

Neural Networks for Machine Learning

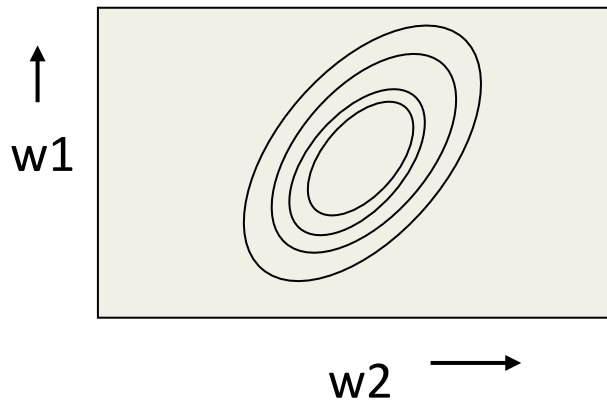
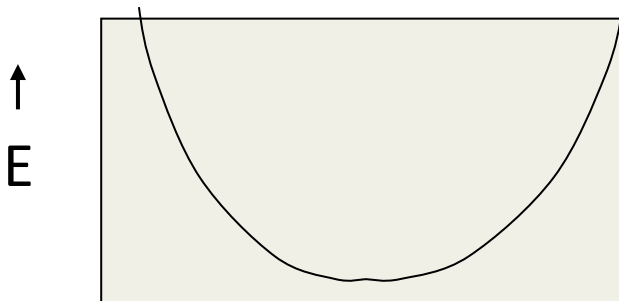
Lecture 6a

Overview of mini-batch gradient descent

Geoffrey Hinton
with
Nitish Srivastava
Kevin Swersky

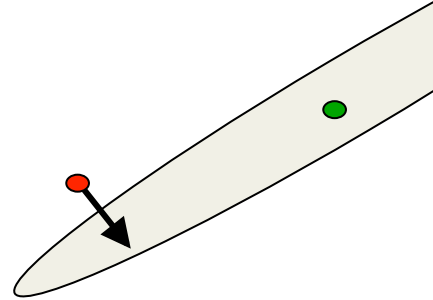
Reminder: The error surface for a linear neuron

- The error surface lies in a space with a horizontal axis for each weight and one vertical axis for the error.
 - For a linear neuron with a squared error, it is a quadratic bowl.
 - Vertical cross-sections are parabolas.
 - Horizontal cross-sections are ellipses.
- For multi-layer, non-linear nets the error surface is much more complicated.
 - But locally, a piece of a quadratic bowl is usually a very good approximation.



Convergence speed of full batch learning when the error surface is a quadratic bowl

- Going downhill reduces the error, but the direction of steepest descent does not point at the minimum unless the ellipse is a circle.
 - The gradient is big in the direction in which we only want to travel a small distance.
 - The gradient is small in the direction in which we want to travel a large distance.



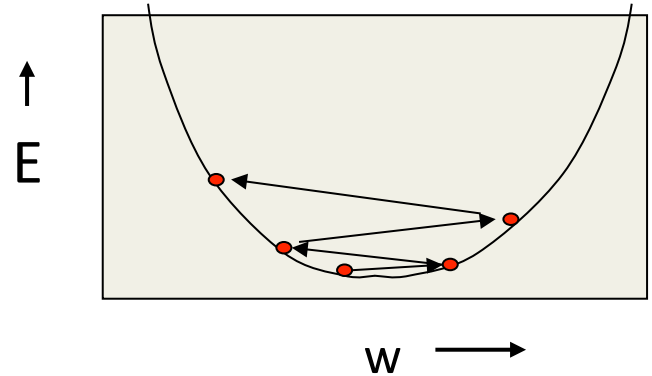
Even for non-linear multi-layer nets, the error surface is locally quadratic, so the same speed issues apply.

globally very curved in one direction and very uncurved in another direction

How the learning goes wrong

- If the learning rate is big, the weights slosh to and fro across the ravine.
 - If the learning rate is too big, this oscillation diverges.
- What we would like to achieve:
 - Move quickly in directions with small but consistent gradients.
 - Move slowly in directions with big but inconsistent gradients.

reverse sign gradient we want to move slowly



Stochastic gradient descent

- If the dataset is highly redundant, the gradient on the first half is almost identical to the gradient on the second half.
 - So instead of computing the full gradient, update the weights using the gradient on the first half and then get a gradient for the new weights on the second half.
 - The extreme version of this approach updates weights after each case. Its called “online”.
- Mini-batches are usually better than online.
 - Less computation is used updating the weights.
 - Computing the gradient for many cases simultaneously uses matrix-matrix multiplies which are very efficient, especially on GPUs
- Mini-batches need to be balanced for classes

avoid min batches that are uncharacteristic of the data
can take random mini batches

Two types of learning algorithm

If we use the full gradient computed from all the training cases, there are many clever ways to speed up learning (e.g. non-linear conjugate gradient).

- The optimization community has studied the general problem of optimizing smooth non-linear functions for many years.
- Multilayer neural nets are not typical of the problems they study so their methods may need a lot of adaptation.

may need lot of modification

For large neural networks with very large and highly redundant training sets, it is nearly always best to use mini-batch learning.

- The mini-batches may need to be quite big when adapting fancy methods.
- Big mini-batches are more computationally efficient.

A basic mini-batch gradient descent algorithm

most people will use this ^{with} big redundant data set.

- Guess an initial learning rate.
 - If the error keeps getting worse or oscillates wildly, reduce the learning rate.
 - If the error is falling fairly consistently but slowly, increase the learning rate.
- Write a simple program to automate this way of adjusting the learning rate.

you expect it to fluctuate on validation set.

- Towards the end of mini-batch learning it nearly always helps to turn down the learning rate.
 - This removes fluctuations in the final weights caused by the variations between mini-batches.
- Turn down the learning rate when the error stops decreasing.
 - Use the error on a separate validation set

turning down learning rate cause smoothing away of fluctuation in caused by end mini batches.

Neural Networks for Machine Learning

Lecture 6b

A bag of tricks for mini-batch gradient descent

Geoffrey Hinton

with

Nitish Srivastava

Kevin Swersky

Initializing the weights

- If two hidden units have exactly the same bias and exactly the same incoming and outgoing weights, they will always get exactly the same gradient.
 - So they can never learn to be different features.
 - We break symmetry by initializing the weights to have small random values.

the can never
become
different from
one another

if we use big weights that it will tend to saturate it

- If a hidden unit has a big fan-in, small changes on many of its incoming weights can cause the learning to overshoot.
 - We generally want smaller incoming weights when the fan-in is big, so initialize the weights to be proportional to $\sqrt{\text{fan-in}}$.
- We can also scale the learning rate the same way.

it small fan in u want to use
bigger weights

Shifting the inputs

Speed of NN learning depend upon

normalizing mean centered

- When using steepest descent, shifting the input values makes a big difference.

- It usually helps to transform each component of the input vector so that it has zero mean over the whole training set.

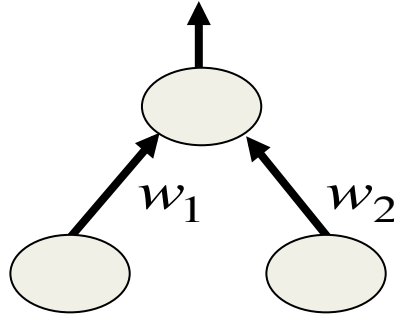
- The hyperbolic tangent (which is $2 \cdot \text{logistic} - 1$) produces hidden activations that are roughly zero mean.

- In this respect its better than the logistic.

fluctuation in the big negative input are ignored by logistic

makes sense to use tanh that goes -1 and 1 true only if inputs are distributed around 0

* for tanh we have out to the end of plateau to ignore negative



color indicates training case

$$101, 101 \rightarrow 2$$

$$101, 99 \rightarrow 0$$

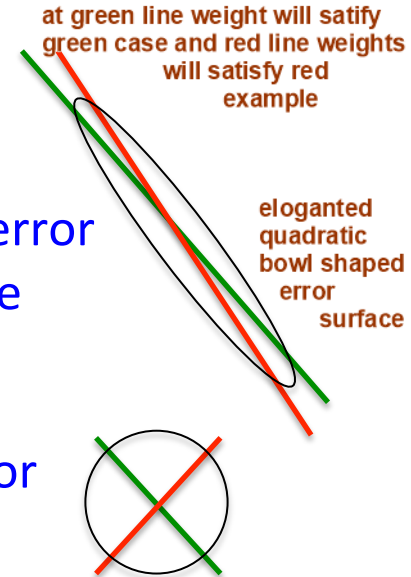
gives error surface

subtract 100 from inputs

$$1, 1 \rightarrow 2$$

$$1, -1 \rightarrow 0$$

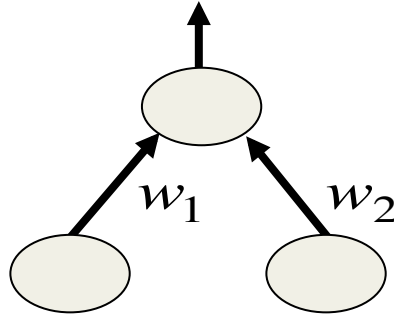
gives error surface



Scaling the inputs

scaling on the basis of variance/ standard deviance
unit variance

- When using steepest descent, scaling the input values makes a big difference.
 - It usually helps to transform each component of the input vector so that it has unit variance over the whole training set.

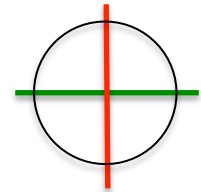
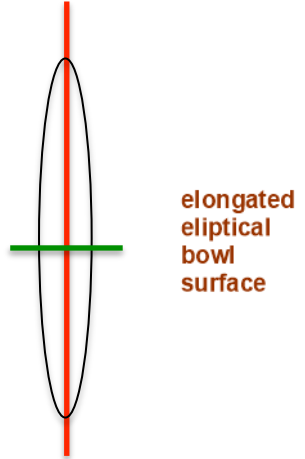


color indicates weight axis

$$\begin{array}{lcl} 0.1, & 10 & \rightarrow 2 \\ 0.1, & -10 & \rightarrow 0 \end{array} \quad \begin{array}{l} \text{gives error} \\ \text{surface} \end{array}$$

x10 x0.1 rescaled

$$\begin{array}{lcl} 1, & 1 & \rightarrow 2 \\ 1, & -1 & \rightarrow 0 \end{array} \quad \begin{array}{l} \text{gives error} \\ \text{surface} \end{array}$$



A more thorough method: Decorrelate the input components

- For a linear neuron, we get a big win by decorrelating each component of the input from the other input components.
- There are several different ways to decorrelate inputs. A reasonable method is to use Principal Components Analysis. PCA
 - Drop the principal components with the smallest eigenvalues.
 - This achieves some dimensionality reduction.
 - Divide the remaining principal components by the square roots of their eigenvalues. For a linear neuron, this converts an axis aligned elliptical error surface into a circular one.
- For a circular error surface, the gradient points straight towards the minimum.

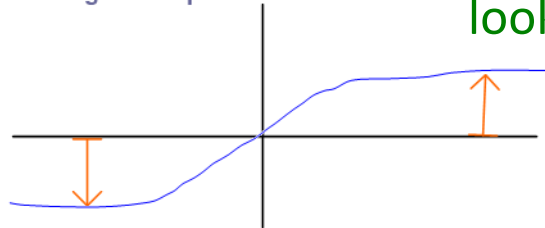
Common problems that occur in multilayer networks

This will cause NN to learn feature from the output without considering the info from input

- If we start with a very big learning rate, the weights of each hidden unit will all become very big and positive or very big and negative.
 - The error derivatives for the hidden units will all become tiny and the error will not decrease. from. outputs
 - This is usually a plateau, but people often mistake it for a local minimum. thus learning will stop
- In classification networks that use a squared error or a cross-entropy error, the best guessing strategy is to make each output unit always produce an output equal to the proportion of time it should be a 1.
 - The network finds this strategy quickly and may take a long time to improve on it by making use of the input. for a network with many layers vanishing gradient problem
 - This is another plateau that looks like a local minimum.

if very big learning rate:

u drive the hidden unit to firmly on or firmly off
There state does not depend on input.



It may take time bcoz it has to get sensible info. from inputs to many hidden layers to output that could take a long time to learn if we start with small weights.

Neural networks often start out not using the inputs. They decrease the weights given to the inputs and set the output to the best constant they can find, which under some loss functions is the average of the targets.

Suppose you are training a neural network to predict the next letter in an English text from the previous 100. One of the first things that tends to happen is that the network learns to say that regardless of the input, there is a 13% chance the next letter is an E, a 9% chance the next letter is a T, 8% chance the next letter is an A, etc. These are simply the character frequencies in English, and this has a lower average loss than the uniform distribution. While this is better than nothing, it makes no use of the inputs. The neural network has not yet learned that if the previous letter was a Q, this increases the chance that the next letter is a U. Eventually, the network should learn to use the inputs, but this may take longer than you are prepared to wait.

If the input is very confusing, then most features constructed from the inputs are useless. These are the initial features found by a neural network. The first thing the network learns is that these random features are not very good, and that trusting them more than an infinitesimal amount makes the squared error or cross-entropy worse. It is easier to learn to ignore these bad features than to learn to fix them, and if you ignore all features you disconnect the inputs from the outputs, and the neural network trains toward the best constant predictor of the labels.

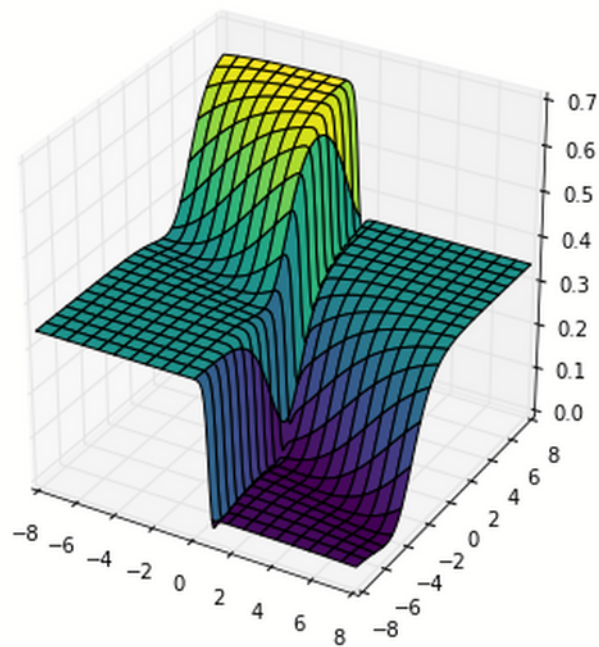
This is true when the hidden unit activations are bounded above and/or below as in the case of the logistic and hyperbolic tangent functions. An update that makes the weights large pushes the activation out to either side where the function flattens out. When the function flattens out, the gradient is small and subsequent learning will be slow.

Consider the logistic activation $\sigma(z)=1/(1+e^{-z})$ whose derivative is $\sigma'(z)=\sigma(z)(1-\sigma(z))$. A large (positive or negative) z will push $\sigma(z)$ close to 0 or 1, and thus one or the other factor in the derivative close to 0. If the initial weights or learning rate is too big, it can take arbitrarily long to converge to a minimum.

If all of the hidden units are close to saturation, the error surface has slowed, a plateau may be indistinguishable from a local minimum

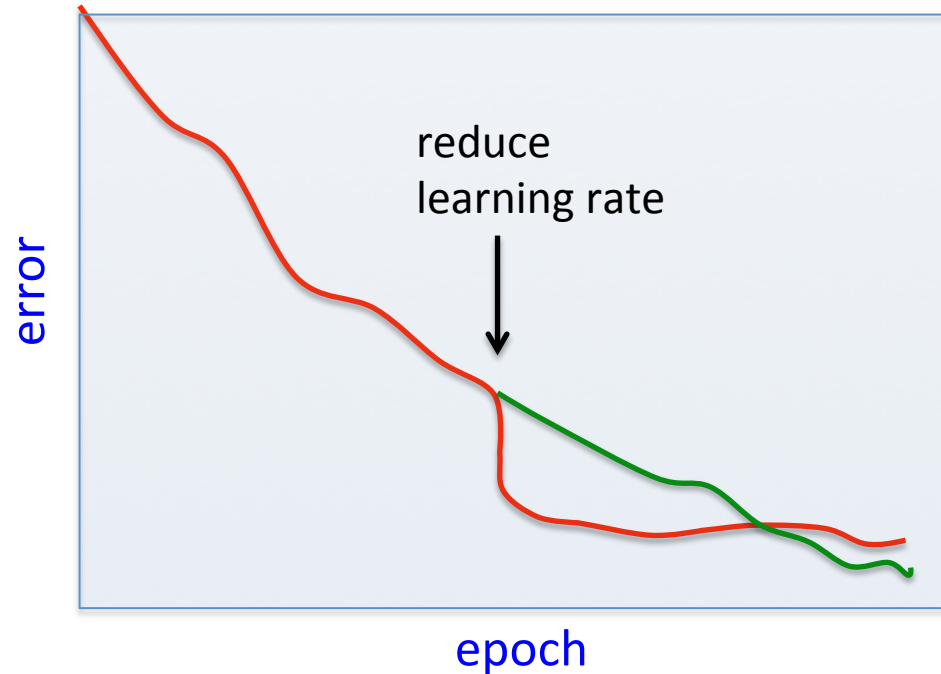
Just to have a visual, here is an example error surface of a two-layer network with two input neurons and a logistic output neuron (using the squared error cost function).

Weight decay helps because it pushes weights towards the origin, but if the learning rate is too big a subsequent update can put you in another bad spot.



Be careful about turning down the learning rate

- Turning down the learning rate reduces the random fluctuations in the error due to the different gradients on different mini-batches.
 - So we get a quick win.
 - But then we get slower learning.
- Don't turn down the learning rate too soon!



Four ways to speed up mini-batch learning

change in weight = weight + gradient x learning rate

this like a ball at a position of weight accelerate downhill according to gradient which changes its velocity than its position.

- 1 Use “momentum”
 - Instead of using the gradient to change the position of the weight “particle”, use it to change the velocity.
- 2 Use separate adaptive learning rates for each parameter
 - Slowly adjust the rate using the consistency of the gradient for that parameter.

1. ball has momentum i.e. it remember previous gradient in its velocity

2. Learning by empirical measurement- are we keep making progress by changing wts in same direction or does gradient keeps oscillating around the sign of gradient keeps changing. if sign keeps changing -> slow down learning rate. If sign keeps same -> increase learning rate.

- 3 rmsprop: Divide the learning rate for a weight by a running average of the magnitudes of recent gradients for that weight.
 - This is the mini-batch version of just using the sign of the gradient.
- 4 Take a fancy method from the optimization literature that makes use of curvature information (not this lecture)
 - Adapt it to work for neural nets
 - Adapt it to work for mini-batches.

3. if gradient are large divide by large number if gradients are small divide by small number
deal very wide range of gradient.

4. use Full batch learning and fancy method

Neural Networks for Machine Learning

Lecture 6c

The momentum method

Geoffrey Hinton
with
Nitish Srivastava
Kevin Swersky

Can be applied to full batch
learning

Best strategy to use stochastic
Gradient descent with mini
batches combined with
momentum method.

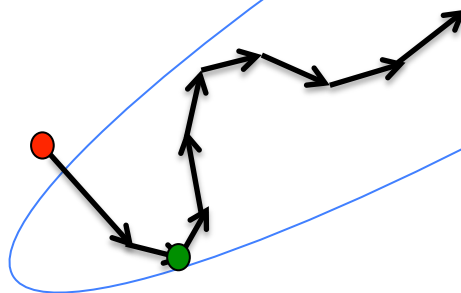
The intuition behind the momentum method

Imagine a ball on the error surface. The location of the ball in the horizontal plane represents the weight vector.

- The ball starts off by following the gradient, but once it has velocity, it no longer does steepest descent.
- Its momentum makes it keep going in the previous direction.

we want to loose energy so we want to have little bit of viscosity/friction velocity die off on each update

- It damps oscillations in directions of high curvature by combining gradients with opposite signs.
- It builds up speed in directions with a gentle but consistent gradient.



The equations of the momentum method

time here is updates of the weights

$$\mathbf{v}(t) = \alpha \mathbf{v}(t-1) - \varepsilon \frac{\partial E}{\partial \mathbf{w}}(t)$$

momentum term

velocity will be decreased a little bit respect to momentum also we want it to go downhill by learning rate ε times the gradient

The effect of the gradient is to increment the previous velocity. The velocity also decays by α which is slightly less than 1.

viscosity or alpha momentum.

$$\Delta \mathbf{w}(t) = \mathbf{v}(t)$$

$$= \alpha \mathbf{v}(t-1) - \varepsilon \frac{\partial E}{\partial \mathbf{w}}(t)$$

$$= \alpha \Delta \mathbf{w}(t-1) - \varepsilon \frac{\partial E}{\partial \mathbf{w}}(t)$$

The weight change is equal to the current velocity.

The weight change can be expressed in terms of the previous weight change and the current gradient.

The behavior of the momentum method

- If the error surface is a tilted plane, the ball reaches a **terminal velocity**.
 - If the momentum is close to 1, this is much faster than simple gradient descent.

Terminal Velocity: the gain in velocity comes from gradient is balanced by multiplied attenuation of velocity due to momentum term

$$\mathbf{v}(\infty) = \frac{1}{1 - \alpha} \left(-\epsilon \frac{\partial E}{\partial \mathbf{w}} \right)$$

1-alpha goes to zero velocity will become high

- ¹ At the beginning of learning there may be very large gradients. so it will go large momentum
 - So it pays to use a small momentum (e.g. 0.5).
0.5 rather than 0 bcoz will avg. out some slashes in obvious ravines.
 - Once the large gradients have disappeared and the weights are stuck in a ravine the momentum can be smoothly raised to its final value (e.g. 0.9 or even 0.99)
- This allows us to learn at a rate that would cause divergent oscillations without the momentum.

1. large initialized weight -> large weight -> big momentum -> quickly change them to make thing better can not take hard problem to just the right value of relative value of different weights to have sensible feature vector.

raising the momentum to 0.99 not smoothly can start oscillation

WHY DON'T USE A BIGGER LEARNING RATE

Using a smaller learning rate and big momentum allows you to get away with an overall learning rate much bigger than that we could have if we have used learning rate alone with no momentum

Using momentum speeds up gradient descent learning because

Directions of consistent change get amplified.

Directions of fluctuations get damped.

Allows using much larger learning rates.

BECAUSE

Momentum accumulates consistent components of the gradient and attenuates the fluctuating ones. It also allows us to use bigger learning rates because the learning is now more stable.

A better type of momentum (Nesterov 1983)

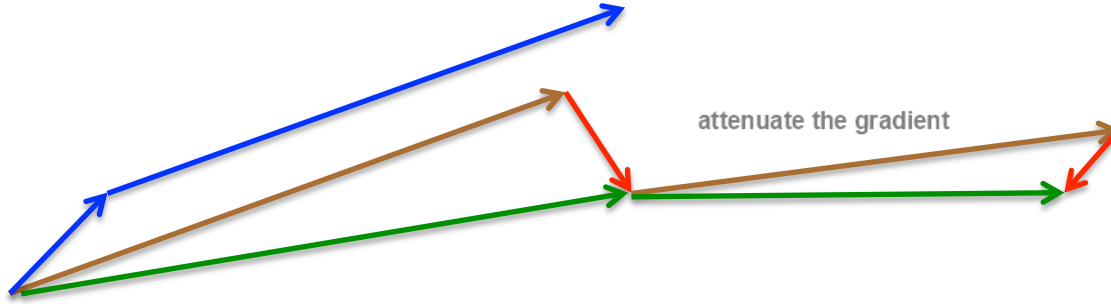
- The standard momentum method first computes the gradient at the current location and then takes a big jump in the direction of the updated accumulated gradient.
- Ilya Sutskever (2012 unpublished) suggested a new form of momentum that often works better.
 - Inspired by the Nesterov method for optimizing convex functions.
- First make a big jump in the direction of the previous accumulated gradient.
- Then measure the gradient where you end up and make a correction.
 - Its better to correct a mistake after you have made it!

In standard method we add in gradient and gamble on the big jump.

In nestrov method use previously accumulated gradient you make the jump and correct then afterwards

A picture of the Nesterov method

- **First** make a big jump in the direction of the previous accumulated gradient.
- **Then** measure the gradient where you end up and make a correction.



brown vector = jump, red vector = correction, green vector = accumulated gradient

blue vectors = standard momentum

Neural Networks for Machine Learning

Lecture 6d

A separate, adaptive learning rate for each connection

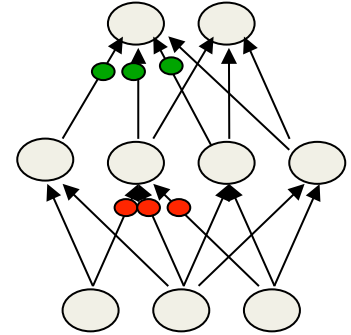
Geoffrey Hinton
with
Nitish Srivastava
Kevin Swersky

ADAPTIVE LEARNING

each weight has its own learning rate when the gradient switches sign frequently we turn the learning rate if not we increase the learning rate

The intuition behind separate adaptive learning rates

- In a ^{Deep} multilayer net, the appropriate learning rates can vary widely between weights:
 - The magnitudes of the gradients are often very different for different layers, especially if the initial weights are small.
 - The fan-in of a unit determines the size of the “overshoot” effects caused by simultaneously changing many of the incoming weights of a unit to correct the same error.
- So use a global learning rate (set by hand) multiplied by an appropriate local gain that is determined empirically for each weight.



Gradients can get very small in the early layers of very deep nets.

The fan-in often varies widely between layers.

It maybe that the unit didn't get enough input, when you change all these weights at the same time to fix up the error, it now gets too much input. Obviously, that effect is going to be bigger if there's a bigger fan-in.

One way to determine the individual learning rates

- Start with a local gain of 1 for every weight.
- Increase the local gain if the gradient for that weight does not change sign.
- Use small additive increases and multiplicative decreases (for mini-batch)
 - This ensures that big gains decay rapidly when oscillations start.
 - If the gradient is totally random the gain will hover around 1 when we increase by **plus** δ half the time and decrease by **times** $1 - \delta$ half the time.

$$\Delta w_{ij} = -\varepsilon g_{ij} \frac{\partial E}{\partial w_{ij}}$$

$$\text{if } \left(\frac{\partial E}{\partial w_{ij}}(t) \frac{\partial E}{\partial w_{ij}}(t-1) \right) > 0$$

$$\text{then } g_{ij}(t) = g_{ij}(t-1) + .05$$

$$\text{else } g_{ij}(t) = g_{ij}(t-1) * .95$$

if gradient hover around 0 then g_{ij} wil hover around 1

if gradient increases +vely $g_{ij} >> 1$

if gradient decreases -vely $g_{ij} << 1$

Tricks for making adaptive learning rates work better

- Limit the gains to lie in some reasonable range
 - *e.g.* [0.1, 10] or [.01, 100]
- Use full batch learning or big mini-batches
 - This ensures that changes in the sign of the gradient are not mainly due to the sampling error of a mini-batch.
- Adaptive learning rates can be combined with momentum.
 - Use the agreement in sign between the current gradient for a weight and the velocity for that weight (Jacobs, 1989).
- Adaptive learning rates only deal with axis-aligned effects.
 - Momentum does not care about the alignment of the axes.

if we don't limit gain we can get instability and they won't die down fast enough and will destroy all the weights

Neural Networks for Machine Learning

Lecture 6e

rmsprop: Divide the gradient by a running average
of its recent magnitude

Geoffrey Hinton
with
Nitish Srivastava
Kevin Swersky

use for large NN with large
redundant dataset

rprop: Using only the sign of the gradient

- The magnitude of the gradient can be very different for different weights and can change during learning.
 - This makes it hard to choose a single global learning rate.
- For **full batch learning**, we can deal with this variation by only using the sign of the gradient.
 - The weight updates are all of the same magnitude.
 - This escapes from plateaus with tiny gradients quickly.
- rprop: This combines the idea of only using the sign of the gradient with the idea of adapting the step size separately for each weight.
 - Increase the step size for a weight **multiplicatively** (e.g. times 1.2) if the signs of its last two gradients agree.
 - Otherwise decrease the step size multiplicatively (e.g. times 0.5).
 - Limit the step sizes to be less than 50 and more than a millionth (Mike Shuster's advice).

if inputs are tiny then we should have big weights in order for inputs to be affected

Why rprop does not work with mini-batches

It violates

- The idea behind stochastic gradient descent is that when the learning rate is small, it averages the gradients over successive mini-batches.
 - Consider a weight that gets a gradient of $+0.1$ on nine mini-batches and a gradient of -0.9 on the tenth mini-batch.
 - We want this weight to stay roughly where it is.
- rprop would increment the weight nine times and decrement it once by about the same amount (assuming any adaptation of the step sizes is small on this time-scale).
 - So the weight would grow a lot.
- Is there a way to combine:
 - The robustness of rprop.
 - The efficiency of mini-batches.
 - The effective averaging of gradients over mini-batches.

rmsprop: A mini-batch version of rprop

- rprop is equivalent to using the gradient but also dividing by the size of the gradient.
 - The problem with mini-batch rprop is that we divide by a different number for each mini-batch. So why not force the number we divide by to be very similar for adjacent mini-batches?
- rmsprop: Keep a moving average of the squared gradient for each weight
$$MeanSquare(w, t) = 0.9 MeanSquare(w, t-1) + 0.1 \left(\frac{\partial E}{\partial w}(t) \right)^2$$
- Dividing the gradient by $\sqrt{MeanSquare(w, t)}$ makes the learning work much better (Tijmen Tieleman, unpublished).

For each connection keep running avg. of the rms gradient and divide gradient by that

Further developments of rmsprop

- Combining rmsprop with standard momentum
 - Momentum does not help as much as it normally does. Needs more investigation.
- Combining rmsprop with Nesterov momentum (Sutskever 2012)
 - It works best if the RMS of the recent gradients is used to divide the correction rather than the jump in the direction of accumulated corrections.
- Combining rmsprop with adaptive learning rates for each connection
 - Needs more investigation.
- Other methods related to rmsprop
 - Yann LeCun's group has a fancy version in "No more pesky learning rates"

Summary of learning methods for neural networks

- For small datasets (e.g. 10,000 cases) or bigger datasets without much redundancy, use a full-batch method.
 - Conjugate gradient, LBFGS ...
 - adaptive learning rates, rprop ...
- For big, redundant datasets use mini-batches.
 - Try gradient descent with momentum.
 - Try rmsprop (with momentum ?)
 - Try LeCun's latest recipe.
- Why there is no simple recipe:
Neural nets differ a lot:
 - Very deep nets (especially ones with narrow bottlenecks).
 - Recurrent nets.
 - Wide shallow nets.Tasks differ a lot:
 - Some require very accurate weights, some don't.
 - Some have many very rare cases (e.g. words).