i, (b^i)_i)$ represent all params & $x \mapsto h(\cdot, \theta)$
is the input/output map of the network.

Let $N$ = # of classes.

We will classify "goodness" of $\theta$ using the cross entropy loss. Given an input/label pair $(x,y)$, $x \in \mathbb{R}^n$, $y \in \{0, \dots, N\}$, the cross entropy loss is given by

$$CE(h(x,\theta), y) = -\sum_{i=1}^{N} (\mathbb{1}_y)_i \ln(h(x,\theta)_i)$$

indicates only
$$=\begin{cases} 1 & \text{if } i = y \\ 0 & \text{else} \end{cases}$$

$$= -\ln(\underbrace{h(x,\theta)_y}_{\text{}})$$

the $y$-th coordinate of $h$.

We want to optimize $\theta$ by using gradient descent on the empirical risk.

Let $D = ((x_i, y_i))_i$ denote the MNIST training set. The empirical risk is given by

$$L(\theta) = \sum_{(x,y) \in D} \underbrace{CE(h(x,\theta), y)}_{=: \, l(\theta, (x,y))}$$

The gradient of $L$ is given by

gradient w.r.t. $\theta$, of course.

$$\nabla L(\theta) = \sum_{(x,y) \in D} \overbrace{\nabla CE(h(x,\theta), y)}^{}$$

$$\underbrace{\qquad}_{(*) \;=: \; l(\theta, (x,y))}$$

Computing $(*)$ for each $(x,y)$ is costly. If we choose a random $(\tilde{x}, \tilde{y}) \in D$ (uniform random in $D$) ~~than~~ & let $g = \nabla CE(h(\tilde{x}, \theta), \tilde{y})$, then

$$\mathbb{E}\, g = \frac{1}{|D|} \nabla L(\theta).$$ We would like to use SGD,

$$\theta_{t+1} = \theta - \alpha g, \qquad g = \nabla CE(\cdot), \; (\tilde{x}, \tilde{y}) \sim U(D)$$

However, drawing a random sample from $D$ each iteration of SGD is costly. Instead, we will Shuffle $D$ randomly & simply iterate through it:

$$\theta_{t+1} = \theta - \alpha g_i, \qquad g_i = ~~\nabla CE(\cdot)~~ \nabla l(\theta, (x_i, y_i)).$$

We need to compute $\nabla_\theta CE(h(x,\theta), y)$. To make this unambiguous, we'll use the following notation.

$$\frac{\partial}{\partial \theta} f(g(\theta))$$

means the derivative of $h = f \circ g$ taken w.r.t. $\theta$.

$$D_z f(z)$$

Means ~~Jacobian of $f$~~ Derivative of $f$ in usual sense

Basically, $\frac{\partial}{\partial \theta}$ means were going to have to use the chain rule to evaluate.

$D_z$ means were not taking the derivative of a composition.

They really, mean the same thing, but I'm hoping this give context. We also still use $\frac{\partial}{\partial x_i}$ for standard problems. (... Maybe use $\partial_\theta$ to suggest chain rule?) Want to compute

$$\frac{\partial}{\partial \theta} CE\left(h(x,\theta), y\right) = \frac{\partial}{\partial \theta} - \ln [h(x,\theta)]_y$$

We'll need to compute ~~the follow derivates~~ derivates of the following funcs

- ~~$\frac{\partial}{\partial}$~~
- $\theta$

- $\sigma$
- linr
- Softmax
- CE.

In all subsequent derivations, we use the following

for $f: \mathbb{R}^n \to \mathbb{R}^m$,

$$D_z f(z)\Big|_{z=w} = \begin{bmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & \frac{\partial f_1}{\partial z_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial z_1} & \cdots & \frac{\partial f_m}{\partial z_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

where each partial is evaluated at the pt $w$.

Sometimes, we use the shorthand $D_z f(w)$,
$\uparrow$ base point

For $\underset{\text{composites}}{f \circ g \circ h}$, the chain rule is given by

$$\frac{\partial}{\partial z} f(g(h(z))) = $$

$$= D_w f(w)\Big|_{w = g(h(z))} \; D_w g(w)\Big|_{w = h(z)} \; D_w h(w)\Big|_{h=z}$$

In shorthand:

$$= D f(g(h(z))) \; D g(h(z)) \; D h(z)$$

To Do: prove chain rule sometime?

## Derivation of $\sigma$

In an abuse of notation, we use $\sigma$ to mean the scalar sigmoid function applied elementwise. To clarify, here (and only here) we'll say,

$$\sigma: \mathbb{R} \to (0,1)$$

is

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

and

$$\vec{\sigma}: \mathbb{R}^n \to \mathbb{R}^n$$

is given elementwise by

$$\vec{\sigma}_i(z) = \sigma(z_i).$$

Note

$$\frac{\partial \vec{\sigma}_i(z)}{\partial z_j} = \begin{cases} \sigma'(z_i) & j=i \\ 0 & \text{else} \end{cases}$$

Hence

$$D\vec{\sigma}(z) = \text{diag}(\sigma'(z_1) \dots, \sigma'(z_n)).$$

To be complete (so we can code it) lets compute $\sigma'$.

$$\frac{d}{dz}\sigma(z) = \frac{d}{dz}\frac{1}{1+e^{-z}} = (1+e^{-z})^{-1} = -\frac{1}{(1+e^{-z})^2}(-e^{-z}) = \frac{e^{z}}{(1+e^{-z})^2}$$

use chain rule.

$$= \frac{e^{z}}{1+e^{-z}}\frac{1}{1+e^{-z}}$$

(*)
Note: $1 - \frac{e^{z}}{1+e^{-z}} = \frac{1}{1+e^{z}} = \sigma(z)$

$$\underbrace{\frac{e^z}{1+e^{-z}}}_{(*) = 1-\sigma(z)} \underbrace{\frac{1}{1+e^{-z}}}_{=\sigma(z)}.$$

Hence $\frac{e^z}{1+e^{-z}} = 1-\sigma(z)$

We should already (stably) from the forward pass. So, now we have a complete $D\vec{\sigma}$.

Derivation of Softmax

$$SM(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Note:

① $\dfrac{\partial SM(x)_i}{\partial x_l} = \dfrac{-e^{x_i}}{(\sum e^{x_j})^2} e^{x_l}$    $l \neq i$

② $\dfrac{\partial SM(x)_i}{\partial x_i} = \dfrac{e^{x_i}}{\sum_j e^{x_j}} + -e^{x_i} \left(\sum_j e^{x_j}\right)^2 e^{x_i}$

Note that ① corresponds to diagonals of $DSM$

Let

$C := \left(\sum_j e^{x_j}\right)^2$

$M := \underbrace{\exp(x) \; \exp(x)^T}_{\text{outer product between the vectors}}$

$v = \exp(x) \left(\sum e^{x_j}\right)^{-1}$

Then

$D\, SM(x) = \cancel{c(M+v)}$

$\quad c\left(M + \text{diag}(v)\right)$

## Derivative of cross entropy

$$CE(p,q) = -\sum q_i \ln p_i$$

We want the derivative w.r.t. $p$.

$$D_p CE(p,q) = -\sum_i \frac{q_i}{p_i} \quad \longleftarrow \quad \frac{d}{dx}\ln(x) = \frac{1}{x}$$

## Derivative of Linear function

Let $f(x) = W$ $\quad f(x,w,b) = Wx + b$

We must compute $D_x f$, $D_w f$, and $D_b f$.

### 1. $D_x f$

Note that $f(x)_i = W_i^T x + b_i$

$\hat{\imath}$ th row of $W$

$$\frac{\partial f(x)_i}{\partial x_j} = W_{ij}$$

$$\Rightarrow D_x f = W$$

### 2. $D_w f$

We will treat $W$ as a big column vector, So $D_w f$ takes the form

Let $m \times n$
$W \in \mathbb{R}^{m \times n}$

$$D_w f = \begin{bmatrix} \frac{\partial f_1}{\partial w_{11}} & \cdots & \frac{\partial f_1}{\partial w_{1n}} & \frac{\partial f_1}{\partial w_{m1}} & \cdots & \frac{\partial f_1}{\partial w_{mn}} \\ \vdots & & & & & \\ \frac{\partial f_m}{\partial w_{11}} & \cdots & & & & \frac{\partial f_1}{\partial w_{mn}} \end{bmatrix}$$

Note that $\quad f(w)_i \overset{\leftarrow \ \text{abuse of notat. omit nonactive args.}}{=} w_i^T x + b_i$.

$$\frac{\partial f(w)_i}{\partial w_{jk}} = \begin{cases} 0 & j \neq i \\ x_k & \text{else} \end{cases}$$

Hence

$$\left[ \frac{\partial f(w)_i}{\partial w_{j1}} \quad \cdots \quad \frac{\partial f(w)_i}{\partial w_{jn}} \right] = \begin{cases} 0 & j \neq i \\ x^T & j = i \end{cases}$$

Hence

$$D_w f = \begin{bmatrix} x^T & 0 & \cdots & 0 \\ 0 & x^T & \cdots & 0 \\ 0 & 0 & \cdots 0 & x^T \end{bmatrix}$$

3. $D_b f$

$$f(b)_i = w_i^T x + b_i$$

$$\frac{\partial f(b)_i}{\partial b_j} = \begin{cases} 1 & i = j \\ 0 & \text{else} \end{cases} \implies D_b f = I_m$$

Finally, how do we compute derivatives w.r.t. $b^i$ and $w^i$.

Recall the previous discussion of reverse-mode auto diff.

Let $J(\cdot)$ denote our cost function (w.r.t. some fixed data).

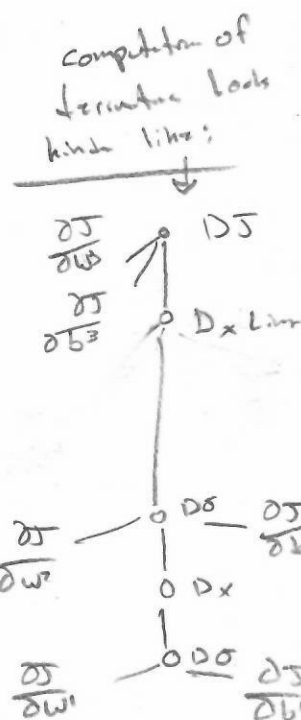For concreteness, Suppose we have a network w/ 3 layers

$$h(x;\theta) = W^3 \sigma \left( W^2 \sigma \left( \underbrace{\underbrace{\underbrace{W^1 x + b^1}_{\hat{z}^1}}_{\hat{z}^2}}_{} \right) + b^2 \right) + b^3$$

$$\underbrace{\phantom{W^3 \sigma ( W^2 \sigma ( W^1 x + b^1 ) + b^2 ) + b^3}}_{\hat{z}^3}$$

We have

$$\frac{\partial J}{\partial W^3} = \underbrace{DJ(\hat{z}^3)}_{} \; D_w \text{Linear} \left( \sigma(\hat{z}^2) \right)$$

$$\frac{\partial J}{\partial W^2} = \underbrace{DJ(\hat{z}^3) \; D_x \text{Linear} \left( \sigma(\hat{z}^2) \right) \; D\sigma(\hat{z}^2) \; D_w \text{Linear} \left( \sigma(\hat{z}^1) \right)}_{=: \; Q, \; \text{for lack of a better name}}$$

$$\frac{\partial J}{\partial W^2} = Q \; D_x \text{Linear} \left( \sigma(\hat{z}^1) \right) \; D\sigma(\hat{z}^1) \; D_w \text{Linear}(x)$$

computation of derivative looks kinda like:

$$\frac{\partial J}{\partial W^3} \nearrow \quad DJ$$

$$\frac{\partial J}{\partial b^3} \quad D_x \text{Linear}$$

$$\frac{\partial J}{\partial W^2} \quad \frac{}{} \quad D\sigma \quad \frac{\partial J}{\partial b^2}$$
$$D_x$$

$$\frac{\partial J}{\partial W^1} \quad DD\sigma \quad \frac{\partial J}{\partial b^1}$$

Note how as you work backwards, computing the derivatives w.r.t. parameters, you can reuse previous derivative computations. Also, the base points for derivatives, $\hat{z}^1$, $\sigma(\hat{z}^1)$, etc. are saved from the forward pass.

# Appendix:   log sum exp   &   Stable Softmax

Suppose you wish to compute the softmax of $X \in \mathbb{R}^n$

$$Sm(x)_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

If any elements of $X$ are very large, you will run into overflow when computing $e^{x_i}$ or $\sum_j e^{x_j}$. To fix this, you can use the following trick:

$$\frac{e^{x_i}}{\sum_j e^{x_j}} = \exp\left( \ln\left( \frac{e^{x_i}}{\sum e^{x_j}} \right) \right)$$

$$= \exp\left( x_i - \underbrace{\ln \sum_j e^{x_j}}_{(*)} \right)$$

$(*)$ can be computed stably by recentering as follows. For any $c \in \mathbb{R}$ we have

$$\ln \sum_j e^{x_j} = \ln e^c \sum e^{x_j - c} = c + \ln \sum e^{x_j - c}$$

If you choose $c = \max\{x_1, \ldots, x_n\}$, this tends to work well. I don't think max vs min vs mean makes a huge difference. Mostly you want to get what goes into the exponent centered close to zero. I think...

observations:

Final form for numerically stable softmax is

$$Sm(x)_i = \exp\left(x_i - c - \ln \sum_j e^{x_j - c}\right), \quad c = \max\{x_1, \ldots, x_n\}$$

Note 3 things:

- Inside the sum, $e^{x_i - c}$ should be less prone to overflow/underflow because $x_i - c$ should hopefully be closer to zero. (Maybe not if $\min x_i \ll 0$ & $\max x_i \gg 0$, but that seems like it shouldn't arise oftn)

- In the outer exp, we have $\exp(\underbrace{x_i - c} + stuff)$
  This should be more numerically stable for the same reason

- In the log we have $\ln \sum_j e^{x_j - c}$.
  We would like the argument to $\ln(\cdot)$ to be $\geq 0$, to avoid the singularity and instability in log near 0. This is guaranteed by our choice of $c$, since for $x_k \in \text{argmax}\{x_1, \ldots, x_n\}$ $e^{x_k - c} = 1$ & $e^{x_j - c} > 0$ for all other $j$.