



A Convolutional Neural Network for Image Classification of Cats and Dogs

Status update



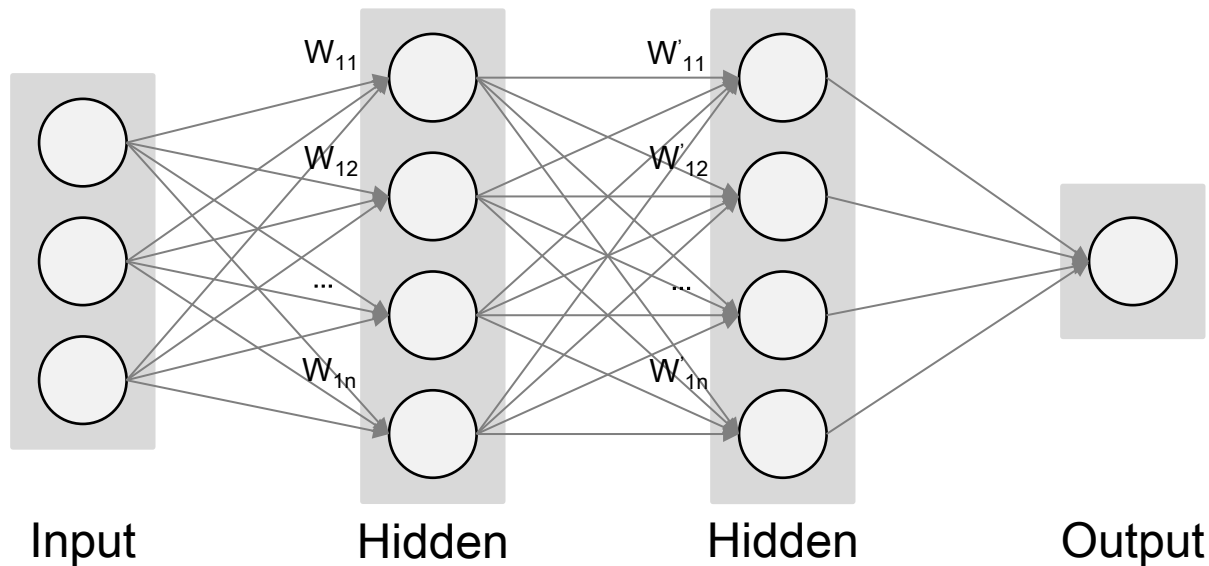
Structure

- Neural Nets (NN)
- Convolution NN (CNN)
- Problem
- Evaluation
- Aims



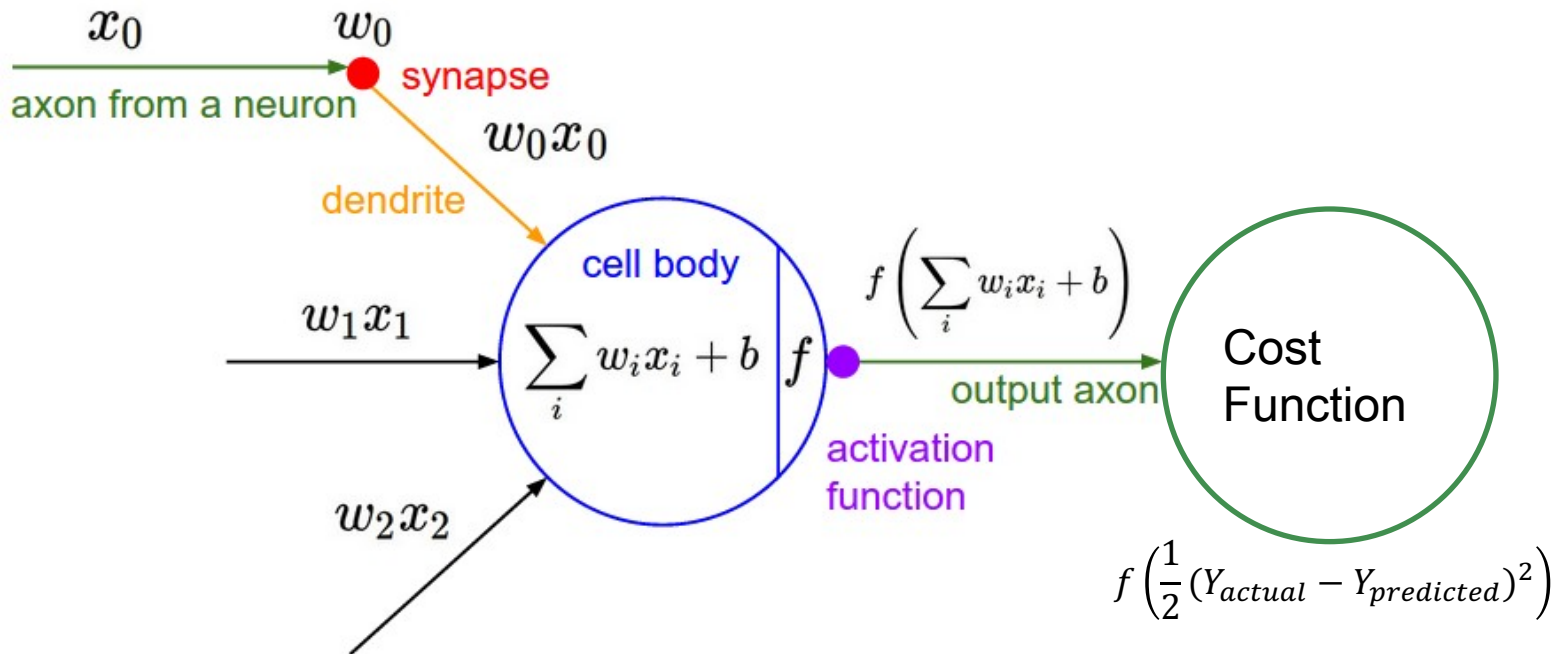
NEURAL NETS

Introduction – Neural Nets



Quelle: <http://cs231n.github.io/convolutional-networks/>

Neuron Model



Mathematical view

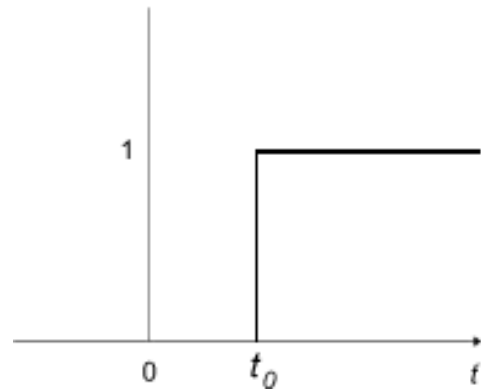
- Input, Weights
- $l_0 = X_i, W = \text{rand}()$
- Compute Sigmoid (Activation Function)
- $l_1 = \text{Sig}(X_i \cdot W)$
- Measure how much we missed (Cost Function)
- $\text{Err} = l_0 - l_1$
- Multiply Err by the Sigmoid slope
- $\Delta l_1 = \text{Err} \cdot \Delta(\text{Sig}(\text{Err}))$
- Update Weights
- $W = W + (l_0 \cdot \Delta l_1)$

Activation Functions

- Non-Linear
 - sigmoid, tanh, elu, softplus, and softsign
- continuous but not everywhere differentiable functions
 - relu, relu6, crelu and relu_x

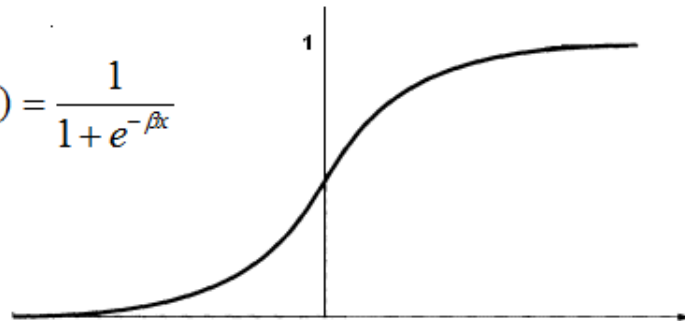
Activation functions

- Why Sigmoid?
 - Not telling in which direction should we move in.
 - Non-differentiability at certain points



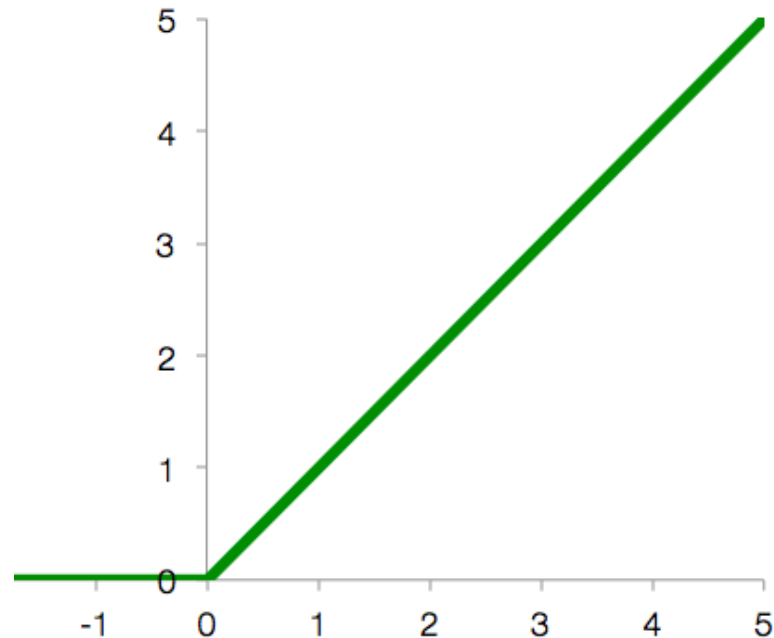
Sigmoid

$$f(x) = \frac{1}{1 + e^{-\beta x}}$$



Activation functions

- Why Relu?
 - Accelerates the convergence rate
 - Simply implementation





Cost Functions

- Squared Error Measure
- Softmax Cross-entropy Function

Squared Error Measure Function

- $Error = \frac{1}{2} (Y_{actual} - Y_{predicted})^2$
- Drawbacks
 - No gradient to get from 0.000...1 to 1.
 - To do so it will take quite longer.
 - Deprives NN of probability information.

Advantages to Squared Error Measure

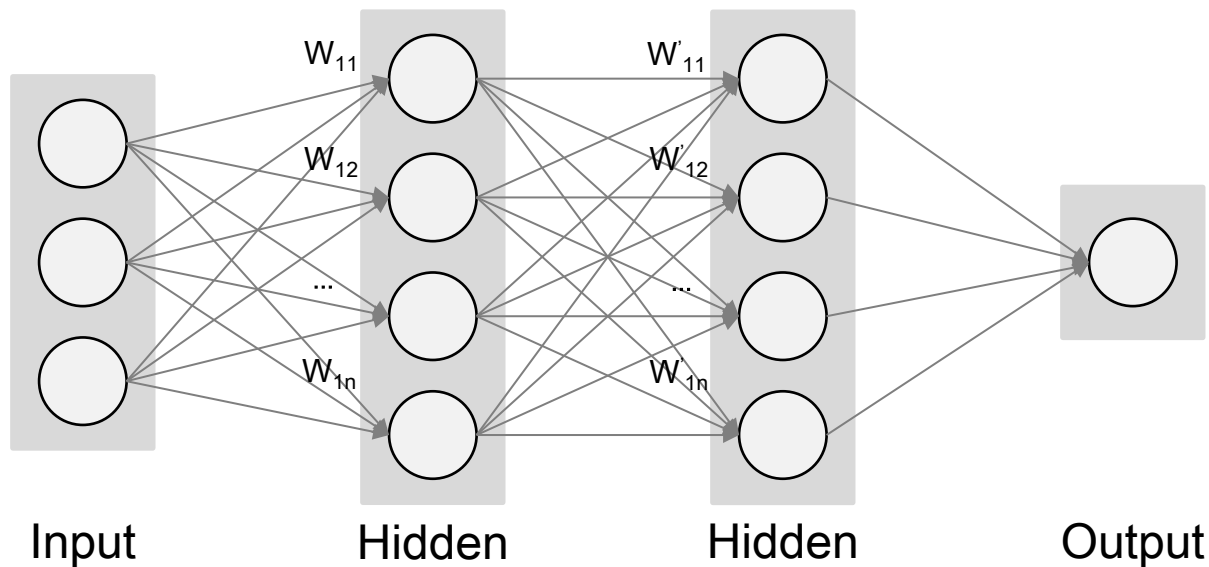
- $C = -\sum_j t_j \log Y_j$
- Very big gradient when:
 - Target value is 1.
 - Actual output is 0.
- Balance between
 - Steepness of $\frac{dC}{dy}$ and
 - Flatness of $\frac{dy}{dZ}$

$$\frac{\partial C}{\partial Z_j} = \sum_j \frac{\partial C}{\partial y_j} \frac{\partial y_j}{\partial Z_j}$$



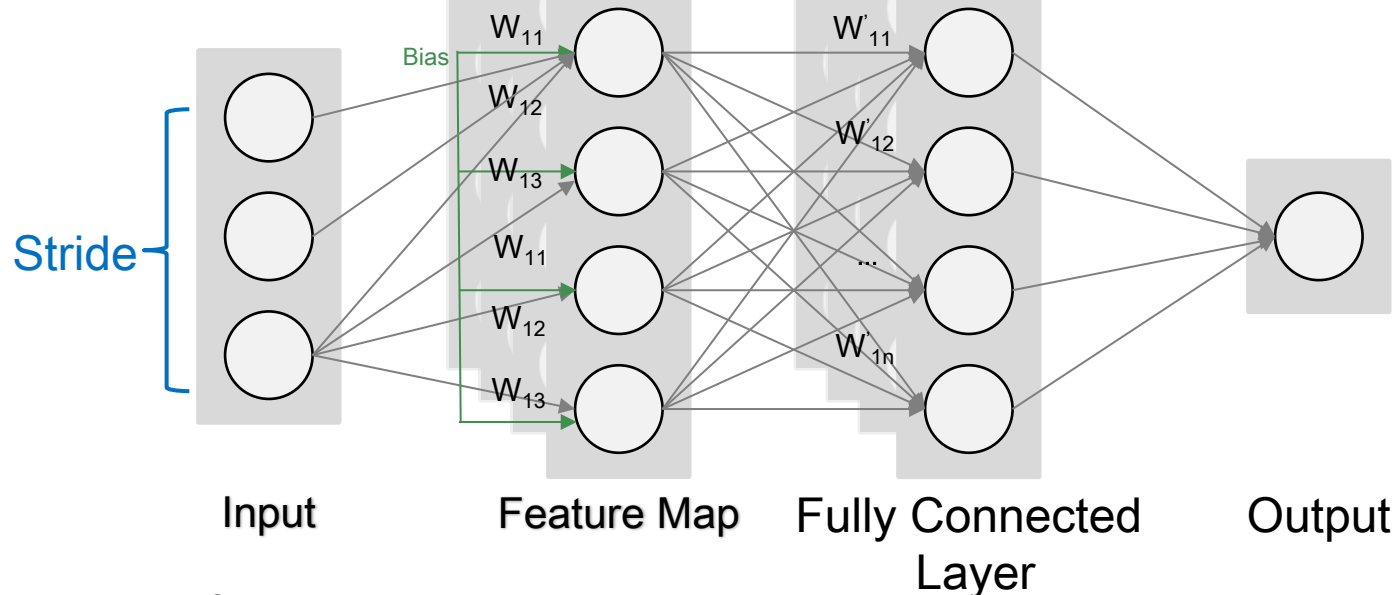
CONVOLUTIONAL NN

Neural Nets to CNN



Quelle: <http://cs231n.github.io/convolutional-networks/>

Introduction – Neural Nets



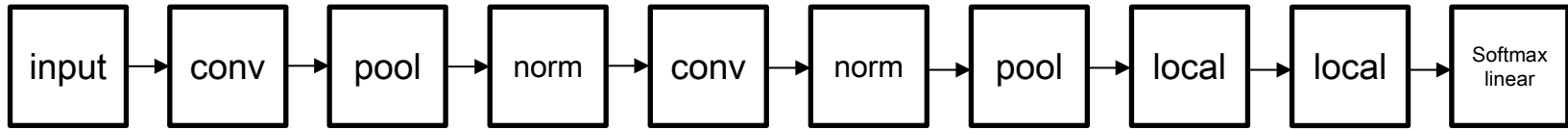
Quelle: <http://cs231n.github.io/convolutional-networks/>

TensorFlow

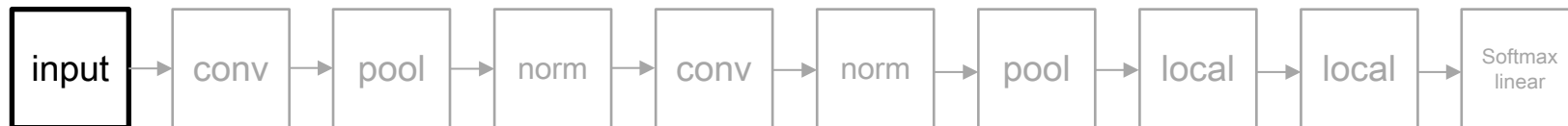


- Developed by Google Brain Team
- Use cases
 - Handwritten patterns, image recognition, Word2Vec
- Input data
 - Audio, image, text
- Used techniques
 - Linear classifiers, NN

Structure of the CNN we used



Input layer



Convolutional layer - Filter

Input

| | | | | | | |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 2 | 1 | 0 | 0 | 2 | 0 |
| 0 | 2 | 2 | 2 | 0 | 2 | 0 |
| 0 | 1 | 2 | 0 | 2 | 0 | 0 |
| 0 | 2 | 1 | 2 | 1 | 2 | 0 |
| 0 | 1 | 2 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Filter

| | | |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 2 | 1 |
| 0 | 2 | 2 |

Bias

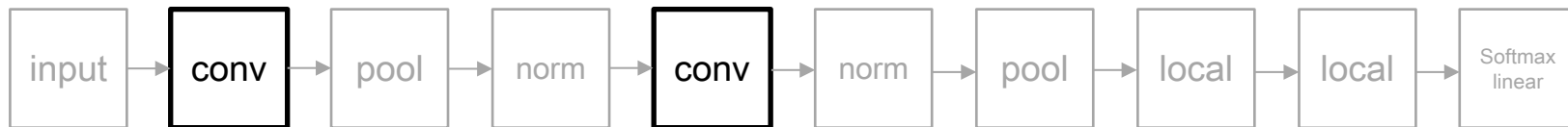
| |
|---|
| 1 |
|---|

Output

| |
|---|
| 1 |
|---|

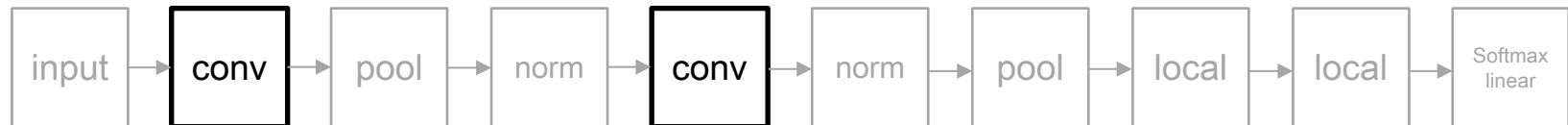
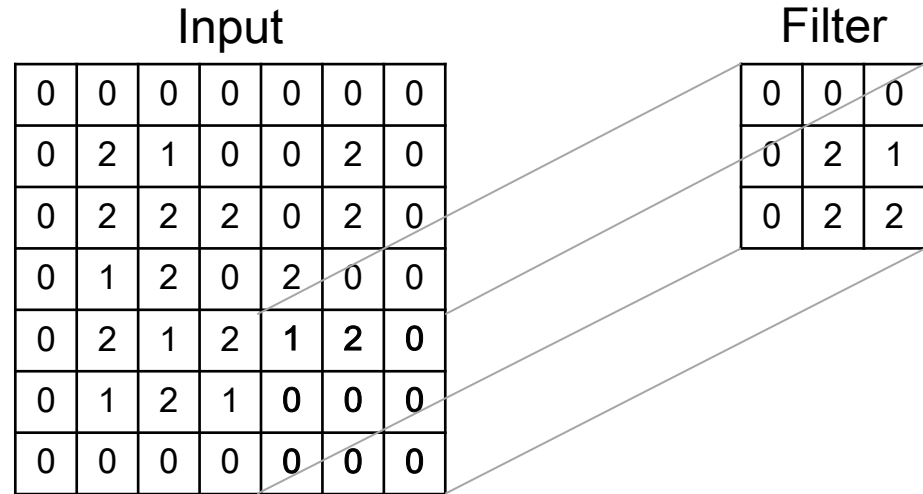
$$\begin{aligned}
 &1*0 + 2*0 + 0*0 + \\
 &0*0 + 0*2 + 0*1 + \\
 &0*0 + 0*2 + 0*2 + \\
 &1 = 1
 \end{aligned}$$

<http://cs231n.github.io/convolutional-networks/>



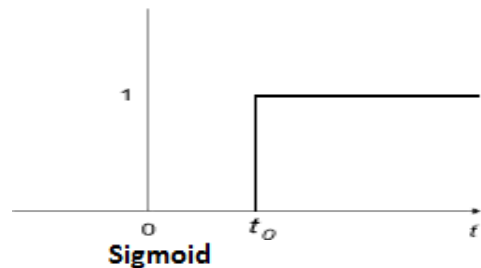
Convolutional layer - Parameters

- Input volume size
- Number of filters
- Filter size
- Zero padding

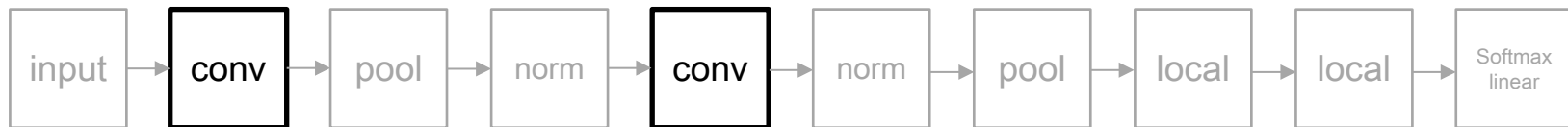
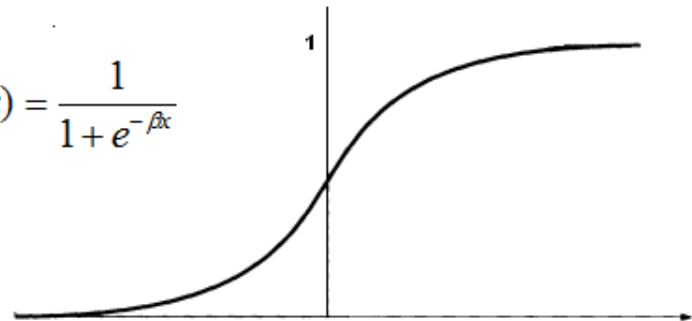


Convolutional layer – Activation function

- Sigmoid
 - Not telling in which direction should we move in.
 - Non-differentiability at certain points

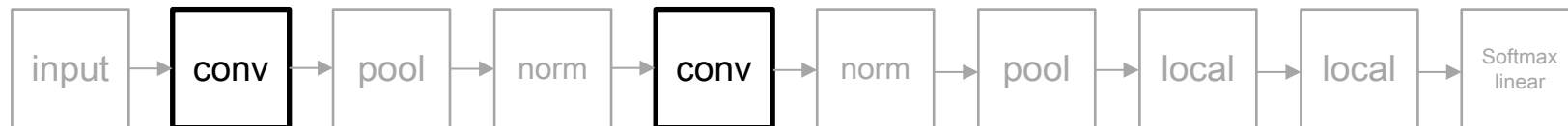


$$f(x) = \frac{1}{1 + e^{-\beta x}}$$

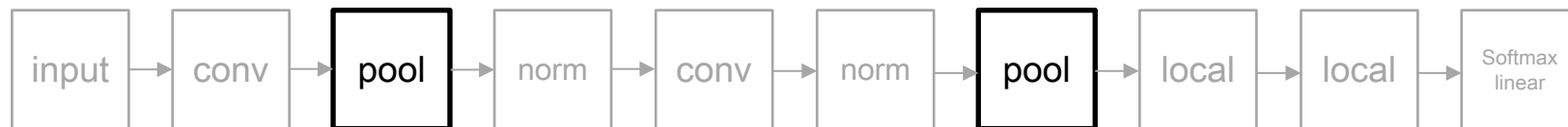


Convolutional layer – Activation function

- Rectified linear
 - *Element Wise*: $\max(0, x)$
 - Leaky ReLU
 - If $x < 0$, Output = $0.01x$.
 - Non-zero gradient when the input is negative

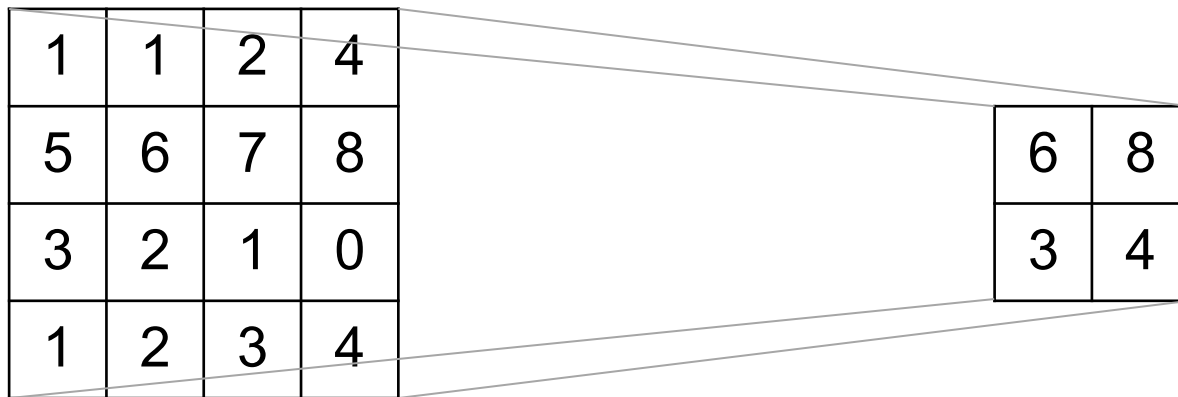


Pool layer

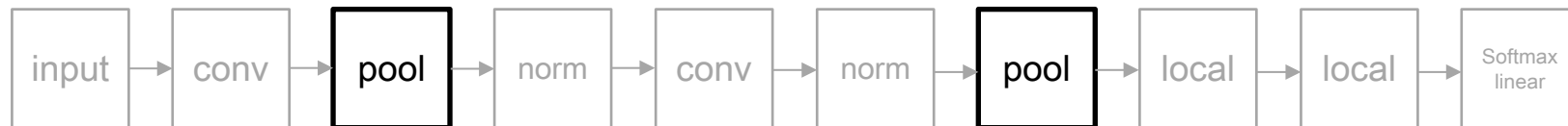


Pool layer – Max pooling

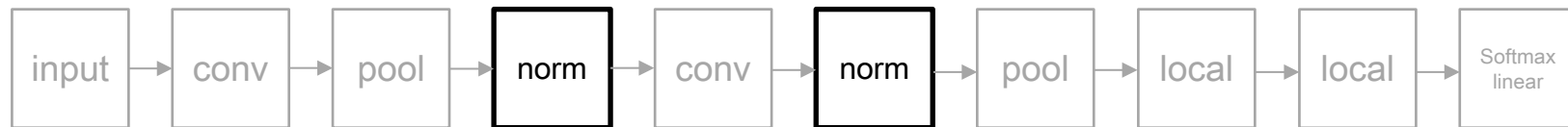
- Reduce the spatial dimension of an image



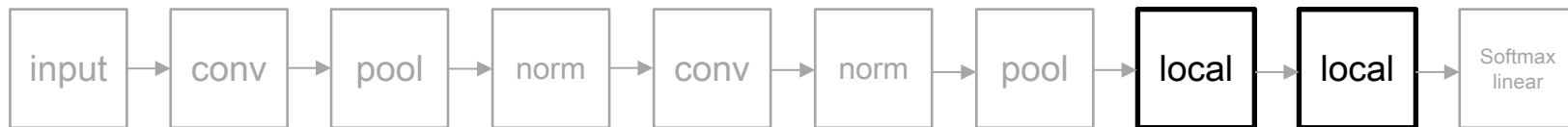
<http://cs231n.github.io/convolutional-networks/>



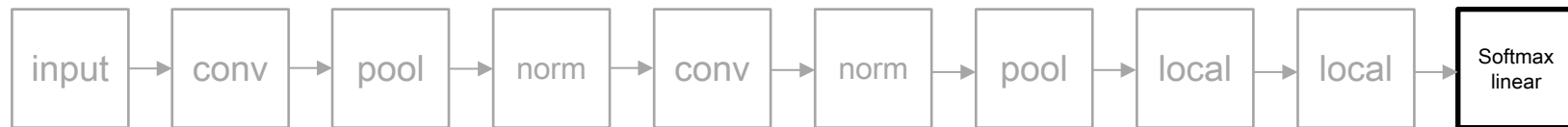
Norm layer



Local layer

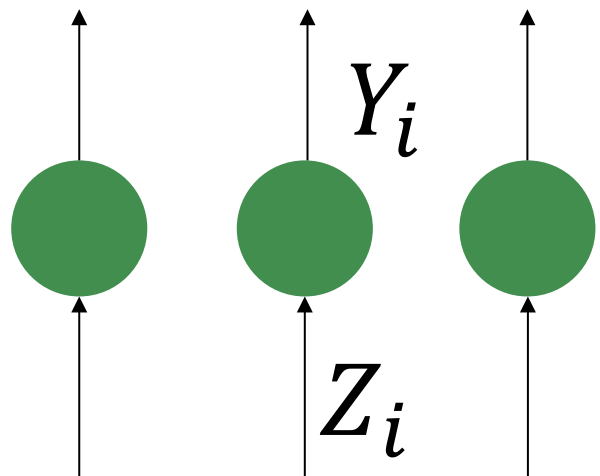


Softmax-linear layer



Softmax Output Function

- Soft continuous version of Max Function
- Forces $\sum(\text{Output of NN}) = 1$.


$$Y_i = \frac{e^{Z_i}}{\sum_{j \in \text{group}} e^{Z_j}}$$

Derivative Softmax

- $\frac{\delta Y_i}{\delta Z_i} = Y_i (1 - Y_i)$
- Nice Simple derivative
- Even though Y_i depends of Z_i ,
 - Derivative
 - for an individual neuron
 - of an I/P in respect to O/P is just $Y_i (1 - Y_i)$

Cost Measure for Softmax Output Function

- $\mathcal{C} = - \sum_j t_j \log Y_j$
 - Negative log probability of correct answer
- Maximise the log probability of getting answer right



HYPERPARAMETERS

Learning Rate

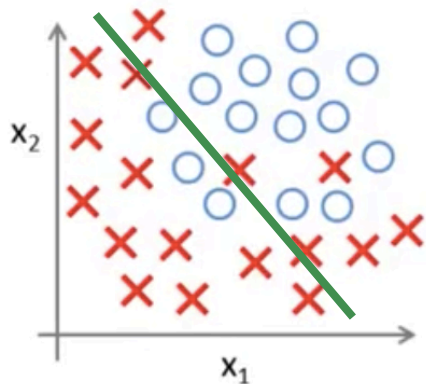
- how fast the network trains
- High learning rate
 - Convergence or global minimum finding is problem
- Low learning rate
 - High training times

Learning Rate decay

- Learning rate decay means the learning rate decreases over time
 - higher learning rate is well suited to get close to the global minimum
 - small learning rate is better at fine tuning the global minimum
- Several ways to do it:
 - Exponential decay, reduction by factor of n
 - GoogLeNet: function to decrease the learning rate by 4%

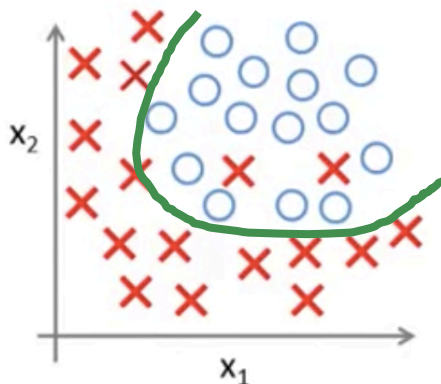
Overfitting or Underfitting

Example: Logistic regression

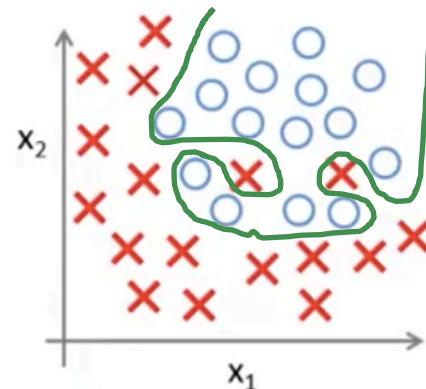


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(g = sigmoid function)



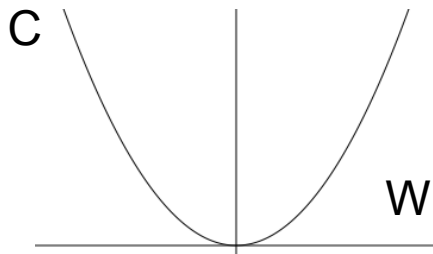
$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

Weight Penalty

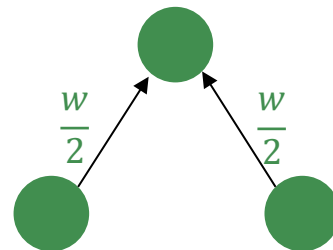
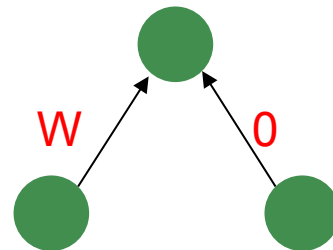
- Adding λ to penalise
 - Keeps weight small
 - Big error derivatives



- $C = E + \frac{\lambda}{2} \sum_{i=1} w_i^2$
- $\frac{\partial C}{\partial w_i} = \frac{\partial E}{\partial w_i} + \lambda w_i$
- When $\frac{\partial C}{\partial w_i} = 0$;
 - $w_i = -\frac{1}{\lambda} \frac{\partial E}{\partial w_i}$
 - So, at minimum of Cost function if $\frac{\partial E}{\partial w_i}$ is big, the weights are big

Weight Penalty - Advantages

- Preventing network from the weights it does not need
 - Don't have a lot of weights not doing anything
 - So output changes more slowly as input changes.
- Putting half the weight on each and not on one





PROBLEM

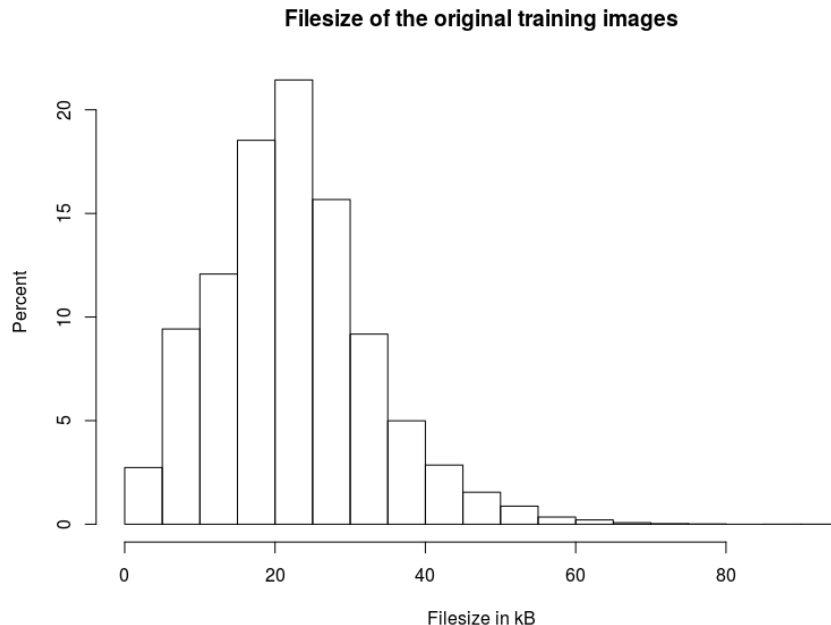
The data

- Images of cats and dogs
- File format is *.jpg
- Color space is RGB



Training data

- 25,000 images
 - 12,500 of dogs
 - 12,500 of cats
- Avg. file size
 - 22.34 kB



Test data

- 12,500 images
 - x of dogs
 - y of cats
 - $x + y = 12,500$

Process images

- Resize to $32 * 32 * 3$
- Convert to array
 - $25,000 * 3,073$



dog1.jpg



cat10.jpg

Process images

- Resize to $32 * 32 * 3$
- Convert to array
 - $25,000 * 3,073$



dog1.jpg



cat10.jpg





EVALUATION



AIMS



Aims

- Removing normalization layer



QUESTIONS

Quellen

- <http://cs231n.github.io/convolutional-networks/>
- https://www.tensorflow.org/tutorials/deep_cnn/
- Maas, Andrew L., Awni Y. Hannun, and Andrew Y. Ng. "Rectifier nonlinearities improve neural network acoustic models." *Proc. ICML*. Vol. 30. No. 1. 2013.