

August 2000

0

EFFICIENT BACKPROP

A great paper that looks at issues that you can run into with Backprop and it then presents solutions and tries to explain why these solutions work.

Authors note that there is no silver bullet for designing NN's as their design is problem and data dependent. But there are useful heuristics.

So the note the bias v. Variance tradeoff where bias is the difference between the network output across all datasets and the desired function's output. And variance is how much the network output varies from dataset to dataset.

obviously as you train more your bias will go down as you learn the underlying function. But if you train for too long the variance will go up as you learn the noise of this specific dataset. (that's bad). So the goal is to minimize both Variance and bias.

So they state that there is no guarantee that the network will 1) converge to a good solution 2) converge fast 3) even converge at all.

Now for ways to converge to better solutions + do it faster

A) Stochastic v.s. Batch learning

In batch learning you compute the avg gradient across the whole train set + then take a step in that direction.

2

But in Stochastic or "online" learning you pick a random sample from X & then step in the direction of that example. (a more noisy step).

They claim that Stochastic learning:

1. Often MUCH faster
2. Often Better solutions are found
3. Can be used for tracking changes.

There is often lots of redundancy in a dataset because it hopefully has patterns across the data, so on-line does not look at all the examples saving time on the redundancy.

The Better solutions occur because In Batch GD it will converge to the minimum of the basin that the weights were initialized in however with on-line it allows it to jump to new basins because of the noisy steps

also on point #3 if the function you are trying to model constantly changes then on-line learning allows you to constantly take in a new sample & train and the model will slowly change if the underlying model also is changing in real life.

But the very noise that helps SGD also hurts it.

- Because of the noise we can't use second order methods for SGD
- Because of the noise the learning may never converge. But just fluctuate around the local minima

However in practice SGD is preferred

B) Example order

- for Batch this does not matter
- for SGD this does matter
- Networks learn the most from examples that are the most unfamiliar.

- we can shuffle Data so that successive examples are not from the same class
- Present inputs that cause larger error more frequently
- *note be careful as this perturbs the importance the network places on different types of examples.*

c) Normalizing Inputs

- If we normalize input values training is usually faster
- Also good to decorrelate inputs or don't use if they are like $x_2 = 2x_1$ then don't use x_2 as this will slow down learning.

D) Sigmoid

- If not normalized input than Learning can often not happen

E) choosing Target values

- Don't use $(-1, 1)$ for classification labels as will cause weight updates to be very small when using Sigmoid.
- Also this would mean you get no indication of confidence when inputs are near the decision boundary.
- So pick target to maximize the second derivate of sigmoid to fix this problem

F) You don't want weights to be too big and too small as both screw up learning with Sigmoid
 So. pick weights from Distribution with mean 0 & $sd = m^{-1/2}$
 where $m = \# \text{ input features}$

G) choosing a Learning Rate

- give each weight its own learning rate
- lower layers lower Learning rate
- Higher layers Higher learning rate

E) Momentum can help if loss surface is nonspherical

They also note Adaptive learning rates & Radial Basis Functions.

NOW on the topic of gradient decent convergence for some of the theory behind the earlier tricks.

if α_{opt} = optimal learning rate

then $\alpha < \alpha_{opt}$ = too small of steps

then $\alpha > \alpha_{opt}$ = too large and will oscillate but eventually converge

then $\alpha > 2\alpha_{opt}$ = it will Diverge

7

They then Note several Methods for improving training speed via second order methods often using the Hessian to account for the curvature + there for take a Better learning Step.

They then discuss and present Several Ways to compute the Hessian (or it's approximate at least)

They then go on to explain why large eigenvalues are a problem.

After all this discussion they present methods on how to Impliment this in practice