

Neural Networks

Saurabh Mathur

Machine Learning

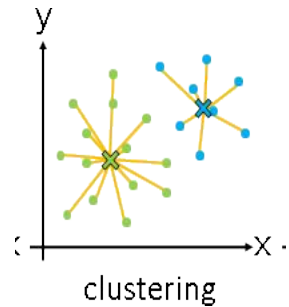
Programming computers to learn from experience should eventually eliminate the need for much of this detailed programming effort.

- Arthur Samuel

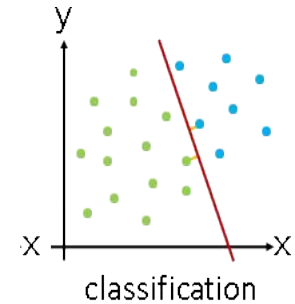
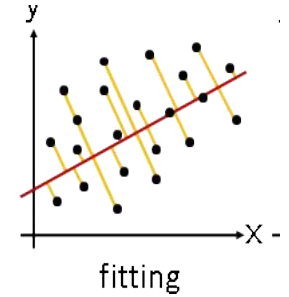
Machine learning research seeks to develop computer systems that automatically improve their performance through experience.

- Tom Mitchell

Problem formulation



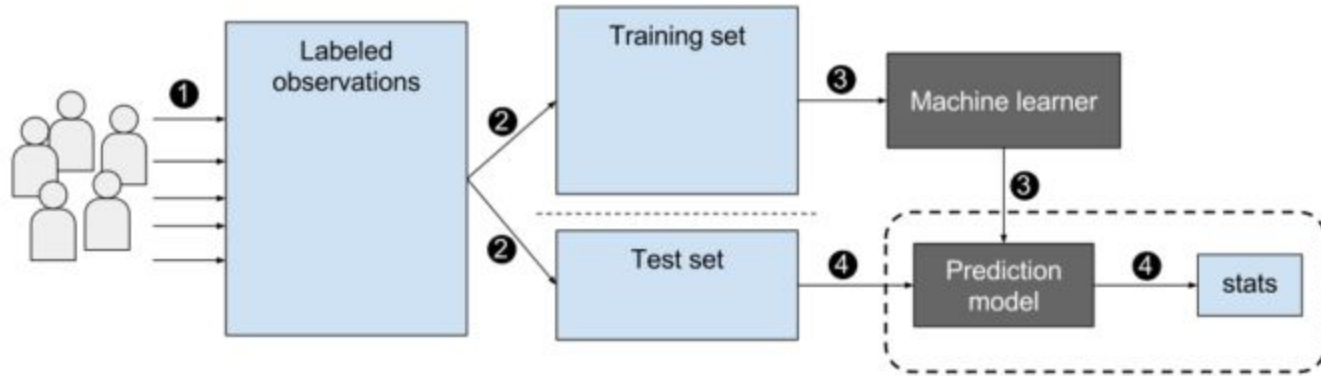
Task	Experience	Performance
Predict house prices	House attributes and prices	Closeness to correct price (RMSE)
Grouping tweets together	Text content of tweets	Coherent groupings (Within group similarity)
Classify new images as either Cat or Dog	Images of cats & dogs	Accuracy of classification, Cross-entropy



Machine Learning Tasks

- Supervised Learning (We have examples and labels)
 - Classification
 - Regression
- Unsupervised Learning (We only have examples)
 - Topic modeling
 - Clustering
- Semi-Supervised Learning (Some examples have labels)

Supervised Learning



Source: <https://blogs.nvidia.com/blog/2018/08/02/supervised-unsupervised-learning/>

Supervised learning

- Model
 - A simplified view of the world.
 - Encodes assumptions about the problem.
 - How was the data generated ?
- Learning Algorithm
 - Estimates specifics of the model from data.
- Measure of correctness
 - Loss Function
 - Log Likelihood

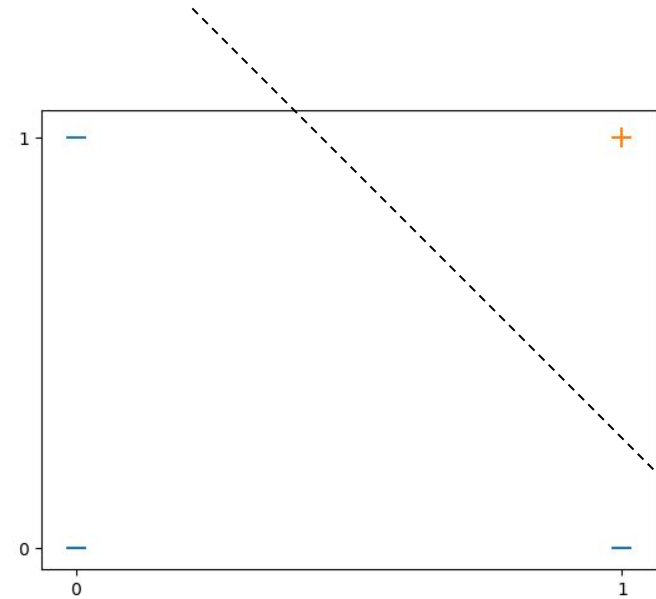
Example: Classification



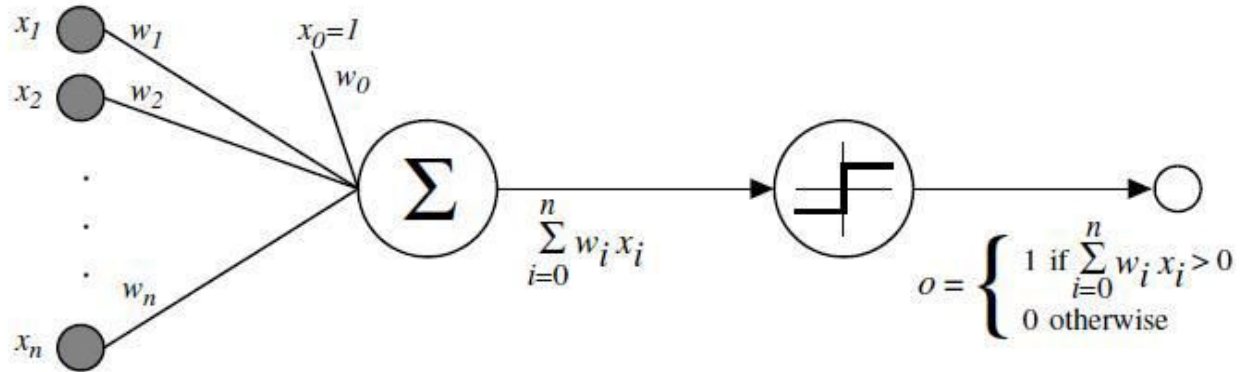
x_1	x_2	y (AND function)
0	0	-1
0	1	-1
1	0	-1
1	1	+1

Example: Classification

- Model
 - A line separates the examples : $y = ax + b$
 - Examples above the line are + others are -
- Measure of correctness
 - Accuracy = #correct / total
 - Error = sum(prediction - correct)
- Learning Algorithm ?

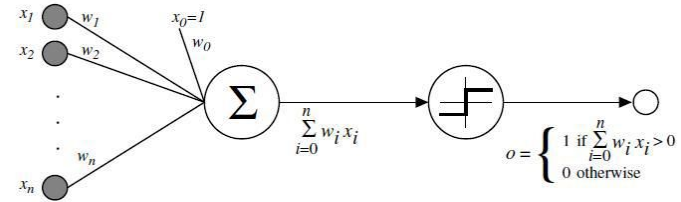


The Perceptron



The Learning Algorithm

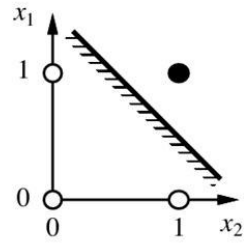
1. Set w to $[0, 0, \dots, 0]$
2. For each example x and output y in dataset:
 - 2.1. Compute prediction, $p = \text{Sign}(w \cdot x)$
 - 2.2. If $p \neq y$
 - 2.2.1. For i in $1, \dots, M$
Update, $w = w + x * y$



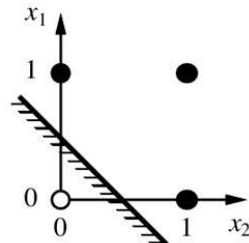
Limitations of perceptrons

- Learning converges only if data is linearly separable.
- Does not model uncertainty - outputs either 0 or 1.
- Input features should be well-defined.

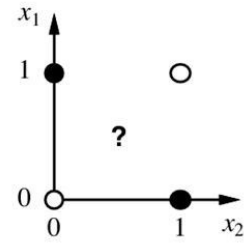
Linear separability



x_1 and x_2

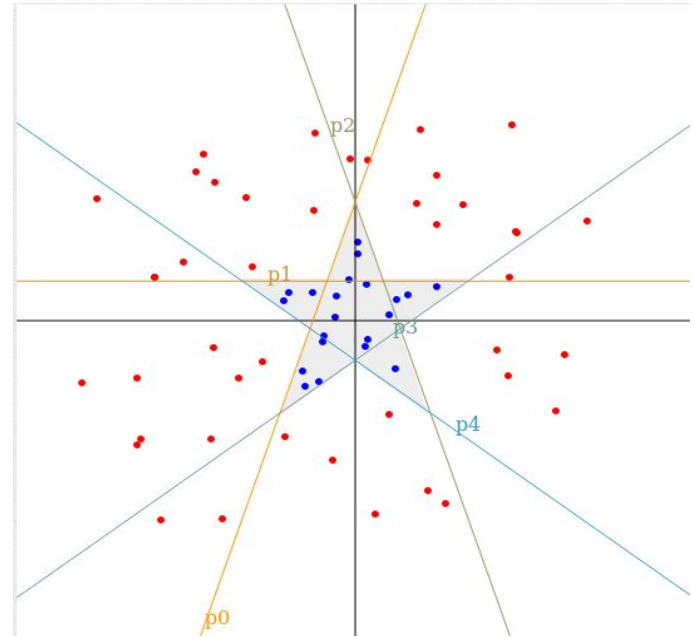
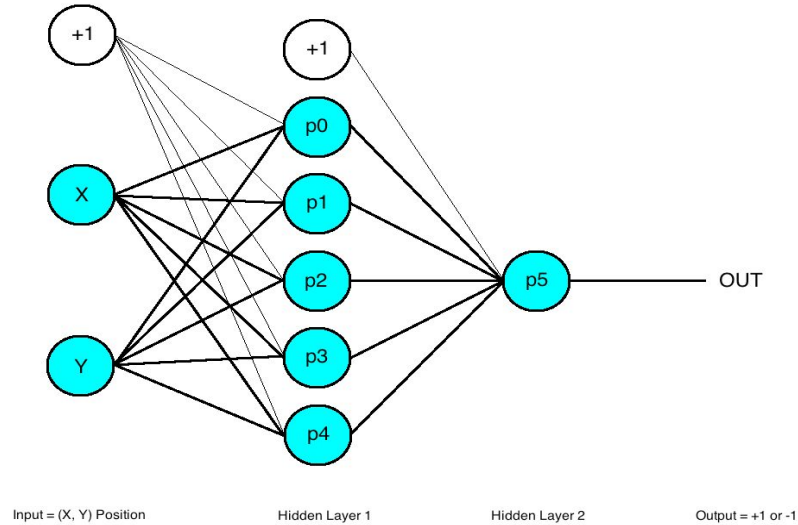


x_1 or x_2



x_1 xor x_2

Multi-layered Perceptron

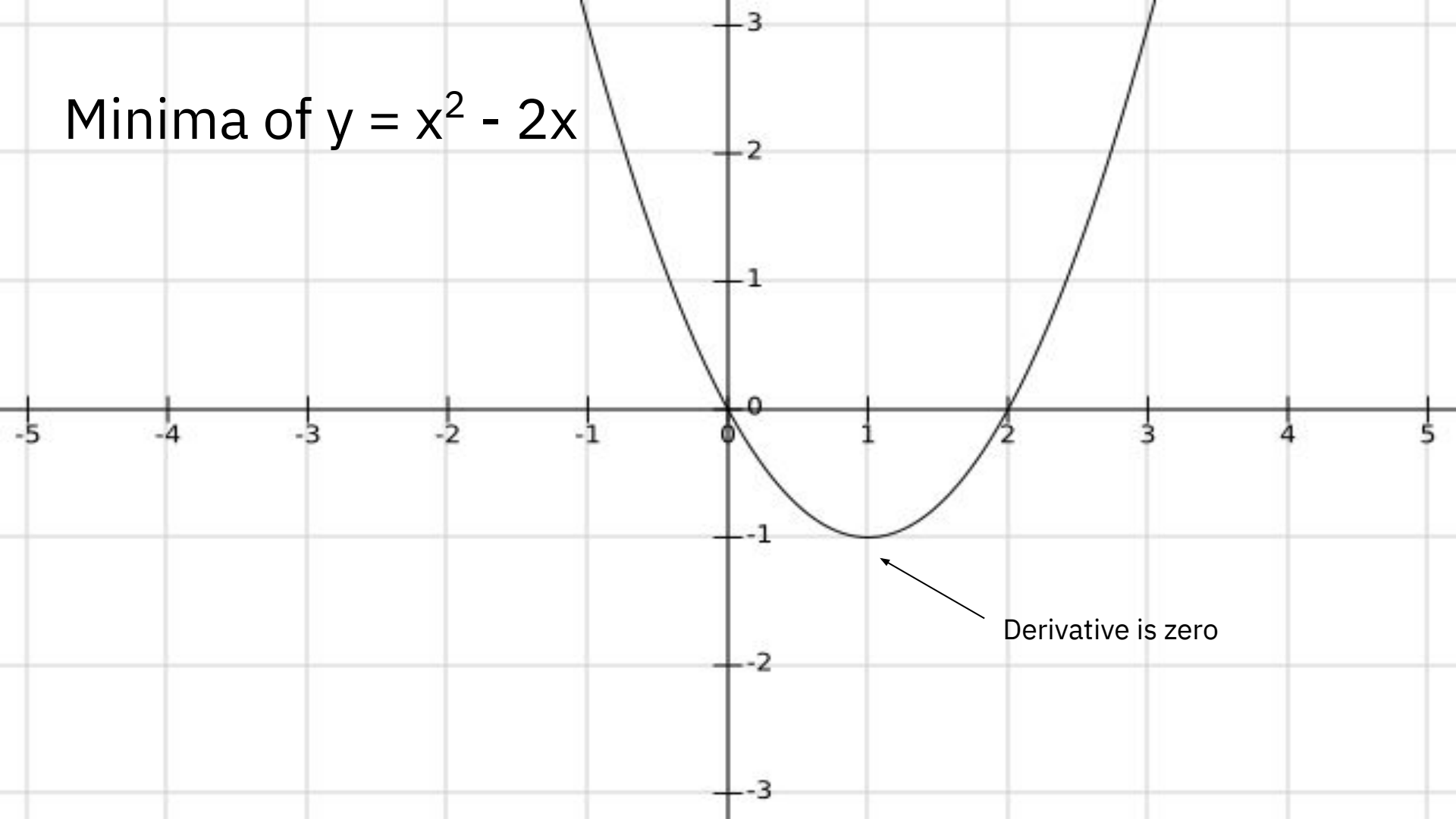


Source: <https://www.cs.utexas.edu/~teammco/misc/mlp/>

Learning as error-minimization

- Empirical risk minimization
- Hypothesis $\mathbf{h} : X \rightarrow y$
- Find \mathbf{h} by minimizing an error function over training data.

Minima of $y = x^2 - 2x$



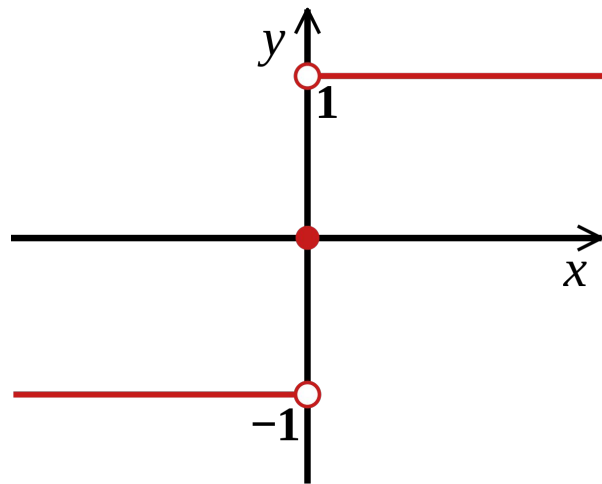
Requirements for direct error-minimization

- Differentiable expression for loss function
- Closed-form solution for minimum (derivative = 0)

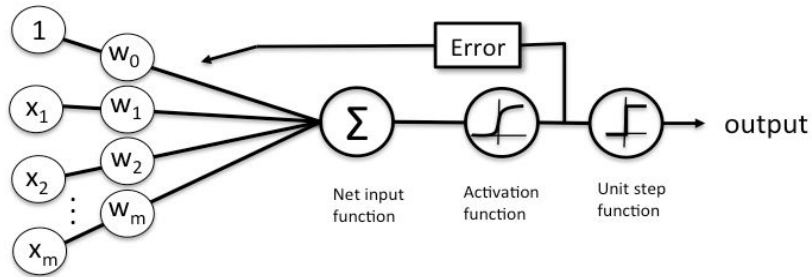
Perceptron Loss function

Perceptron uses the Sign function for prediction.

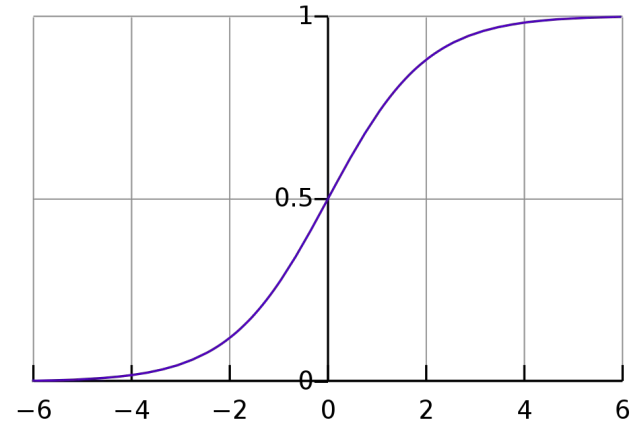
Is this differentiable ?



The Sigmoid Function



Schematic of a logistic regression classifier.



$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}.$$

Logistic Regression

Replace Sign with Sigmoid.

$$Loss(w) = \sum_{i=1}^m y^{(i)} \log P(y = 1) + (1 - y^{(i)}) \log P(y = 0)$$

Differentiate and solve for w

$$Loss(w) = \sum_{i=1}^m y^{(i)} \log P(y = 1) + (1 - y^{(i)}) \log P(y = 0)$$

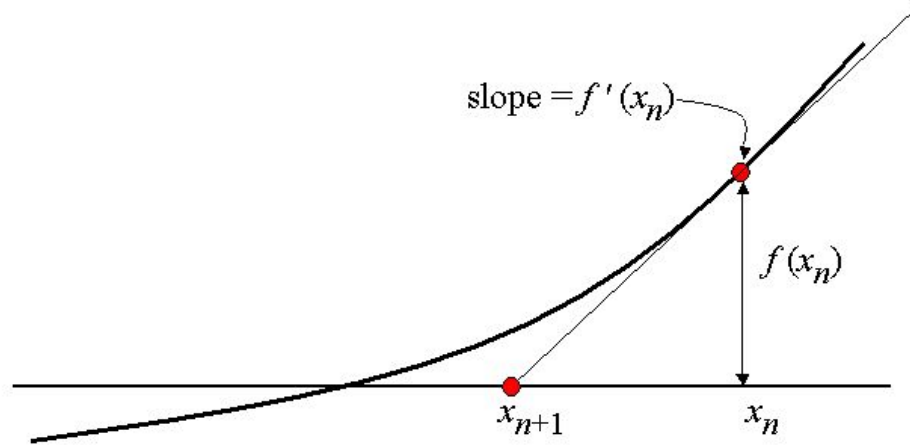
$$Loss(w) = \sum_{i=1}^m y^{(i)} \log \left(\frac{1}{1 + e^{-w^T x}} \right) + (1 - y^{(i)}) \log \left(1 - \frac{1}{1 + e^{-w^T x}} \right)$$

$$\frac{\partial Loss}{\partial w} = \sum_i \left[y_i - \frac{1}{1 + \exp(x'_i w)} \right] x_i = 0$$

There is no closed-form solution

$$\frac{\partial Loss}{\partial w} = \sum_i \left[y_i - \frac{1}{1 + \exp(x'_i w)} \right] x_i = 0$$

Newton's Method to find zero-point

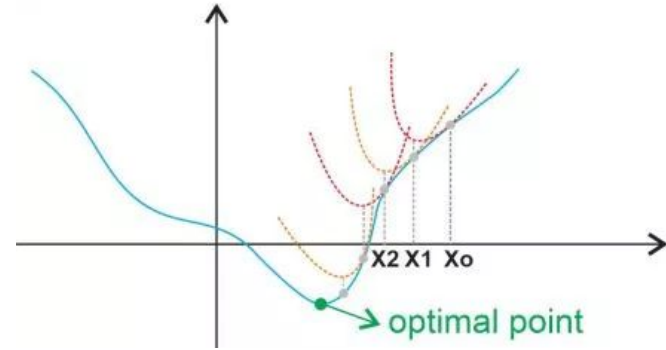


$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Newton's Method for minimization

- Find zero-point of derivative.

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

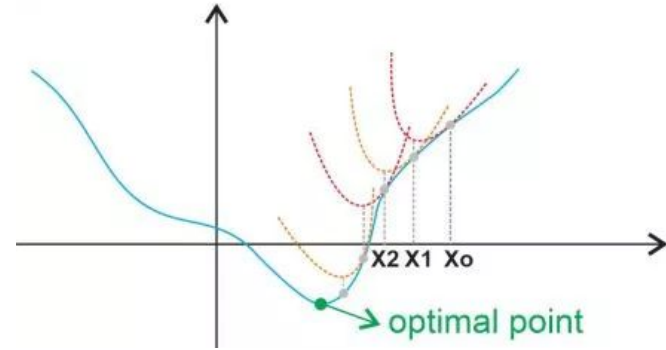


Gradient Descent for minimization

- Second derivative is expensive

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

$$x_{n+1} = x_n - \eta f'(x_n)$$

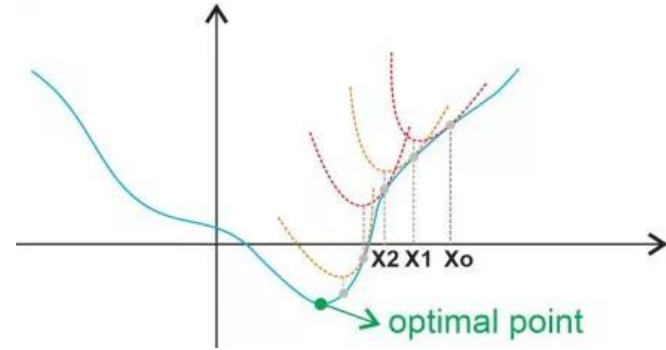


Stochastic Gradient Descent

- Approximate derivative using single example

$$x_{n+1} = x_n - \eta f'(x_n)$$

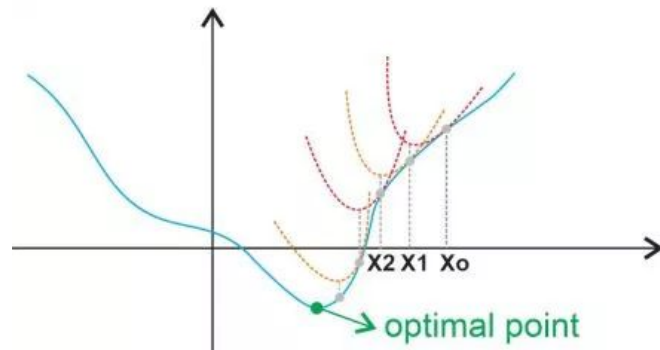
$$\sum_{i=1}^N \frac{\partial \text{Loss}(x_i, y, w_n)}{\partial w_n} \simeq N * \frac{\partial \text{Loss}(x_i, y, w_n)}{\partial w_n}$$



Extensions of SGD : Selecting LR

- ADAM
- RMSProp
- Adagrad
- Adadelta

<http://ruder.io/optimizing-gradient-descent/index.html>



"Shallow" Learning on images

Pixel values → Feature extraction → Shallow learning → Prediction

Histograms (HOG)
Feature Points (SIFT)
Visual words

Logistic Regression
Perceptron
SVM

Object class
Object bounding box
Depth

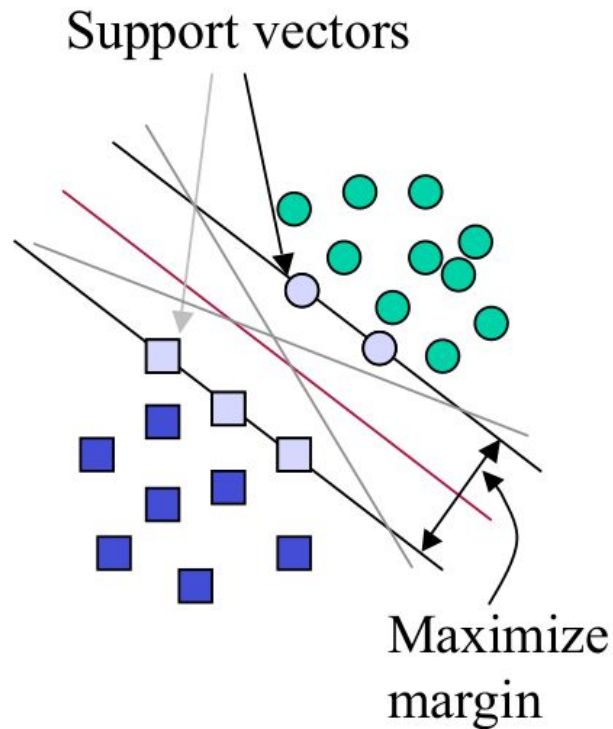
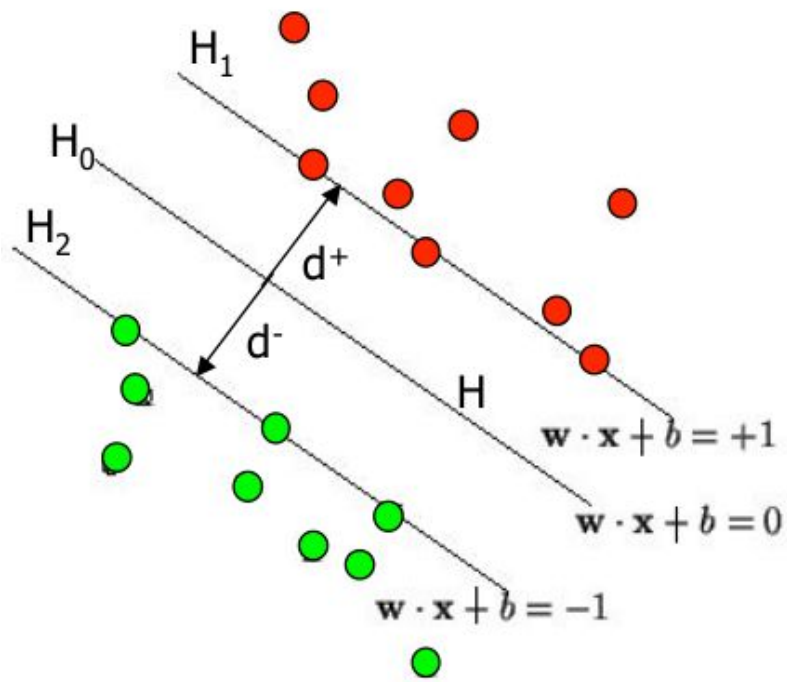
Activation Functions

	Propagation	Back-propagation
Sigmoid	$y_s = \frac{1}{1+e^{-x_s}}$	$\left[\frac{\partial E}{\partial x}\right]_s = \left[\frac{\partial E}{\partial y}\right]_s \frac{1}{(1+e^{x_s})(1+e^{-x_s})}$
Tanh	$y_s = \tanh(x_s)$	$\left[\frac{\partial E}{\partial x}\right]_s = \left[\frac{\partial E}{\partial y}\right]_s \frac{1}{\cosh^2 x_s}$
ReLU	$y_s = \max(0, x_s)$	$\left[\frac{\partial E}{\partial x}\right]_s = \left[\frac{\partial E}{\partial y}\right]_s \mathbb{I}\{x_s > 0\}$

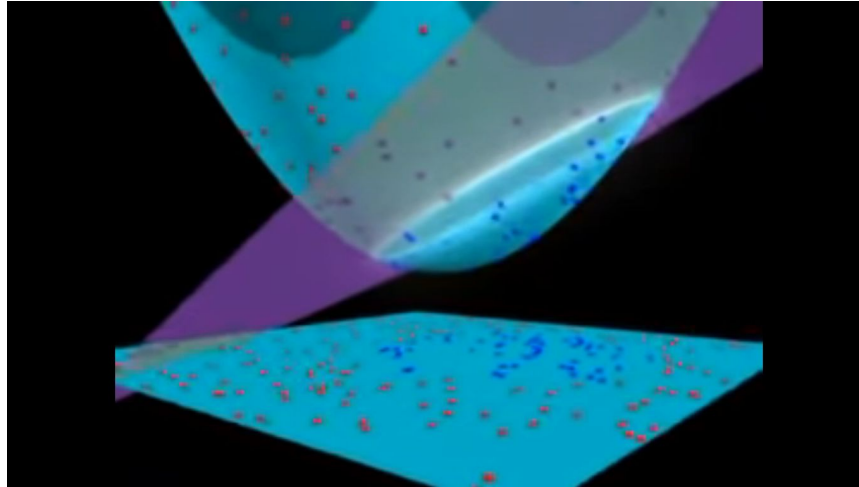
Loss Functions

Problem Type	Output Type	Final Activation Function	Loss Function
Regression	Numerical value	Linear	Mean Squared Error (MSE)
Classification	Binary outcome	Sigmoid	Binary Cross Entropy
Classification	Single label, multiple classes	Softmax	Cross Entropy
Classification	Multiple labels, multiple classes	Sigmoid	Binary Cross Entropy

Support Vector Machines

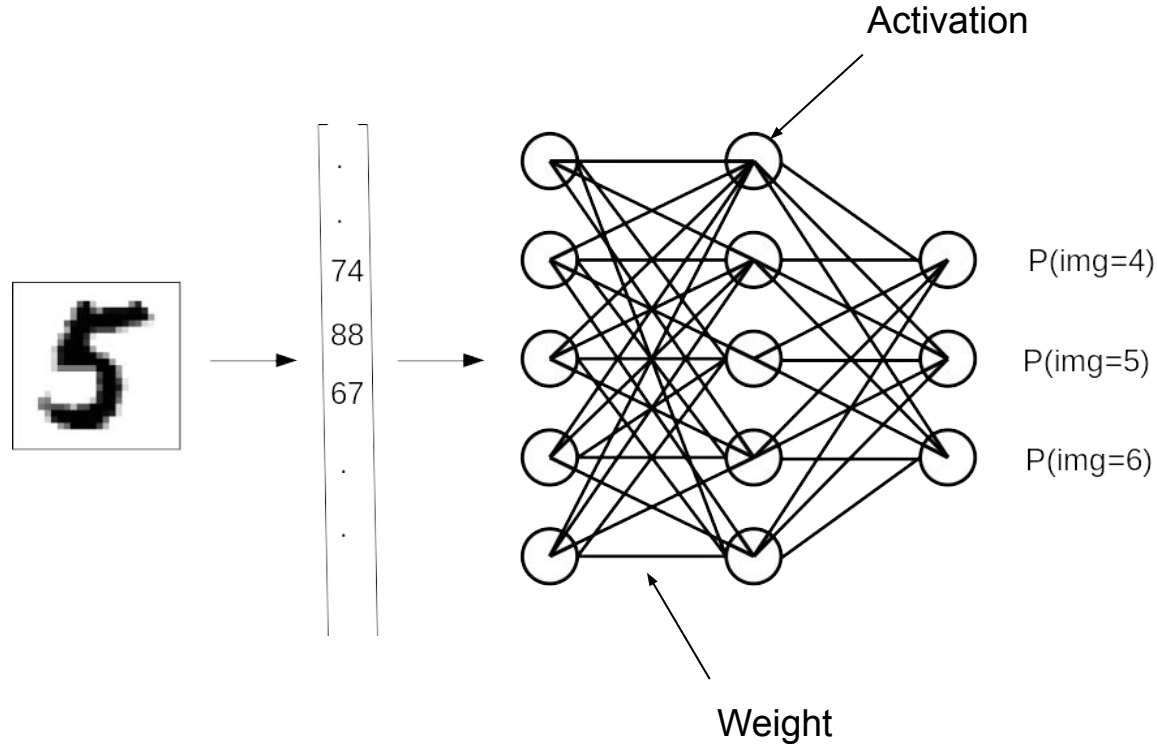


Support Vector Machines: Kernels



Source: <https://www.youtube.com/watch?v=3liCbRZPrZA>

Neural Network for image classification



Backpropagation - Key Ideas

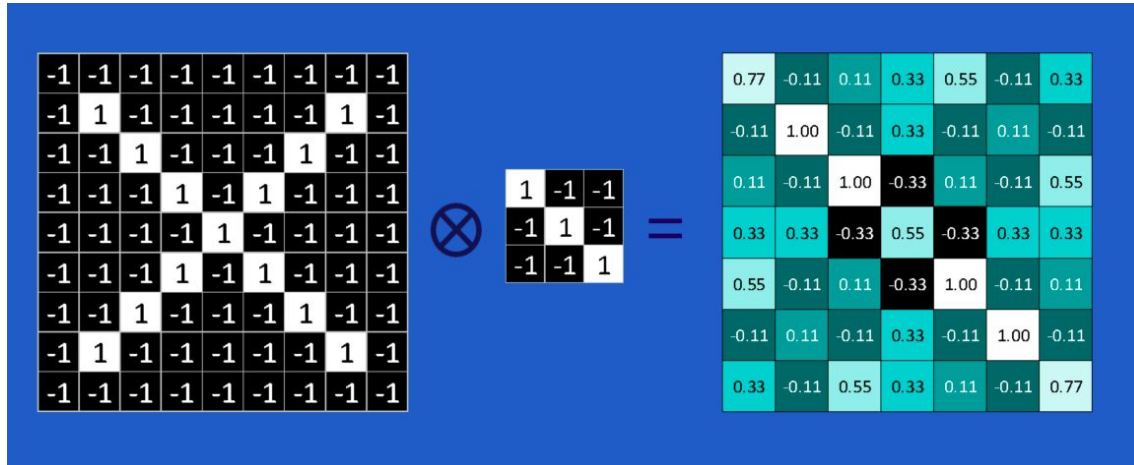
- Forward-propagation
- Use chain-rule to expand derivatives.
- Save intermediate derivatives (dynamic programming)

<https://google-developers.appspot.com/machine-learning/crash-course/backprop-scroll/>

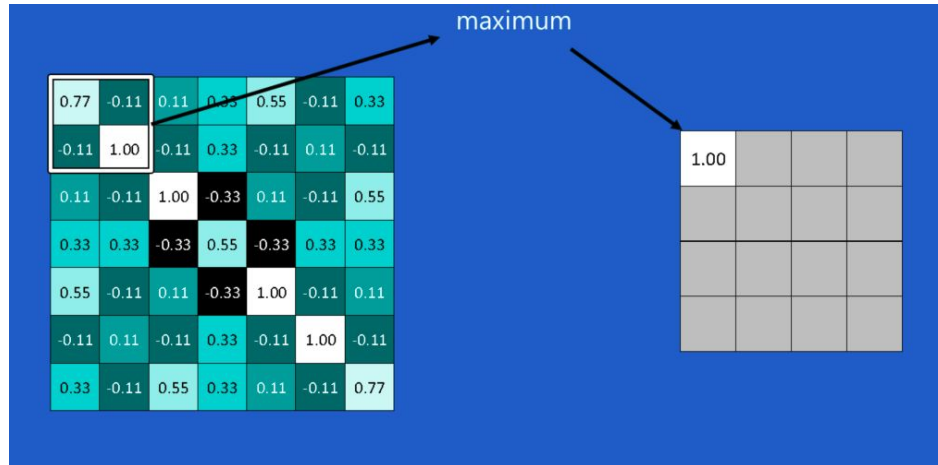
Convolutional Neural Networks



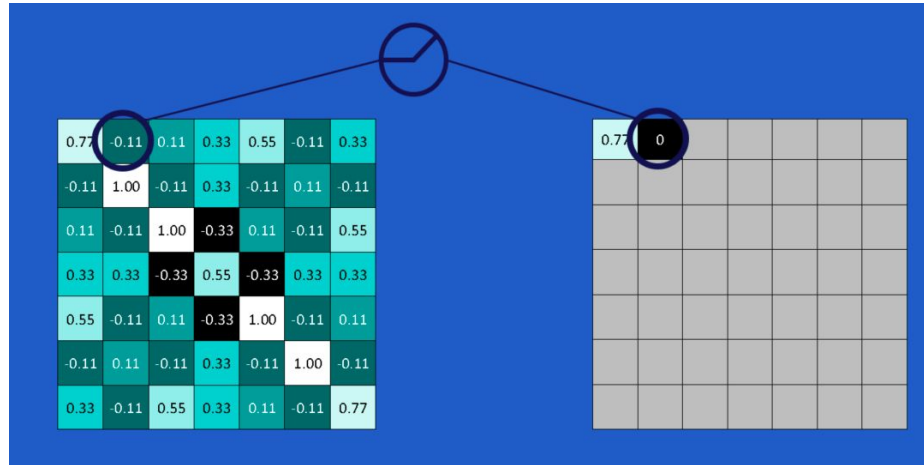
Convolution Layer



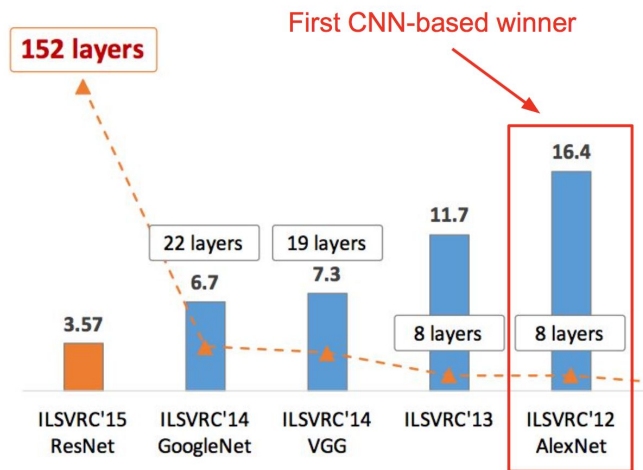
Pooling Layer



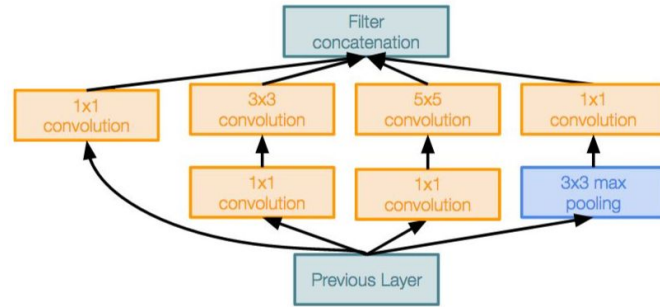
Rectified Linear Unit



Key Idea - More Layers

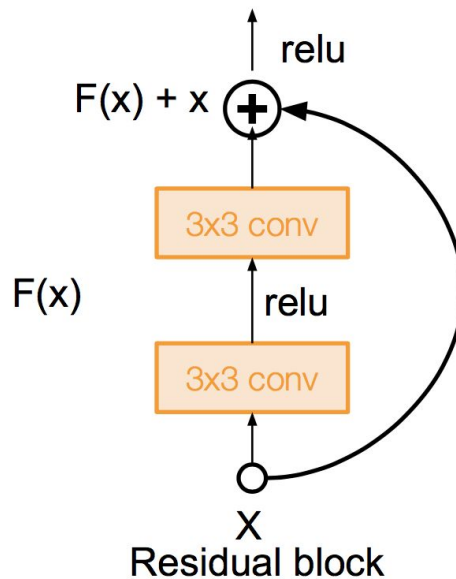


Key Idea - Parallel Layers

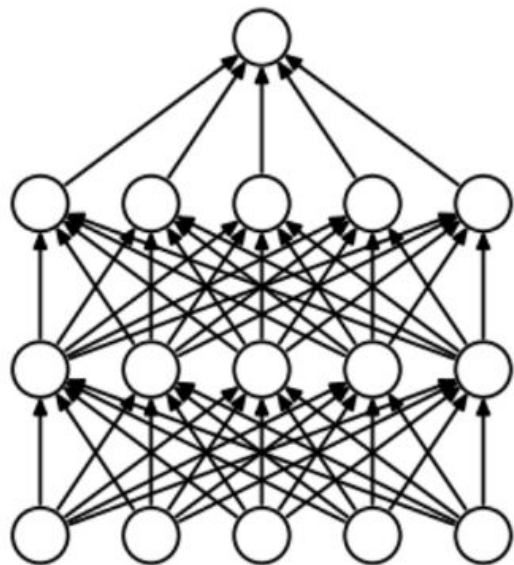


Inception module

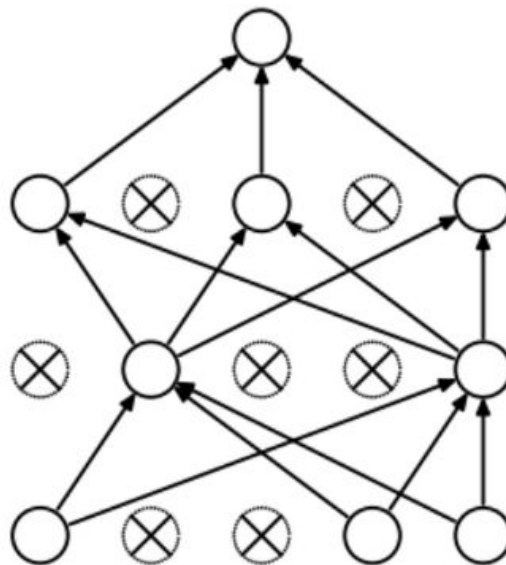
Key Idea - Skip Connections



Key Idea - Dropout

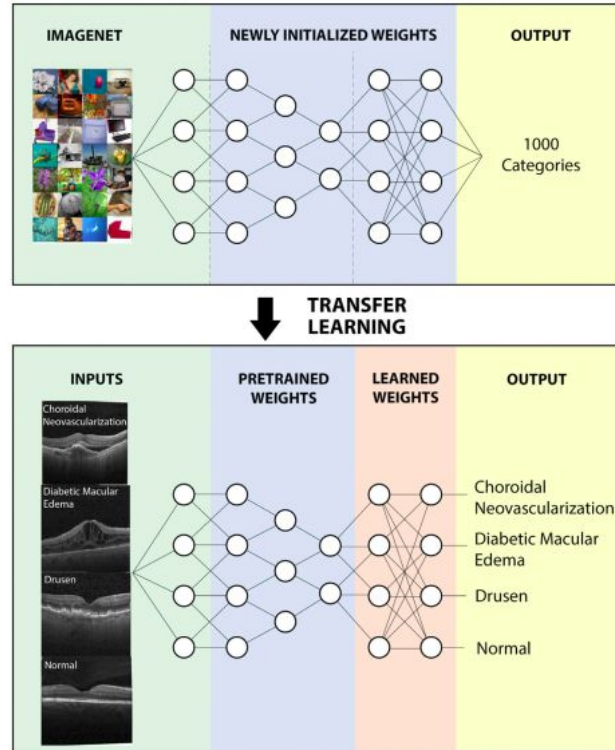


(a) Standard Neural Net

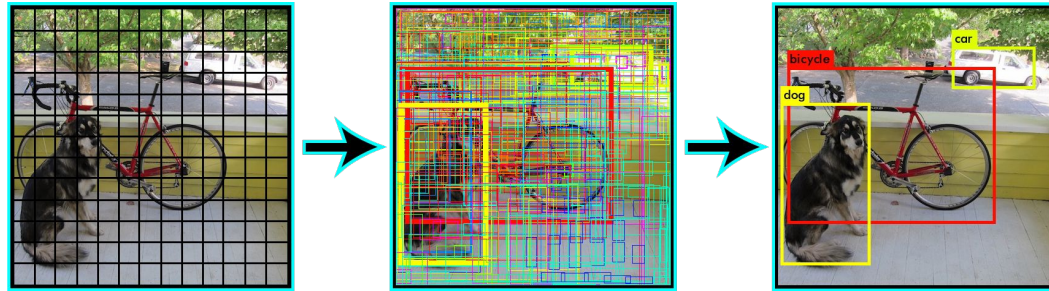


(b) After applying dropout.

Key Idea - Transfer Learning



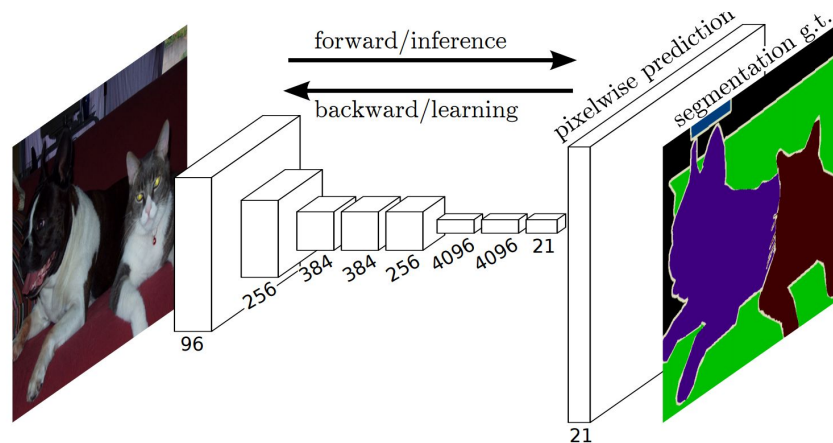
YOLO - Very Fast Object Detection



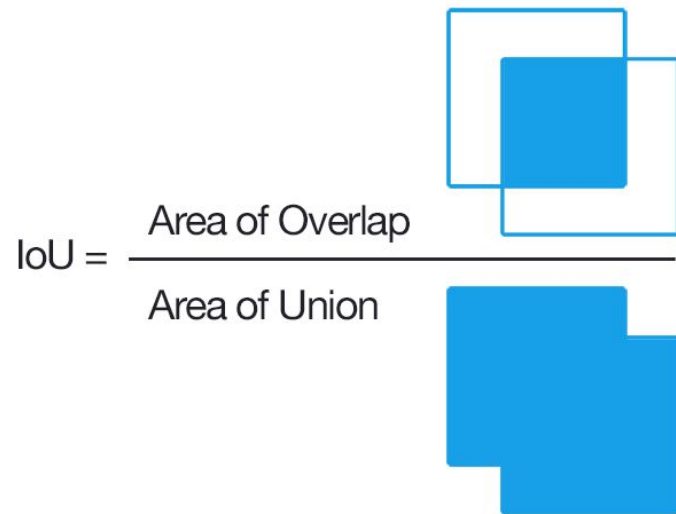
Quantization

Non-Maximal Suppression

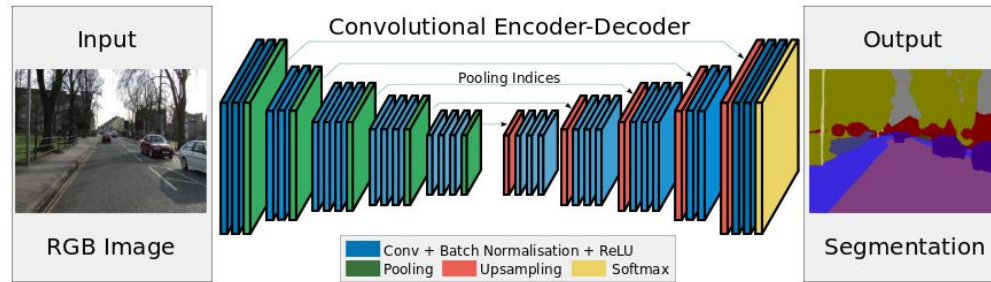
FCN - Semantic Segmentation



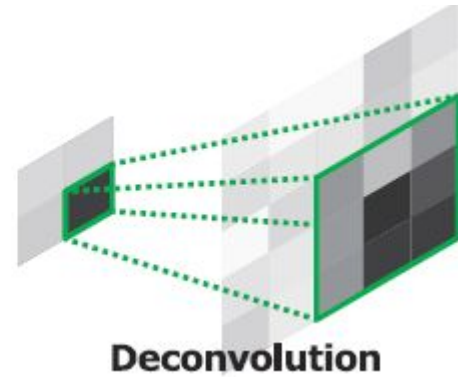
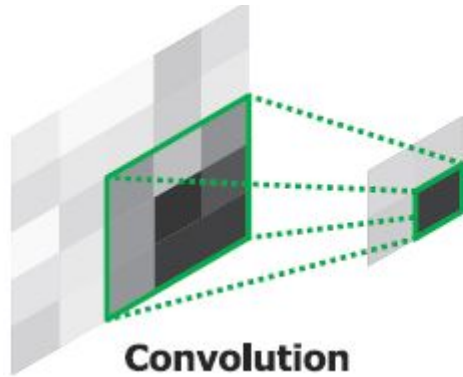
Intersection over Union



SegNet : Trainable Encoder-Decoder Architecture



Key Idea - Transposed Convolutions for Upsampling



Key Idea - Transposed Convolutions for Upsampling

$$\begin{pmatrix} w_{0,0} & 0 & 0 & 0 \\ w_{0,1} & w_{0,0} & 0 & 0 \\ w_{0,2} & w_{0,1} & 0 & 0 \\ 0 & w_{0,2} & 0 & 0 \\ w_{1,0} & 0 & w_{0,0} & 0 \\ w_{1,1} & w_{1,0} & w_{0,1} & w_{0,0} \\ w_{1,2} & w_{1,1} & w_{0,2} & w_{0,1} \\ 0 & w_{1,2} & 0 & w_{0,2} \\ w_{2,0} & 0 & w_{1,0} & 0 \\ w_{2,1} & w_{2,0} & w_{1,1} & w_{1,0} \\ w_{2,2} & w_{2,1} & w_{1,2} & w_{1,1} \\ 0 & w_{2,2} & 0 & w_{1,2} \\ 0 & 0 & w_{2,0} & 0 \\ 0 & 0 & w_{2,1} & w_{2,0} \\ 0 & 0 & w_{2,2} & w_{2,1} \\ 0 & 0 & 0 & w_{2,2} \end{pmatrix}^T$$

Summary

