Topics

- Review of gradient descent
- Loss functions
- Multivariable Chain Rule
- Basic neural networks (preview)
- Backpropagation

These notes are loosely based on ILA Ch 9.2. Reviewing CalcBLUE 2 Ch. 5 on the chain rule is recommended.

## Backpropagation — Motivation

Last class, we studied the unconstrained optimization

$$\text{minimize } f(\underline{x}) \qquad (P)$$

over $\underline{x} \in \mathbb{R}^n$, where we look for the $\underline{x} \in \mathbb{R}^n$ that makes the value of the cost function $f : \mathbb{R}^n \to \mathbb{R}$ as small as possible. We saw that one way to find either a local or global minimum $\underline{x}^*$ is gradient descent. Starting at an initial guess $\underline{x}^{(0)}$, we iteratively update our guess via

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - s \nabla f(\underline{x}^{(k)}), \qquad k = 0, 1, 2, \dots \qquad (GD)$$

where $\nabla f(\underline{x}^{(k)}) \in \mathbb{R}^n$ is the gradient of $f$ evaluated at the current guess, and $s > 0$ is a step size chosen large enough to make progress towards $\underline{x}^*$, but not so big as to overshoot.

Today, we'll focus our attention on optimization problems $(P)$ for which the cost function takes the following special form

$$f(\underline{x}) = \sum_{i=1}^{N} f_i(\underline{x}), \qquad (S)$$

i.e, cost functions $f$ that decompose into a sum of $N$ "sub-costs" $f_i$. Problems with cost functions of the form $(S)$ are particularly common in machine learning.

For example, a typical problem setup in machine learning is as follows (we saw an example of this when we studied least squares for data-fitting). We are given a set training data $(\underline{z}_i, \underline{y}_i), i = 1, \dots, N$, composed of "inputs" $\underline{z}_i \in \mathbb{R}^p$ "outputs" $\underline{y}_i \in \mathbb{R}^m$. Our goal is to find a set of weights $\underline{x} \in \mathbb{R}^n$ which parameterize a model such that $m(\underline{z}_i; \underline{x}) \approx \underline{y}_i$ on our training data. A common way of doing this is to minimize a loss function of the form

$$\text{loss}\left((\underline{z}_i, \underline{y}_i); \underline{x}\right) = \frac{1}{N} \sum_{i=1}^{N} \ell\left(m(\underline{z}_i; \underline{x}) - \underline{y}_i\right), \qquad (L)$$

Where each term $\ell(m(\underline{z}_i; \underline{x}) - \underline{y}_i)$ is a term penalizing the difference between our model prediction $m(\underline{z}_i; \underline{x})$ on input $\underline{z}_i$ and the observed output $\underline{y}_i$. In this setting, the loss function $(L)$ takes the form $(S)$, with $f_i = \frac{1}{N} \ell(m(\underline{z}_i; \underline{x}) - \underline{y}_i)$ the error between our prediction $\hat{\underline{y}}_i = m(\underline{z}_i; \underline{x})$ and the true output $\underline{y}_i$.

A common
choice for the "sub-loss" function is $\ell(\underline{e}) = \|\underline{e}\|^2$, leading to a least-squares regression problem, but note that most other choices of loss function are compatible with the following discussion.

Now, suppose that we want to implement gradient descent (GD) on the loss function (L). Our first step is to compute the gradient $\nabla_{\underline{x}} \text{loss}((\underline{z}_i, y_i); \underline{x})$. Because of the sum structure of (L), we have that:

$$\nabla_{\underline{x}} \text{loss}((\underline{z}_i, y_i); \underline{x}) = \frac{1}{N} \sum_{i=1}^{N} \nabla_{\underline{x}} \ell(m(\underline{z}_i; \underline{x}) - y_i) ,$$

i.e., the gradient of the loss function is the sum of the gradients of the "sub-losses" on each of the $i = 1, ..., N$ data points.

Our task now is therefore to compute the gradient $\nabla_{\underline{x}} \ell(m(\underline{z}_i; \underline{x}) - y_i)$. This requires the multivariate chain rule, as $f_i(\underline{x}) = \ell(m(\underline{z}_i; \underline{x}) - y_i)$ is a composition of the functions $\ell(e)$, $e = w - y_i$, and $w = m(\underline{z}_i; \underline{x})$.

## The Multivariate Chain Rule ((CalcBLUE 2 Ch. 5)

We begin with a reminder of the chain rule for scalar functions. Let $f : \mathbb{R} \to \mathbb{R}$ and $g : \mathbb{R} \to \mathbb{R}$ be differentiable functions. Then for $h(x) = g(f(x))$, we have that

$$h'(x) = g'(f(x)) f'(x). \qquad (C1)$$

If we define $g = g(f)$, and $f = f(x)$, then we can rewrite (C1) as $\frac{dh}{dx} = \frac{dg}{df} \cdot \frac{df}{dx}$. This is a useful way of writing things as we can "cancel" $df$ on the RHS to check that our formula is correct.

WARNING: $\frac{dh}{dx} = \frac{dg}{df} \cdot \frac{df}{dx}$ is shorthand for $\frac{dh}{dx}(x) = \frac{dg}{df}(f(x)) \frac{df}{dx}(x)$. The evaluation points matter!

Generalizing slightly, suppose now that $f : \mathbb{R}^n \to \mathbb{R}$ maps a vector $\underline{x} \in \mathbb{R}^n$ to $f(\underline{x}) \in \mathbb{R}$. Then for $h(\underline{x}) = g(f(\underline{x}))$, we have:

$$\nabla_{\underline{x}} h(\underline{x}) = g'(f(\underline{x})) \nabla f(\underline{x}), \qquad (C2)$$

which we see is a natural generalization of equation (C1). It will be convenient for us later to define $\frac{df}{d\underline{x}} = \nabla f(\underline{x})^T$ and $\frac{dh}{d\underline{x}} = \nabla_{\underline{x}} h(\underline{x})^T$. Again defining $g = g(f)$ and $f = f(\underline{x})$, we can rewrite (C2) as $\frac{dh}{d\underline{x}} = \frac{dg}{df} \cdot \frac{df}{d\underline{x}}$, which looks exactly the same as before!

WARNING: $\frac{dh}{d\underline{x}} = \frac{dg}{df} \cdot \frac{df}{d\underline{x}}$ is shorthand for $\frac{dh}{d\underline{x}}(\underline{x}) = \frac{dg}{df}(f(\underline{x})) \frac{df}{d\underline{x}}(\underline{x})$. The evaluation points matter!

Now, let's apply these ideas to computing the gradient of $h(\underline{x}) = \ell(m(\underline{z}_i; \underline{x}) - y_i)$, where we'll assume for now that $m(\underline{z}_i; \underline{x}), y_i \in \mathbb{R}$. Applying (C2), we get

$$\nabla_{\underline{x}} h(\underline{x}) = \ell'(m(\underline{z}_i; \underline{x}) - y_i) \cdot \nabla_{\underline{x}} (m(\underline{z}_i; \underline{x}) - y_i) = \ell'(m(\underline{z}_i; \underline{x}) - y_i) \cdot \nabla_{\underline{x}} m(\underline{z}_i; \underline{x})$$

where we use that $\partial_x y_i = 0$ (since its a constant). Without knowing more about the functions $\ell$ and $m$, this is all we can say.

Example: Suppose $\ell(e) = \frac{1}{2}e^2$ and $m(\underline{z}_i; \underline{x}) = \underline{x}^T \underline{z}_i$. Then

$$\ell(m(\underline{z}_i; \underline{x}) - y_i) = \frac{1}{2}(\underline{x}^T \underline{z}_i - y_i)^2 \text{ and } \partial_{\underline{x}} \ell(m(\underline{z}_i; \underline{x}) - y_i) = \underbrace{(\underline{x}^T \underline{z}_i - y_i)}_{\ell'(m - y_i)} \cdot \underbrace{\underline{z}_i}_{\partial_{\underline{x}} m}$$

Next class we will have a brief introduction to deep learning. In deep learning, the function $m(\underline{z}_i; \underline{x})$ is often parameterized as a chain of function compositions:

$$m(\underline{z}_i; \underline{x}) = \underline{m}_L(\underline{m}_{L-1}(\cdots(\underline{m}_2(\underline{m}_1(\underline{z}_i)))\cdots)) \qquad (DN)$$

$$= \underline{m}_L \circ \underline{m}_{L-1} \circ \cdots \circ \underline{m}_2 \circ \underline{m}_1 (\underline{z}_i).$$

A more suggestive way of writing this parameterization (that also highlights the dependence on $\underline{x}$) is

$$\underline{o}_0 = \underline{z}_i$$
$$\underline{o}_1 = \underline{m}_1(\underline{o}_0; \underline{x}_1) \qquad \underline{o}_1 \in \mathbb{R}^{P_1}, \ \underline{o}_0 \in \mathbb{R}^{P_0} \qquad (DNN)$$
$$\underline{o}_2 = \underline{m}_2(\underline{o}_1; \underline{x}_2) \qquad \underline{o}_2 \in \mathbb{R}^{P_2}, \ \underline{o}_1 \in \mathbb{R}^{P_1}$$
$$\vdots$$
$$\underline{o}_L = \underline{m}_L(\underline{o}_{L-1}; \underline{x}_L) \quad \underline{o}_L \in \mathbb{R}^{P_L}, \ \underline{o}_{L-1} \in \mathbb{R}^{P_{L-1}}$$

Here the model parameters $\underline{x} = (\underline{x}_1, \ldots, \underline{x}_L)$ we split across the layers $1, \ldots, L$. The intermediate outputs $\underline{o}_i$ can be of different dimension, as can the layer parameters $\underline{x}_i$. Writing (DN) as (DNN) highlights why these functions are called deep neural networks as the number of layers $L$ increases. Our goal is then to compute $\partial_{\underline{x}} \ell(m(\underline{z}_i; \underline{x}) - \underline{y}_i)$ for $m$ of the form (DN), and where $m, y_i \in \mathbb{R}^{P_L}$ are now also possibly vector-valued. To do this, we need the fully general multivariate chain rule.

For $h(\underline{x}) = g(f(\underline{x}))$, with vector-valued $\underline{f}: \mathbb{R}^n \to \mathbb{R}^p$ and $\underline{g}: \mathbb{R}^p \to \mathbb{R}^m$, we need to define the Jacobian matrices for $\underline{f}$ and $\underline{g}$:

$$\frac{d\underline{f}}{d\underline{x}} = \begin{bmatrix} \frac{df_1}{d\underline{x}} \\ \vdots \\ \frac{df_p}{d\underline{x}} \end{bmatrix} = \begin{bmatrix} \partial f_1/\partial x_1 & \cdots & \partial f_1/\partial x_n \\ \vdots & \ddots & \\ \partial f_p/\partial x_1 & \cdots & \partial f_p/\partial x_n \end{bmatrix} \text{ and } \frac{d\underline{g}}{d\underline{f}} = \begin{bmatrix} \frac{dg_1}{d\underline{f}} \\ \vdots \\ \frac{dg_m}{d\underline{f}} \end{bmatrix} = \begin{bmatrix} \partial g_1/\partial f_1 & \cdots & \partial g_1/\partial f_p \\ \vdots & \ddots & \\ \partial g_m/\partial f_1 & \cdots & \partial g_m/\partial f_p \end{bmatrix}$$

$$(J)$$

as the $p \times n$ and $m \times p$ matrices of partial derivatives, respectively.

We'll use our sure intuition of "cancelling" to derive the expression:

$$\frac{dh}{dx} = \frac{dg}{df} \cdot \frac{df}{dx} \qquad (C3)$$

Note that (C3) is defined by a matrix-matrix multiplication of an $m \times p$ and $p \times n$ matrix, meaning $\frac{dh}{dx} \in \mathbb{R}^{m \times n}$. The claim is that $(i,j)^{th}$ entry of $\frac{dh}{dx}$ is the rate of change of $h_i(x) = g_i(f(x))$ with respect to $x_j$. From (5) and (C3), we have

$$\left(\frac{dh}{dx}\right)_{ij} = \underbrace{\frac{dg_i}{df}}_{\substack{i^{th} \text{ row} \\ \text{of } \frac{dg}{df}}} \cdot \underbrace{\begin{bmatrix} \partial f_1/\partial x_j \\ \vdots \\ \partial f_p/\partial x_j \end{bmatrix}}_{\substack{j^{th} \text{ column} \\ \text{of } \frac{df}{dx}}} = \frac{dg_i}{df_1} \cdot \frac{\partial f_1}{\partial x_j} + \cdots + \frac{\partial g_i}{\partial f_p} \frac{\partial f_p}{\partial x_j},$$

which is precisely the expression we were looking for. The "cancellation rule" tells us each term in the sum is computing the part of $\frac{\partial g_i}{\partial x_j}$ in the "$f_i$" channel.

We can apply this formula recursively to our function class (DN) to obtain the formula:

$$\frac{dm}{dx} = \frac{dm_L}{dm_{L-1}} \cdot \frac{\partial m_{L-1}}{\partial m_{L-2}} \cdots \frac{dm_2}{dm_2} \cdot \frac{dm_1}{dx} \qquad (MC),$$

which is a fully general matrix chain rule. We'll use (MC) next to explore the key ideas behind back propagation, which has been a key technical enabler of contemporary deep learning.

## Backpropagation

We are going to work out how to efficiently compute the gradient of

$$\ell(m(z_i; x) - y_i)$$

when $m$ takes the form $m$ (DNN). We'll further assume, as is often the case in deep learning, that each layer function $m_i$ takes the following form:

$$m_i(o_{i-1}; x_i) = \sigma\left(x_i \begin{bmatrix} o_i \\ 1 \end{bmatrix}\right)$$

where $X_i$ is a $\mathbb{R}^{p_i \times (n_i+1)}$ matrix with entries given by $x_i \in \mathbb{R}^{p_i(n_i+1)}$, and $\sigma$ is a pointwise nonlinearity $\sigma(x) = (\sigma(x_1), \ldots, \sigma(x_n))$ called an activation function (more on these next class)

Applying our matrix chain rule to $\ell(\underline{m}(\underline{x}) - \underline{y}_i)$   (we won't write $z_i$ to save space),
we get the expression

$$\frac{d\ell}{d\underline{x}} = \frac{\partial \ell}{\partial \underline{m}} \cdot \frac{d\underline{m}}{d\underline{x}} = \frac{\partial \ell}{\partial \underline{m}_L} \frac{d\underline{m}_L}{d\underline{m}_{L-1}} \cdots \frac{d\underline{m}_2}{d\underline{m}_1} \cdot \frac{d\underline{m}_1}{d\underline{x}}.$$

Here, $\frac{\partial \ell}{\partial \underline{m}}$ is a $P_L$ dimensional row vector, and $\frac{d\underline{m}_i}{d\underline{m}_{i-1}}$ is a $P_i \times P_{i-1}$ matrix.

In modern applications, the layer dimensions, also called layer widths, $P_i$ can be very large (on the order of 100s of thousands or even millions), meaning the $\frac{d\underline{m}_i}{d\underline{m}_{i-1}}$ matrices are **very very** large! Too large to store in memory actually.

Fortunately, since $\frac{\partial \ell}{\partial \underline{m}}$ is a row vector, we can build $\frac{d\ell}{d\underline{x}}$ by sequentially computing inner products. For example, if $\frac{d\underline{m}_L}{d\underline{m}_{L-1}} = [\underline{a}_1 \cdots \underline{a}_{P_{L-2}}]$

$$\frac{\partial \ell}{\partial \underline{m}_L} \cdot \frac{d\underline{m}_L}{d\underline{m}_{L-2}} = \underbrace{[\longleftarrow \frac{\partial \ell}{\partial \underline{m}_L} \longrightarrow]}_{1 \times P_L} \underbrace{\begin{bmatrix} \underline{a}_1 & \cdots & \underline{a}_{P_{L-2}} \end{bmatrix}}_{P_L \times P_{L-2}}$$

$$= [\frac{\partial \ell}{\partial \underline{m}_L} \underline{a}_1 \cdots \frac{\partial \ell}{\partial \underline{m}_L} \underline{a}_{P_{L-2}}],$$

meaning we only ever need to store $\frac{\partial \ell}{\partial \underline{m}_L}$ and $\underline{a}_i$ in memory at any given time, which is only $2P_L$ numbers, as opposed to $P_L \times P_{L-2}$ #s! Then once we've computed $\frac{d\ell}{d\underline{m}_2} \frac{d\underline{m}_L}{d\underline{m}_{L-1}}$, which is now a $P_{L-2}$ dim. row vector, we can continue our way down the chain.

What's left to do is compute the partial derivatives! Let's break down $\frac{d\ell}{d\underline{x}}$ into partial derivatives with respect to a layer's parameters $\underline{x}_i$. For layer $L$, we have:

$$\frac{\partial \ell}{\partial \underline{x}_L} = \frac{\partial \ell}{\partial \underline{m}_L} \cdot \frac{d\underline{m}_L}{d\underline{x}_L} + \frac{\partial \ell}{\partial \underline{m}_L} \cdot \frac{d\underline{m}_L}{d\underline{m}_{L-1}} \frac{d\underline{m}_{L-1}}{d\underline{x}_L} = \frac{\partial \ell}{\partial \underline{m}_L} \cdot \frac{d\underline{m}_L}{d\underline{x}_L}$$

Since $\underline{x}_L$ appears in the last layer, it shows up right away in the first term above, which is the derivative of $\underline{m}_L(\underline{m}_{L-1}; \underline{x}_L)$ with respect to $\underline{x}_L$ (the 2nd argument). The second term

$$\frac{\partial \ell}{\partial \underline{m}_L} \cdot \frac{d\underline{m}_L}{d\underline{m}_{L-1}} \frac{d\underline{m}_{L-1}}{d\underline{x}_L} = 0$$

which measures how $\underline{m}_L$ changes with respect to changes in $\underline{m}_{L-1}$ caused by changes in $\underline{x}_L$ is zero because $\underline{m}_{L-1}$ does not depend on $\underline{x}_L$ at all! This is a key observation in the back-propagation algorithm!

Let's proceed to computing the derivative with respect to the parameters $\underline{s}_{L-2}$:

$$\frac{\partial l}{\partial \underline{s}_{L-2}} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-1}} \cdot \left( \frac{\partial \underline{m}_{L-2}}{\partial \underline{s}_{L-2}} + \frac{\partial \underline{m}_{L-2}}{\partial \underline{m}_{L-2}} \cdot \frac{\partial \underline{m}_{L-2}}{\partial \underline{s}_{L-2}} \nearrow 0 \right)$$

$$= \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-2}} \cdot \frac{\partial \underline{m}_{L-2}}{\partial \underline{s}_{L-2}}$$

We see again that we can "stop" once we hit the layer that depends explicitly on $\underline{s}_{L-2}$. In general, we have:

$$\frac{\partial l}{\partial \underline{s}_L} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{s}_L}$$

$$\frac{\partial l}{\partial \underline{s}_{L-1}} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-1}} \cdot \frac{\partial \underline{m}_{L-1}}{\partial \underline{s}_{L-1}} \qquad \left( \frac{\partial l}{\partial \underline{m}_{L-1}} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-2}} \right)$$

$$\frac{\partial l}{\partial \underline{s}_{L-2}} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-1}} \cdot \frac{\partial \underline{m}_{L-1}}{\partial \underline{m}_{L-2}} \cdot \frac{\partial \underline{m}_{L-2}}{\partial \underline{s}_{L-2}} \qquad \left( \frac{\partial l}{\partial \underline{m}_{L-2}} = \frac{\partial l}{\partial \underline{m}_{L-1}} \cdot \frac{\partial \underline{m}_{L-1}}{\partial \underline{m}_{L-2}} \right)$$

$$\vdots$$

$$\frac{\partial l}{\partial \underline{s}_j} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-1}} \cdot \frac{\partial \underline{m}_{L-1}}{\partial \underline{m}_{L-2}} \cdot \frac{\partial \underline{m}_{L-2}}{\partial \underline{s}_{L-2}} \cdots \frac{\partial \underline{m}_{j+1}}{\partial \underline{m}_j} \cdot \frac{\partial \underline{m}_j}{\partial \underline{s}_j} \qquad \left( \frac{\partial l}{\partial \underline{m}_{j-1}} = \frac{\partial l}{\partial \underline{m}_j} \cdot \frac{\partial \underline{m}_j}{\partial \underline{m}_{j-1}} \right)$$

Notice that there is a lot of reuse of expressions, which means we don't have to recompute things over and over. In particular

$$\frac{\partial l}{\partial \underline{m}_{L-1}} = \frac{\partial l}{\partial \underline{m}_L} \cdot \frac{\partial \underline{m}_L}{\partial \underline{m}_{L-1}} \quad , \quad \frac{\partial l}{\partial \underline{m}_{L-2}} = \frac{\partial l}{\partial \underline{m}_{L-1}} \cdot \frac{\partial \underline{m}_{L-1}}{\partial \underline{m}_{L-2}} \quad ,$$

and in general

$$\frac{\partial l}{\partial \underline{m}_{j-1}} = \frac{\partial l}{\partial \underline{m}_j} \cdot \frac{\partial \underline{m}_j}{\partial \underline{m}_{j-1}} \quad )$$

where $\frac{\partial l}{\partial \underline{m}_j}$ will have been computed at the layer above. This is another key piece of backpropagation!

The only thing left to compute is $\frac{\partial \underline{m}_j}{\partial \underline{s}_j}$ — this is now just an exercise in calculus, so we'll not work it out in class, but the online notes will provide links to pages with further information for those interested.

Optional:
We apply our chain rule (with $\underline{w} = X_j \left[ \frac{a_{j-1}}{2} \right]$) to get

$$\frac{\partial \underline{m}_j}{\partial x_j} = \frac{\partial}{\partial x_j} \sigma \left( X_j \left[ \frac{a_{j-1}}{2} \right] \right) = \frac{\partial \underline{\sigma}}{\partial \underline{w}} \cdot \frac{\partial \underline{w}}{\partial x_j}.$$

Now for $\underline{\sigma}(\underline{w}) = \begin{bmatrix} \sigma(w_1) \\ \vdots \\ \sigma(w_{p_{j-1}}) \end{bmatrix}$, $\frac{\partial \underline{\sigma}}{\partial \underline{w}} = \begin{bmatrix} \sigma'(w_1) & & \\ & \ddots & \\ & & \sigma'(w_{p_{j-1}}) \end{bmatrix}$. Next, we need to find

$$\frac{\partial \underline{w}}{\partial x_j} = \frac{\partial}{\partial x_j} \left( X_j \left[ \frac{a_{j-1}}{2} \right] \right).$$ This can be computed using multilinear algebra (tensors).
We won't work it out, but note that it can be found efficiently.