

Random Ponderings

Tuesday, January 13, 2015

A Brief Overview of Deep Learning

(This is a guest post by [Ilya Sutskever](#) on the intuition behind deep learning as well as some very useful practical advice. Many thanks to Ilya for such a heroic effort!)

Deep Learning is really popular these days. Big and small companies are getting into it and making money off it. It's hot. There is some substance to the hype, too: large deep neural networks achieve the best results on speech recognition, visual object recognition, and several language related tasks, such as machine translation and language modeling.

But why? What's so special about deep learning? (from now on, we shall use the term **Large Deep Neural Networks --- LDNN ---** which is what the vaguer term "Deep Learning" mostly refers to). Why does it work now, and how does it differ from neural networks of old? Finally, suppose you want to train an LDNN. Rumor has it that it's very difficult to do so, that it is "black magic" that requires years of experience. And while it is true that experience helps quite a bit, the amount of "trickery" is surprisingly limited --- one needs be on the lookout for only a small number well-known pitfalls. Also, there are many open-source implementations of various state-of-the-art neural networks (c.f. [Caffe](#), [cuda-covnet](#), [Torch](#), [Theano](#)), which makes it much easier to learn all the details needed to make it work.

Why Does Deep Learning Work?

It is clear that, to solve hard problems, we must use powerful models. This statement is obvious. Indeed, if a model is not powerful, then there is absolutely no chance that it can succeed in solving a hard problem, no matter how good the learning algorithm is.

The other necessary condition for success is that our model is trainable. That too is obvious, for if we cannot train our model, then its power is useless --- it will never amount to anything, and great results will not be achieved. The model will forever remain in a state of unrealized potential.

Fortunately, LDNNs are both trainable and powerful.

Why Are LDNNs Powerful?

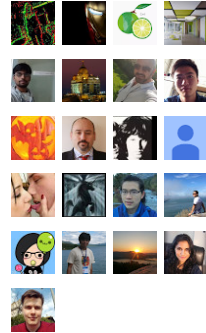
When I talk about LDNNs, I'm talking about 10-20 layer neural networks (because this is what can be trained with today's algorithms). I can provide a few ways of looking at LDNNs that will illuminate the reason they can do as well as they do.

- Conventional statistical models learn simple patterns or clusters. In contrast, LDNNs learn computation, albeit a massively parallel computation with a modest number of steps. Indeed, this is the key difference between LDNNs and other statistical models.
- To elaborate further: it is well known that any algorithm can be implemented by an appropriate very deep circuit (with a layer for each timestep of the algorithm's execution -- [one example](#)). What's more, the deeper the circuit, the more expensive are the algorithms that can be implemented by the circuit (in terms of runtime). And given that neural networks are circuits as well, deeper neural networks can implement algorithms with more steps ---- which is why depth = more power.
 - N.B.: It is easy to see that a single neuron of a neural network can compute the conjunction of its inputs, or the disjunction of its inputs, by simply setting their connections to appropriate values.
- Surprisingly, neural networks are actually more efficient than [boolean circuits](#). By more efficient, I mean that a fairly shallow DNN can solve problems that require many more layers of boolean circuits. For a specific example, consider the highly surprising fact that a DNN with 2 hidden layer and a modest number of units can sort N N -bit numbers! I found the result shocking when I heard about it, so I implemented a small neural network and trained it to sort 10 6-bit numbers,

- [my website](#)
- [my photostream](#)
- [my Quora answers](#)

Followers

关注者 (65 人) [下一步](#)



[关注](#)

Labels

- [adventures](#)
- [announcements](#)
- [computer science](#)
- [education](#)
- [humor](#)
- [information retrieval](#)
- [internet / networks](#)
- [life extension](#)
- [machine learning](#)
- [math](#)
- [random ponderings](#)
- [rationality / intelligence](#)
- [science / technology](#)

Subscribe To

- [Posts](#) ▾
- [Comments](#) ▾

My Blog List

- The Wronger Box**
Facts of the Day (8/7)
4 hours ago
- Slate Star Codex**
Highlights From The
Comments On Billionaire
Philanthropy
6 hours ago
- Turing Machine**
collatz microscope
14 hours ago
- Computational Complexity**
Obstacles to improving
Classical Factoring
Algorithms
15 hours ago
- Gödel's Lost Letter
and P=NP**
Code It Up
1 day ago

which was easy to do to my surprise. It is impossible to sort N N -bit numbers with a boolean circuit that has two hidden layers and that are not gigantic.

- The reason DNNs are more efficient than boolean circuits is because neurons perform a threshold operation, which cannot be done with a tiny boolean circuit.
- Finally, human neurons are slow yet humans can perform lots of complicated tasks in a fraction of a second. More specifically, it is well-known that a human neuron fires no more than 100 times per second. This means that, if a human can solve a problem in 0.1 seconds, then our neurons have enough time to fire only 10 times --- definitely not much more than that. It therefore follows that a large neural network with 10 layers can do anything a human can in 0.1 seconds.
- This is not scientific fact since it is conceivable that real neurons are much more powerful than artificial neurons, but real neurons may also turn out to be much less powerful than artificial neurons. In any event, the above is certainly a plausible hypothesis.
- This is interesting because humans can solve many complicated perception problems in 0.1 seconds --- for example, humans can recognize the identity of an object that's in front of them, recognize a face, recognize an emotion, and understand speech in a fraction of a second. In fact, if there exists even just one person in the entire world who has achieved an uncanny expertise in performing a highly complex task of some sort in a fraction of a second, then this is highly convincing evidence that a large DNN could solve the same task --- if only its connections are set to the appropriate values.
- But won't the neural network need to be huge? Maybe. But we definitely know that it won't have to be exponentially large ---- simply because the brain isn't exponentially large! And if human neurons turn out to be noisy (for example), which means that many human neurons are required to implement a single real-valued operation that can be done using just one artificial neuron, then the number of neurons required by our DNNs to match a human after 0.1 seconds is greatly diminished.

These four arguments suggest (strongly, in my opinion), that for a very wide variety of problems, there exists a setting of the connections of a LDNN that basically solves the problem. Crucially, the number of units required to solve these problems is far from exponential --- on the contrary, the number of units required is often so "small" that it is even possible, using current hardware, to train a network that achieves super-high performance on the task of interest. It is this last point which is so important, and requires additional elaboration:

- We know that most machine learning algorithms are consistent: that is, they will solve the problem given enough data. But consistency generally requires an exponentially large amount of data. For example, the **nearest neighbor algorithm** can definitely solve any problem by memorizing the correct answer to every conceivable input. The same is true for a support vector machine --- we'd have a support vector for almost every possible training case for very hard problems. The same is also true for a neural network with a single hidden layer: if we have a neuron for every conceivable training case, so that neuron fires for that training case and but not for any other, then we could also learn and represent every conceivable function from inputs to outputs. Everything can be done given exponential resources, but it is never ever going to be relevant in our limited physical universe.
- And it is in this point that LDNNs differ from previous methods: we can be reasonably certain that a large but not huge LDNN will achieve good results on a surprising variety of problems that we may want to solve. If a problem can be solved by a human in a fraction of a second, then we have a very non-exponential super-pessimistic upper bound on the size of the smallest neural network that can achieve very good performance.
- But I must admit that it is impossible to predict whether a given problem will be solvable by a deep neural network ahead of time, although it is often possible to tell whenever we know that a similar problem can be solved by an LDNN of a manageable size.

So that's it, then. Given a problem, such as visual object recognition, all we need is to train a giant **convolutional neural network** with 50 layers. Clearly a giant convnet with 50 layers can be configured to achieve human-level performance on object recognition --- right? So we simply need to find these weights. Once once we do, the problem is solved.

Learning.

The Gates Notes

The hidden costs of unreliable electricity
3 days ago

Study Hacks

"It's Like I Couldn't Stop": A Digital Minimalism Case Study
3 days ago

Overcoming Bias

Advice Wiki
1 week ago

Shtetl-Optimized

Links, proofs, talks, jokes
1 week ago

What's new

Sharp bounds for multilinear curved
Kakeya, restriction and oscillatory integral estimates away from the endpoint
1 week ago

Show All

Science News

SA Scientific American

The Death of Social Reciprocity in the Era of Digital Distraction
10 minutes ago

NYT > Science

This Bread Recipe Starts With 4,000-Year-Old Yeast
53 minutes ago

PhysOrg.com - latest science and technology news

Zooming in on top-quark production
2 hours ago

T TIME.com: Top Science and Health Stories

It's a Moment of Reckoning For How We Use the Planet to Halt Climate Change, Warns U.N. Report
2 hours ago

Wired Science

New IPCC Report Shows How Our Abuse of Land Drives Climate Change
3 hours ago

Show All

What is learning? Learning is the problem of finding a setting of the neural network's weights that achieves the best possible results on our training data. In other words, we want to "push" the information from the labelled data into the parameters so that the resulting neural network will solve our problem.

The success of Deep Learning hinges on a very fortunate fact: that well-tuned and carefully-initialized **stochastic gradient descent (SGD)** can train LDNNs on problems that occur in practice. It is not a trivial fact since the training error of a neural network as a function of its weights is highly non-convex. And when it comes to non-convex optimization, we were taught that all bets are off. Only convex is good, and non-convex is bad. And yet, somehow, SGD seems to be very good at training those large deep neural networks on the tasks that we care about. The problem of training neural networks is NP-hard, and in fact there exists a family of datasets such that the problem of finding the best neural network with three hidden units is NP-hard. And yet, SGD just solves it in practice. This is the main pillar of deep learning.

We can say fairly confidently that successful LDNN training relies on the "easy" correlation in the data, which allows learning to bootstrap itself towards the more "complicated" correlations in the data. I have done an experiment that seems to support this claim: I found that training a neural network to solve the **parity problem** is hard. I was able to train the network to solve parity for 25 bits, 29 bits, but never for 31 bits (by the way, I am not claiming that learning parity is impossible for over 30 bits --- only that I didn't succeed in doing so). Now, we know that parity is a highly unstable problem that doesn't have any linear correlations: every linear function of the inputs is completely uncorrelated with the output, which is a problem for neural networks since they are mostly linear at initialization time (so perhaps I should've used larger initial weights? I will discuss the topic of weight initialization later in the text). So my hypothesis (which is shared by many other scientists) is that neural networks start their learning process by noticing the most "blatant" correlations between the input and the output, and once they notice them they introduce several hidden units to detect them, which enables the neural network to see more complicated correlations. Etc. The process goes on. I imagine some sort of a "spectrum" of correlations --- both easy and hard, and the network jumps from a correlation to a more complicated correlation, much like an opportunistic mountain climber.

Generalization.

While it is very difficult to say anything specific about the precise nature of the optimization of neural networks (except near a local minimum where everything becomes convex and uninteresting), we can say something nontrivial and specific about generalization.

And the thing we can say is the following: in his famous 1984 **paper** called "A Theory of the Learnable", **Valiant** proved, roughly speaking, that if you have a finite number of functions, say N , then every training error will be close to every test error once you have more than $\log N$ training cases by a small constant factor. Clearly, if every training error is close to its test error, then **overfitting** is basically impossible (overfitting occurs when the gap between the training and the test error is large). (I am also told that this result was given in **Vapnik's book** as small exercise). This theorem is easy to prove but I won't do it here.

But this very simple result has a genuine implication to any implementation of neural networks. Suppose I have a neural network with N parameters. Each parameter will be a float32. So a neural network is specified with $32N$ bits, which means that we have no more than 2^{32N} distinct neural networks, and probably much less. This means that we won't overfit much once we have more than $32N$ training cases. Which is nice. It means that it's theoretically OK to count parameters. What's more, if we are quite confident that each weight only requires 4 bits (say), and that everything else is just noise, then we can be fairly confident that the number of training cases will be a small constant factor of $4N$ rather than $32N$.

The Conclusion:

If we want to solve a hard problem we probably need a LDNN, which has many parameters. So we need a large high-quality labelled training set to make sure that it has enough information to specify all the network's connections. And once we get that training set, we should run SGD on it until the network solves the problem. And it probably will, if our neural network is large and deep.

What Changed Since the 80s?

In the old days, people believed that neural networks could "solve everything". Why couldn't they do it in the past? There are several reasons.

- **Computers were slow.** So the neural networks of past were tiny. And tiny neural networks cannot achieve very high performance on anything. In other words, small neural networks are not powerful.

- **Datasets were small.** So even if it was somehow magically possible to train LDNNs, there were no large datasets that had enough information to constrain their numerous parameters. So failure was inevitable.
- **Nobody knew how to train deep nets.** Deep networks are important. The current best object recognition networks have between 20 and 25 successive layers of convolutions. A 2 layer neural network cannot do anything good on object recognition. Yet back in the day everyone was very sure that deep nets cannot be trained with SGD, since that would've been too good to be true!

It's funny how science progresses, and how easy it is to train deep neural networks, especially in retrospect.

Practical Advice.

Ok. So you're sold. You're convinced that LDNNs are the present and the future and you want to train it. But rumor has it that it's so hard, so difficult... or is it? The reality is that it used to be hard, but now the community has consolidated its knowledge and realized that training neural networks is easy as long as you keep the following in mind.

Here is a summary of the community's knowledge of what's important and what to look after:

- **Get the data:** Make sure that you have a high-quality dataset of input-output examples that is large, representative, and has relatively clean labels. Learning is completely impossible without such a dataset.
- **Preprocessing:** it is essential to center the data so that its mean is zero and so that the variance of each of its dimensions is one. Sometimes, when the input dimension varies by orders of magnitude, it is better to take the $\log(1 + x)$ of that dimension. Basically, it's important to find a faithful encoding of the input with zero mean and sensibly bounded dimensions. Doing so makes learning work much better. This is the case because the weights are updated by the formula: change in $w_{ij} \propto x_i dL/dy_j$ (w denotes the weights from layer x to layer y , and L is the loss function). If the average value of the x 's is large (say, 100), then the weight updates will be very large and correlated, which makes learning bad and slow. Keeping things zero-mean and with small variance simply makes everything work much better.
- **Minibatches:** Use **minibatches**. Modern computers cannot be efficient if you process one training case at a time. It is vastly more efficient to train the network on minibatches of 128 examples, because doing so will result in massively greater throughput. It would actually be nice to use minibatches of size 1, and they would probably result in improved performance and lower overfitting; but the benefit of doing so is outweighed the massive computational gains provided by minibatches. But don't use very large minibatches because they tend to work less well and overfit more. So the practical recommendation is: use the smaller minibatch that runs efficiently on your machine.
- **Gradient normalization:** Divide the gradient by minibatch size. This is a good idea because of the following pleasant property: you won't need to change the learning rate (not too much, anyway), if you double the minibatch size (or halve it).
- **Learning rate schedule:** Start with a normal-sized learning rate (LR) and reduce it towards the end.
 - A typical value of the LR is **0.1**. Amazingly, 0.1 is a good value of the learning rate for a large number of neural networks problems. Learning rates frequently tend to be smaller but rarely much larger.
 - Use a **validation set** ---- a subset of the training set on which we don't train --- to decide when to lower the learning rate and when to stop training (e.g., when error on the validation set starts to increase).
 - A practical suggestion for a learning rate schedule: if you see that you stopped making progress on the validation set, divide the LR by 2 (or by 5), and keep going. Eventually, the LR will become very small, at which point you will stop your training. Doing so helps ensure that you won't be (over-)fitting the training data at the detriment of validation performance, which happens easily and often. Also, lowering the LR is important, and the above recipe provides a useful approach to controlling via the validation set.
- But most importantly, worry about the **Learning Rate**. One useful idea used by some researchers (e.g., **Alex Krizhevsky**) is to monitor the ratio between the update norm and the weight norm. This ratio should be at around 10^{-3} . If it is much smaller then learning will probably be too slow, and if it is much larger then learning will be unstable and will probably fail.

- **Weight initialization.** Worry about the random initialization of the weights at the start of learning.
 - If you are lazy, it is usually enough to do something like $0.02 * \text{randn}(\text{num_params})$. A value at this scale tends to work surprisingly well over many different problems. Of course, smaller (or larger) values are also worth trying.
 - If it doesn't work well (say your neural network architecture is unusual and/or very deep), then you should initialize each weight matrix with the $\text{init_scale} / \sqrt{\text{layer_width}} * \text{randn}$. In this case init_scale should be set to 0.1 or 1, or something like that.
 - Random initialization is super important for deep and **recurrent nets**. If you don't get it right, then it'll look like the network doesn't learn anything at all. But we know that neural networks learn once the conditions are set.
 - Fun story: researchers believed, for many years, that SGD cannot train deep neural networks from random initializations. Every time they would try it, it wouldn't work. Embarrassingly, they did not succeed because they used the "small random weights" for the initialization, which works great for shallow nets but simply doesn't work for deep nets at all. When the nets are deep, the many weight matrices all multiply each other, so the effect of a suboptimal scale is amplified.
 - But if your net is shallow, you can afford to be less careful with the random initialization, since SGD will just find a way to fix it.

You're now informed. Worry and care about your initialization. Try many different kinds of initialization. This effort will pay off. If the net doesn't work at all (i.e., never "gets off the ground"), keep applying pressure to the random initialization. It's the right thing to do.

- If you are training **RNNs** or **LSTMs**, use a **hard constraint** over the norm of the gradient (remember that the gradient has been divided by batch size). Something like 15 or 5 works well in practice in my own experiments. Take your gradient, divide it by the size of the minibatch, and check if its norm exceeds 15 (or 5). If it does, then shrink it until it is 15 (or 5). This one little trick plays a huge difference in the training of RNNs and LSTMs, where otherwise the exploding gradient can cause learning to fail and force you to use a puny learning rate like $1e-6$ which is too small to be useful.
- **Numerical gradient checking:** If you are not using **Theano** or **Torch**, you'll be probably implementing your own gradients. It is easy to make a mistake when we implement a gradient, so it is absolutely critical to use numerical gradient checking. Doing so will give you a complete peace of mind and confidence in your code. You will know that you can invest effort in tuning the hyperparameters (such as the learning rate and the initialization) and be sure that your efforts are channeled in the right direction.
- If you are using **LSTMs** and you want to train them on problems with very long range dependencies, you should initialize the biases of the forget gates of the LSTMs to large values. By default, the forget gates are the sigmoids of their total input, and when the weights are small, the forget gate is set to 0.5, which is adequate for some but not all problems. This is the one non-obvious caveat about the initialization of the LSTM.
- **Data augmentation:** be creative, and find ways to algorithmically increase the number of training cases that are in your disposal. If you have images, then you should translate and rotate them; if you have speech, you should combine clean speech with all types of random noise; etc. Data augmentation is an art (unless you're dealing with images). Use common sense.
- **Dropout.** **Dropout** provides an easy way to improve performance. It's trivial to implement and there's little reason to not do it. Remember to tune the dropout probability, **and to not forget to turn off Dropout and to multiply the weights by** (namely by $1 - \text{dropout probability}$) **at test time**. Also, be sure to train the network for longer. Unlike normal training, where the validation error often starts increasing after prolonged training, dropout nets keep getting better and better the longer you train them. So be patient.
- **Ensembling.** Train 10 neural networks and average their predictions. It's a fairly trivial technique that results in easy, sizeable performance improvements. One may be mystified as to why averaging helps so much, but there is a simple reason for the effectiveness of averaging. Suppose that two classifiers have an error rate of 70%. Then, when they agree they are right. But when they disagree, one of them is often right, so now the average prediction will place much more weight on the correct answer. The effect will be especially strong whenever the network is confident when it's right and unconfident when it's wrong.

I am pretty sure that I haven't forgotten anything. The above 13 points cover literally everything that's needed in order to train LDNNs successfully.

So, to Summarize:

- LDNNs are powerful.
- LDNNs are trainable if we have a very fast computer.
- So if we have a very large high-quality dataset, we can find the best LDNN for the task.
- Which will solve the problem, or at least come close to solving it.

The End.

But what does the future hold? Predicting the future is obviously hard, but in general, models that do even more computation will probably be very good. The **Neural Turing Machine** is a very important step in this direction. Other problems include **unsupervised learning**, which is completely mysterious and incomprehensible in my opinion as of 8 Jan 2015. Learning very complicated "things" from data without supervision would be nice. All these problems require extensive research.

Posted by Unknown at 10:22 PM



Labels: computer science, machine learning, science / technology

675 comments:

1 – 200 of 675 [Newer](#) [Newest](#)



Unknown said...

Very nice and useful summary. However, I would like to argue that the reasons stated for the success of deep neural nets would apparently work as well for kernel machines (if you were to train them using an efficient algorithm that is not quadratic in the number of examples, and it is not hopeless to achieve that). In particular, I believe that deep learning offers something specific in terms of generalization, something that a general-purpose kernel such as the Gaussian kernel does not offer. I have written about the expressive power of deep nets (see recent Montufar et al NIPS paper, for example) and about their ability to generalize far from the training examples (NIPS'2005 and my 2009 book and 2013 PAMI review), which a Gaussian kernel machine or a decision tree could not do. The latter may need an exponentially large number of training examples to get the same generalization error as some deep net (exponential in depth). Now this is associated with priors that are coming with the deep net (there is no free lunch and it will not work for *any* function), such as the priors associated with distributed representations (the existence of underlying factors that generate the data, so that one can learn about each factor without having to observe all the configurations of values of all the others) and with depth (the assumption of a hierarchy of factors). To summarize, I believe that the reasons for deep nets to work go beyond what was stated in this nice post.

1/14/2015 7:32 PM



Yoshua Bengio said...

I forgot to sign the above comment: Yoshua Bengio

1/14/2015 7:37 PM



Ilya Sutskever said...

Thanks!

I think that we are expressing similar concepts in different words, and I suspect that we use the term generalization in slightly different ways.

I don't see a particular difference between a shallow net with a reasonable number of neurons and a kernel machine with a reasonable number of support vectors (its not useful to consider Kernel machines with exponentially many support vectors just like there isn't a point in considering the universal approximation theorem as both require exponential resources) --- both of these models are nearly identical, and thus equally unpowerful. Both of these models will be inferior to an LDNN with a comparable number of parameters precisely because the LDNN can do computation and the shallow models cannot. The LDNN can sort, do integer-multiplication, compute analytic functions, decompose an input into small pieces and recombine it later in a higher level representation, partition the input space into an exponential number of non-arbitrary tiny regions, etc. Ultimately, if the LDNN has 10,000 layers, then it can, in principle, execute any parallel algorithm that runs in fewer than 10,000 steps, giving this LDNN an incredible expressive power. Thus, I don't think that the argument in the article suggests that huge kernel machines should be able to solve these hard problems --- they would need to have exponentially many support vectors.

Although I didn't define it in the article, generalization (to me) means that the gap between the training and the test error is small. So for example, a very bad model that has similar training and test errors does not overfit, and hence generalizes, according to the way I use these concepts. It follows that generalization is easy to achieve whenever the capacity of the model (as measured by the number of parameters or its VC-dimension)