

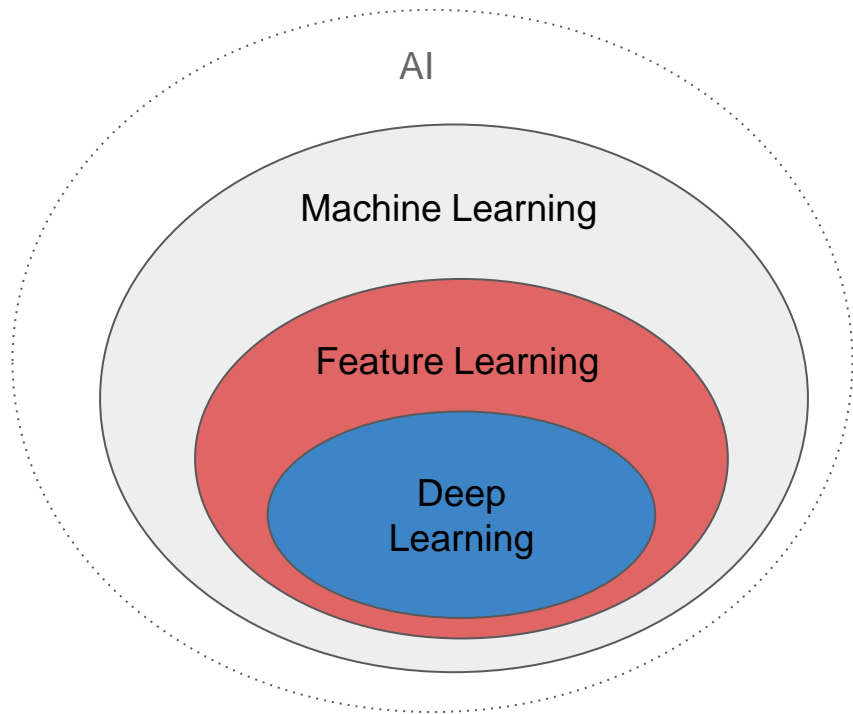
Deep Neural Networks

Thanks to: [Deep Learning by Google - Take machine learning to the next level](http://www.boredpanda.com/inceptionism-neural-network-deep-dream-art/)



Example of Inceptionism
<http://www.boredpanda.com/inceptionism-neural-network-deep-dream-art/>

What is Deep Learning



Deep Learning (DL) has emerged around the '10 as a general tool to solve recognition problems in:

- computer vision
- speech recognition
- robotics
- *discovering* new medicines
- *understanding* natural language
- *understanding* documents
- ranking
- ... and many other applications!

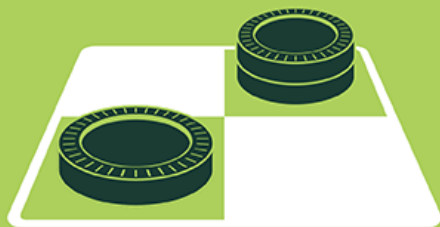
Overview

- History
- Preliminaries: logistic classification
- Training
- Deep networks
- Regularization
- Architectures
 - Convolutional networks
 - Embeddings
 - Recurrent models

History

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's

1960's

1970's

1980's

1990's

2000's

2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

Everything can be optimized in Computer Science

Given a problem to solve P , it can be formalized as $\{P, C, F\}$

P := the problem formulation

$C = \{c_1, c_1, \dots, c_n\}$:= set of configurations, each one of them representing a possible solution to P

$f: C \rightarrow R$:= function which provides a goodness measure of the configuration w.r.t. the problem to be solved

Casting the problem via minimization means to maximize or minimize the function f in the C space, *independently on the implied meaning of P*

Minimization: to be used *always*?

Problem P_1 : *sort numbers x_1, x_2, \dots, x_N in increasing order*

In this case, minimization could be left apart

In facts, there is at least one algorithm (e.g., quicksort) which brings directly to the best (in sense of the function f) configuration

Problem P_2 : *foresee the stocks' trend*

Much more difficult to formalize into an algorithm

Minimization comes to help [Yong et al. 2015]

[Yong et al. 2015] Hu, Yong, et al. "Application of evolutionary computation for rule discovery in stock algorithmic trading: A literature review." *Applied Soft Computing* 36 (2015): 534-551.

Inside the minimization approach

The main goal of an optimization approach is that of exploring the configuration space C looking for the best configuration given the function f (obviously avoiding the brute force way!)

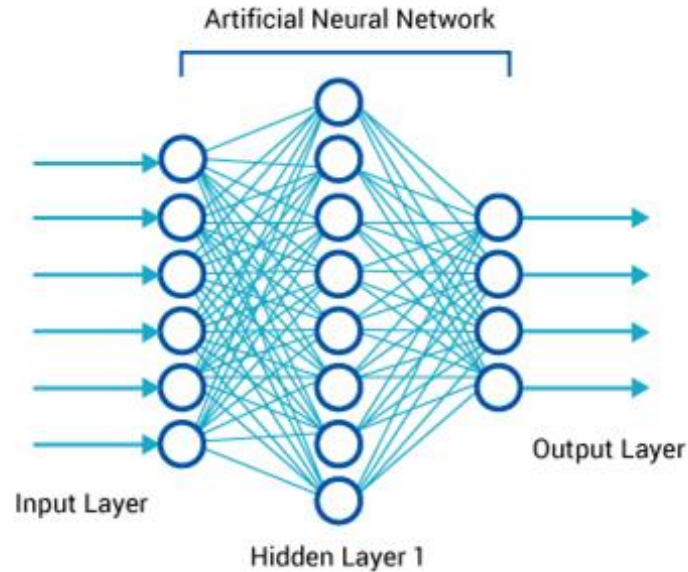
The set of configurations C give a space to explore (very often, a manifold)

Optimizing means to explore the manifold by iterative approaches (e.g., the gradient descent family of strategies)

The more the manifold is complex (non concave, multimodal), the more often local minima are met

Neural Networks [1943 - McCulloch & Pitts]

- Optimization approaches which scale *very well* with data
- We are talking about *artificial neurons* and *layered computation*

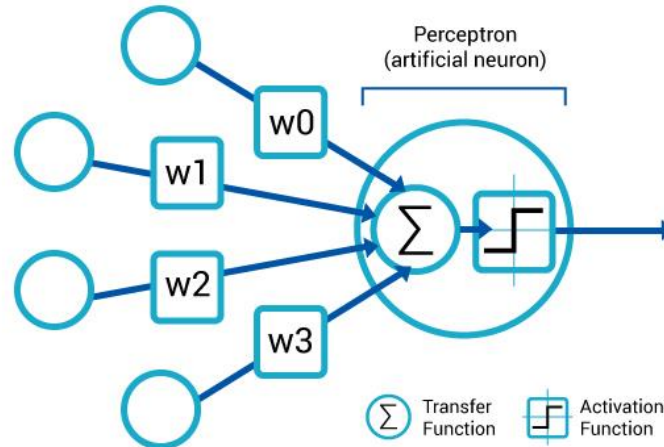
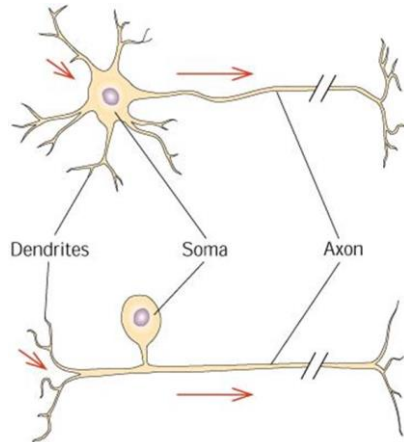


Neural Networks - Neurons

NN are composed by artificial neurons (1943 - McCulloch & Pitts).

Each neuron has:

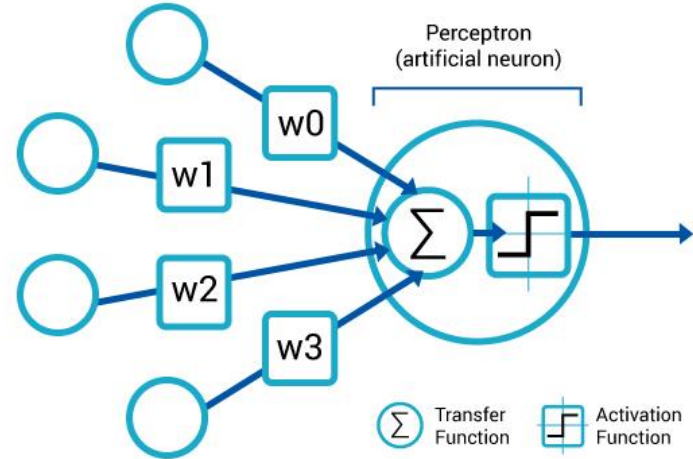
- *dendrites* (**inputs**)
- a *nucleus/soma/perceptron* (**transfer fuction + activation function**)
- an *axon* (**output**)

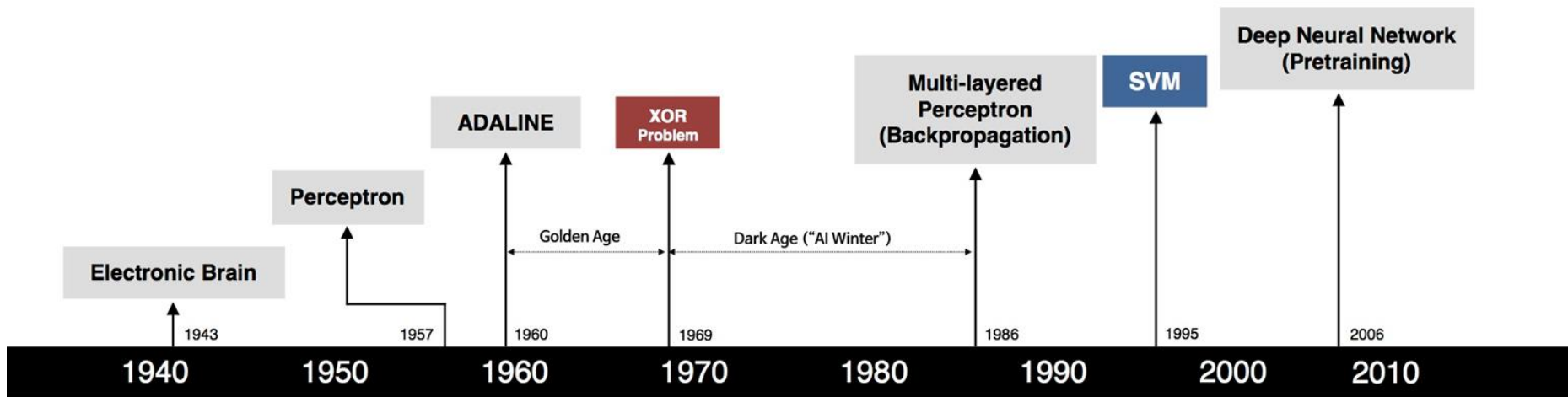


Neural Networks - Neurons (2)

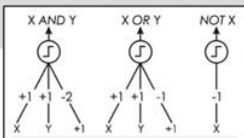
The information flow is **unidirectional**:

- The neuron get **inputs** (electric potentials) from the *dendrites*, that weight them (w_i 's)
- In the *nucleus*, the weighted inputs are summed together (the **transfer** Σ of the whole information coming from the dendrites)
- In the *nucleus*, the summation flows into an **activation function**, which may inhibit, diminish or amplify it
- The **output** of the activation function is transmitted through the *axon*





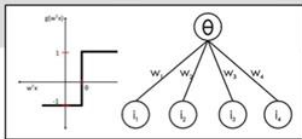
S. McCulloch – W. Pitts



- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



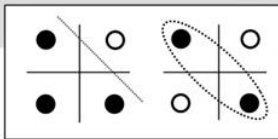
- Learnable Weights and Threshold



B. Widrow – M. Hoff



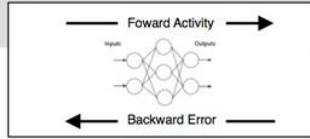
M. Minsky – S. Papert



- XOR Problem



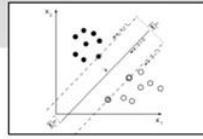
D. Rumelhart – G. Hinton – R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



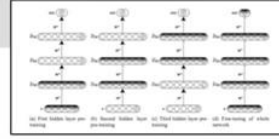
V. Vapnik – C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention

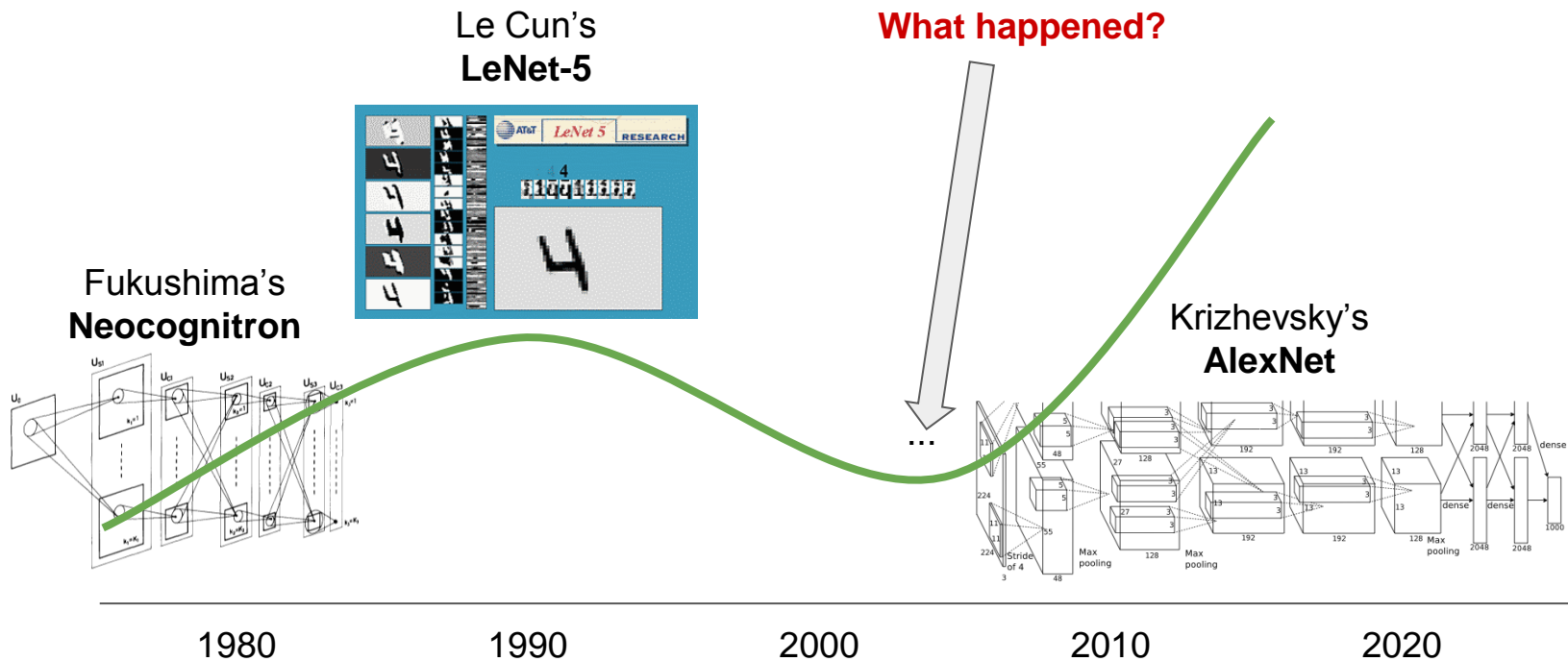


G. Hinton – S. Ruslan

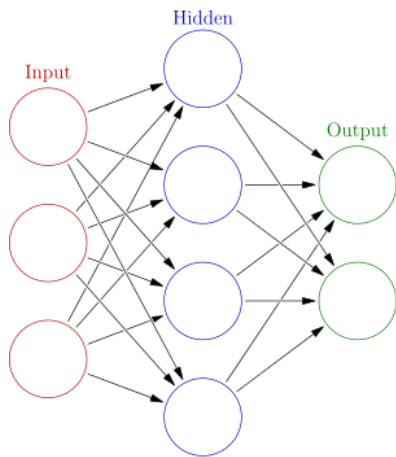


- Hierarchical feature Learning

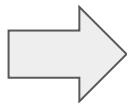
Neural Networks - the renaissance



Neural Networks and Deep Learning



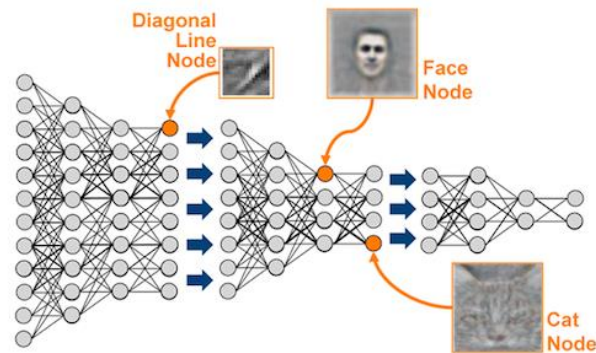
Neural
Networks



GPUs



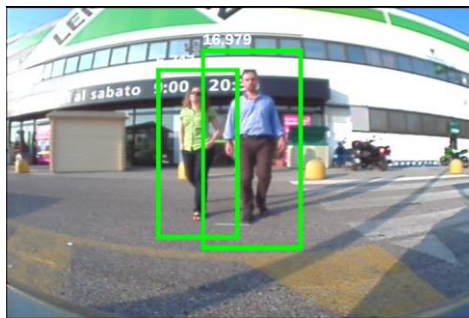
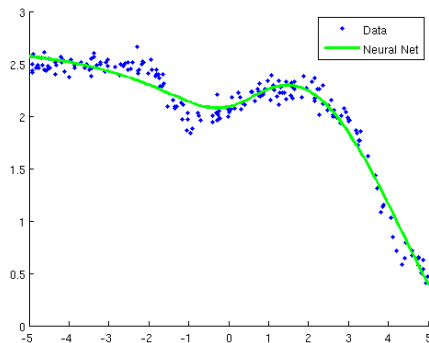
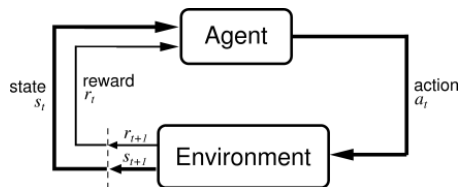
Data



Deep
Neural
Networks

Supervised Classification

- Traditional kind of problem the NN do solve
 - Regression
 - Ranking
 - Reinforcement Learning
 - Detection

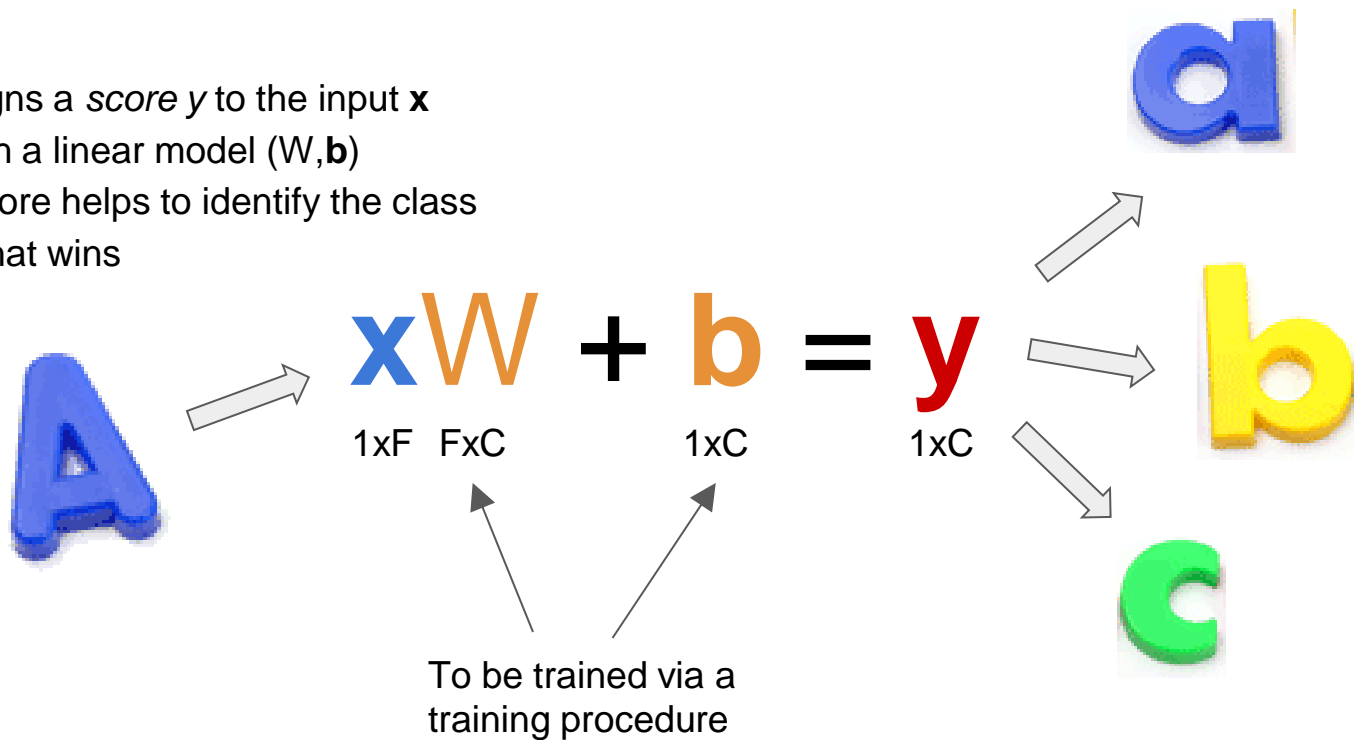


Labels {'a', 'b', 'c', 'd', 'e'}

Preliminaries: logistic classification

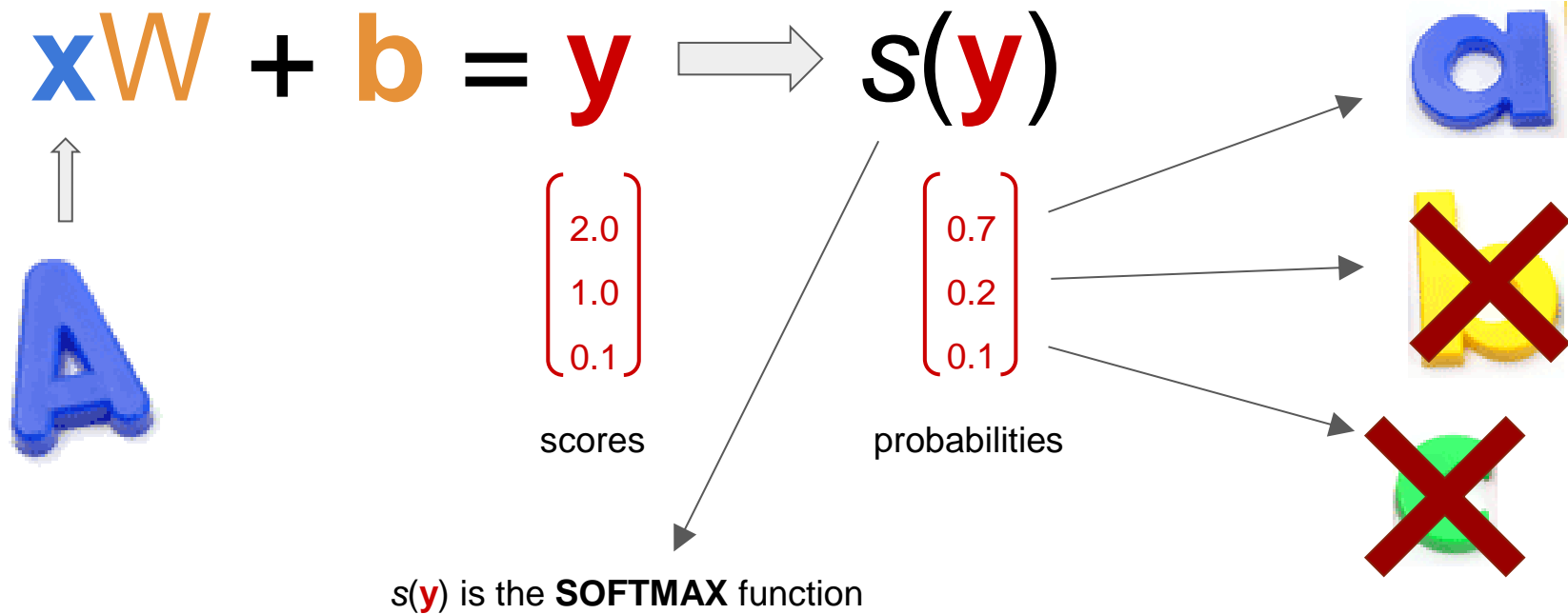
Logistic Classifier

- It assigns a *score* y to the input x through a linear model (W, \mathbf{b})
- The score helps to identify the class label that wins



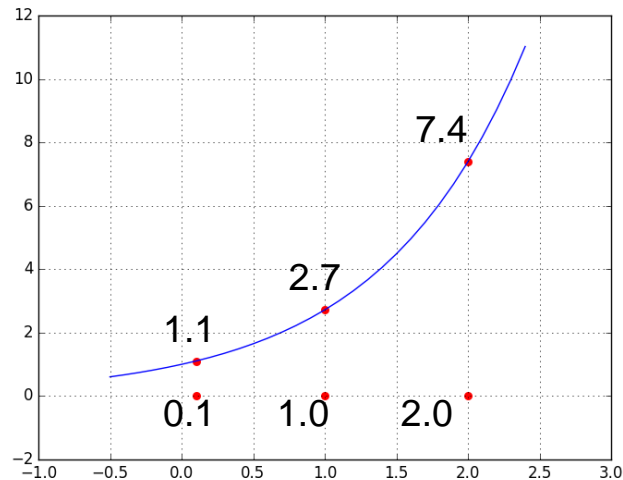
x = input or feature vector; F = number of features; W = weights matrix;
 C = number of classes b = biases; y = output or logits/scores vector

Logistic Classifier - the score is not enough



Softmax function

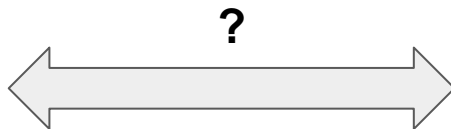
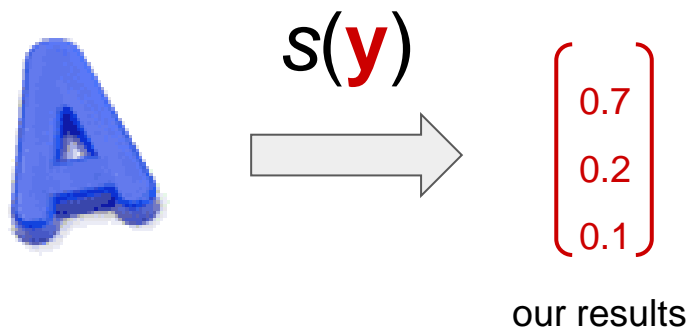
- Converts scores into probability distributions
 - $\mathbb{R} \rightarrow (0,1)$
 - *Open* codomain!
- The softmax function highlights the largest values and suppress values which are significantly below the maximum value



$$\begin{matrix} \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \\ \text{scores} \end{matrix} \Rightarrow s(y_i) = \frac{e^{y_i}}{\sum_{j=1}^C e^{y_j}} \Rightarrow \begin{matrix} \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix} \\ \text{probabilities} \end{matrix}$$

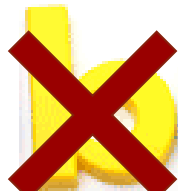
One-Hot Encoding (OHE)

- → *There is only one correct label for each input sample*
- → *There is the need to evaluate the classification result*
- OHE encodes labels for a C-class problem in \mathbb{R}^C , indicating the c-th class label with 1, the rest with 0's
- OHE needs sparse representation which is inefficient especially in the case C is very big (thousands or millions of classes...!)



Which distance measure?

MSE
Cross-Entropy



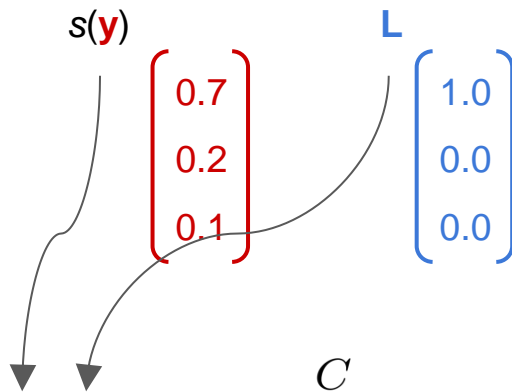
Cross-Entropy

GOAL: computes the distance between two probability vectors

- Non symmetric function

$$D(\mathbf{s}, \mathbf{L}) \neq D(\mathbf{L}, \mathbf{s})$$

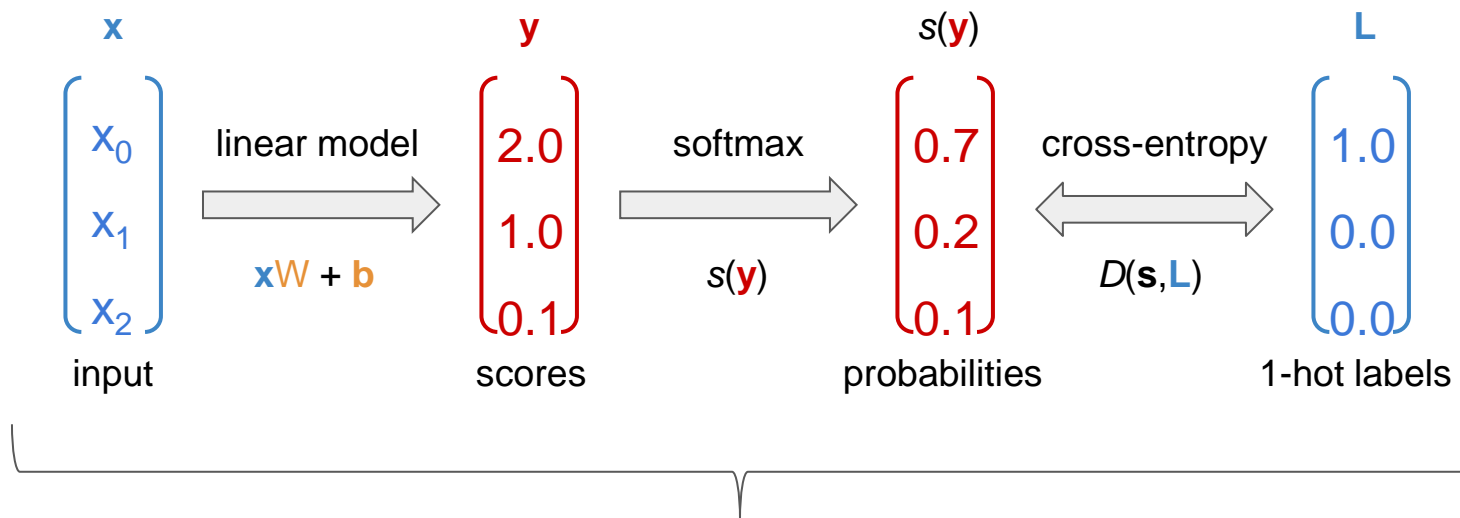
$$D(\mathbf{s}, \mathbf{L}) \begin{cases} \rightarrow 0 & \text{high similarity} \\ \rightarrow \infty & \text{low similarity} \end{cases}$$



$$D(\mathbf{s}, \mathbf{L}) = - \sum_{j=1}^C L_j \log s_j = 0.36$$

- The $\log(\cdot)$ makes the training faster to converge than the MSE function ($\sum (\mathbf{s}(\mathbf{y}) - \mathbf{l})^2$)

Résumé



Multinomial Logistic Classification

$$D(s(xW + b), L)$$

Training

Gradient Descent

GOAL: search for the nearest local minimum of a function F

IDEA: iterate on the parameter set proportionally to the negative of the function gradient

$$\theta_{t+1} = \theta_t - \alpha \nabla F(\theta_t), t \geq 0$$

such that

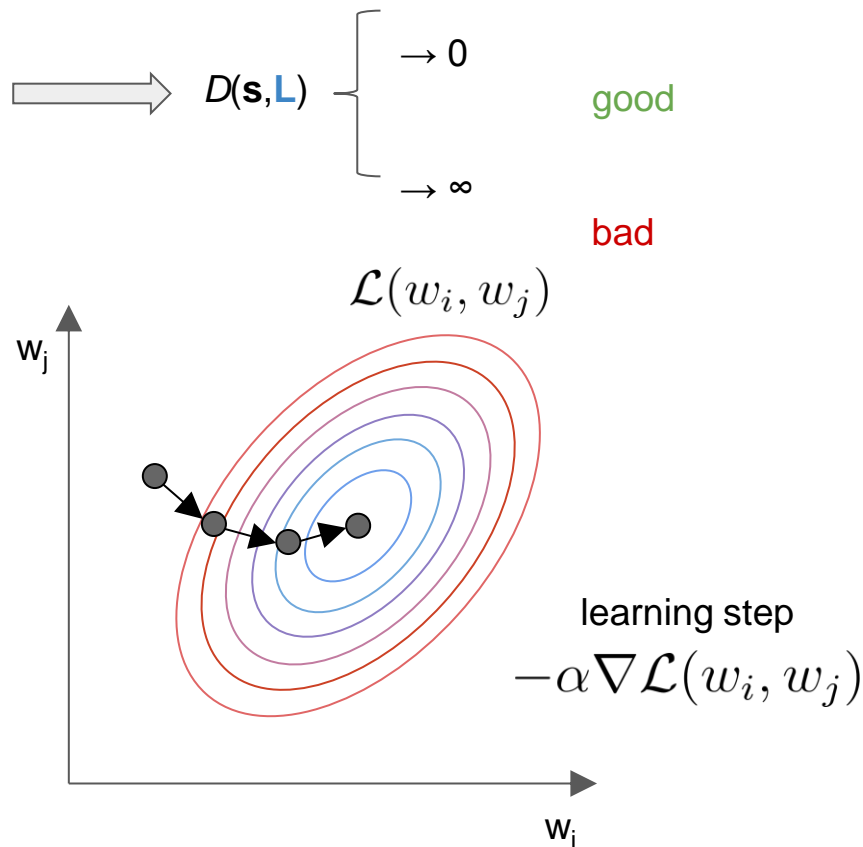
$$F(\theta_0) \geq F(\theta_1) \geq F(\theta_2) \geq \dots$$

Gradient Descent

- GOAL: minimize a loss function
- Needs to compute the entire training set performance of our linear model, that consists in N inputs (which is, in general, very big)
- Needs to minimize a loss function, which depends on W (big matrix) and \mathbf{b}

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N D(s(\mathbf{x}_i W + \mathbf{b}), \mathbf{L}_i)$$

Loss = average cross-entropy



Stochastic Gradient Descent (SGD)

- IDEA: use a random subset (*batch*) of the data (of a given *size*) to compute an approximation of the gradient of the loss function

Iterative implementation of the GD algorithm

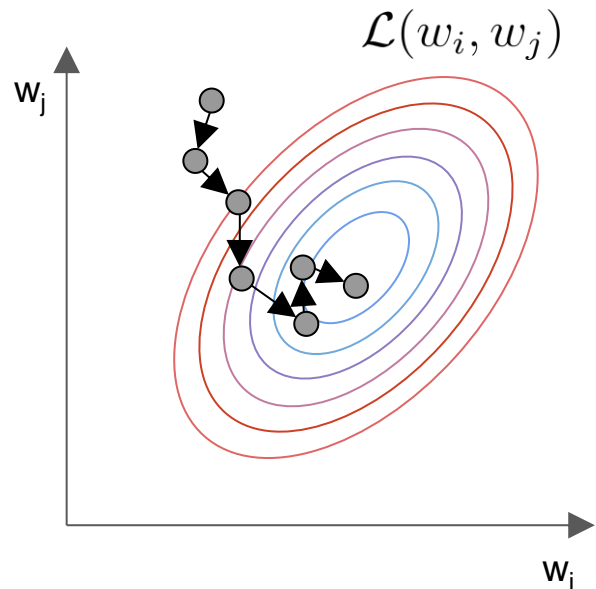
At each step, a new batch is extracted

- Pros:

- simple but sufficiently effective
- fast (depending on batch size)
- scale the problem with data and model size

- Cons:

- needs more iterations to converge
- bad approximation of the loss
- suffers from local minima
- requires data preprocessing to avoid numerical instability



... but tricks to ameliorate SGD are present!

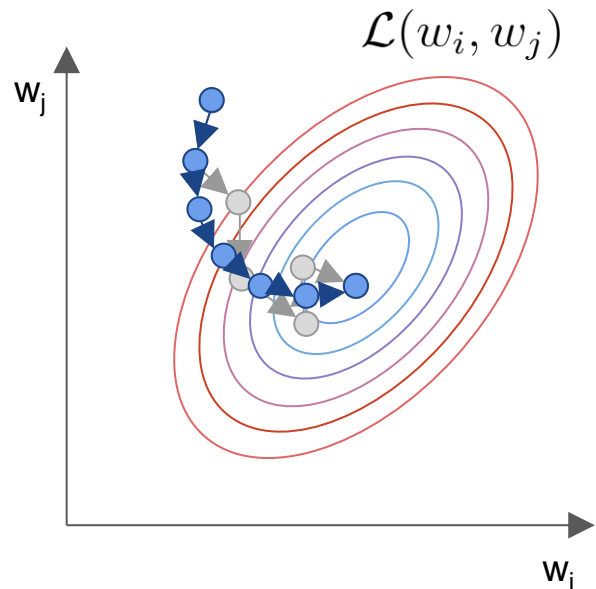
SGD trick 1: momentum

- GOAL: improve the convergence of the optimizer exploiting the accumulated knowledge from previous steps
- IDEA: add a fraction of the previous update vector to the current update vector

$$M_0 = \nabla \mathcal{L}_0(w_i, w_j) = 0$$

$$M_t = \alpha \nabla \mathcal{L}_t(w_i, w_j) + \beta M_{t-1}$$

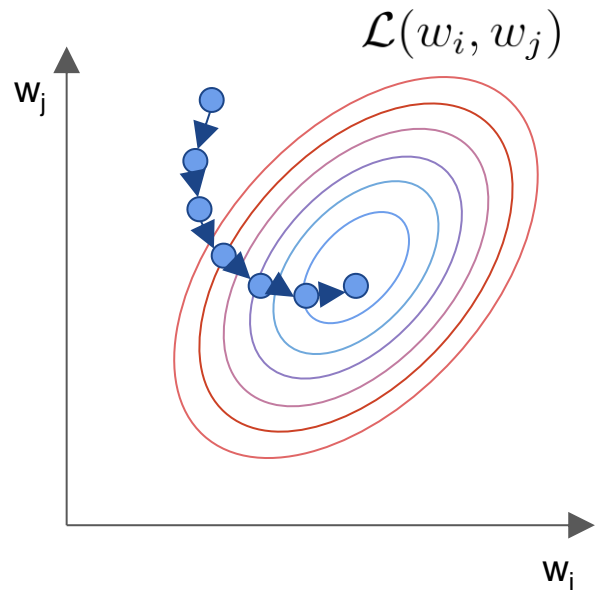
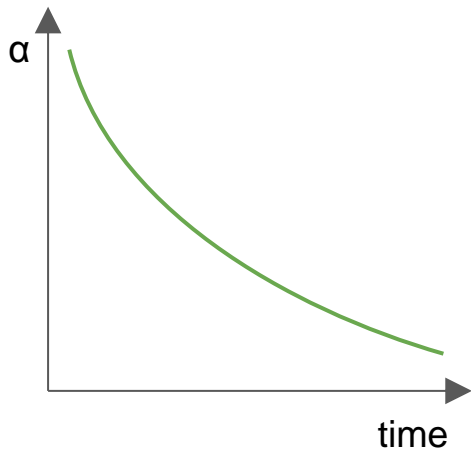
$$(w_i, w_j)_t = (w_i, w_j)_t - M_t$$



faster convergence and oscillation reduction

SGD trick 2: learning-rate decay

- GOAL: make the optimization more robust and accurate over time
- IDEA: apply a decay function to the learning rate or reduce it when the loss function reaches a plateau



SGD trick 3: z-normalization

- GOAL: avoid numerical instability

The values involved in the calculation of the gradient descent never get too big or too small

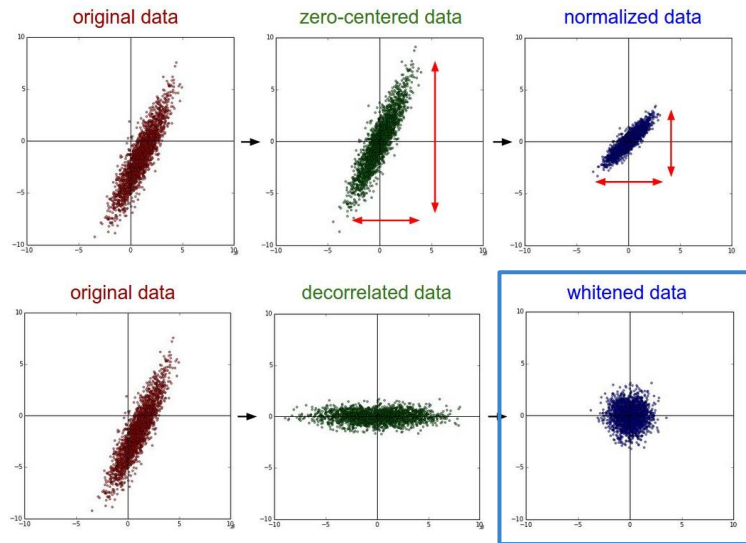
- IDEA: remove the mean and normalize over the variance of the i -th feature of the input vector \mathbf{x}

z-normalization (*whitening*)

$$\mathbf{n} = \frac{\mathbf{x} - E[\mathbf{x}]}{Var[\mathbf{x}]}$$

$$E[\mathbf{x}] = 0$$

$$\forall(i, j), Var[\mathbf{x}_i] = Var[\mathbf{x}_j]$$



SGD trick 4: initialization

A random initialization of the weights and the zero-init of the biases is critical to get a good starting point for the training phase and the convergence of the SGD algorithm.

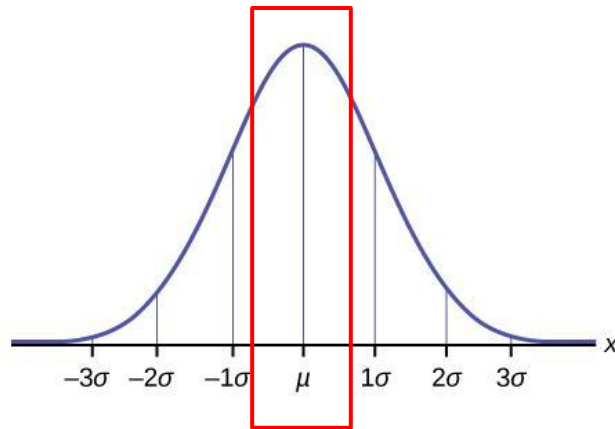
$$w_i = \mathcal{N}(\mu = 0, \sigma \rightarrow 0)$$

$$\mathbf{b} = [0_1 \quad 0_2 \quad \dots \quad 0_C]$$

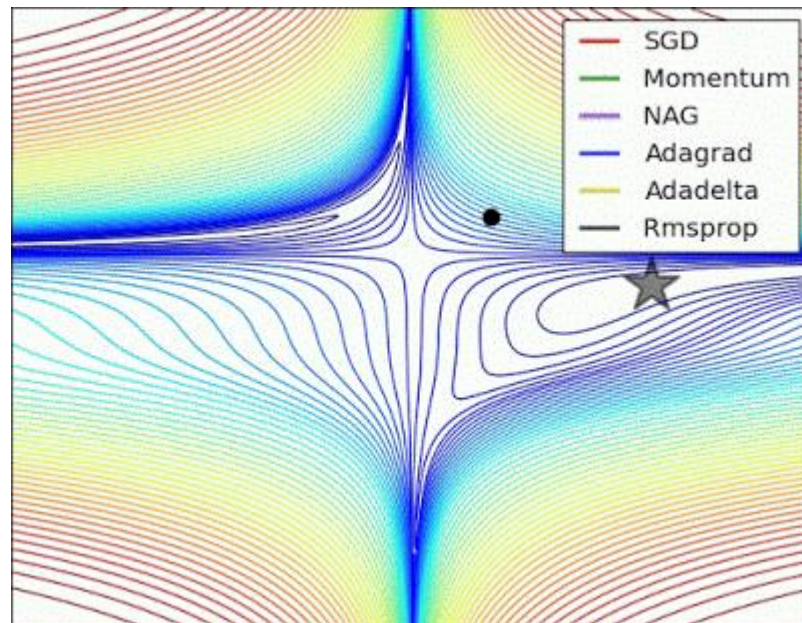
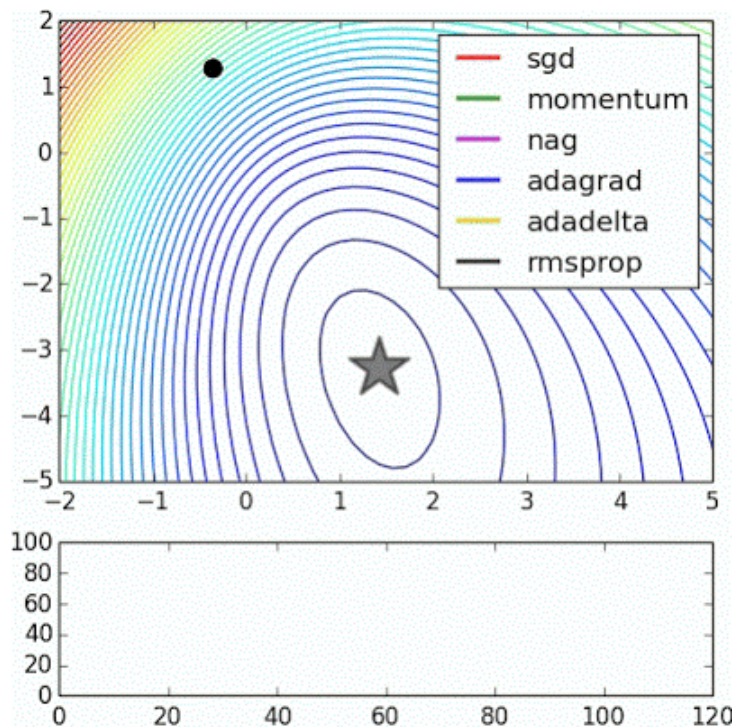
$$\mu = 0 \wedge \sigma \rightarrow 0 \implies \forall(i, j), w_i = w_j \pm \epsilon$$



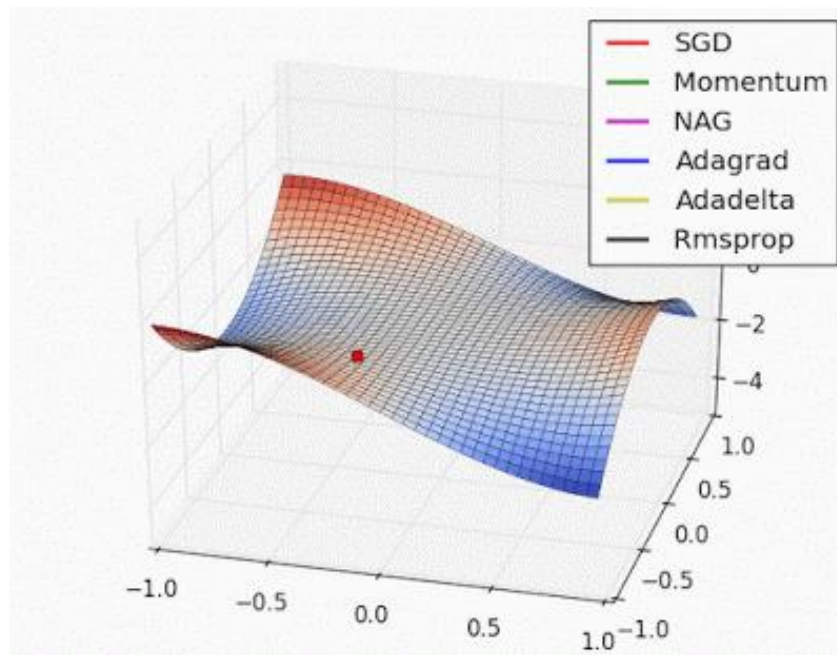
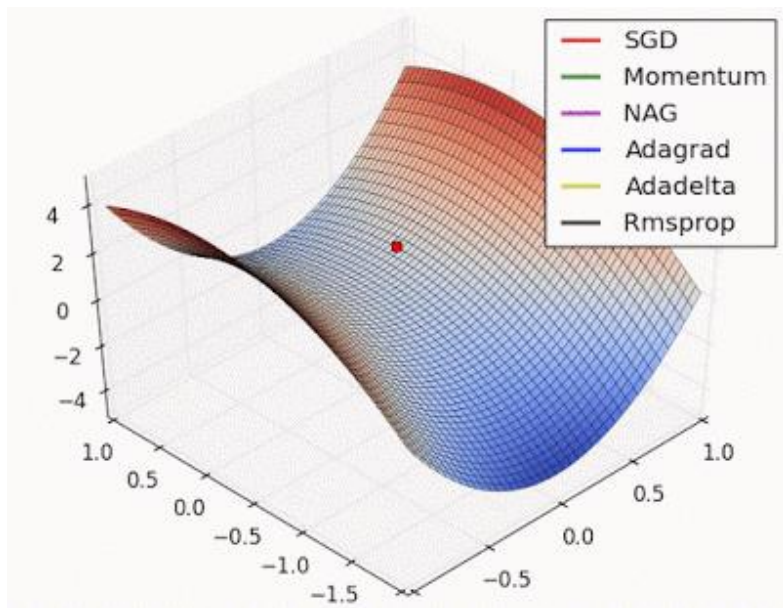
equal probability
of the weights
(no prior)



Gradient Descent: graphical representation (2D)



Gradient Descent: graphical representation (3D)



SGD: tuning

SGD

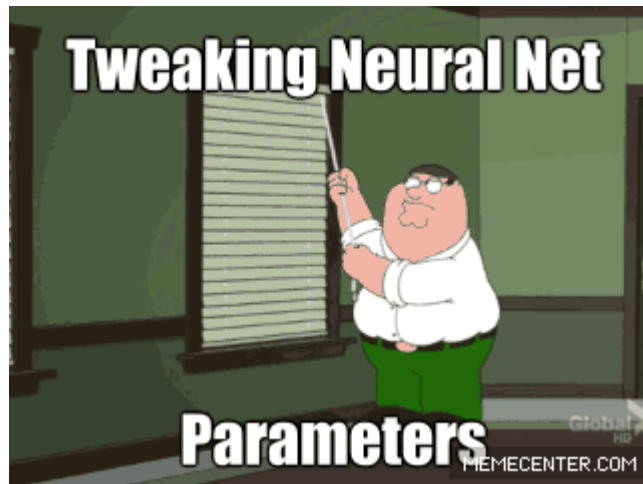
Many hyperparameters:

- initial learning rate
- learning rate decay
- momentum
- batch size
- weight initialization

AdaGrad

SGD modification which implicitly applies momentum and learning rate decay. It uses fewer parameters:

- batch size
- weight initialization



SGD: H. Robbins and S. Monro, "A stochastic approximation method," Annals of Mathematical Statistics, vol. 22, pp. 400–407, 1951.

AdaGrad: J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," in COLT, 2010.