

# Deep Learning

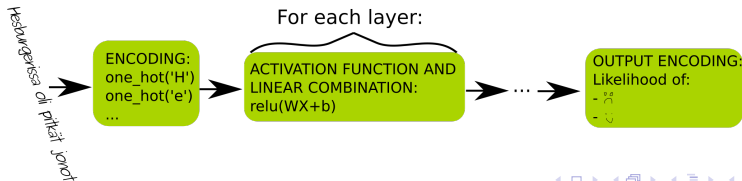
Learning High-degree Non-linear Models:  
Enabling factors, tips and tricks

Tero Keski-Valkama



2016-08-24

- Deep learning refers to deep neural networks. It took about 10 years for neural network technology to surpass the issues they stalled with in the 1990s. In 00s, if an academic paper mentioned neural networks, it was less likely to be published than if not.
- 1990s networks were shallow, and hence relatively useless: It requires a bunch of non-trivial tricks to make deeper networks possible. This presentation will present the most important of these tricks.
- Neural networks are just complex non-linear models with lots of parameters. They are formed by layers of neurons, successive operations with an a linear combination of inputs from the previous layer and an activation function.
- In general, you always need a training set, a test set and a validation set. Typically you would divide your data into three parts, train the system with the training set, test if the system overfits with the test set, and finally for the final system you can validate that your metaparameters didn't just learn the test set using the validation set.
- Underfitting = high bias, overfitting = high variance



- Networks can be trained in a supervised fashion, with target outputs(labels), using backpropagation of error. The output error is known, and it is propagated backwards, layer by layer, estimating an error for each weight parameter. The weights are then updated using this error multiplied by the learning factor ( $\ll 1$ ). Lots of learning steps makes the network learn the target function (*input*  $\rightarrow$  *target\_output*)
- Non-linearities (activation functions) between the layers make the model interesting, compared to other statistical models.
- When the layers are getting smaller than the previous layer, the information is being compressed (or filtered), and the layer activations represent a higher abstraction level, a higher semantic level information about the signal.

- Unsupervised learning (no target output) can be done using backpropagation and autoassociativity ( $input \rightarrow input$ ), (or prediction:  $delayed\_input \rightarrow input$ ), or using Deep Belief Nets and Contrastive Divergence <sup>1</sup> (minimizes the energy for generating real data-related distributions while maximizing the energy for "confabulations") for pairs of subsequent layers one by one.
- Metaoptimizing the learning hyperparameters and neural network structure should always be performed.
- A good platform to use for neural network experimentation is Google's TensorFlow, based on Python.

Supervised Backpropagation Non-energy-based Methods	Unsupervised Contrastive Divergence Energy-based Methods
---	--

<sup>1</sup>[http://deeplearning.cs.cmu.edu/notes/yuxiong\\_CD.pdf](http://deeplearning.cs.cmu.edu/notes/yuxiong_CD.pdf)

## Problems

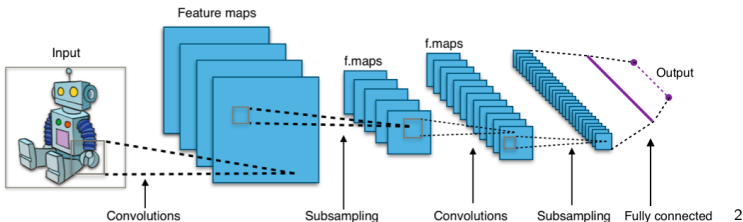
- Overlearning
- Getting stuck in local minima
- Exploding/diminishing gradients

## Solutions

- 1 Weight-sharing – less parameters
- 2 Unsupervised pre-training – lots of data
- 3 Near-linear activation functions
- 4 Drop-out
- 5 Gradient clipping
- 6 Metaoptimization
- 7 Bonus: Reservoir computing

- Reducing the number of parameters or degrees of freedom of the model is regularization. Regularization prevents overfitting.
- In Convolutional Neural Networks, the weights are identical for each window, and window outputs are max-pooled (sub-sampling) to the next layer. This exploits the translational symmetry of the domain, and is often used for images.
- Regularization can also be done by introducing a loss term for weights, L2 norm is often used.
- As always in multiobjective optimization, the factors and slopes of the loss function components are important.

$$\text{RegularizedLoss} = \text{Loss} + \alpha \cdot \text{SquaredWeights}^\beta$$



<sup>2</sup>Image from: Aphex34, Wikipedia

- Unsupervised pre-training makes use of a lot of unlabeled data, learning the overall structure of the domain, which can be fine-tuned with backpropagation for specific labels at a later stage. More data, less labels. <sup>3</sup>
- This could be thought of as "pre-training as regularization", or approaching neural reinforcement learning.
- Free pre-trained models are available from Google and Oxford University. Mostly convolutional nets trained on images and video. <sup>4</sup>
- The models are taught using for example greedy DBN Contrastive Divergence or other autoencoding methods.
- In general, neural networks get better when trained on a different task with similar data. <sup>5</sup>

---

<sup>3</sup><http://www.jmlr.org/papers/volume11/erhan10a/erhan10a.pdf>

<sup>4</sup><http://www.vlfeat.org/matconvnet/pretrained/>

<sup>5</sup><http://blogs.microsoft.com/next/2015/12/10/microsoft-researchers-win-imagenet-computer-vision-challenge/>

- Deep neural networks are effective because of non-linear activation functions – Stacking linear layers only reduces to a single linear layer.
- However, non-linearities are challenging to learn. Minimizing the non-linearity has been shown to be really useful. Deep neural networks can be trained without unsupervised pre-training if rectifier activation functions are used. (otherwise cannot)

Figure: Rectifier

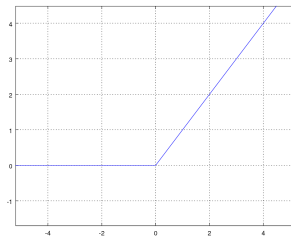
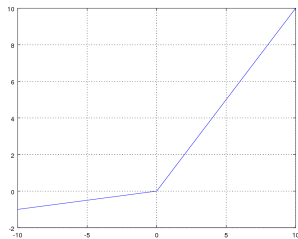
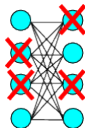


Figure: Leaky Rectifier



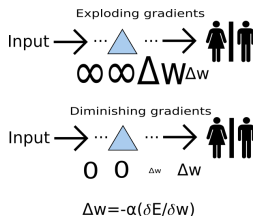


- Dropout is a very effective strategy for preventing overfitting.<sup>6</sup>
- In practice it means randomly dropping out for example half of the internal representation neurons and their output weights for each training step, and scaling the outputs accordingly.
- Learning is done by the randomly thinned networks, and the final use is done using the whole network.
- This prevents neurons from co-adapting with each other, and learn features more independently. I.e. one neuron learns one feature, not multiple neurons learning one feature.
- It is a very strong method against overfitting, and should be almost always used.
- Analogous to Random Forests or Ensemble Averaging Committee Machine.
- The expense is more learning iterations.



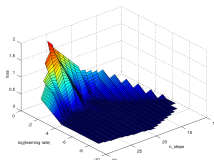
<sup>6</sup><https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>

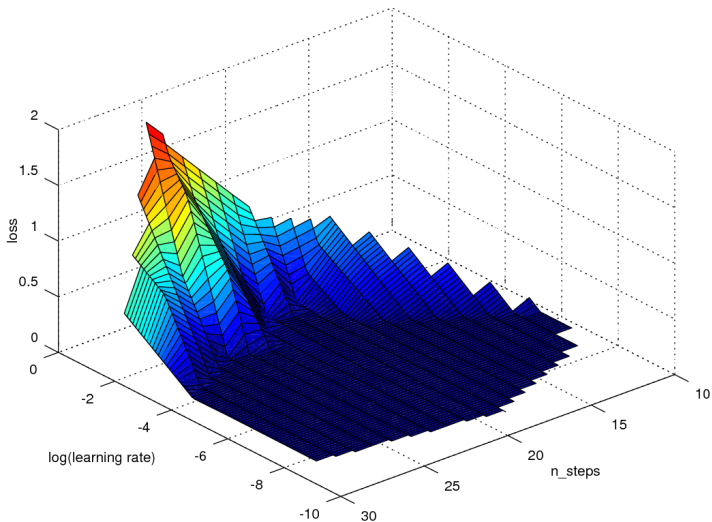
- Gradient explosion happens in backpropagation when the squared error gradients are amplified through activation functions of multiple layers. This typically causes infinities and divergence. Limited activation functions (sigmoid, tanh) are prone to this behavior.
- Gradients are often clipped to some maximum absolute value to prevent infinities.
- Gradients can also diminish to zero, when a small change in internal representation means a too large change in output. If a gradient goes to zero, it means it continues to be zero in backpropagation towards input. Adding noise might help the system in escaping local plateaus.
- Microsoft uses special one-bit gradients only signifying the direction of change, which is always non-zero. <sup>7</sup>



<sup>7</sup><https://www.microsoft.com/en-us/research/publication/1-bit-stochastic-gradient-descent-and-application-to-data-parallel-distributed-training-of-speech-dnn/>

- With increased computational capacity we have better capability for metaoptimization.
- Metaoptimization means finding optimal network structures and learning parameters by running the learning multiple times with varying parameters.
- Two variants: Automatic Bayesian optimization and graphical optimization.
- Automatic Bayesian optimization should be used for stand-alone purposes, it finds good parameters without human aid. Parameter values to test are picked automatically.
- Graphical optimization is often better when human is available, and it also helps in understanding the system. You should sample parameter values randomly, instead of using an uniform grid.
- Process:
  - 1 Pick random values for 1-2 parameters
  - 2 Run the training for a while (i.e. 1 hour)
  - 3 Save the latest test set loss (not the best one)
  - 4 Draw a graph
- Note that sometimes random noise causes apparent optimums. Always note that the other measured points are in agreement.





- A common name for methods that utilize a large, pre-initialized network, and then train a simple linear model with this large reservoir.
- Based on the observation that neural networks tend to learn mostly on the final layers anyhow.
- For example: Echo-State Networks <sup>8</sup> , Extreme Learning Machine <sup>9</sup> .
- Fast to train, work very well in domains where semantic depth is not required.
- Each neuron in the recurrent echo-state network reservoir represents a single non-linear signal response. The echo-state network learns an arbitrary response by fitting a linear combination of these responses to the desired output signal. Analogous to Fourier decomposition where any linear response can be decomposed in frequency plane, as a weighted sum of discrete cosine responses.
- Extreme Learning Machine fits a linear combination of a large reservoir of non-linear responses in a similar way but for feed-forward networks. Activation functions used in the ELM network used can be any mixture of infinitely differentiable functions, e.g. sigmoid, radial basis, sine, cosine, exponential, and many nonregular functions.

---

<sup>8</sup><http://minds.jacobs-university.de/sites/default/files/uploads/papers/PracticalESN.pdf>

<sup>9</sup>[https://www.researchgate.net/profile/Chee\\_Siew/publication/4116697\\_Extreme\\_learning\\_machine\\_A\\_new\\_learning\\_scheme\\_of\\_feedforward\\_neural\\_networks/links/00b4952f8672c66db1000000.pdf](https://www.researchgate.net/profile/Chee_Siew/publication/4116697_Extreme_learning_machine_A_new_learning_scheme_of_feedforward_neural_networks/links/00b4952f8672c66db1000000.pdf)

- To make deep neural networks actually do anything useful, you need a bag of tricks:
- Regularization: Managing the number of parameters
- Pre-training with massive datasets
- Rectifier activation functions
- Dropout
- Gradient clipping
- These should be always used where relevant. In addition, the network structure and learning parameters should always be meta-optimized.
- The next lectures will be about:
  - Data encoding / representation
  - TensorFlow and meta-optimization