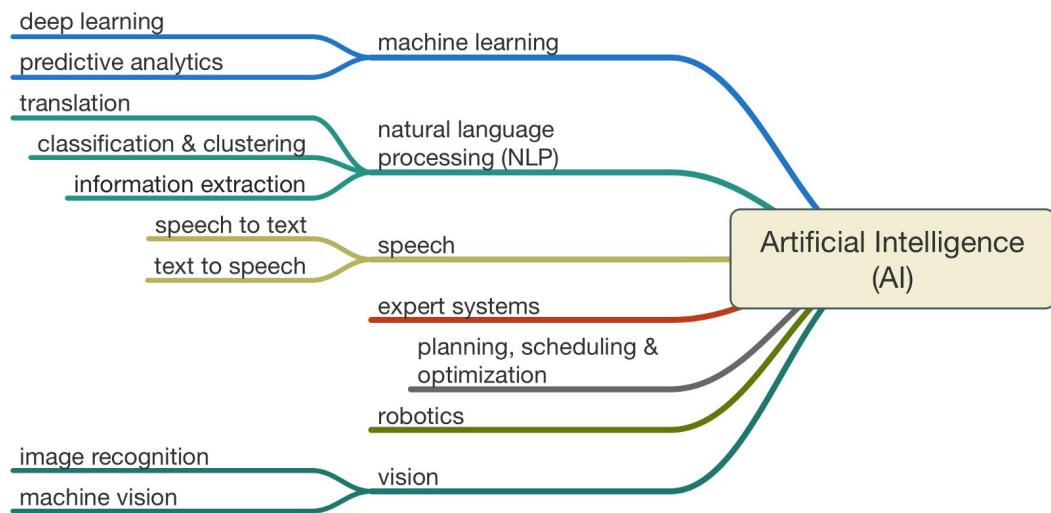


Deep Learning: Foundations and Practicing

Introduction to Deep Learning

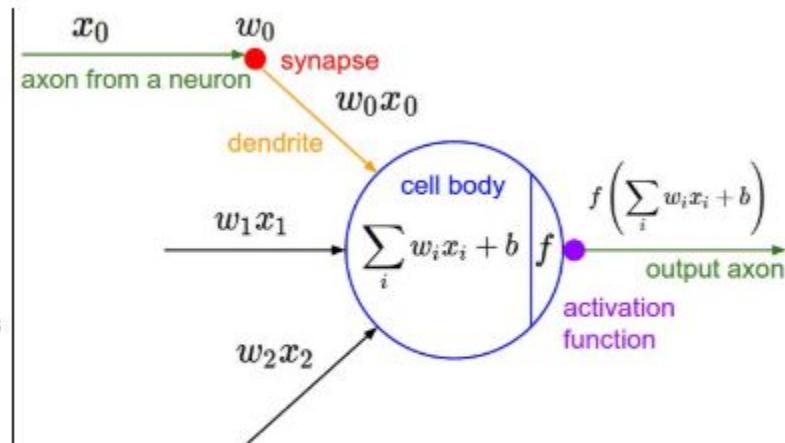
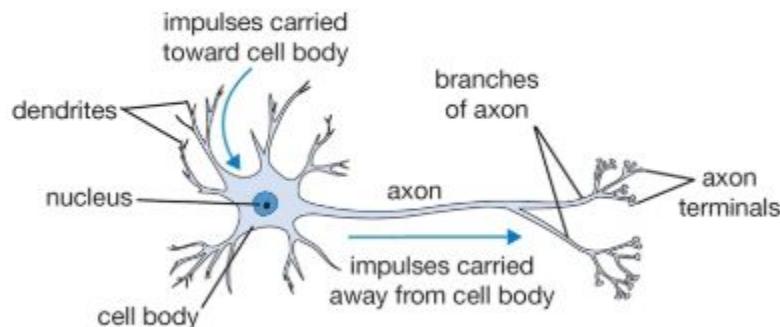
- One of the threads of IA and Machine Learning
- Automatic feature extractor





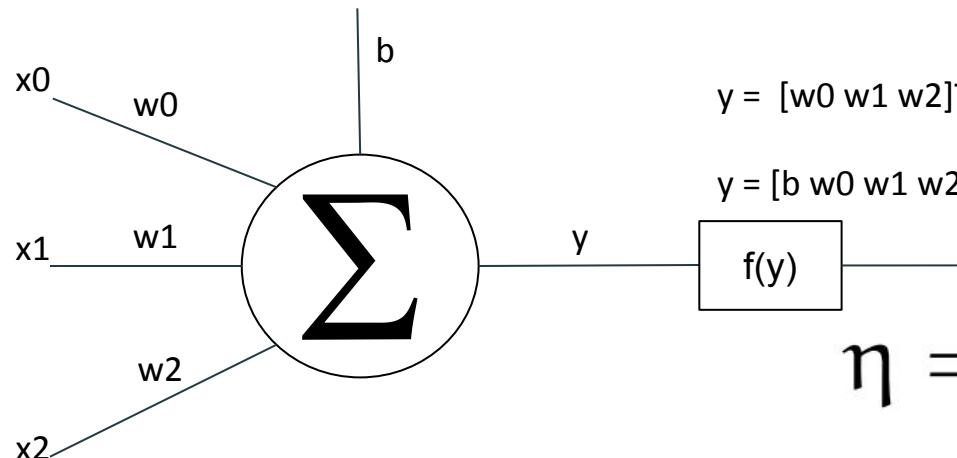


Perceptron (NN's basic unit)



A cartoon drawing of a biological neuron (left) and its mathematical model (right).

Neuron

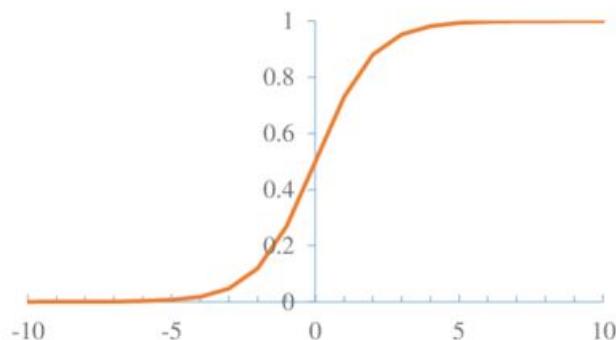


$$y = [w_0 \ w_1 \ w_2]^T X [x_0 \ x_1 \ x_2] + b$$

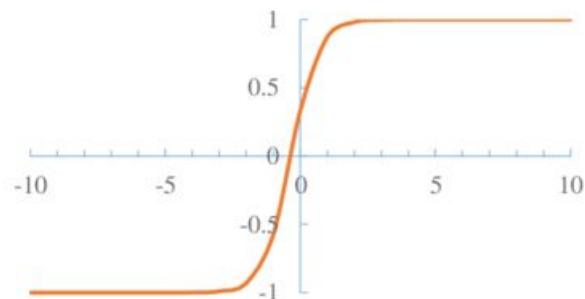
$$y = [b \ w_0 \ w_1 \ w_2]^T X [1 \ x_0 \ x_1 \ x_2]$$

$$\eta = \beta^T x + \beta_0$$

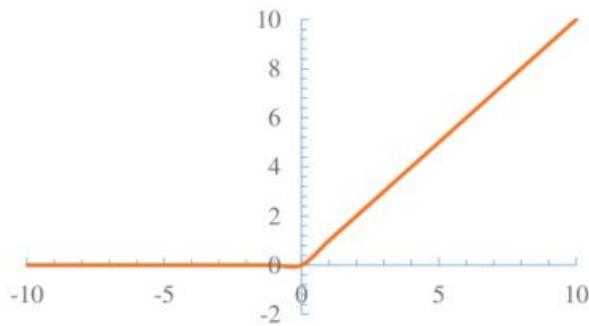
$$y = \eta + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$



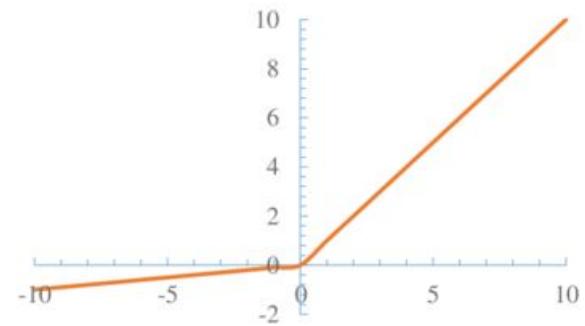
(a)



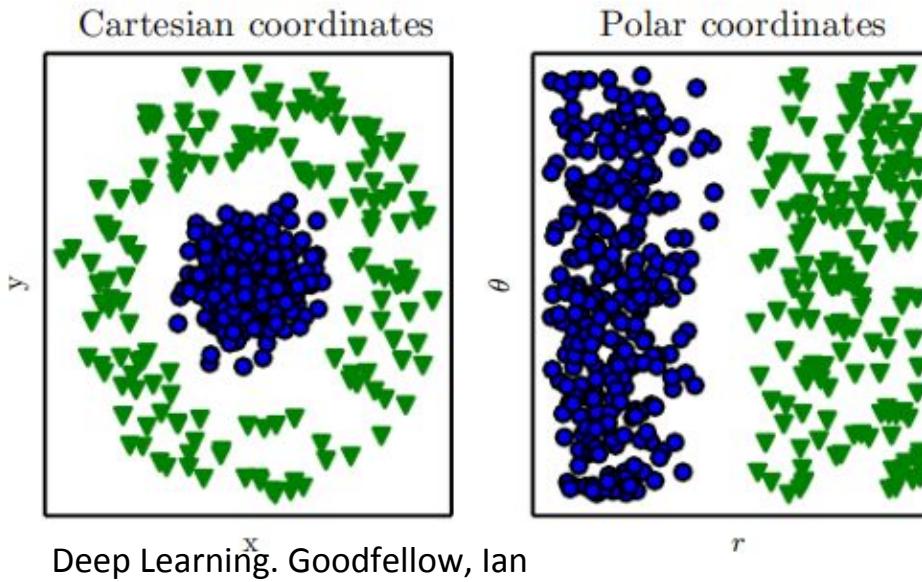
(b)

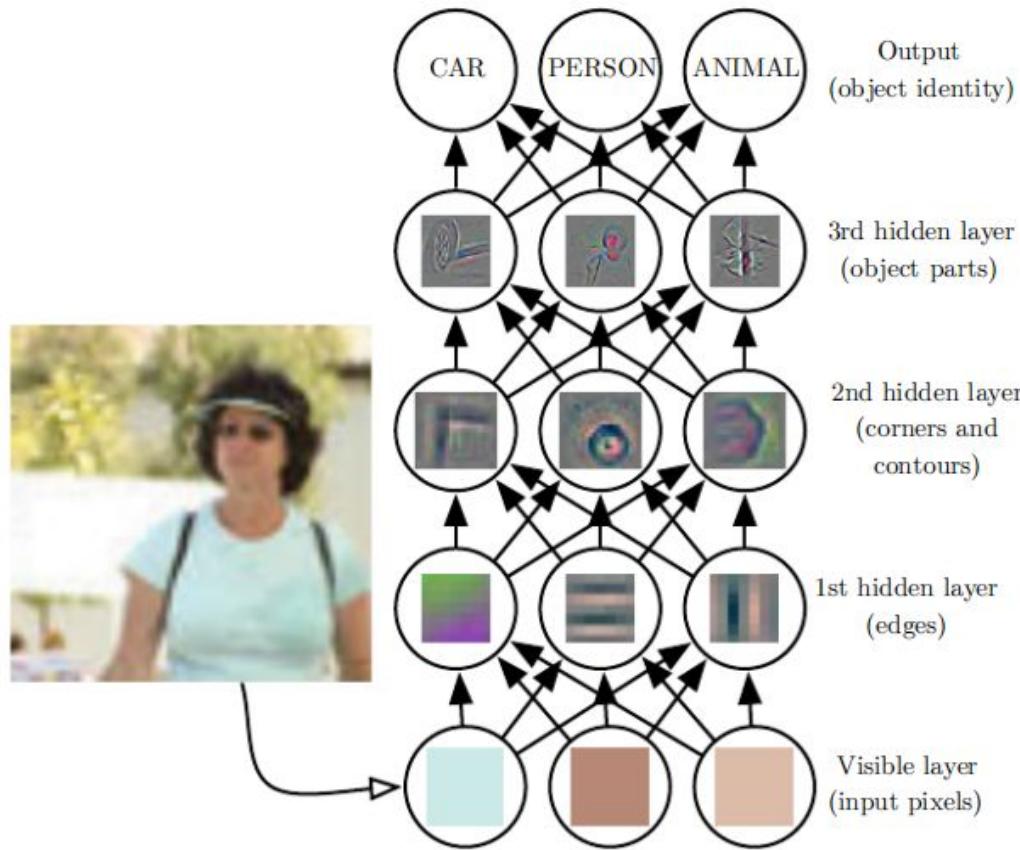


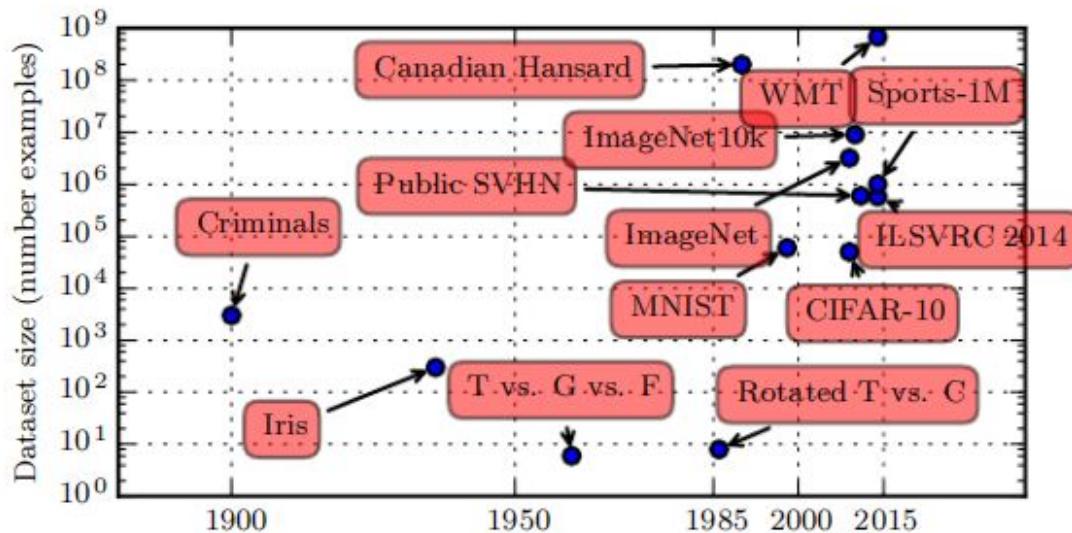
(c)



(d)

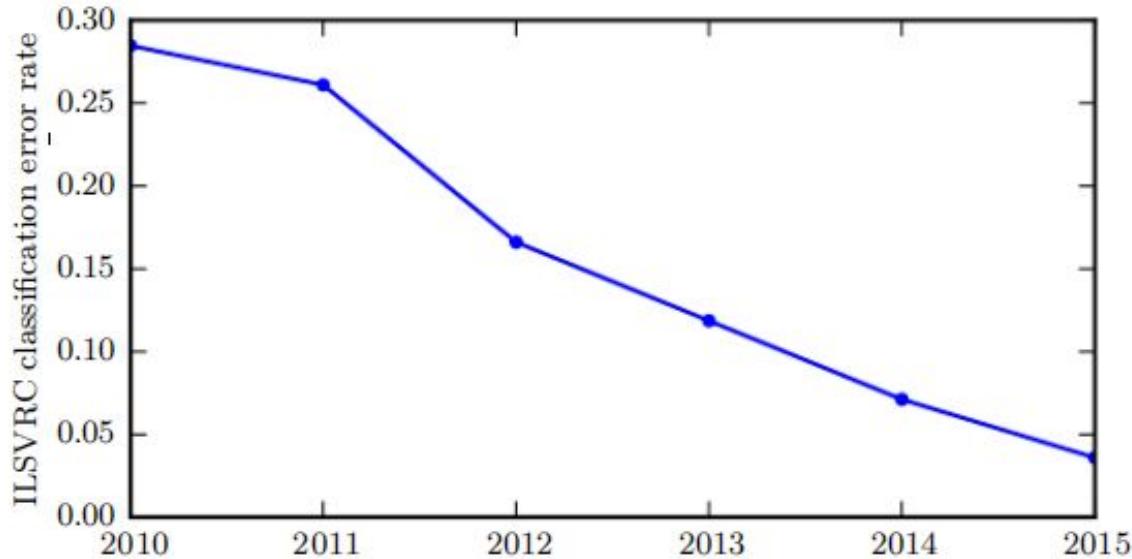




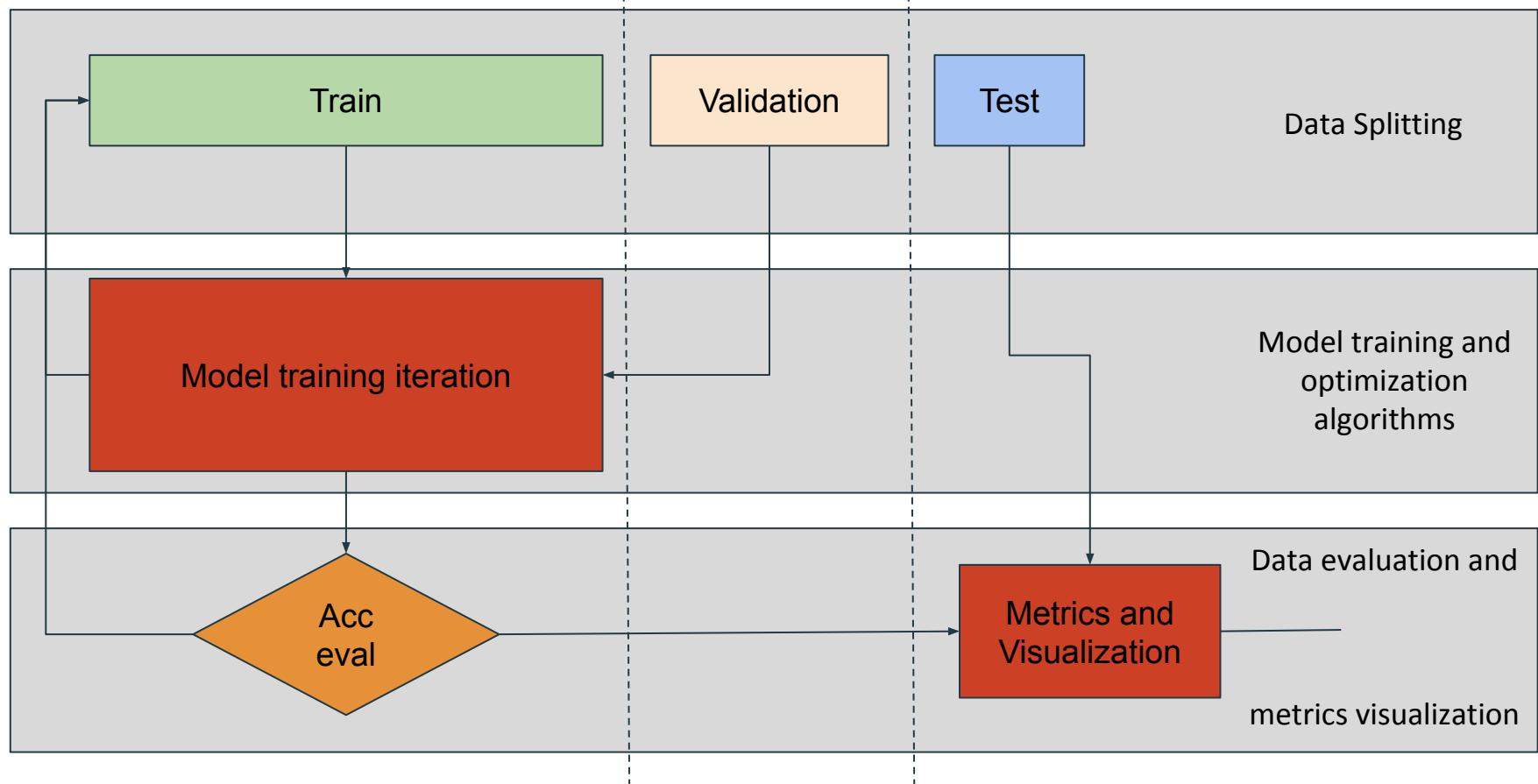




IMAGENET

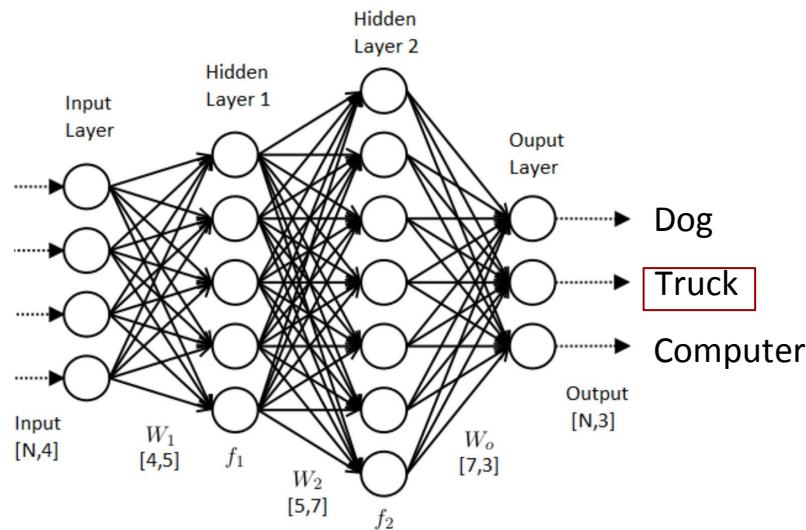


Overview

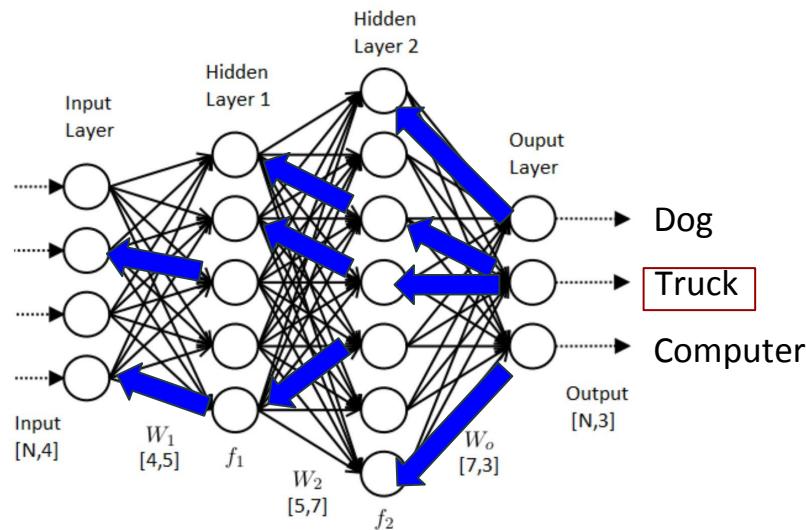


How does it learn?

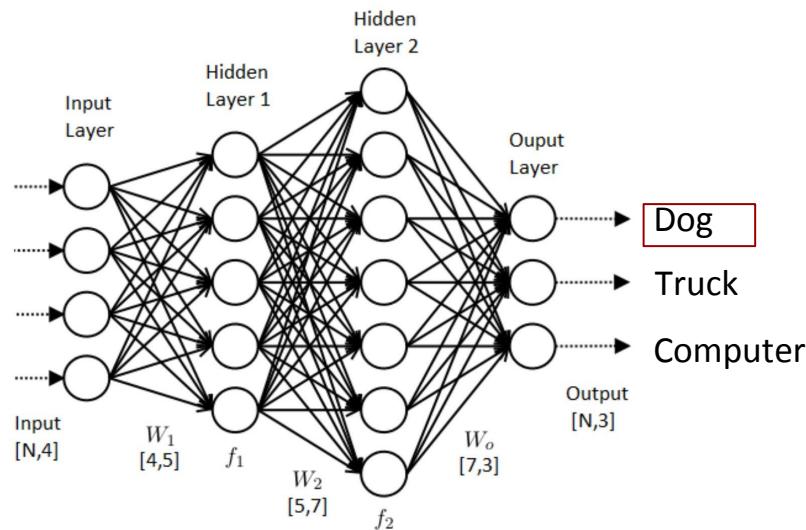
How does it learn?



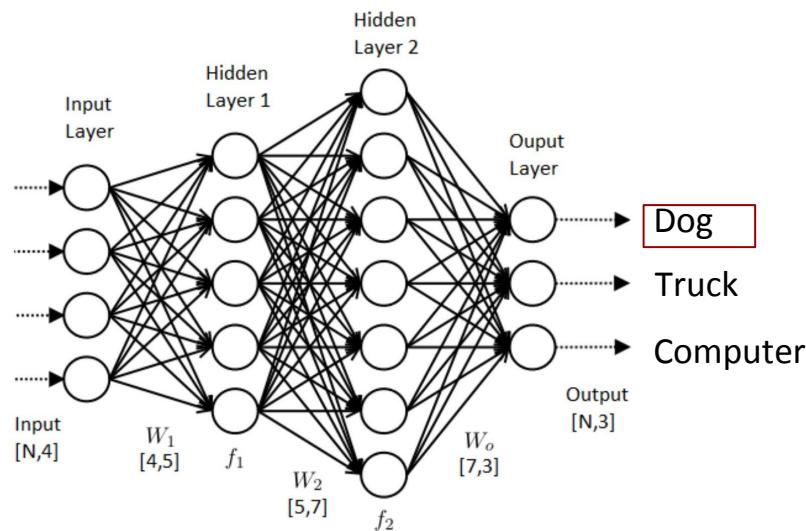
How does it learn?



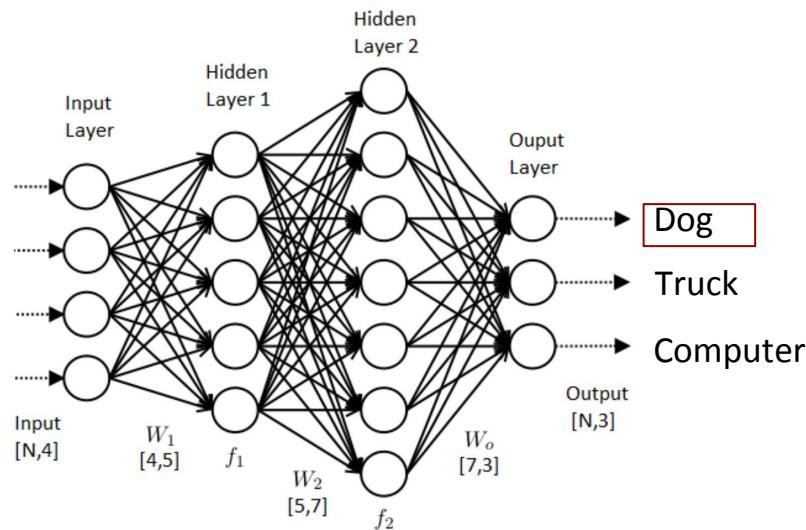
How does it learn?



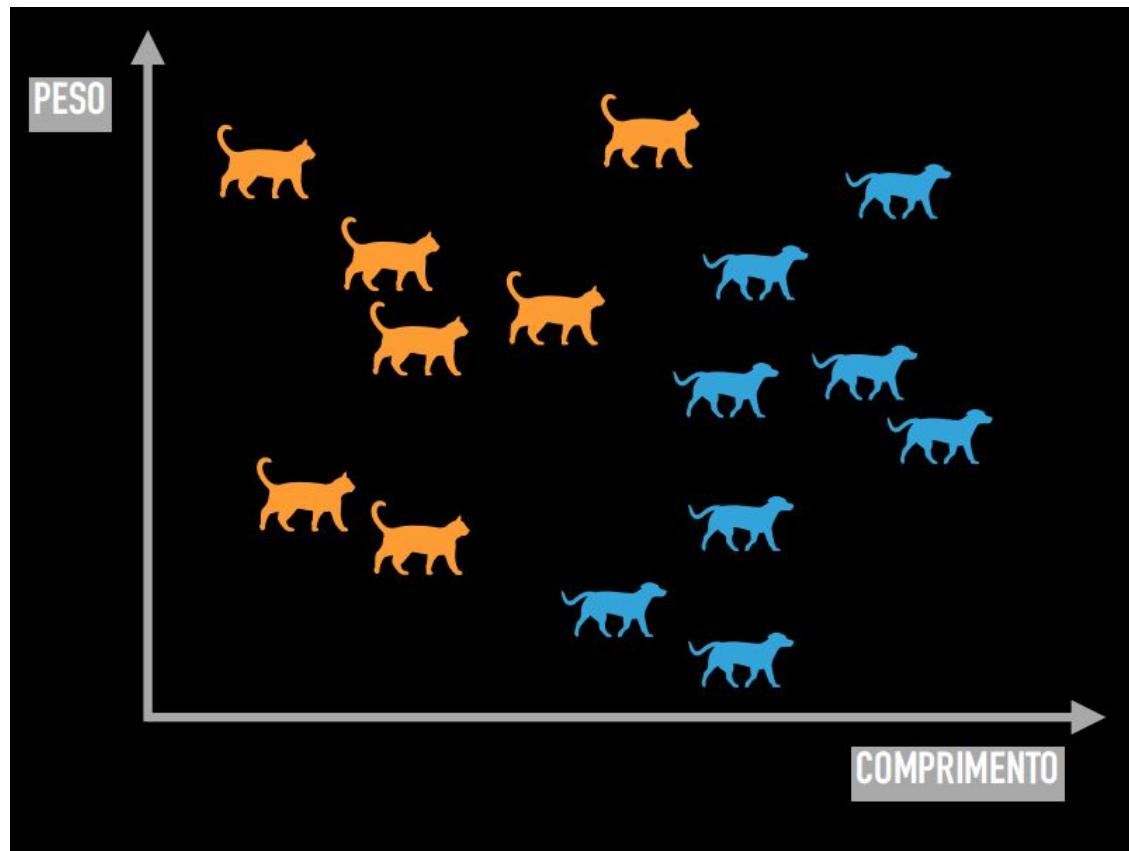
How does it learn?

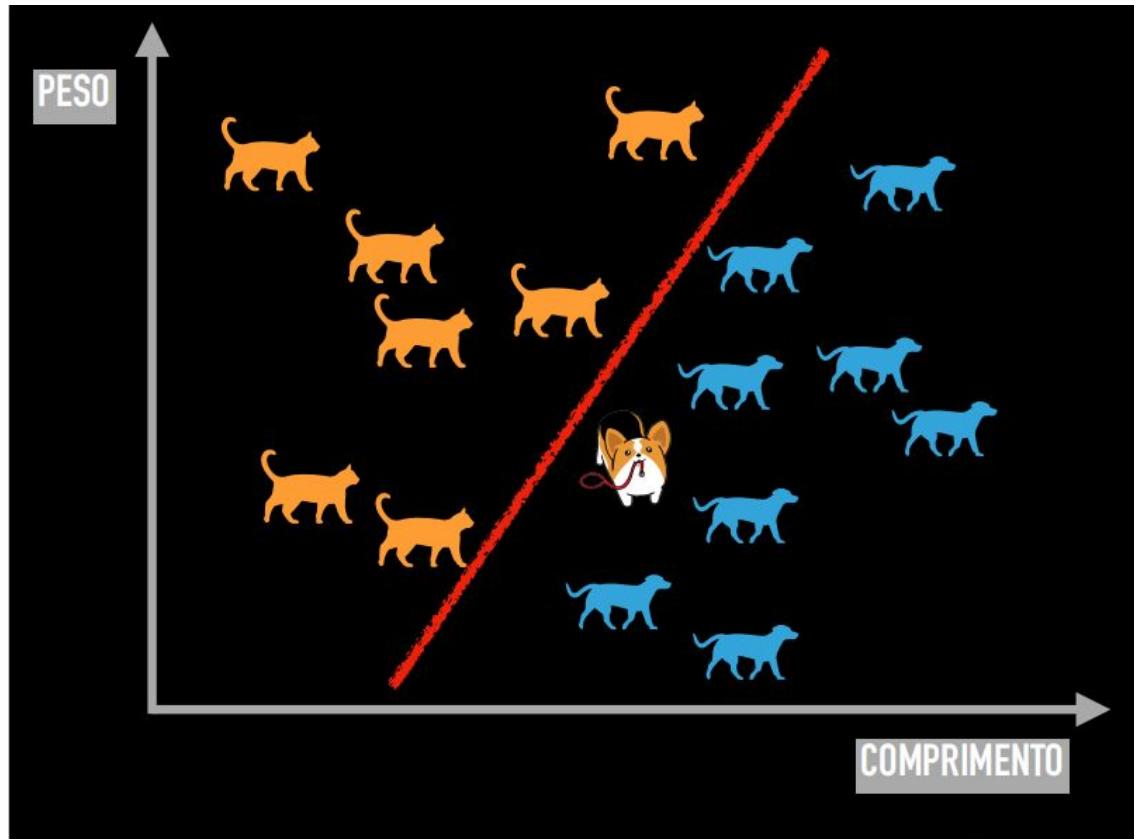


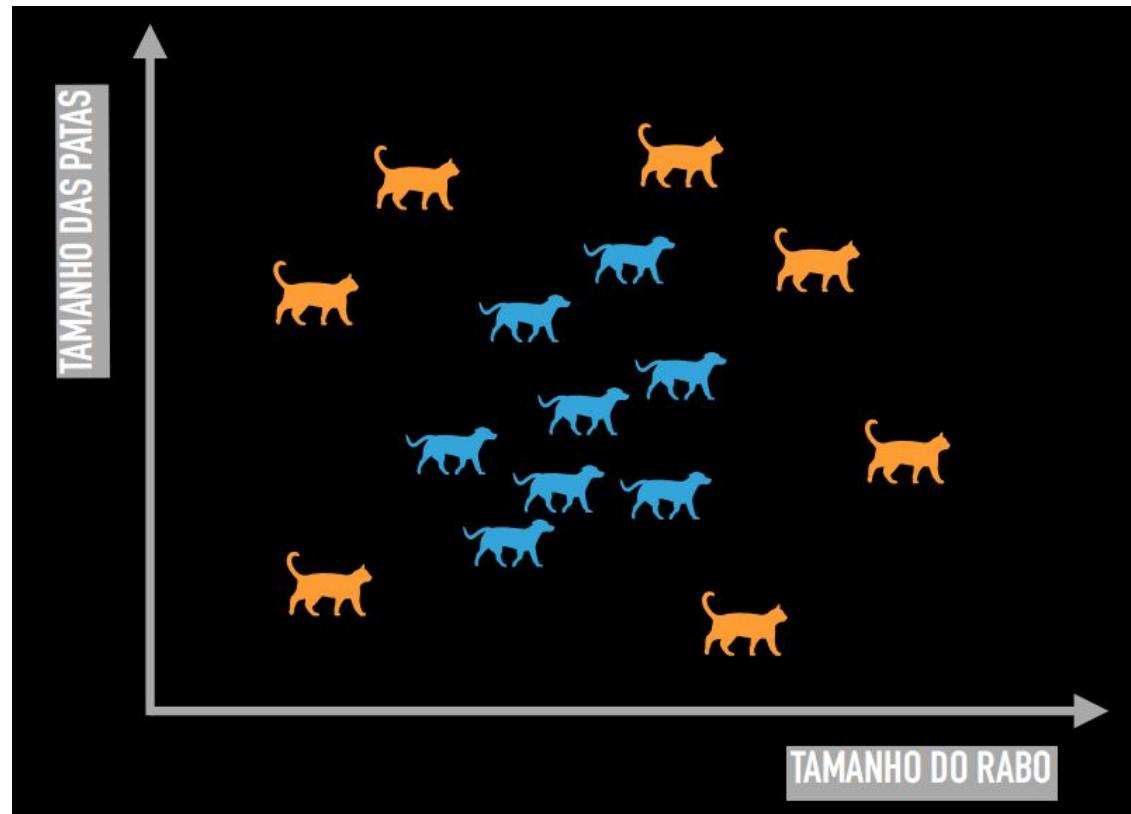
How does it learn?

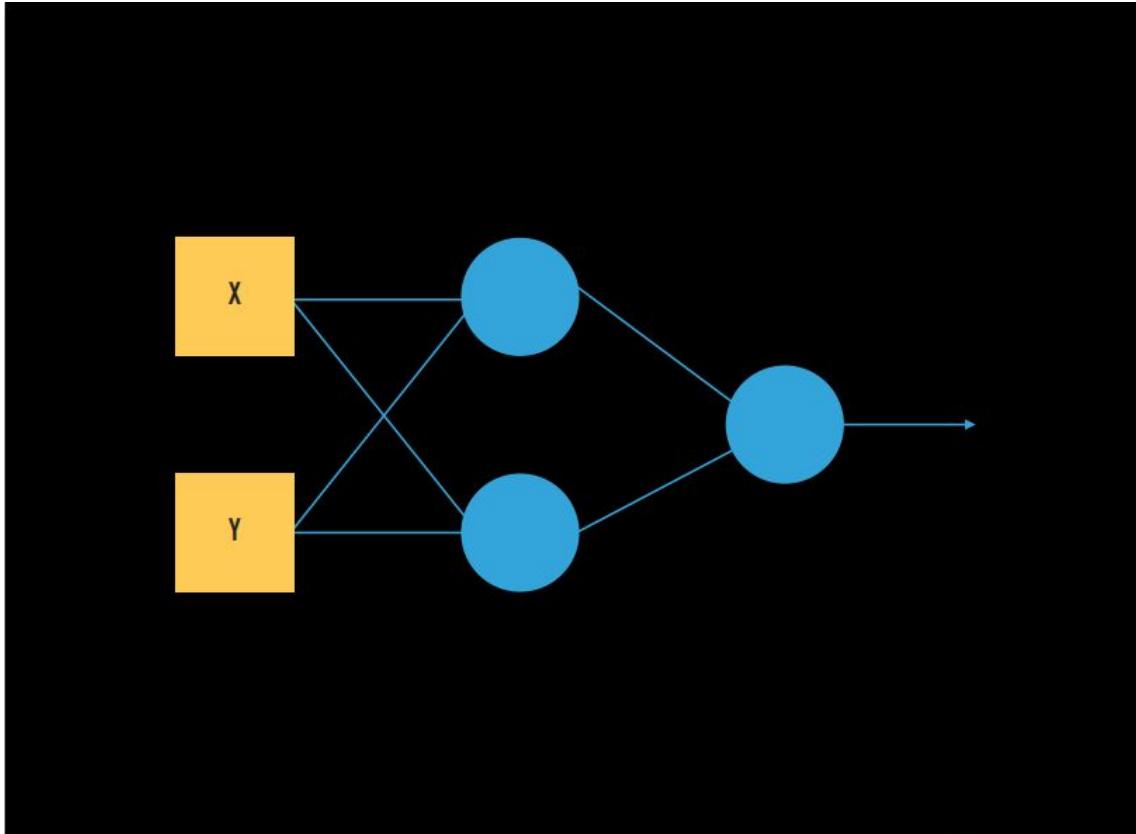


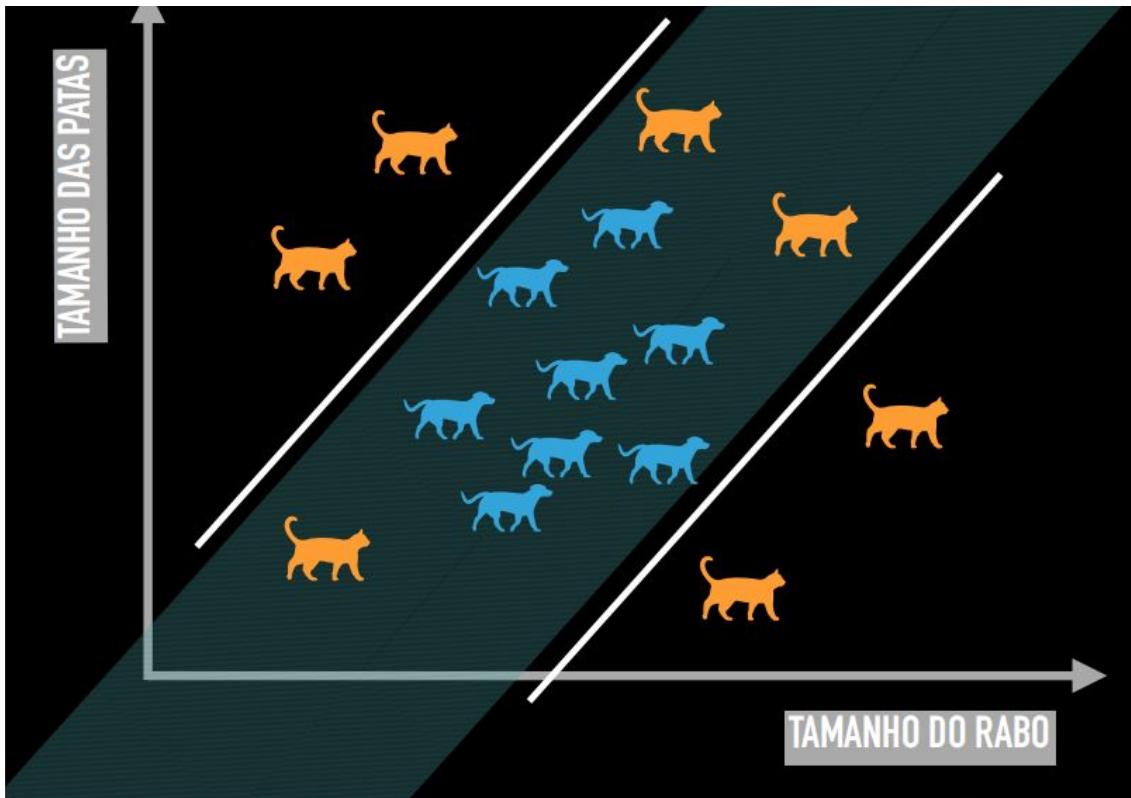
Why is it Deep?

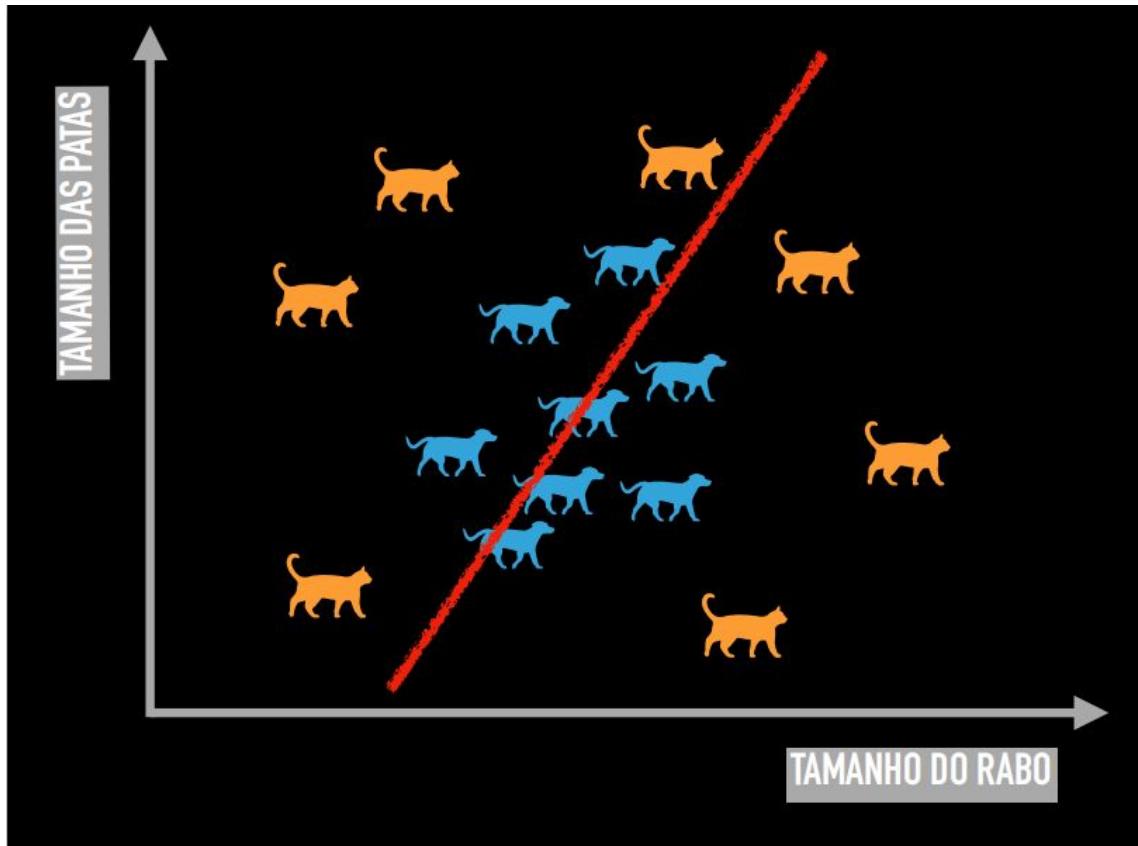


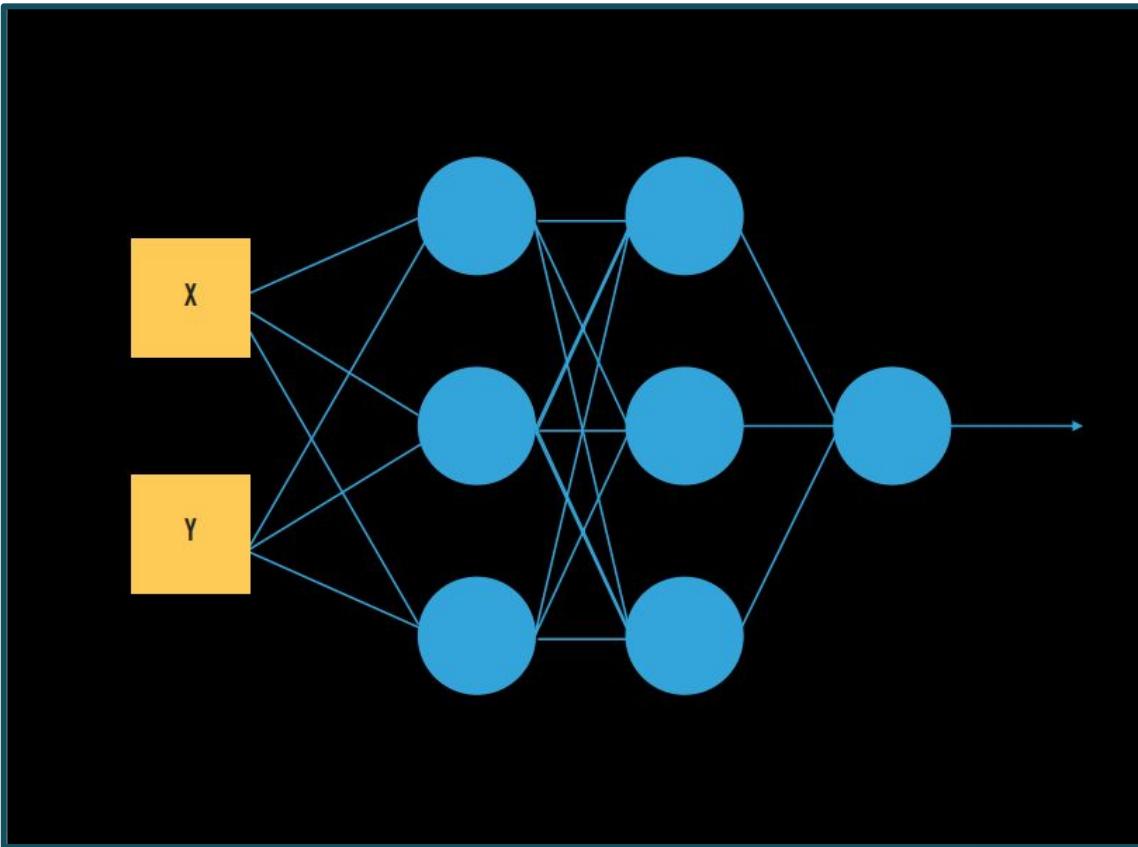


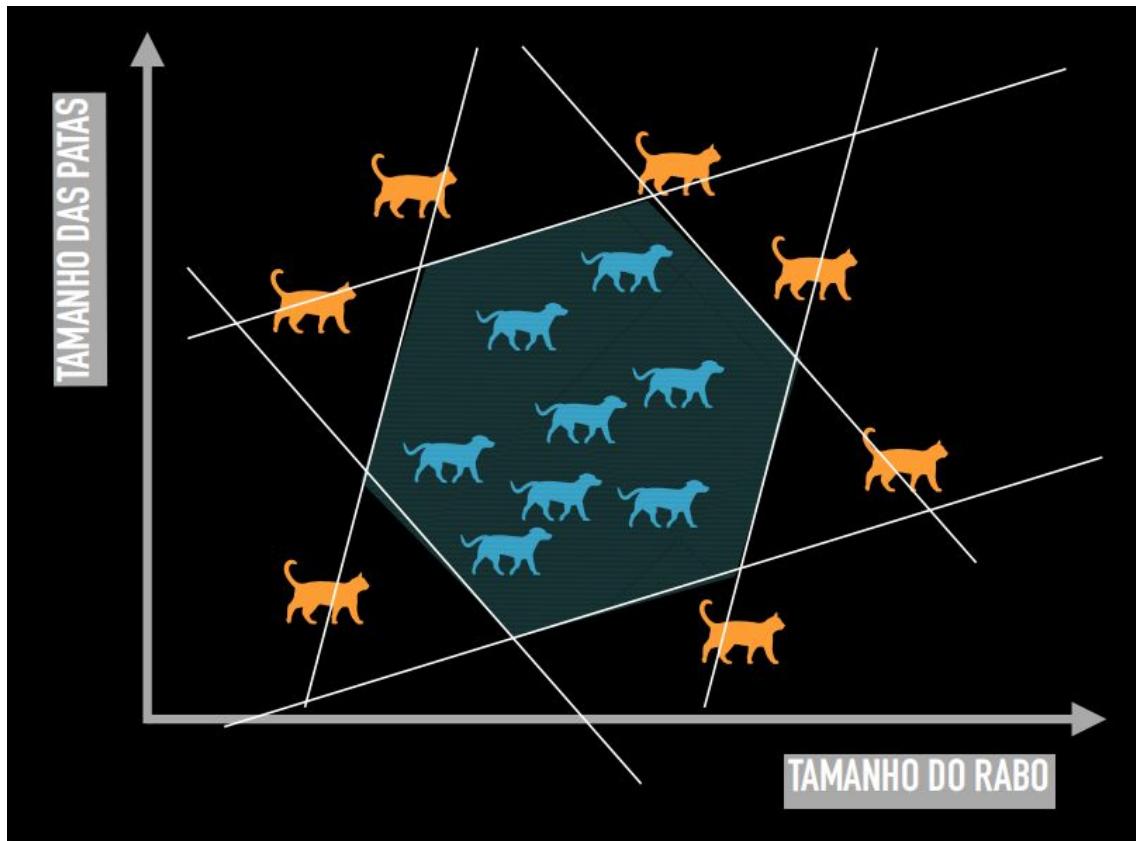


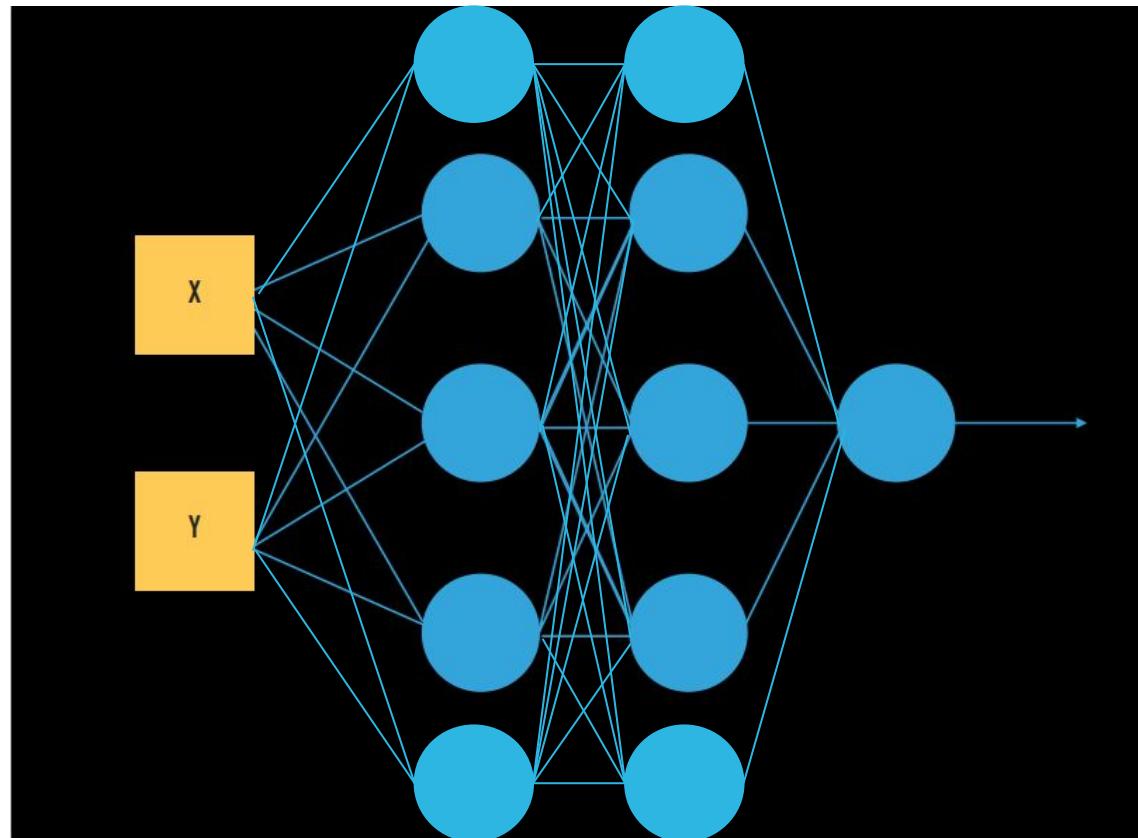


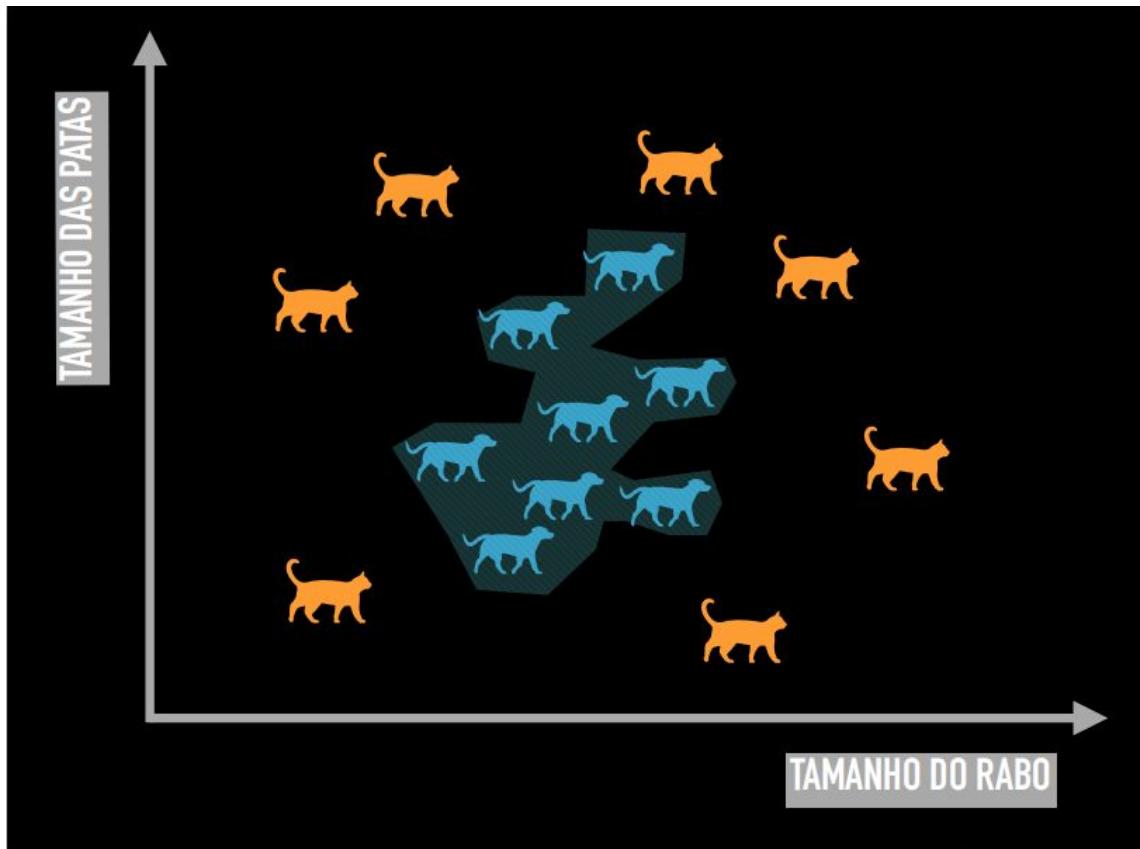


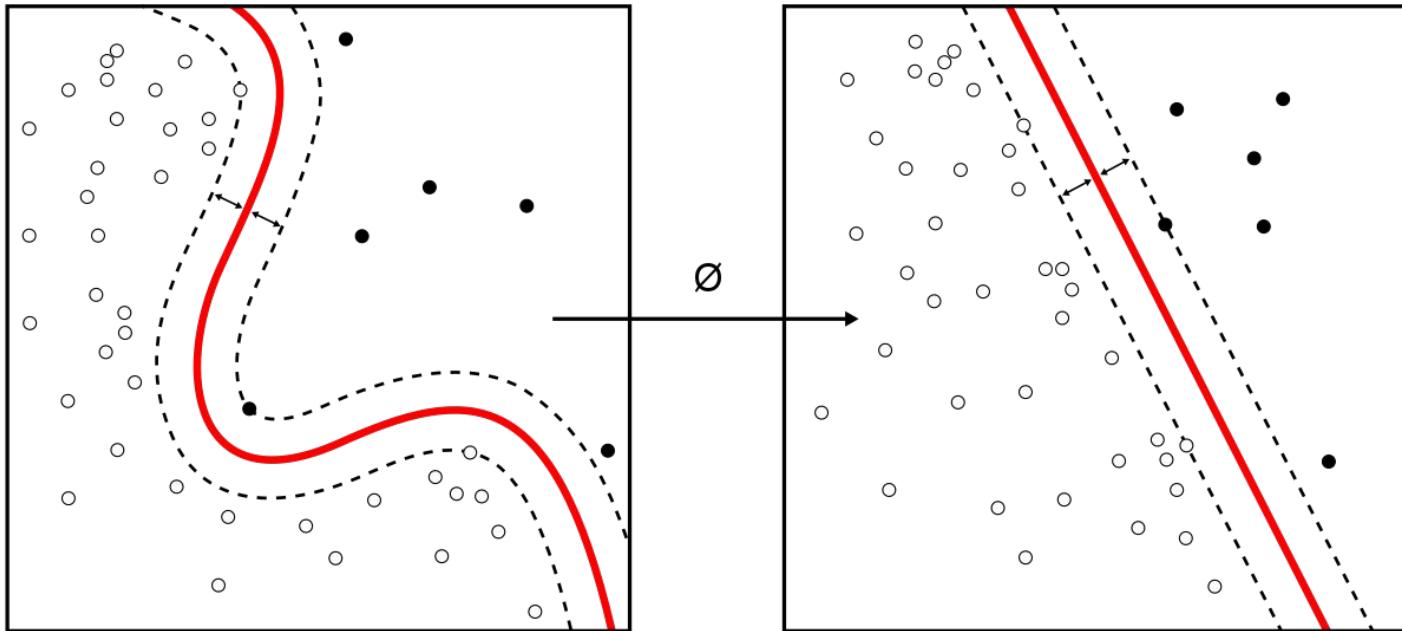










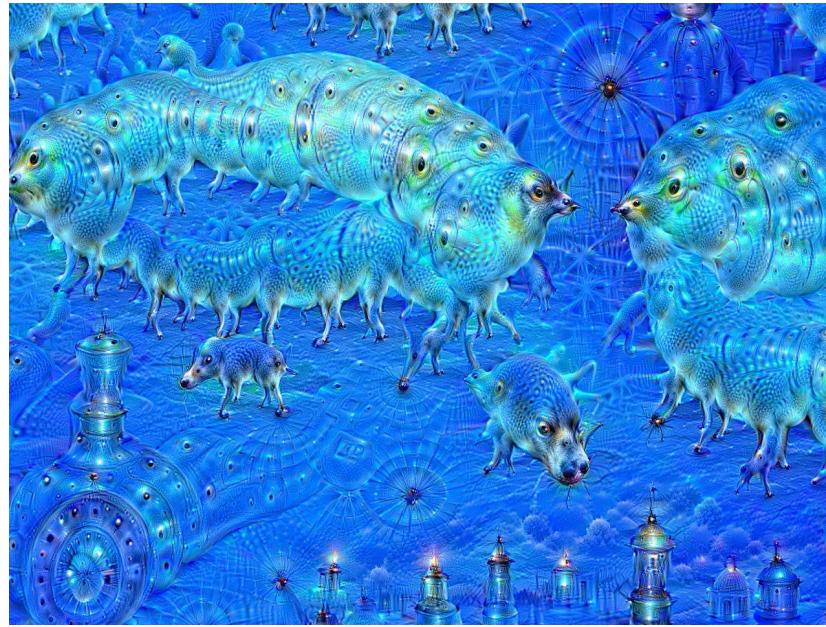


Interesting stuffs in Deep Learning

Deep Dream

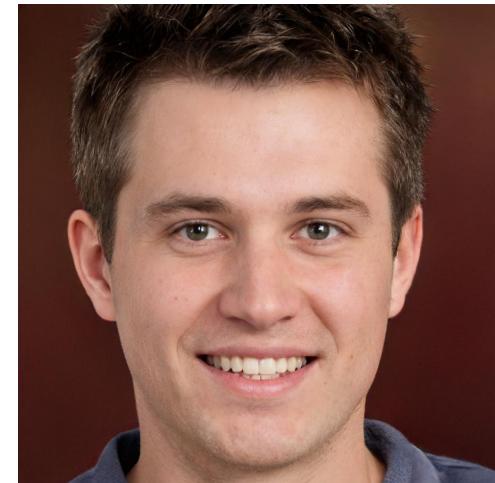


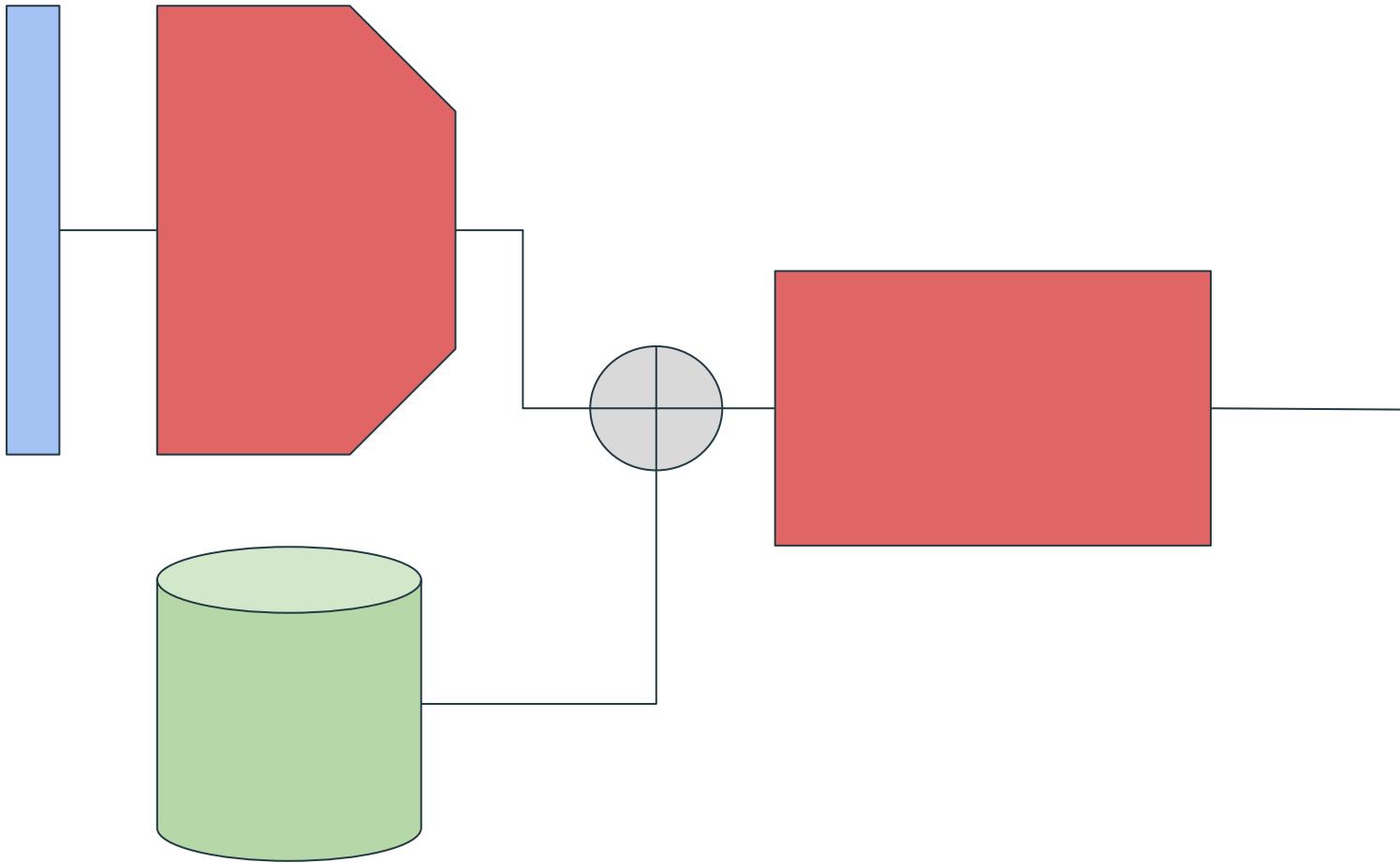


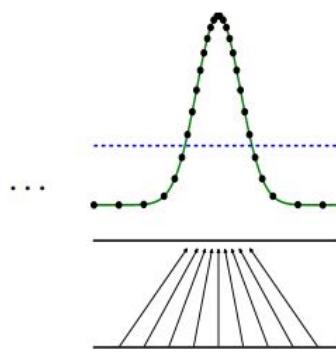
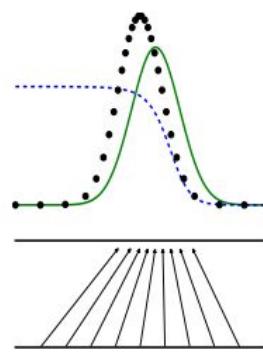
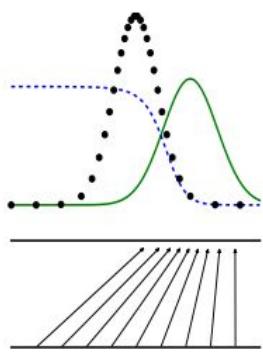
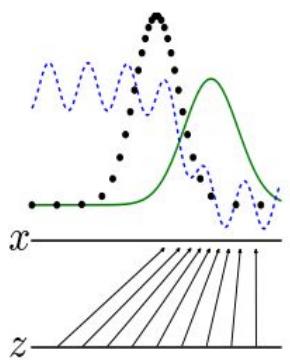




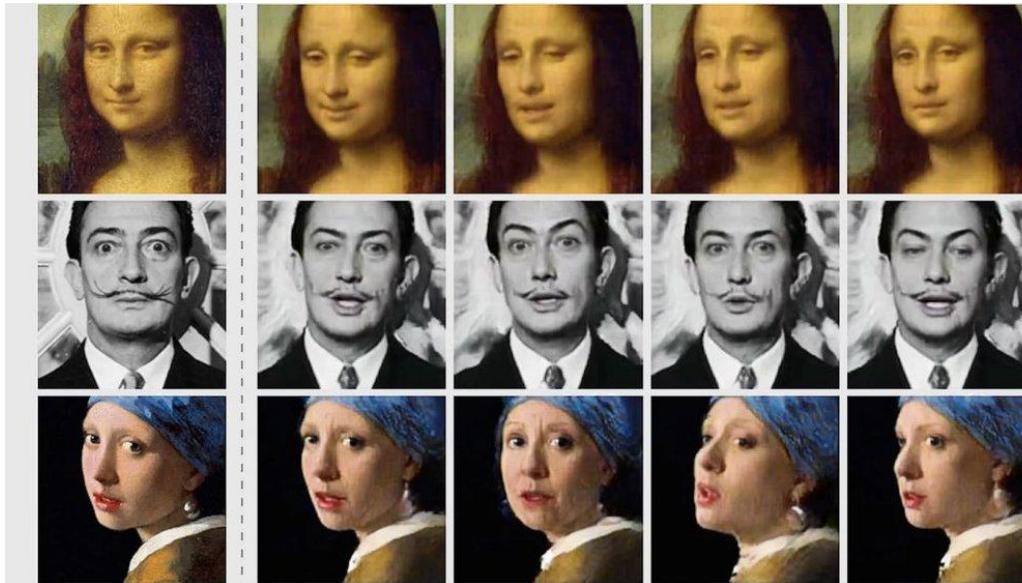
Generative Adversarial Networks



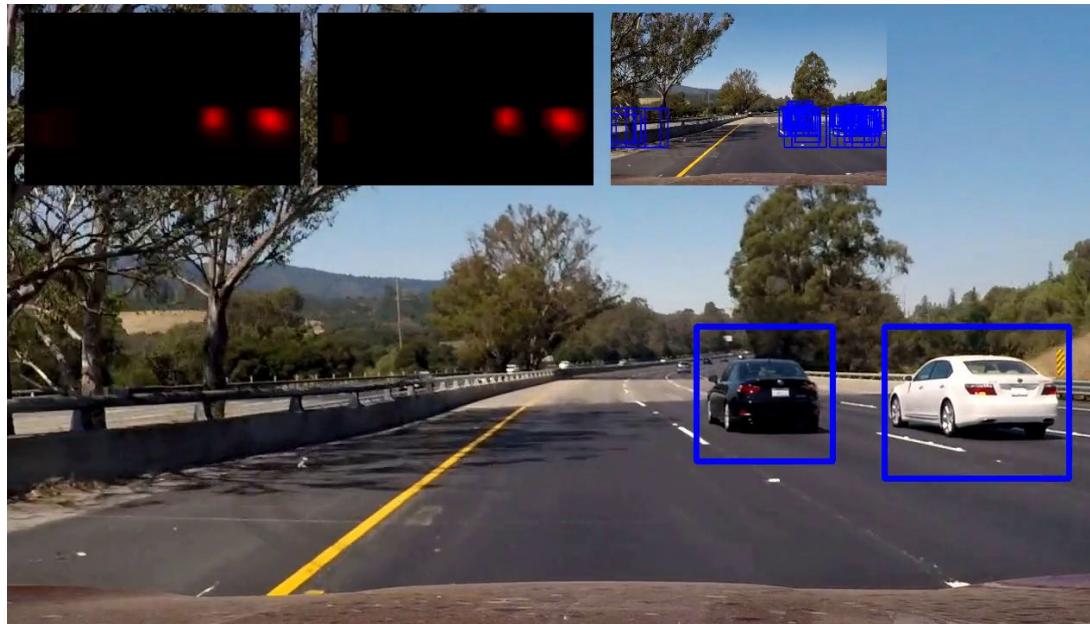




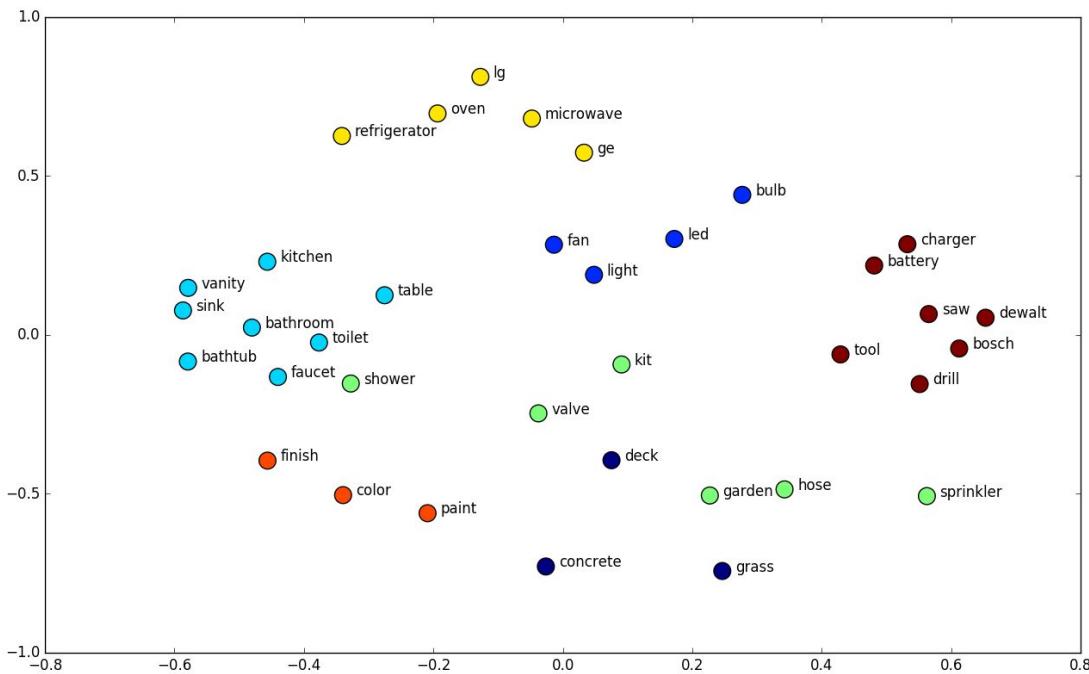
Deep Fake

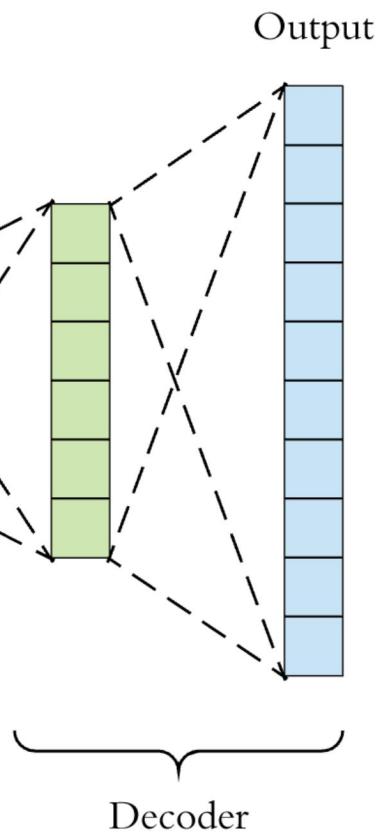
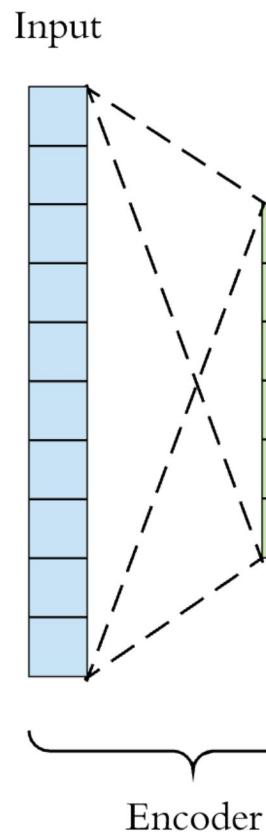


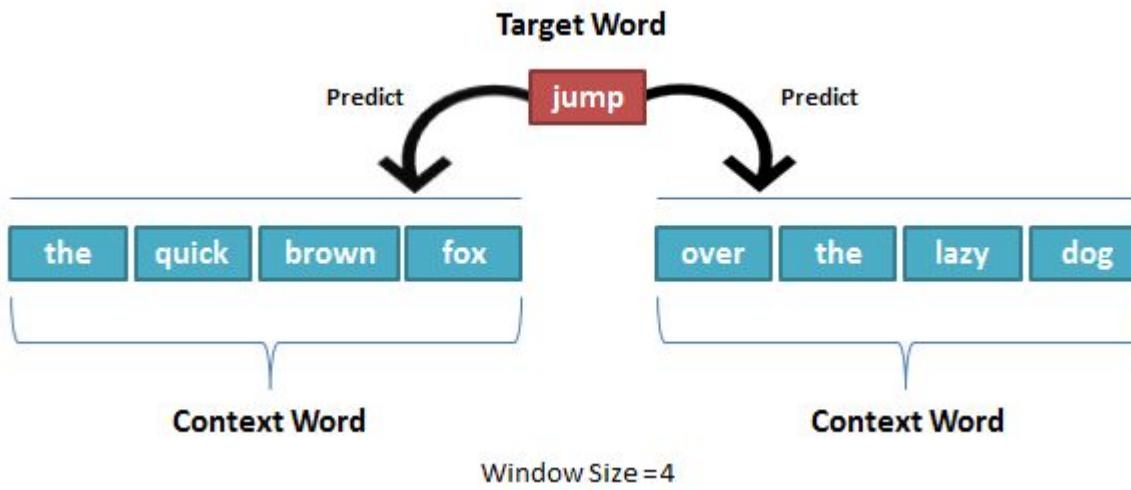
Self-driving Cars



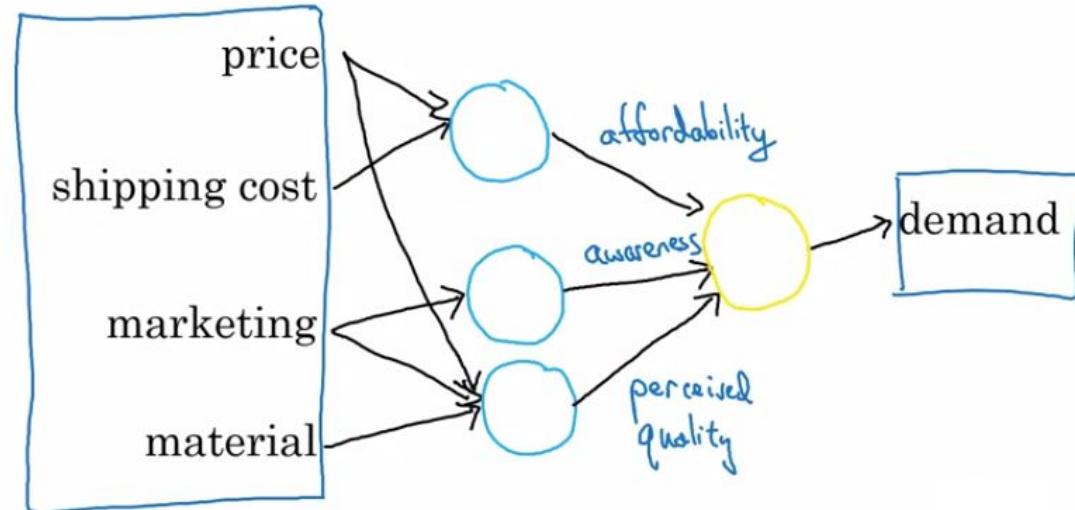


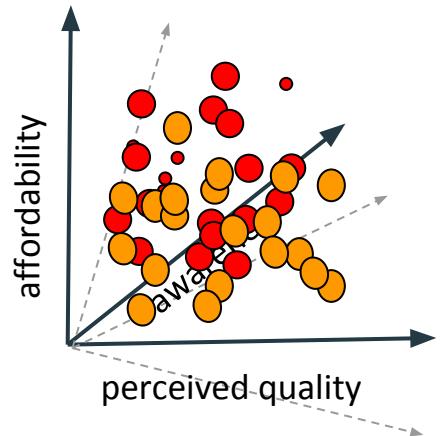






Non-technical explanation of DL





Some problems found in building models

- Unbalanced data (weighted loss, augmentation, etc.)
- Data bias (the ML engineer needs to analyse)
- Small dataset (transfer learning, k-fold cross validation, other ML models rather Deep Learning)
- Constrained resources, i.e., embedded systems (distilling networks, embedded models, network compression)
- Real problem data too distant from the training dataset (augmentation, engineering solutions)
- The need of a massive pre-processing on the data
- Lack of visualization and explanation

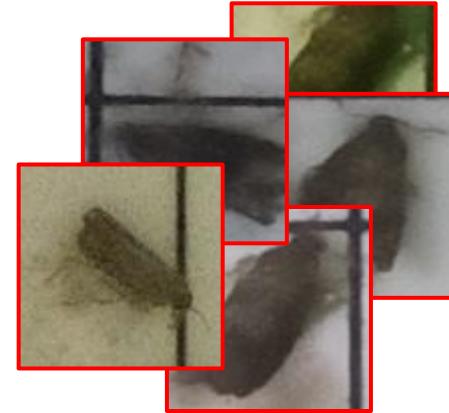
Unbalanced Data



~58%



%22%



~20%

Data Bias

Portuguese

English

Turkish

Detect language



English

Portuguese

Russian

Translate

O bir hemşire
O bir doktor



26/5000

She is a nurse
He is a doctor ✓





(a) Three samples in criminal ID photo set S_c .



(b) Three samples in non-criminal ID photo set S_n

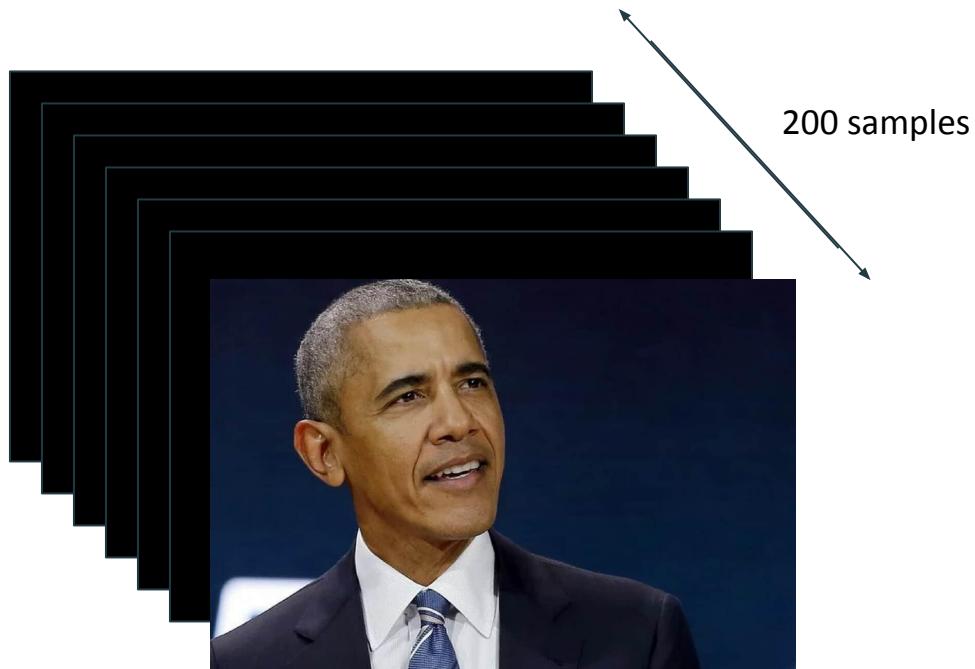
CRIMINAL



NON-CRIMINAL

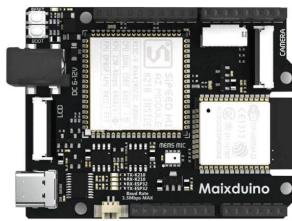


Small Dataset



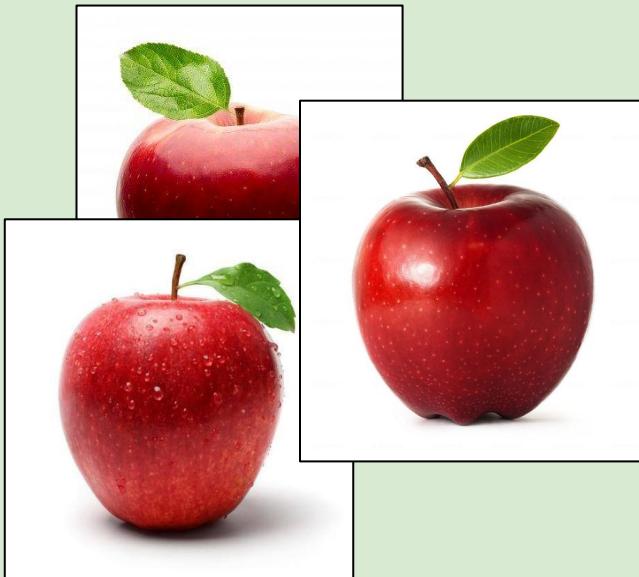
200 samples

Constrained Resources



The representation learned from the data
diverges from the real data distribution

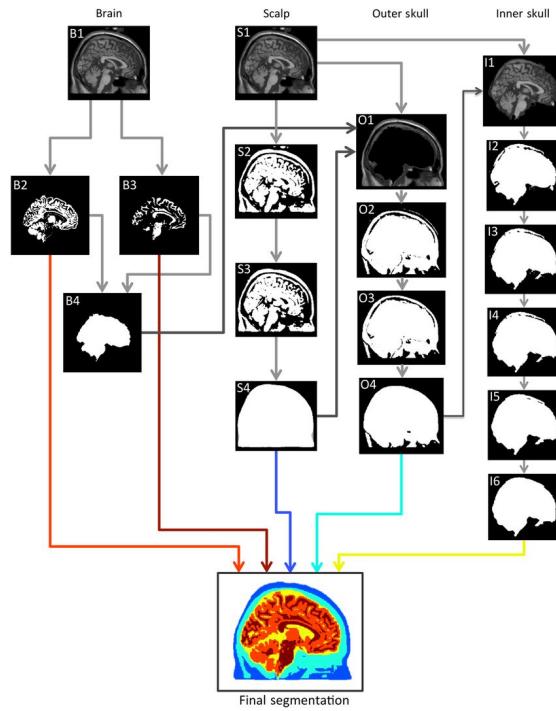
Training and Test data



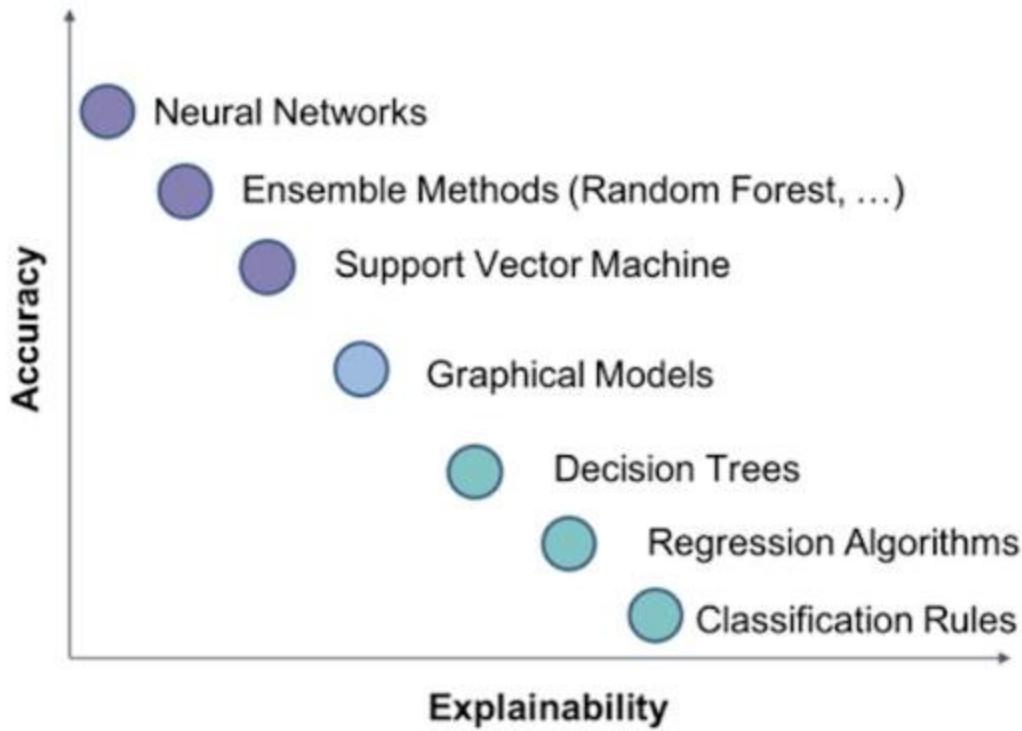
Real data when deployed



Massive Pre-processing



Lack of explainability



Techniques used to improve generalization and accuracy

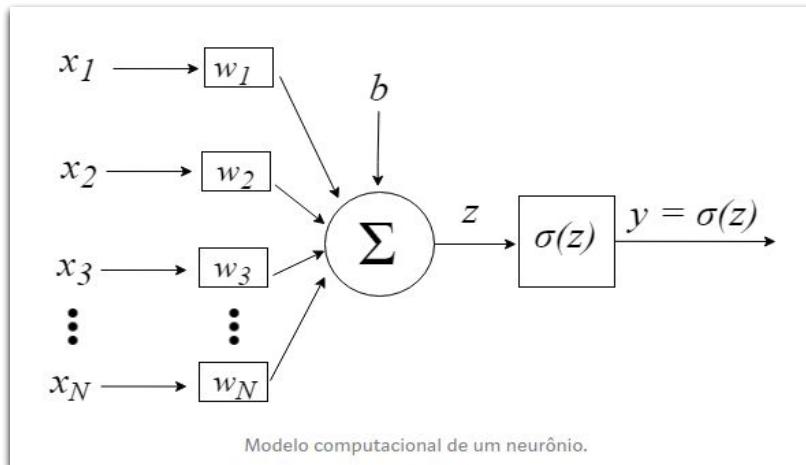
- Normalized inputs (even between layers)
- DropOut (like our brain, sparse networks performs better than dense ones)
- Augmentation (synthetically augmenting samples can help in generalization)
- Transfer learning (small datasets combined to large models improve the risks of overfitting. With transfer learning is needed less epochs and mitigates the risk of overfitting).
- Avoiding temporal correlation between samples

Where and when should you apply deep learning models?

- Large amount of data (big data problem)
- Often in images classification, detection, segmentation Deep Learning outperforms other models
- spacing-correlated, time-series data (images, videos, audio, text, e.g., sequencing)
- non-straight data
- hard-to-describe patterns (how does a cat looks like? Can we describe exactly?)

Deeper into Deep Learning

Mathematical model of a Neuron

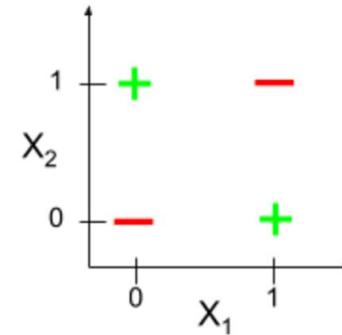
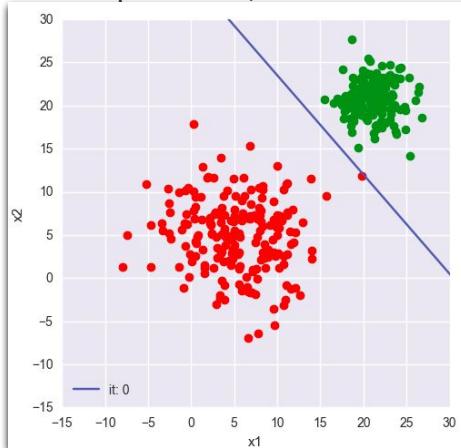


$$z = \sum_{i=1}^N x_i w_i + b$$

$$w_{t+1} = w_t + f(x, y) - f(x, \hat{y}).$$

Regressor

- Classification Problem: it performs well in linearly separable problems .
- Regression: it performs well for linear data.
- In real problems, it doesn't work so well where there are multidimensional data (many covariants)



Loss Functions

- mean square error
- log loss (cross-entropy)
- L1 (absolute difference)
- NLLoss

Mean Square Error

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

Cross-Entropy

$$H(p, q) = - \sum_x p(x) \log q(x)$$

L1 Loss

$$L = \sum_{i=1}^n |y_i - f(x_i)|$$

NLLLoss

$$\text{loss}(x, \text{class}) = -\log \left(\frac{\exp(x[\text{class}])}{\sum_j \exp(x[j])} \right) = -x[\text{class}] + \log \left(\sum_j \exp(x[j]) \right)$$

Evaluation

How to evaluate a machine learning model?

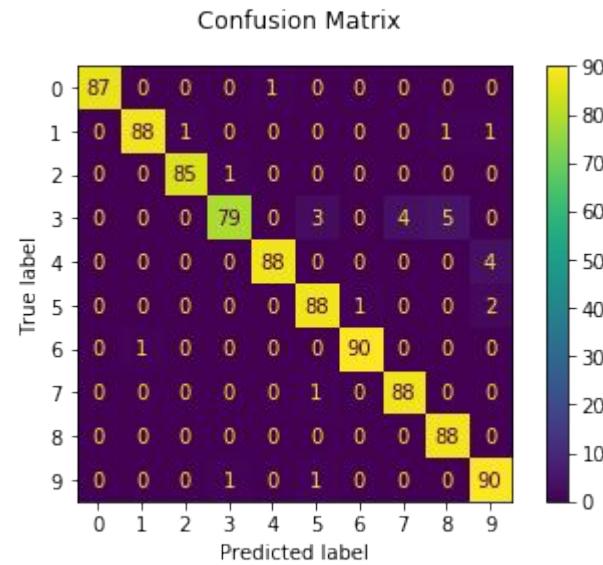
What each metric means?

There's an overall metric?

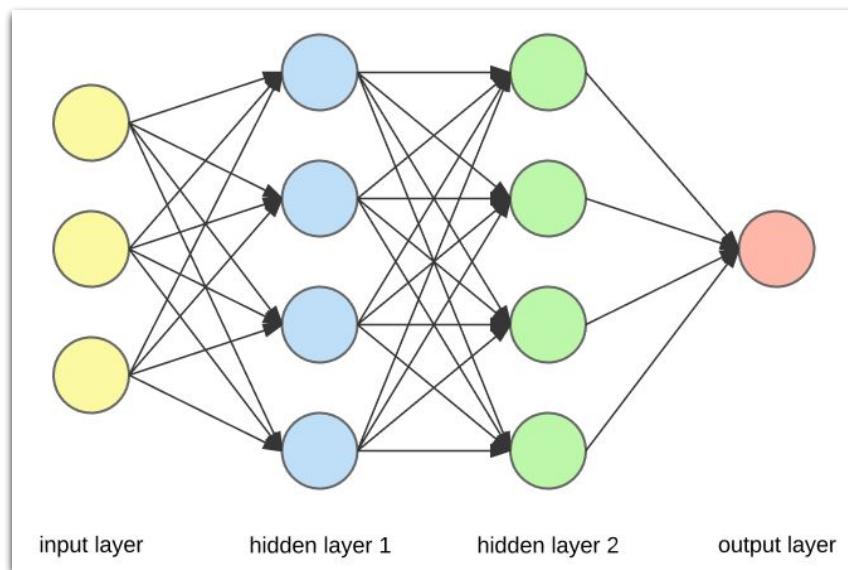
Metric Visualization

Evaluation Metrics

- Loss ??
- accuracy
- true positive rate
- true negative rate
- precision
- recall
- sensitivity
- specificity
- F1-score

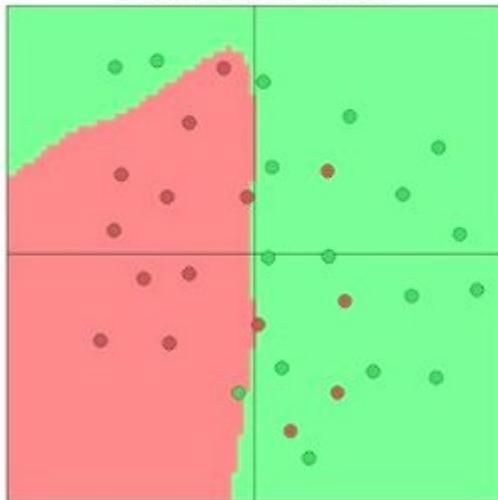


Multi Layer Perceptron

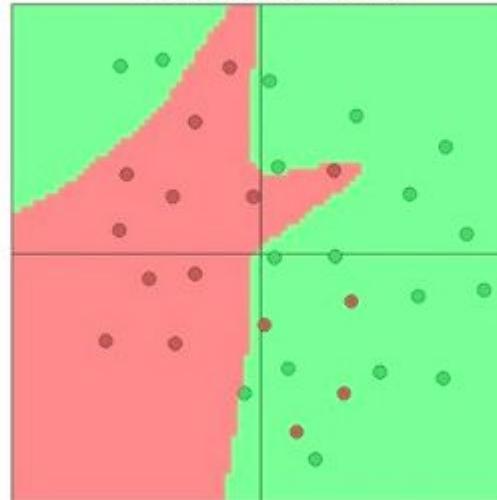


Multi Layer Perceptron

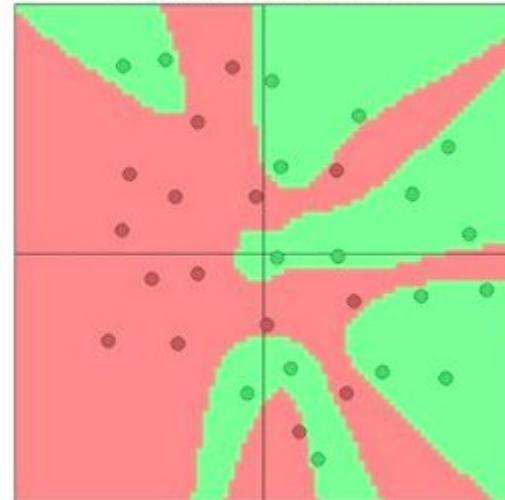
3 hidden neurons



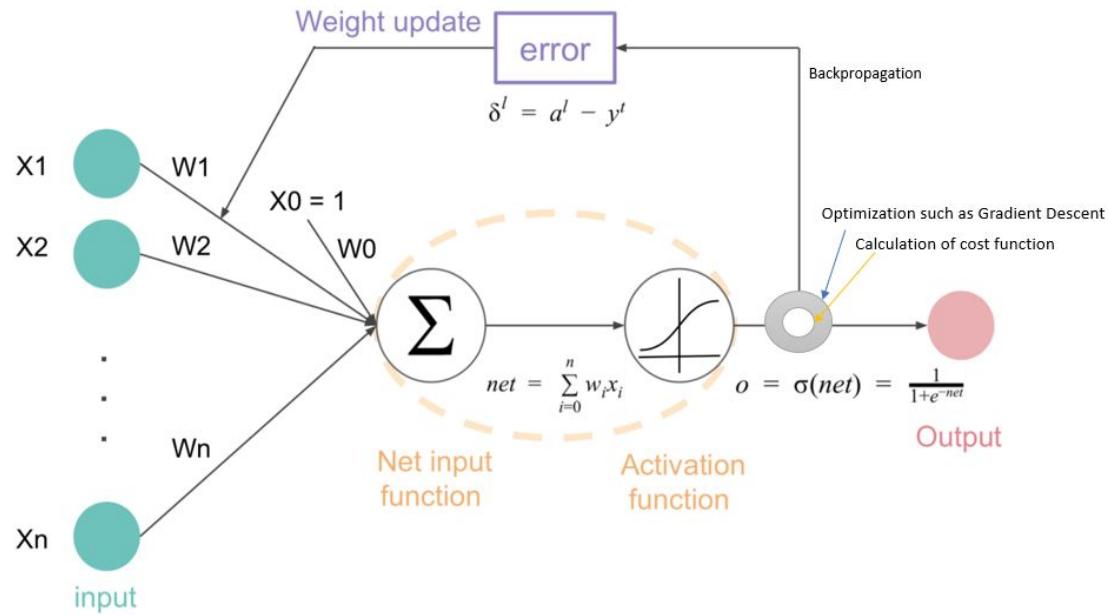
6 hidden neurons



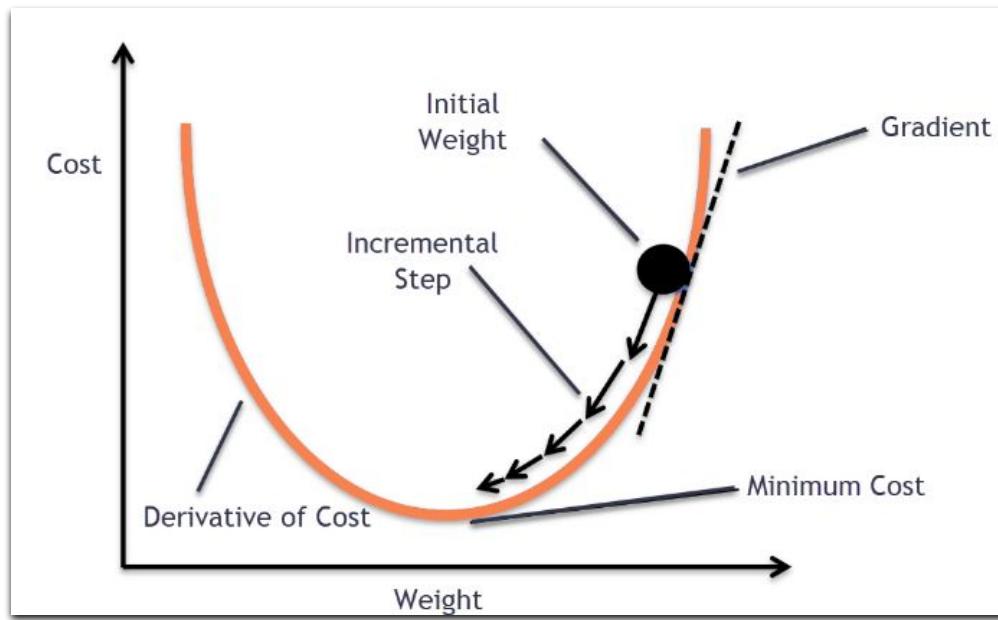
20 hidden neurons



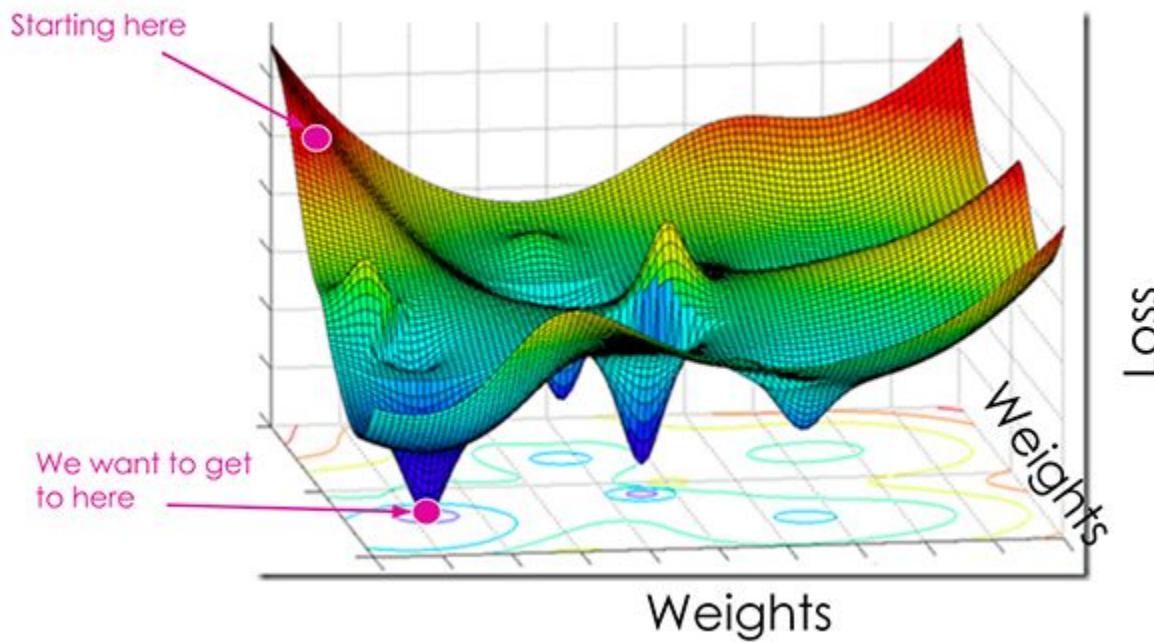
Multi Layer Perceptron

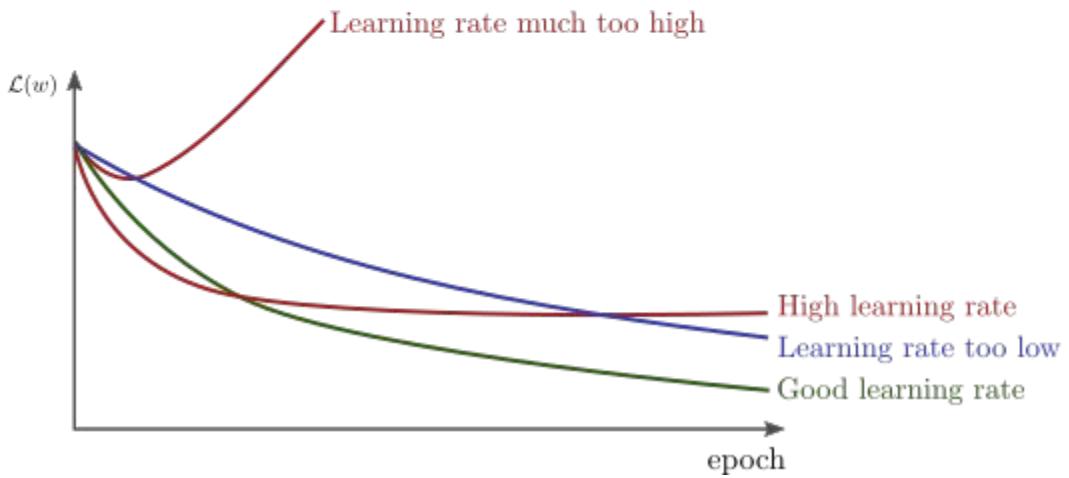
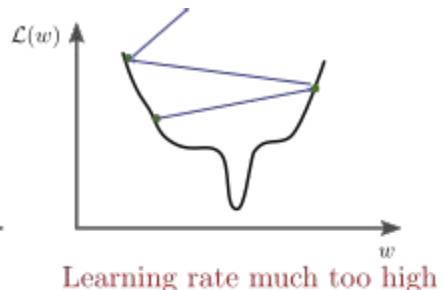
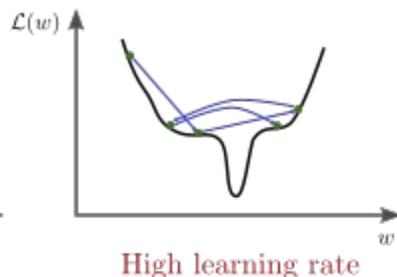
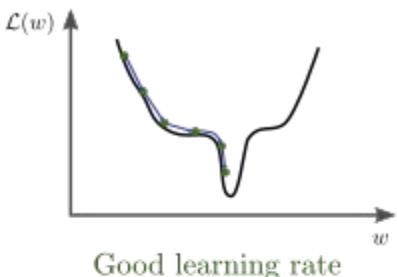
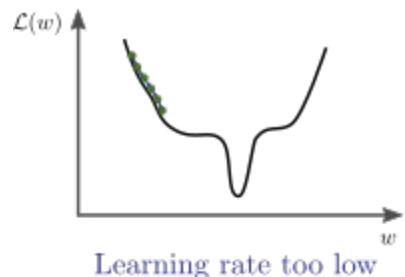


Gradient Descent

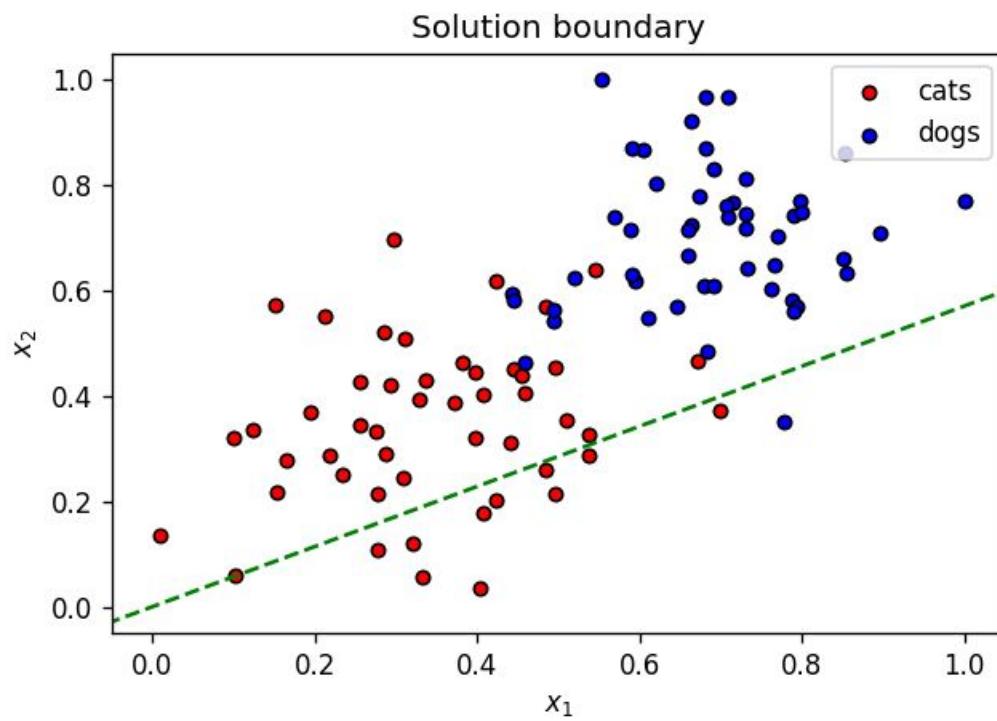


Weight Initialization



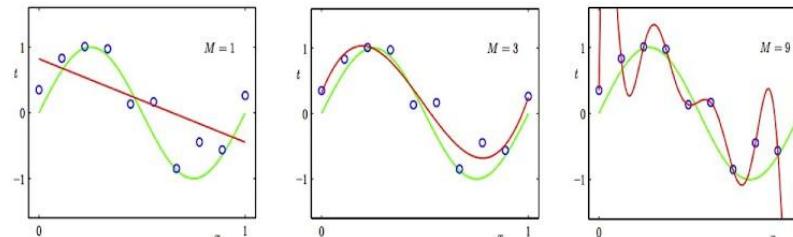


Vanishing Gradient



Overfitting and Underfitting

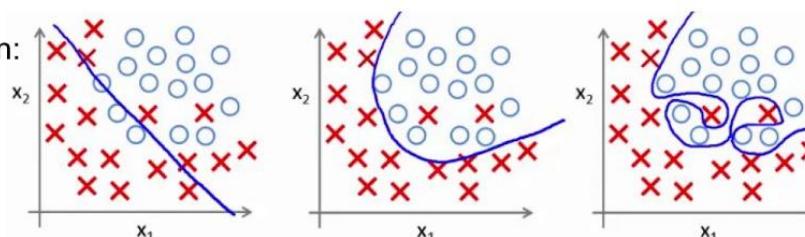
Regression:



predictor too inflexible:
cannot capture pattern

predictor too flexible:
fits noise in the data

Classification:



Copyright © 2014 Victor Lavrenko

Regularizers

- DropOut: sparse nets often generalizes better

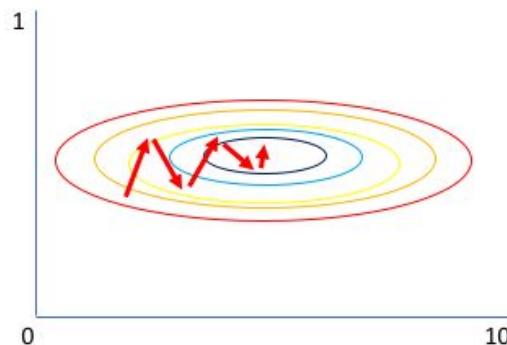
Regularizers

- DropOut
- L1 e L2: large weights are penalized, restricting them to close to zero numbers.

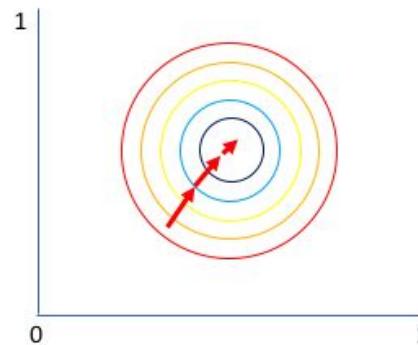
Regularizers

- DropOut
- L1 e L2
- Augmentation: the dataset is virtually enlarged with digital imaging processing

Normalization and Pre-training



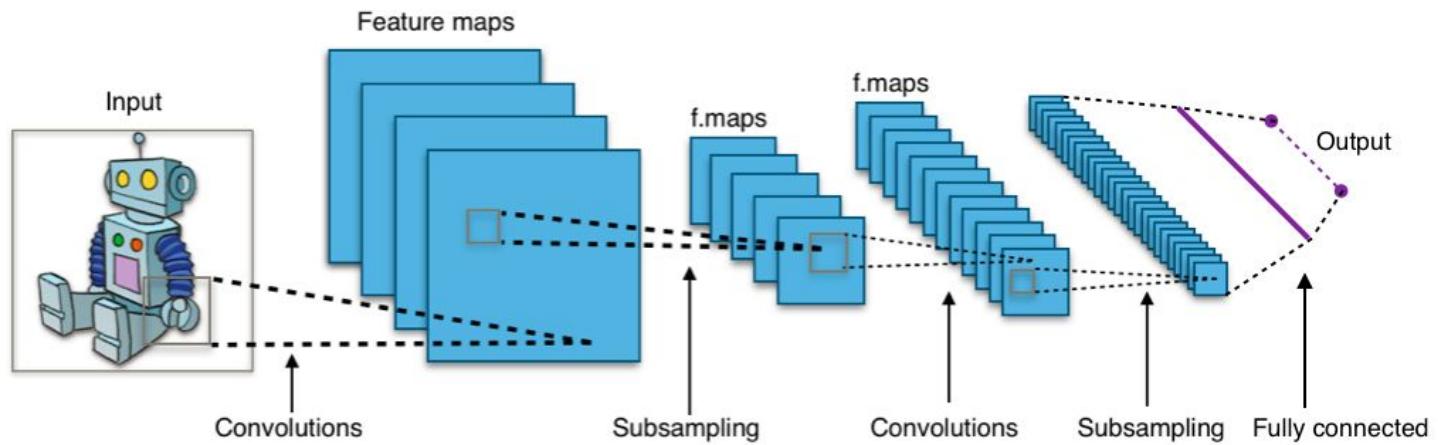
Gradient of larger parameter
dominates the update

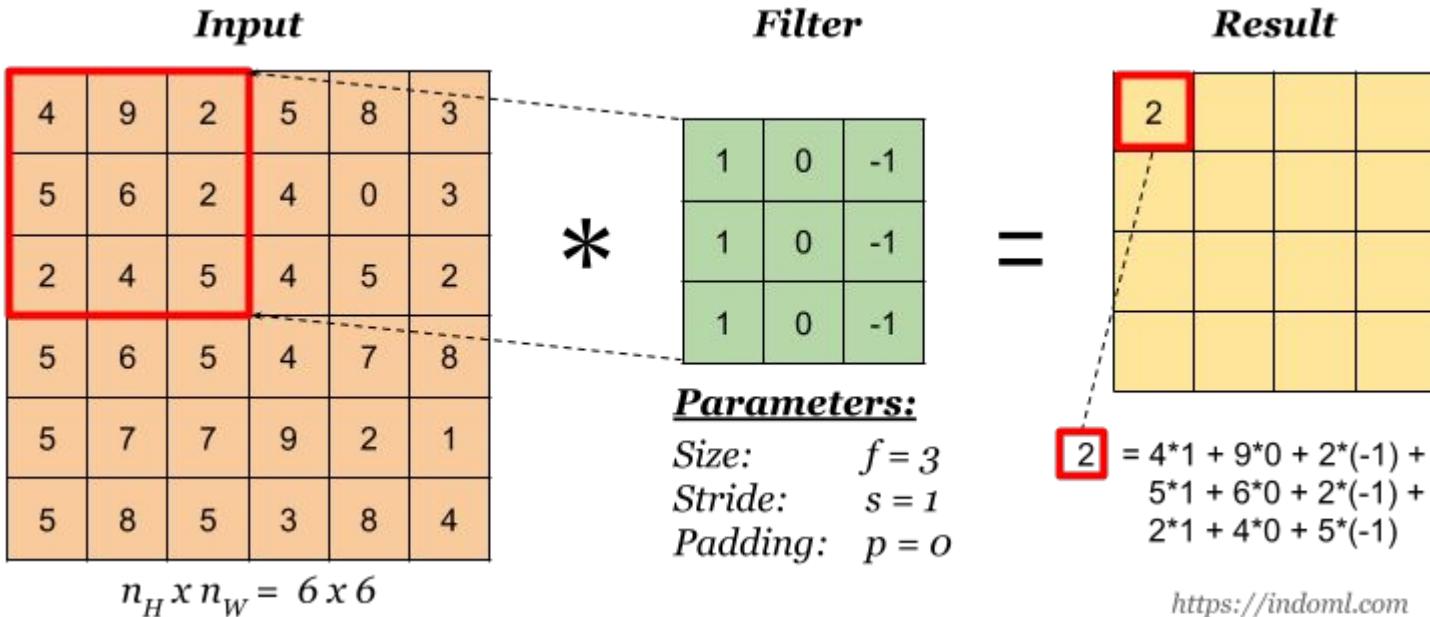


Both parameters can be
updated in equal proportions

Convolutional Neural Networks

- Constrained application domain
- Embedded knowledge (constrained domain) helps to outperforms generic networks such as MLP
- kernels and feature maps
- Automatic feature extractor (it can handle any classifier on top of layers stack)





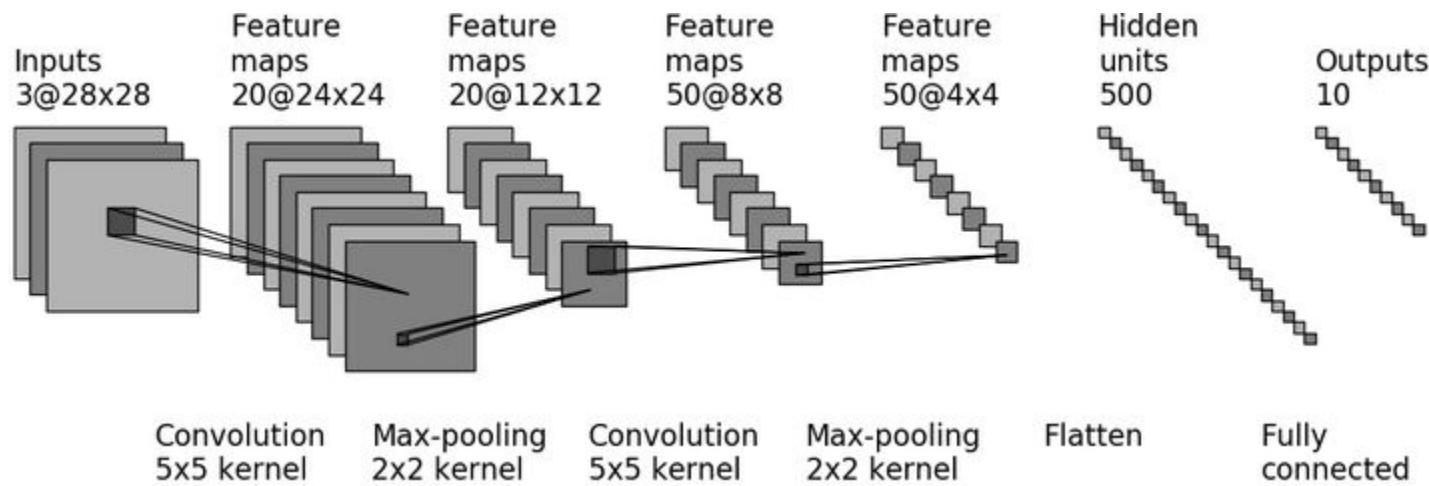
<https://indoml.com>

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

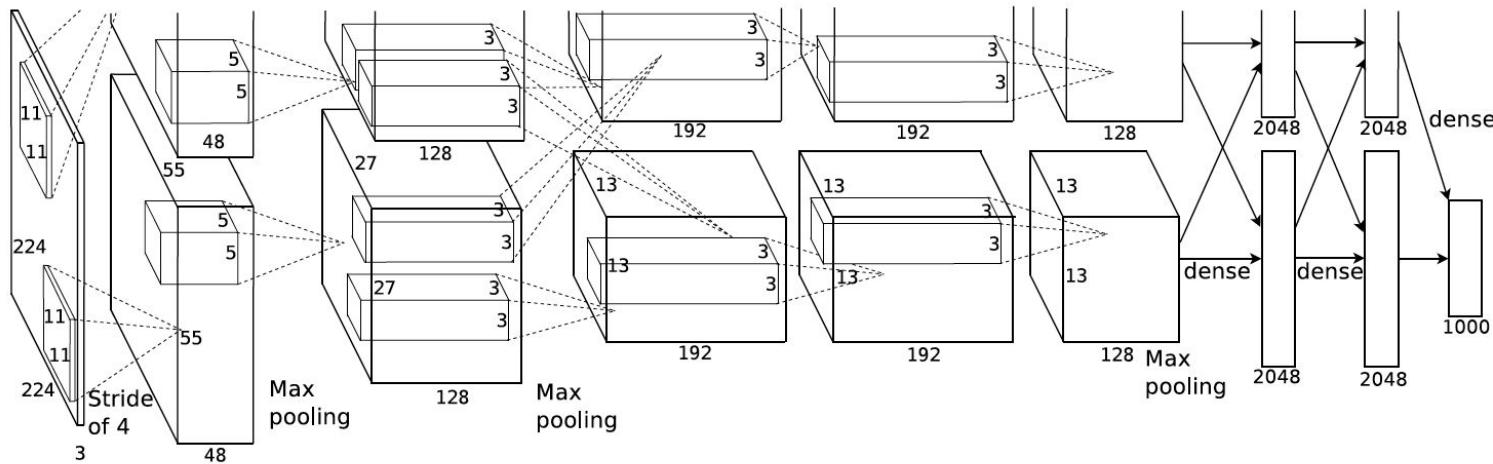
2×2 Max-Pool

20	30
112	37

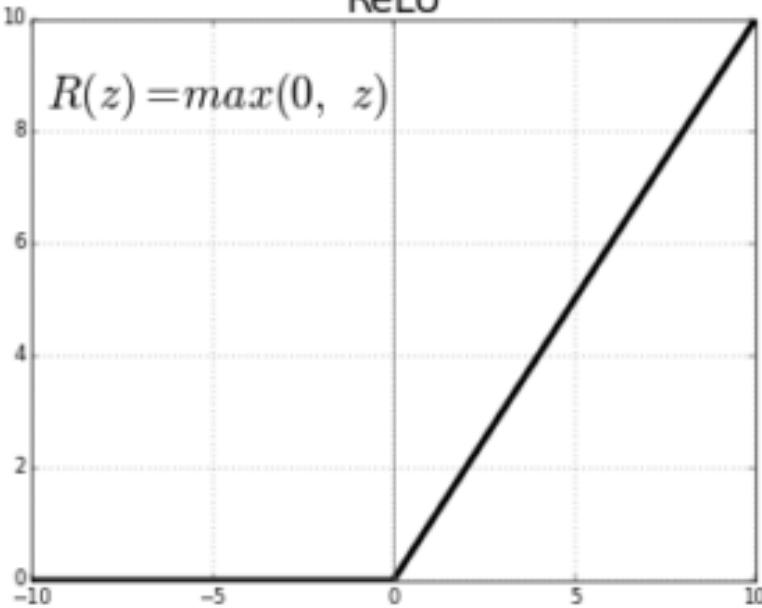
LeNet



AlexNet



ReLU



ResNet

Inception

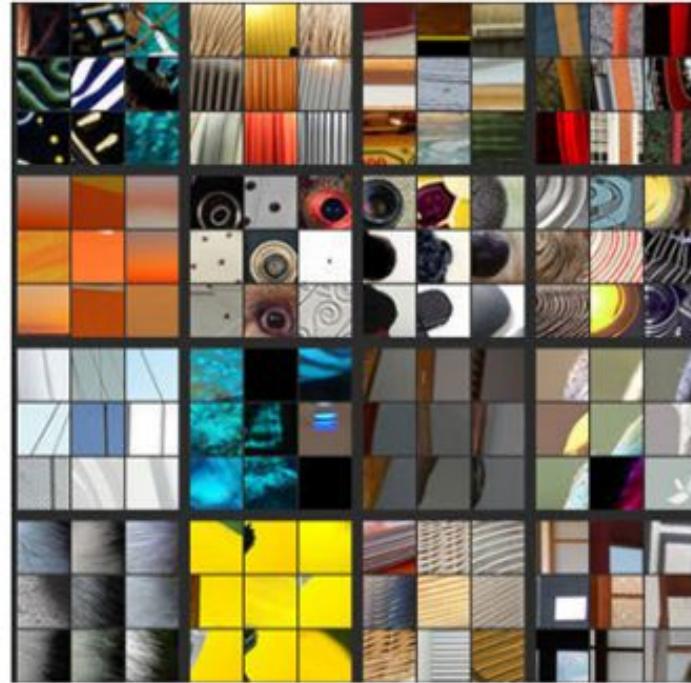
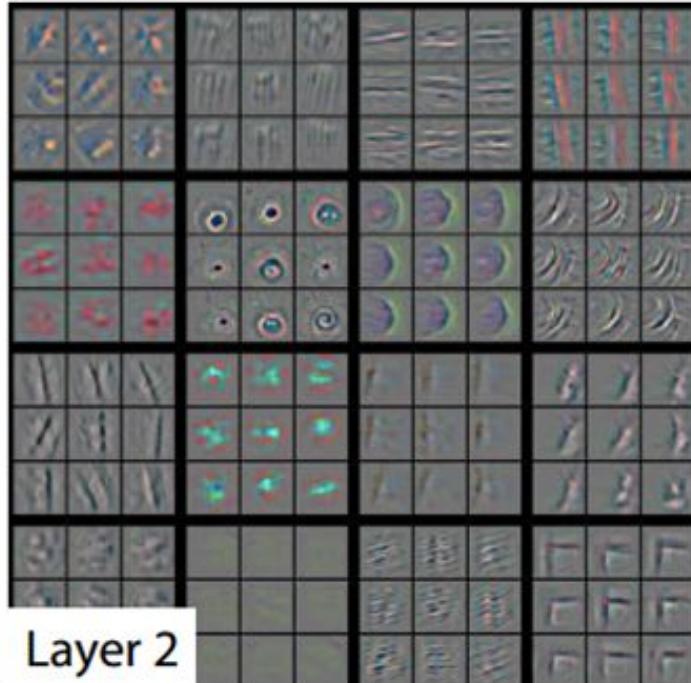
ResNext



Layer 1



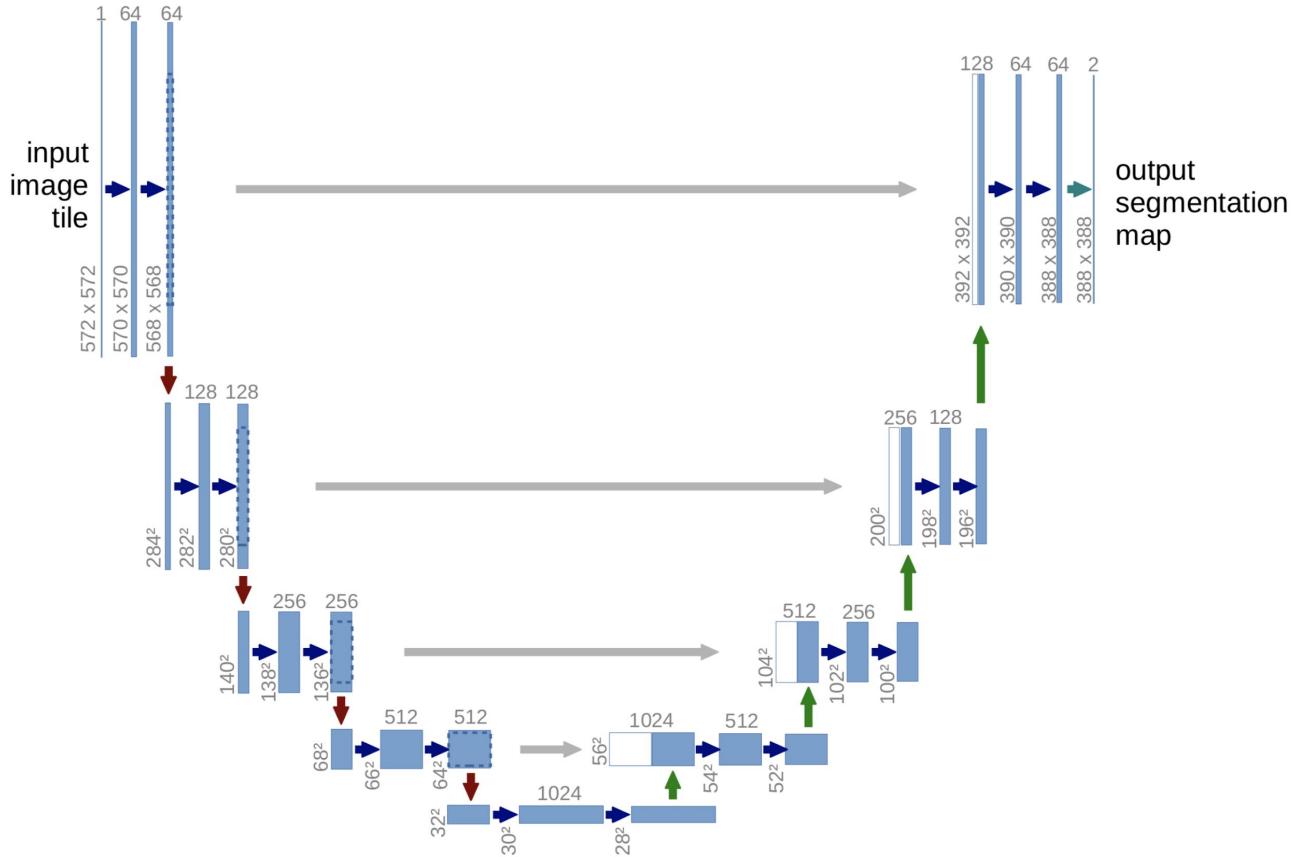
Layer 2

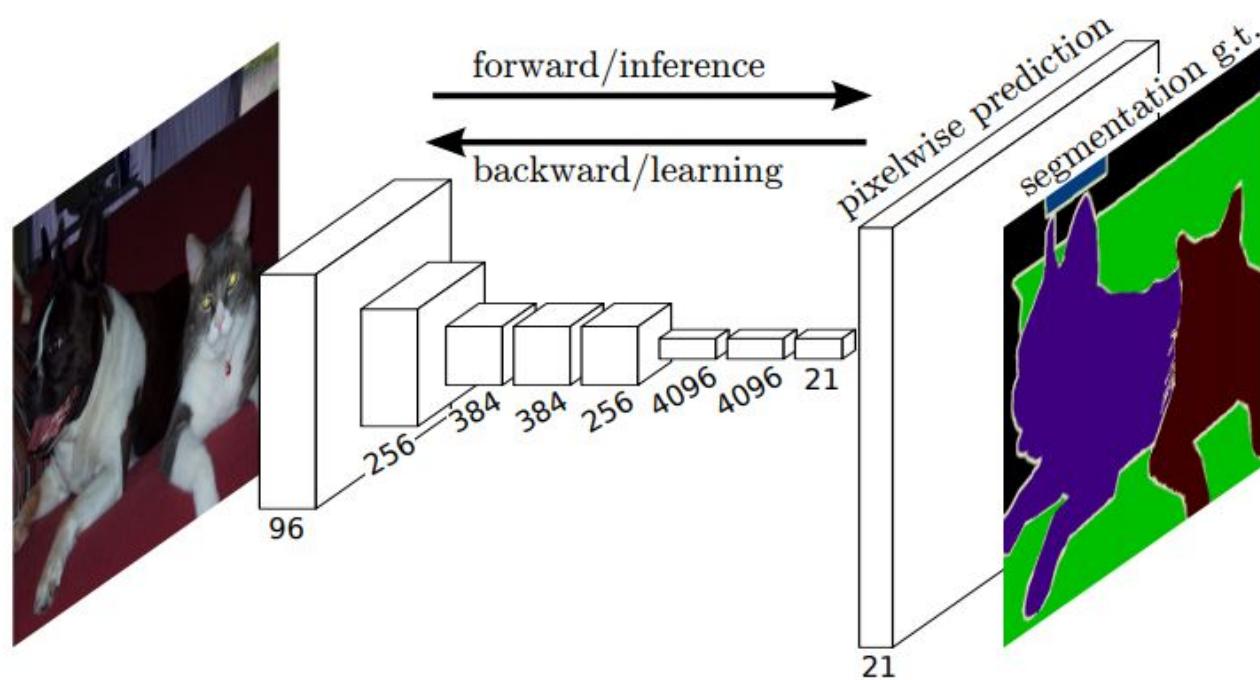


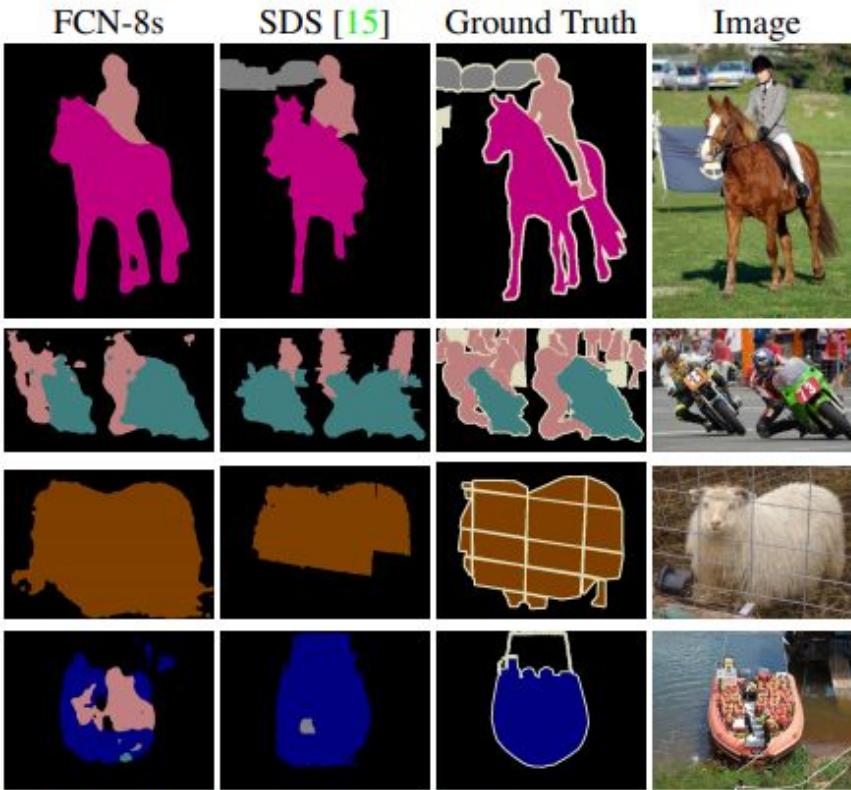
Visualizations of Layer 1 and 2. Each layer illustrates 2 pictures, one which shows the filters themselves and one that shows what part of the image are most strongly activated by the given filter. For example, in the space labeled Layer 2, we have representations of the 16 different filters (on the left)

Segmentadores usando Fully convolutional NN

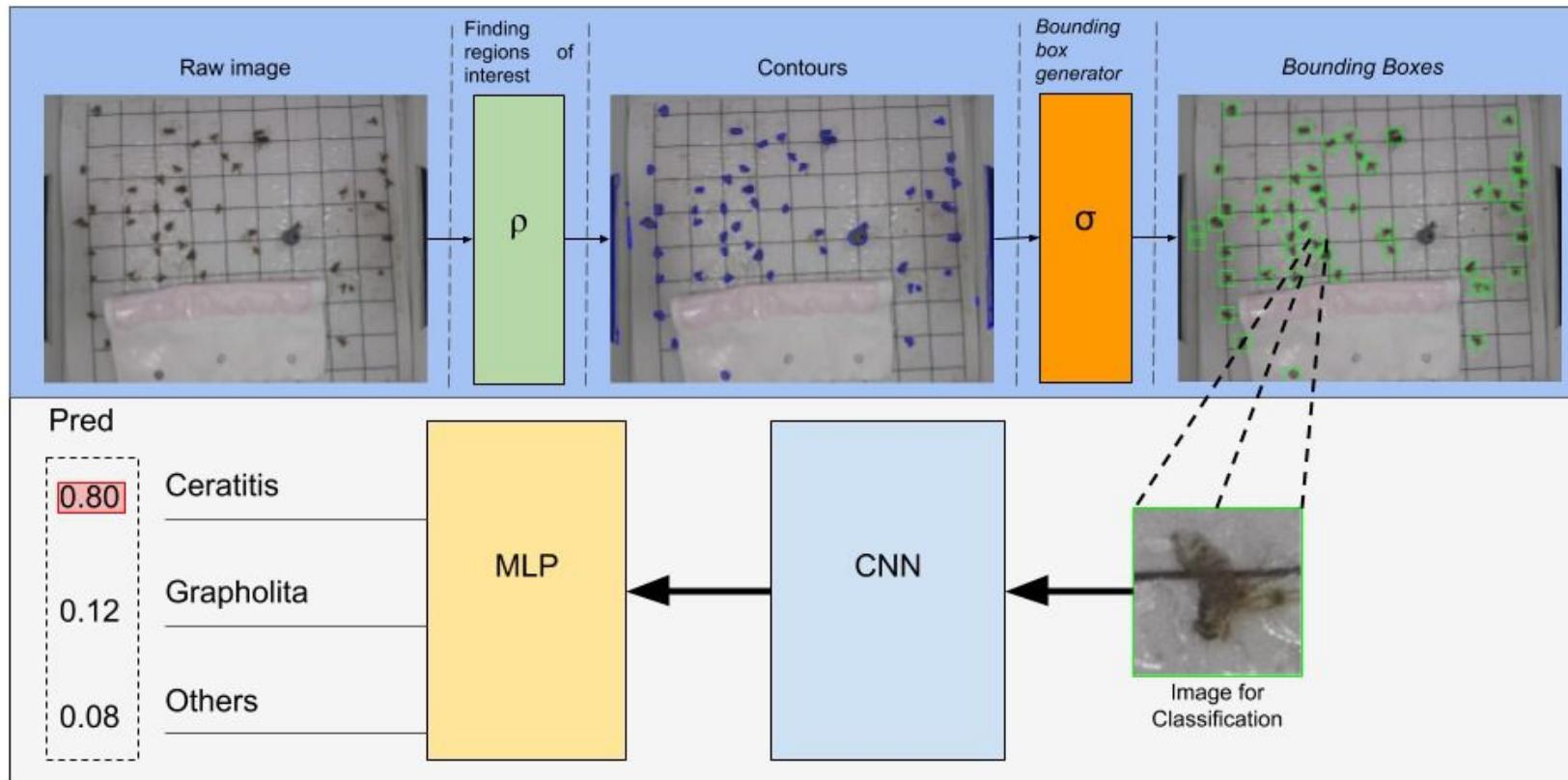
- Semantic segmentation (multiclass classification in a single image)
- It is a labor work to build a dataset for this task





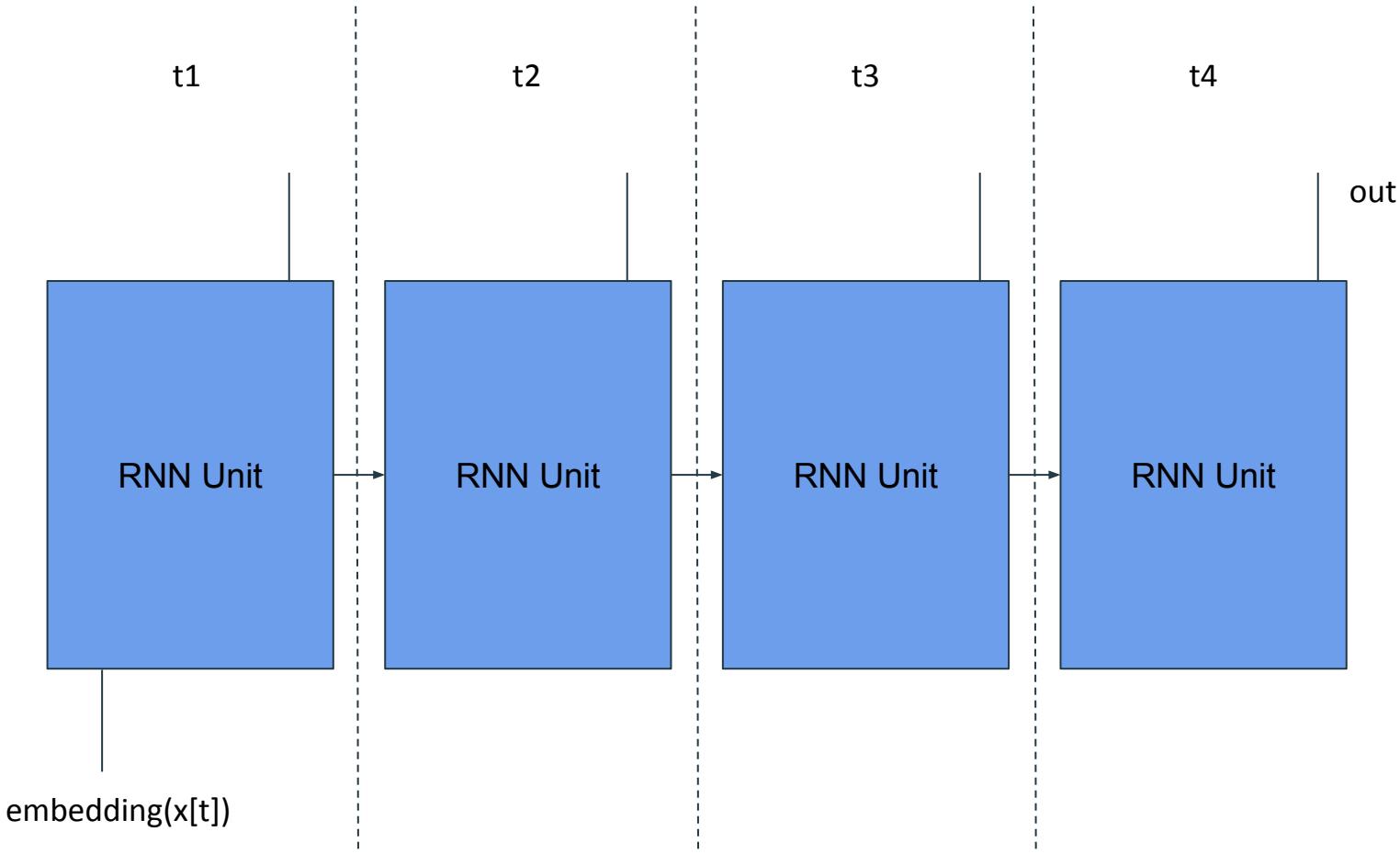


R-based CNN

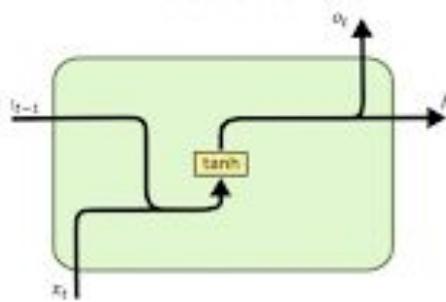


Recurrent Neural Networks

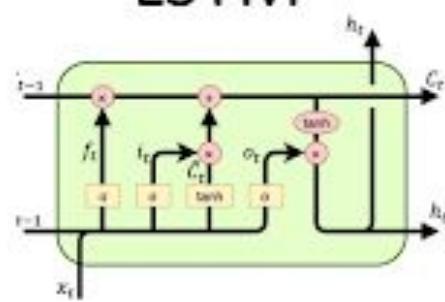
- applied to time-series data



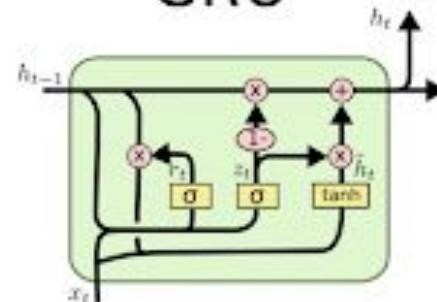
RNN



LSTM

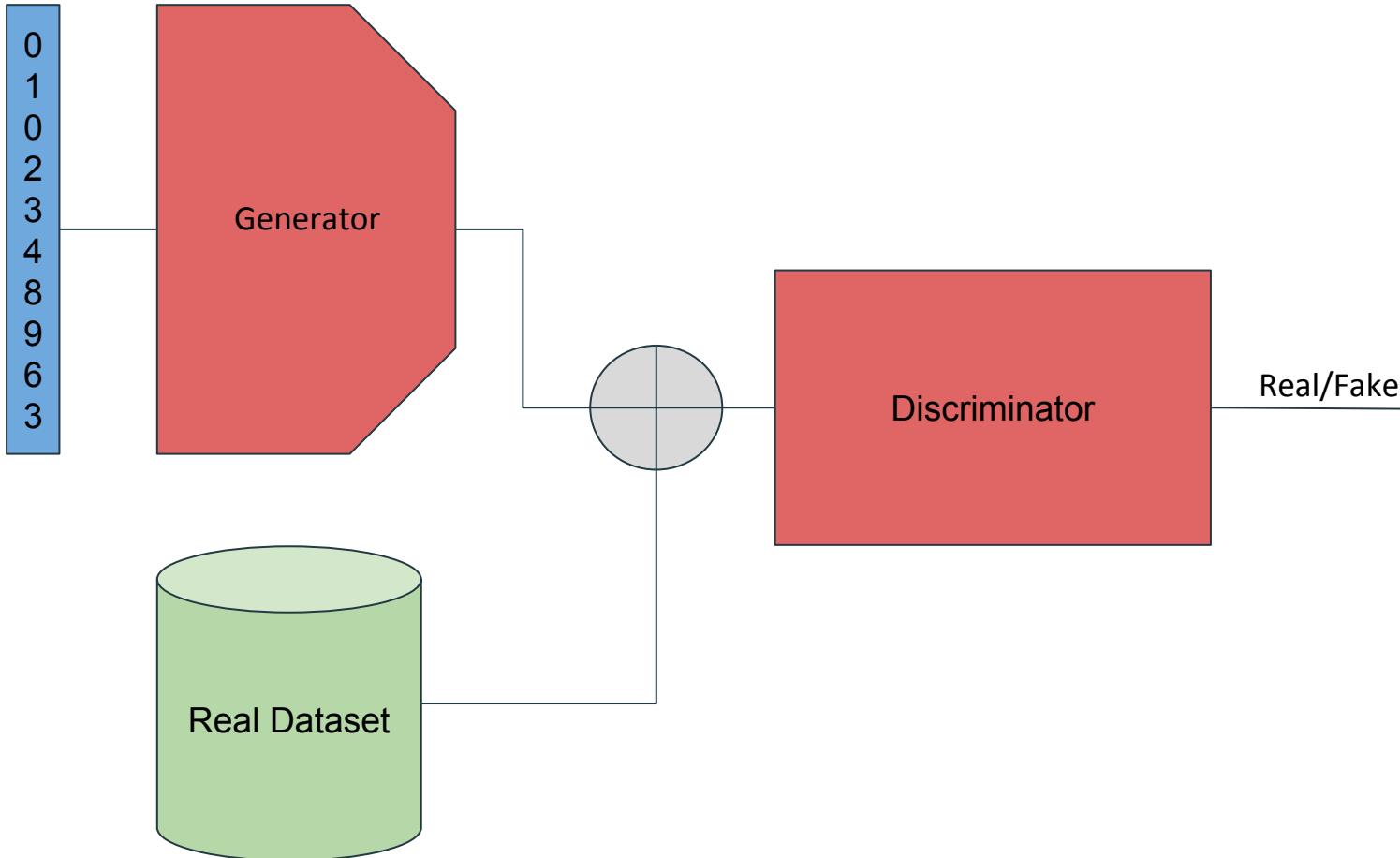


GRU



Generative Adversarial Networks

- Applications
- Hard convergence
- Generates artificial datasets





a)



b)

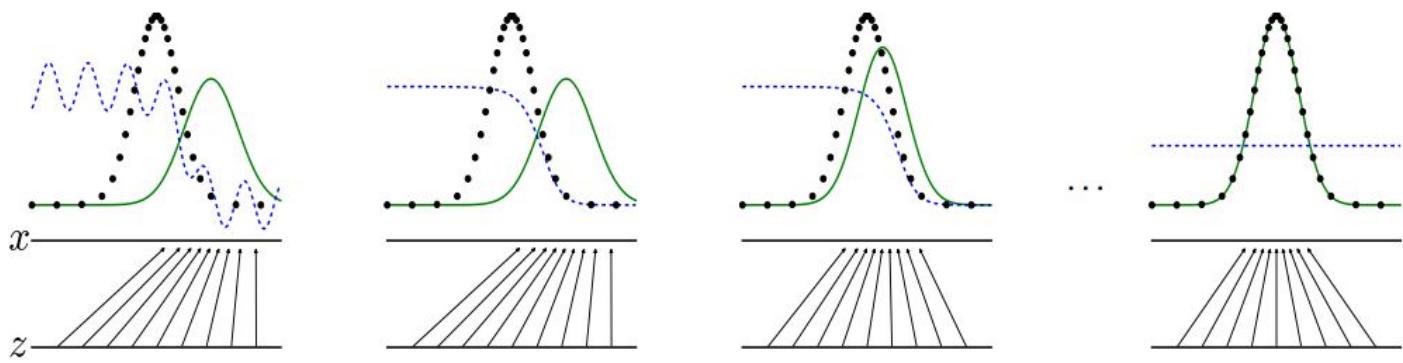


c)



d)

Generative Adversarial Networks; Goodfellow, Ian.



Generative Adversarial Networks; Goodfellow, Ian.



thispersondoesnotexit.com

Interesting Links

<https://affinelayer.com/pixsrv/> : adding fancy to your sketch

https://colab.research.google.com/github/tensorflow/hub/blob/master/examples/colab/biggan_generation_with_tf_hub.ipynb#scrollTo=HuCO9tv3IKT2 : interpolation of a dog and a bath

<https://www.thispersondoesnotexist.com/> : sort of people who don't exist

<https://www.youtube.com/watch?v=MwtVkPKx3RA> : Death Metal generated by a RNN

<https://www.youtube.com/watch?v=cQ54GDm1eL0> : Fake Obama

Interesting Links

<https://machinelearningmastery.com/how-to-configure-the-number-of-layers-and-nodes-in-a-neural-network/>

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

<https://www.nature.com/articles/s41467-019-10168-2.pdf>

<https://arxiv.org/pdf/1811.12808.pdf>

<https://jalammar.github.io/illustrated-word2vec/>

References

- RUDER, Sebastian. An overview of gradient descent optimization algorithms. **arXiv preprint arXiv:1609.04747**, 2016.
- LONG, Jonathan; SHELHAMER, Evan; DARRELL, Trevor. Fully convolutional networks for semantic segmentation. In: **Proceedings of the IEEE conference on computer vision and pattern recognition**. 2015. p. 3431-3440.LONG, Jonathan; SHELHAMER, Evan; DARRELL, Trevor. Fully convolutional networks for semantic segmentation. In: **Proceedings of the IEEE conference on computer vision and pattern recognition**. 2015. p. 3431-3440.
- Goodfellow, Ian, et al. "Generative adversarial nets." *Advances in neural information processing systems*. 2014.
- LECUN, Yann et al. Gradient-based learning applied to document recognition. **Proceedings of the IEEE**, v. 86, n. 11, p. 2278-2324, 1998.
- LECUN, Yann et al. Backpropagation applied to handwritten zip code recognition. **Neural computation**, v. 1, n. 4, p. 541-551, 1989.

References

GOODFELLOW, Ian; BENGIO, Yoshua; COURVILLE, Aaron. **Deep learning**. MIT press, 2016.