

General Gradient of Neural Network

tomas.ukkonen@iki.fi, 2017

Backpropagation is a commonly known algorithm for computing the gradient of error function which arises when we know target values and the loss, cost or error function is one dimensional. Generalizing this to general gradient calculation when we seek to find the maximum or minimum value of a neural network (thought often ill-fated because of local optimas produced by an over-fitted neural network) is then important. This is needed, for example, when implementing certain reinforcement learning methods.

Consider a two-layer neural network

$$y(\mathbf{x}) = f(\mathbf{W}^{(2)} g(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$$

The gradients of the final layer are (non-zero terms are at the j :th row):

$$\frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(2)}} = \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}\right) \begin{pmatrix} 0 \\ g_i \\ 0 \end{pmatrix}$$

$$\frac{\partial y(\mathbf{x})}{\partial b_j^{(2)}} = \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}\right) \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

The derivation chain-rule can be used to calculate the second (and more deep layers' gradients):

$$\frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(1)}} = \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}\right) \frac{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}{\partial \mathbf{g}} \frac{\partial \mathbf{g}(\mathbf{x})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})} \frac{\partial (\mathbf{W}^{(2)} \mathbf{x} + \mathbf{b}^{(2)})}{\partial w_{ji}^{(1)}}$$

$$\frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(1)}} = \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}\right) \mathbf{W}^{(2)} \text{diag}\left(\frac{\partial \mathbf{g}(\mathbf{x})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})}\right) \begin{pmatrix} 0 \\ x_i \\ 0 \end{pmatrix}$$

$$\frac{\partial y(\mathbf{x})}{\partial b_j^{(1)}} = \text{diag}\left(\frac{\partial f(\mathbf{x})}{\partial (\mathbf{W}^{(2)} \mathbf{g} + \mathbf{b}^{(2)})}\right) \mathbf{W}^{(2)} \text{diag}\left(\frac{\partial \mathbf{g}(\mathbf{x})}{\partial (\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})}\right) \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

By analysing the chain rule we can derive generic backpropagation formula for the full gradient. Let $\mathbf{v}^{(k)}$ be a k :th layers local field, $\mathbf{v}^{(k)} = \mathbf{W}^{(k)} f(\mathbf{v}^{(k-1)}) + \mathbf{b}^{(k)}$. Then local gradient matrices $\delta^{(k)}$ are

$$\delta^{(L)} = \text{diag}\left(\frac{\partial f(\mathbf{v}^{(L)})}{\partial \mathbf{v}^{(L)}}\right)$$

$$\delta^{(k-1)} = \delta^{(k)} \mathbf{W}^{(k)} \text{diag}\left(\frac{\partial f(\mathbf{v}^{(k-1)})}{\partial \mathbf{v}^{(k-1)}}\right)$$

And network's parameter gradient matrices for each layer are (only j :th element of each row is non-zero):

$$\frac{\partial y(\mathbf{x})}{\partial w_{ji}^{(k)}} = \delta^{(k)} \begin{pmatrix} 0 \\ f(v_i^{(k-1)}) \\ 0 \end{pmatrix}$$

$$\frac{\partial y(\mathbf{x})}{\partial b_j^{(k)}} = \delta^{(k)} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

To test that gradient matrix is correctly computed it can be compared with normal squared error calculations (normal backpropagation).

$$\varepsilon(\mathbf{x}|\mathbf{w}) = \frac{1}{2} \|y_i - y(\mathbf{x}|\mathbf{w})\|^2$$

$$\frac{\partial \varepsilon(\mathbf{x}|\mathbf{w})}{\partial \mathbf{w}} = (y(\mathbf{x}|\mathbf{w}) - y_i)^T \frac{\partial y(\mathbf{x}|\mathbf{w})}{\partial \mathbf{w}}$$

Sometimes also needs gradient with respect to \mathbf{x} and not weights parameters \mathbf{w} . This can be calculated using the chain rule again. For simplicity, let's consider two-layer case initially.

$$\mathbf{g}(\mathbf{x}) = \mathbf{f}(\mathbf{W}^{(2)} \mathbf{h}(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$$

The gradient is:

$$\begin{aligned}\frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} &= \frac{\partial \mathbf{f}(\mathbf{v}^{(2)})}{\partial (\mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)})} \frac{\partial (\mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)})}{\partial \mathbf{h}} \frac{\partial \mathbf{h}(\mathbf{v}^{(1)})}{\partial (\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})} \frac{\partial (\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})}{\partial \mathbf{x}} \\ \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} &= \frac{\partial \mathbf{f}(\mathbf{v}^{(2)})}{\partial \mathbf{v}^{(2)}} \mathbf{W}^{(2)} \frac{\partial \mathbf{h}(\mathbf{v}^{(1)})}{\partial \mathbf{v}^{(1)}} \mathbf{W}^{(1)}\end{aligned}$$

This results into following formula (diag() entries are square matrices which diagonal is nonzero):

$$\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x}|\mathbf{w}) = \text{diag}(\nabla_{\mathbf{v}^{(L)}} \mathbf{f}(\mathbf{v}^{(L)})) \mathbf{W}^{(L)} \dots \text{diag}(\nabla_{\mathbf{v}^{(2)}} \mathbf{f}(\mathbf{v}^{(2)})) \mathbf{W}^{(2)} \text{diag}(\nabla_{\mathbf{v}^{(1)}} \mathbf{f}(\mathbf{v}^{(1)})) \mathbf{W}^{(1)}$$

In continuous reinforcement learning, we need to maximize the given policy's μ average \mathbf{Q} -value $\mathbf{Q}(\mathbf{x}, \mu(\mathbf{x}))$ which gradient can be computed by using the chain-rule but there is additionally linear pre- and postprocessings $\mathbf{W}\mathbf{x} + \mathbf{b}$ in μ and \mathbf{Q} which makes the calculation of gradient more complicated.

$$\nabla_{(\mathbf{W}_{\mu}\mathbf{x} + \mathbf{b}_{\mu})} \mathbf{W}'_Q \mathbf{Q}(\mathbf{W}_Q \mathbf{z} + \mathbf{b}_Q) + \mathbf{b}'_Q, \mathbf{z} = [\mathbf{x}, \mathbf{W}'_{\mu} \mu(\mathbf{W}_{\mu}\mathbf{x} + \mathbf{b}_{\mu}) + \mathbf{b}'_{\mu}] =$$

$$\mathbf{W}'_Q \nabla \mathbf{Q}(\mathbf{W}_Q \mathbf{z} + \mathbf{b}_Q) \mathbf{W}_Q^{(\mu)} \mathbf{W}'_{\mu} \nabla \mu$$

But in practice we don't have post-processing for μ so the gradient becomes

$$\mathbf{W}'_Q \nabla \mathbf{Q}(\mathbf{W}_Q \mathbf{z} + \mathbf{b}_Q) \mathbf{W}_Q^{(\mu)} \nabla \mu$$

Recurrent Neural Networks and Backpropagation (Similar to RTRL)

The basic learning algorithm for recurrent neural networks (RNN) is BPTT but I use modified RTRL instead. (RTRL - real time recurrent learning). This is done by unfolding neural net in time and computing the gradients. The recurrent neural network is

$$\mathbf{u}(n+1) = \begin{pmatrix} \mathbf{y}(n+1) \\ \mathbf{r}(n+1) \end{pmatrix} = \mathbf{f}(\mathbf{x}(n), \mathbf{r}(n))$$

The error function to minimize is:

$$E(N) = \frac{1}{2} \sum_{n=1}^N \|\mathbf{d}(n+1) - \mathbf{\Gamma}_y \mathbf{f}(\mathbf{x}(n), \mathbf{\Gamma}_r \mathbf{u}(n-1)|\mathbf{w})\|^2$$

In which $\mathbf{\Gamma}$ matrices are used to select $\mathbf{y}(n)$ and $\mathbf{r}(n)$ vectors from generic output vector and the initial input to feedforward neural network is zero $\mathbf{u}(0) = \mathbf{0}$.

It is possible to calculate gradient of \mathbf{f} using the chain rule

$$\frac{\partial E(N)}{\partial \mathbf{w}} = \sum_{n=0}^N (\mathbf{\Gamma}_y \mathbf{f}(\mathbf{x}(n), \mathbf{\Gamma}_r \mathbf{u}(n-1)|\mathbf{w}) - \mathbf{d}(n+1))^T \mathbf{\Gamma}_y \nabla_{\mathbf{w}} \mathbf{f}(\mathbf{x}(n), \mathbf{\Gamma}_r \mathbf{u}(n-1)|\mathbf{w})$$

To calculate the gradient $\nabla_{\mathbf{w}} \mathbf{f}$ one must remember that $\mathbf{u}(n)$ now also depends on \mathbf{w} resulting into eq:

$$\nabla_{\mathbf{w}} \mathbf{f}(\mathbf{x}(n), \mathbf{\Gamma}_r \mathbf{u}(n-1)) = \frac{\partial \mathbf{f}}{\partial \mathbf{w}} + \frac{\partial \mathbf{f}}{\partial \mathbf{r}} \mathbf{\Gamma}_r \frac{\partial \mathbf{u}(n-1)}{\partial \mathbf{w}}$$

To further compute gradients we get a generic update rule

$$\frac{\partial \mathbf{u}(n)}{\partial \mathbf{w}} = \frac{\partial \mathbf{f}}{\partial \mathbf{w}} + \frac{\partial \mathbf{f}}{\partial \mathbf{r}} \mathbf{\Gamma}_r \frac{\partial \mathbf{u}(n-1)}{\partial \mathbf{w}}$$

The computation of gradients can be therefore bootstrapped by setting $\frac{\partial \mathbf{u}(0)}{\partial \mathbf{w}} = \frac{\partial \mathbf{f}(\mathbf{0}, \mathbf{0})}{\partial \mathbf{w}}$ and iteratively updating \mathbf{u} gradient while computing the current error for the timestep.

RNN-RBM

RNN-RBM was described on the web to create learn creating "music" by using BPTT but I use extended RTRL approach instead.

(<http://danshiebler.com/2016-08-17-musical-tensorflow-part-two-the-rnn-rbm/>)

In RNN-RBM we have a standard RBM model but RBM's biases $\mathbf{a}(n+1)$ and $\mathbf{b}(n+1)$ are generated by a recurrent neural network $\begin{pmatrix} \mathbf{y}^{(n+1)} \\ \mathbf{r}^{(n+1)} \end{pmatrix} = \mathbf{f}(\mathbf{x}(n), \mathbf{r}(n) | \mathbf{w})$, $\mathbf{y}(n+1) = \begin{pmatrix} \mathbf{a}^{(n+1)} \\ \mathbf{b}^{(n+1)} \end{pmatrix}$ and visible units (MIDI notes) are fed to be inputs of recurrent neural network $\mathbf{x}(n) = \mathbf{v}(n)$.

One can then compute RBM's log-likelihood gradient with respect to recurrent neural networks weights \mathbf{w} maximizing probability of "semi" independent MIDI notes observations

$$-\log[p(\mathbf{v}(1), \mathbf{v}(2) \dots \mathbf{v}(N))] \approx \sum_n -\log(p(\mathbf{v}_n))$$

We want to calculate gradient with respect to \mathbf{w} where only elements \mathbf{a} and \mathbf{b} depend on \mathbf{w} .

$$\frac{\partial -\log(p(\mathbf{v}_i))}{\partial \mathbf{w}} = \frac{\partial -\log(p(\mathbf{v}_i))}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{w}} + \frac{\partial -\log(p(\mathbf{v}_i))}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{w}}$$

This can be rewritten using free-energy $p(\mathbf{v}) = \frac{1}{Z} e^{-F(\mathbf{v})} = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$ and the gradient formula for $\frac{\partial -\log(p(\mathbf{v}))}{\partial \theta} = \frac{\partial F(\mathbf{v})}{\partial \theta} - E_{\mathbf{v}} \left[\frac{\partial F(\mathbf{v})}{\partial \theta} \right]$:

$$\frac{\partial -\log(p(\mathbf{v}_i))}{\partial \mathbf{w}} = \left(\frac{\partial F(\mathbf{v})}{\partial \mathbf{a}} - E_{\mathbf{v}} \left[\frac{\partial F(\mathbf{v})}{\partial \mathbf{a}} \right] \right) \frac{\partial \mathbf{a}}{\partial \mathbf{w}} + \left(\frac{\partial F(\mathbf{v})}{\partial \mathbf{b}} - E_{\mathbf{v}} \left[\frac{\partial F(\mathbf{v})}{\partial \mathbf{b}} \right] \right) \frac{\partial \mathbf{b}}{\partial \mathbf{w}}$$

We assume GB-RBM model (see `RBM_notes.tm`) with the following energy function

$$E_{\text{GB}}(\mathbf{v}, \mathbf{h}) = \frac{1}{2}(\mathbf{v} - \mathbf{a})^T \Sigma^{-1}(\mathbf{v} - \mathbf{a}) - (\Sigma^{-0.5} \mathbf{v})^T \mathbf{W} \mathbf{h} - \mathbf{b}^T \mathbf{h}$$

And we extend RNN to also output/predict variance $\mathbf{z}(n) = \log(\text{diag}(\Sigma(n)))$ [here we use $\mathbf{z}(n)$ for two different things, one for variance parameter of GB-RBM and other for recurrent neural networks output]. Our gradients therefore are (see `RBM_notes.tm`):

$$\frac{\partial \mathbf{u}(n)}{\partial \mathbf{w}} = \frac{\partial \mathbf{f}(\mathbf{v}(n))}{\partial \mathbf{w}} + \frac{\partial \mathbf{f}}{\partial \mathbf{r}} \mathbf{\Gamma}_{\mathbf{r}} \frac{\partial \mathbf{u}(n-1)}{\partial \mathbf{w}}, \quad \mathbf{u}(n) = \begin{pmatrix} \mathbf{a}(n) \\ \mathbf{b}(n) \\ \mathbf{z}(n) \\ \mathbf{r}(n) \end{pmatrix}$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{w}} = \mathbf{\Gamma}_{\mathbf{a}} \frac{\partial \mathbf{u}(n)}{\partial \mathbf{w}}$$

$$\frac{\partial \mathbf{b}}{\partial \mathbf{w}} = \mathbf{\Gamma}_{\mathbf{b}} \frac{\partial \mathbf{u}(n)}{\partial \mathbf{w}}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}} = \mathbf{\Gamma}_{\mathbf{z}} \frac{\partial \mathbf{u}(n)}{\partial \mathbf{w}}$$

$$\frac{\partial F}{\partial \mathbf{a}} = -\Sigma^{-1}(\mathbf{v} - \mathbf{a})$$

$$\frac{\partial F}{\partial \mathbf{b}} = -\text{sigmoid}(\mathbf{W}^T \Sigma^{-1/2} \mathbf{v} + \mathbf{b})$$

$$\frac{\partial F}{\partial \mathbf{z}_i} = -\frac{1}{2}(v_i - a_i)^2 e^{-z_i} + \frac{1}{2}e^{-z_i/2} v_i \sum_j w_{ij} \text{sigmoid}(\mathbf{W}^T \Sigma^{-1/2} \mathbf{v} + \mathbf{b})_j$$

$$\frac{\partial F}{\partial \mathbf{W}} = -\text{sigmoid}(\mathbf{W}^T \Sigma^{-1/2} \mathbf{v} + \mathbf{b}) (\Sigma^{-1/2} \mathbf{v})^T$$

When using these gradients it is important to remember that, in general, one must update variance terms \mathbf{z} independently from other parameters or the GB-RBM doesn't converge. Initially, however, I will make an attempt to learn both variance \mathbf{z} , \mathbf{a} and \mathbf{b} because they now all depend on common weight vector parameter \mathbf{w} . Pseudocode for RNN-RBM optimization:

1. randomly initialize \mathbf{W} and \mathbf{w} using small values.
2. Calculate parameters \mathbf{a} , \mathbf{b} and \mathbf{z} using current RNN (initially with zero recurrent input parameters including previous step's visible MIDI notes $\mathbf{v}(n-1)$)
3. Use RBM and CD-k to calculate (\mathbf{v}, \mathbf{h}) parameters for input sample(s) and calculate contrastive divergence samples in order to calculate gradient of free energy ∇F .
4. Calculate the gradient of the recurrent neural network $\nabla_{\mathbf{w}} \mathbf{u}(n)$ and use ∇F to calculate gradient of the probability $p(\mathbf{v})$ with respect to \mathbf{W} and \mathbf{w} .

5. Repeat steps 2-4 for each song i (time series) of visible notes $\{\mathbf{v}_i(n)\}$ and use the sum of all songs gradients to move parameters \mathbf{W} and \mathbf{w} towards (hopefully) higher probability of data. (IMPLEMENTATION NOTE: concatenate all songs as a single time-serie and try to learn it).

BB-RBM

Instead of GB-RBM which variance learning is complicated. I initially (and also) implement and test BB-RBM implementation as the RBM part of the RNN-RBM. In this case there is no variance terms \mathbf{z}_i to worry about.

$$\frac{\partial F}{\partial \mathbf{a}} = -\mathbf{v}$$

$$\frac{\partial F}{\partial \mathbf{b}} = -\text{sigmoid}(\mathbf{W}\mathbf{v} + \mathbf{b})$$

$$\frac{\partial F}{\partial \mathbf{W}} = -\text{sigmoid}(\mathbf{W}\mathbf{v} + \mathbf{b})\mathbf{v}^T$$

NOTE: initial use of RNN-RBM outlined here seem to diverge quickly to chaos (many random notes played at once) when applying RNN-RBM to classical MIDI notes data (note range: C-4 .. B-6). It seems problem should be regularized somehow to limit number of on notes played at once.

In practice it seems to be difficult to regularize RNN-RBM because of the special form of the error function (log probability). I suggest the use of “negative gradient”. In addition to training samples which probability should be maximized, artificial songs are created where each note has random probability $p > 0.50$ of being in on position and gradient of these additional training samples is calculated normally but the calculated gradient is substracted from the positive gradient so that probability of those “random songs” is greatly reduced.