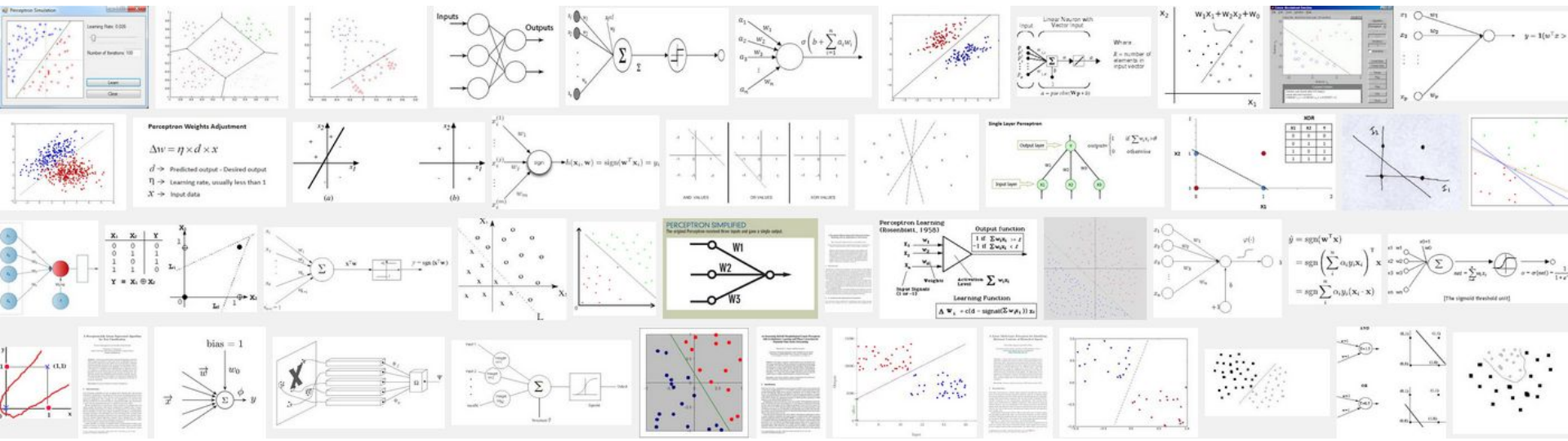
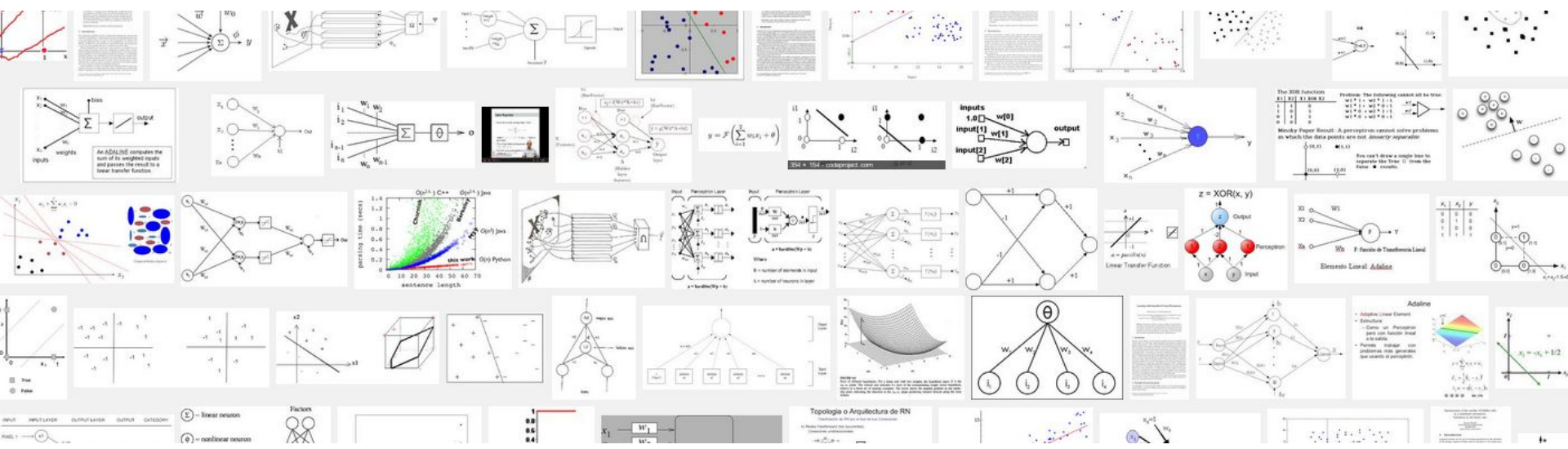


Unidad 2: Redes Neuronales Artificiales

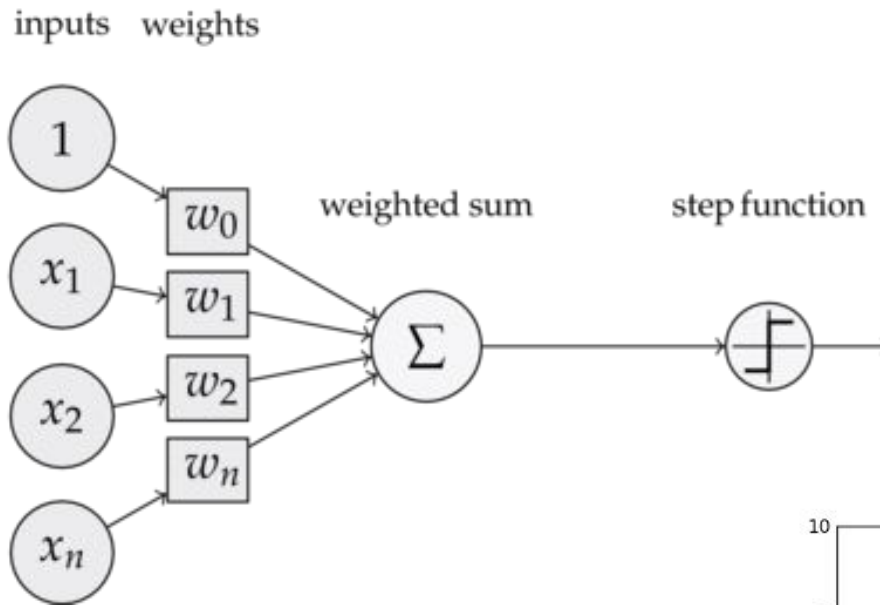
Curso: Redes Neuronales Profundas



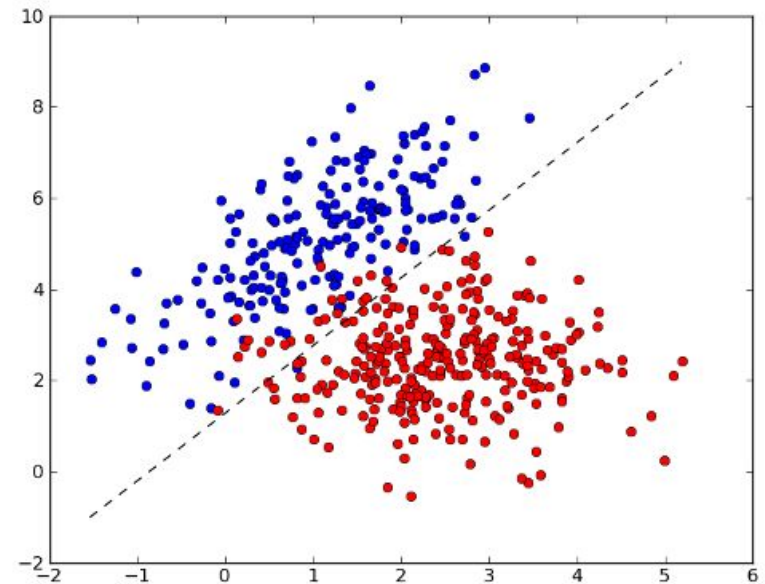
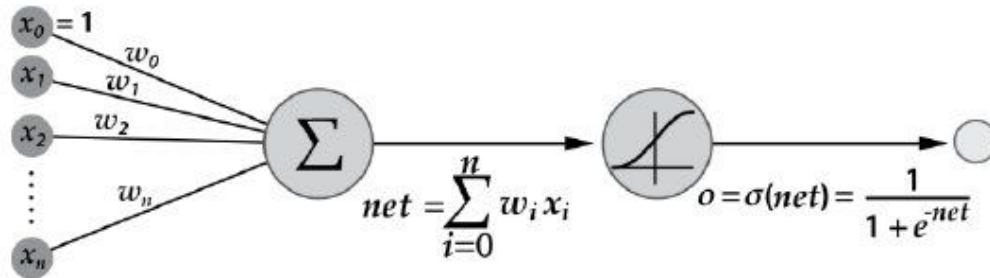
Linear Perceptron



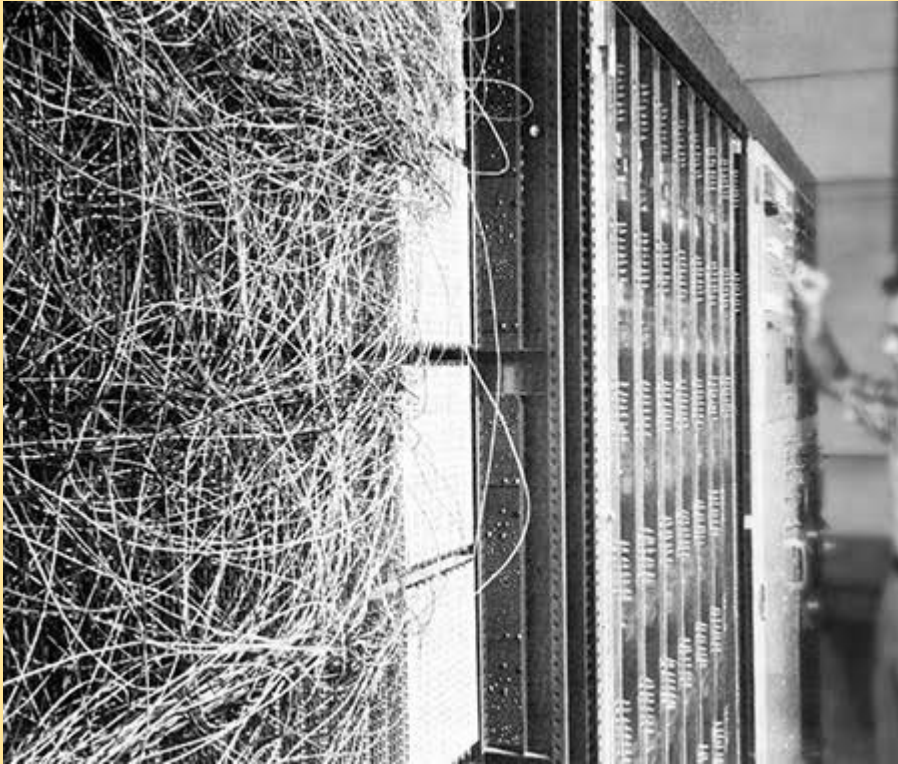
Linear Perceptron



$$g_0(x) = \sigma(x^T W^{(0)} + b^{(0)})$$



Mark I Perceptron (1957)



Input: Imágenes de 20x20 pixeles
(photocells)

Pesos adaptativos: arreglo de
potenciómetros

Updates: motores eléctricos

Frank
Rosenblatt
Cornell
Aeronautical
Laboratory



THE NEW YORK TIMES
REPORTED THE PERCEPTRON TO
BE "THE EMBRYO OF AN
ELECTRONIC COMPUTER THAT
[THE NAVY] EXPECTS WILL BE
ABLE TO

- WALK
- TALK
- SEE
- WRITE
- REPRODUCE ITSELF
- AND BE CONSCIOUS OF ITS
EXISTENCE



Softmax Regression

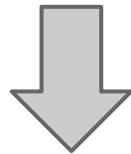
$$\log p(y|x) = \underbrace{x^T W + b}_{\text{affine transform.}} + \underbrace{c(x)}_{\text{normalization}}$$

class probability vector

parametric model

output (multiple categories)

input

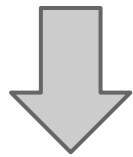


$$p(y|x) = \frac{\exp(x^T W + b)}{\sum_i \exp(x^T W + b)_i} = \text{softmax}(x^T W + b)$$

Softmax Regression Training

$$J(\mathcal{D}, W, b) = \prod_{x, y \in \mathcal{D}} p(y|x) \quad \text{Maximum Likelihood Estimation (MLE)}$$

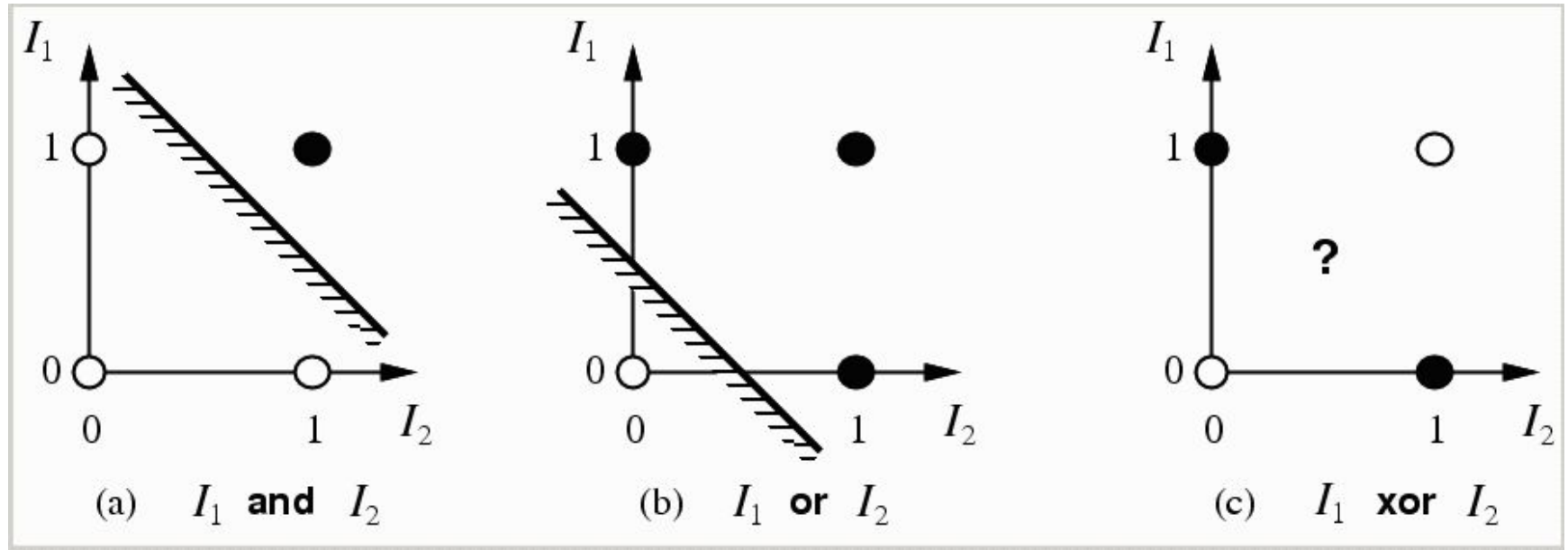
↑
training
dataset

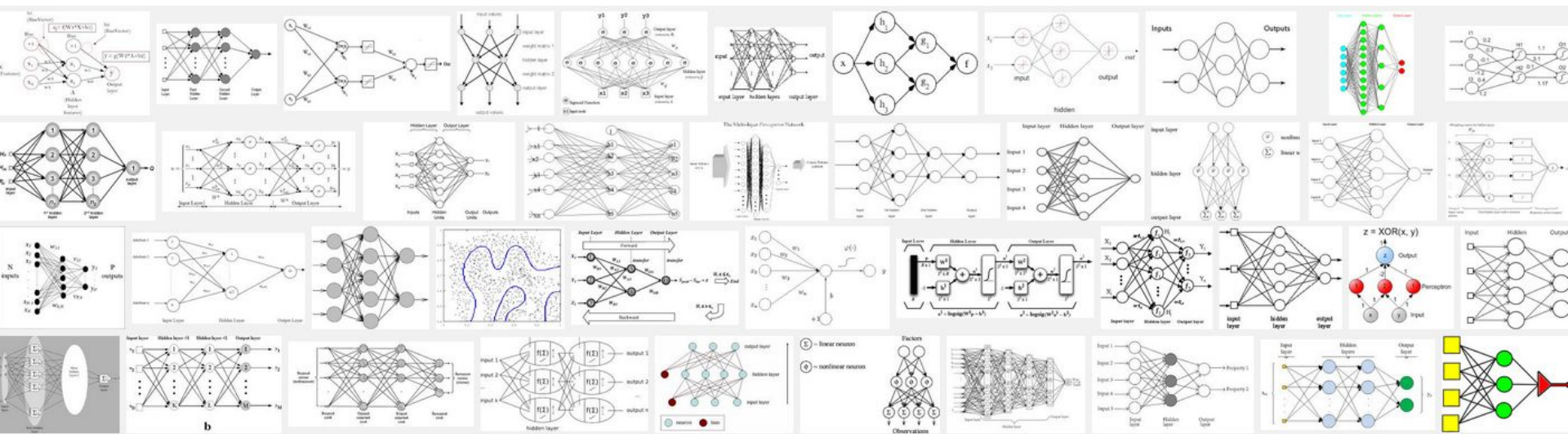


$$J(\mathcal{D}, W, b) = \sum_{x, y \in \mathcal{D}} \log p(y|x) \quad \text{Log-Likelihood}$$

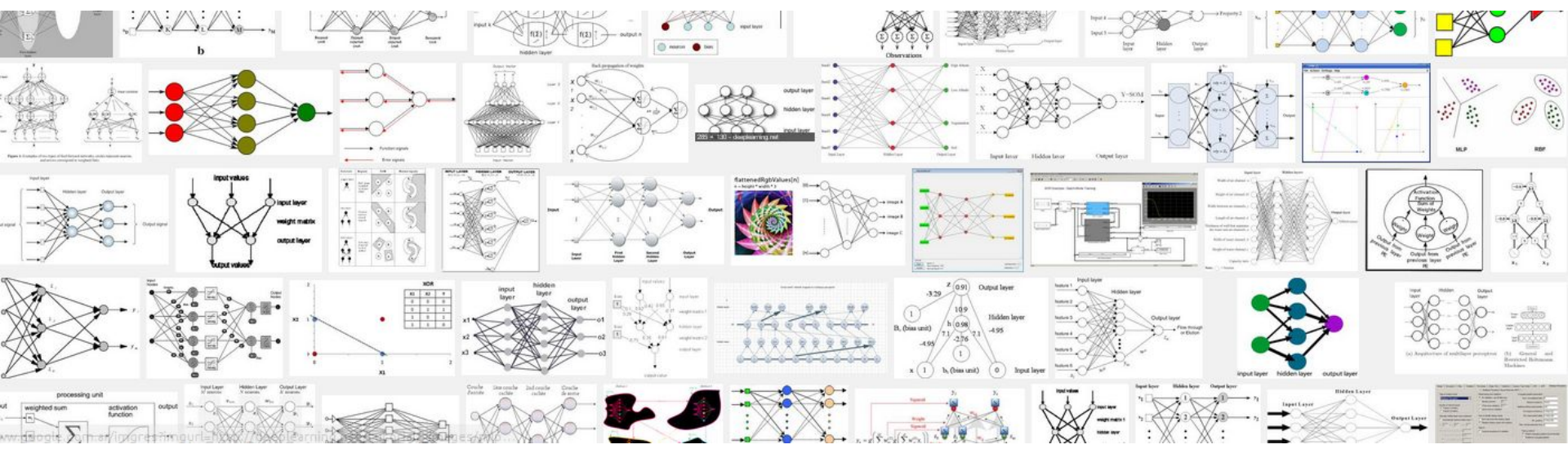
$$NLL(\theta, \mathcal{D}) = - \sum_{i=0}^{|\mathcal{D}|} \log P(Y = y^{(i)} | x^{(i)}, \theta) \quad \text{Negative Log-Likelihood}$$

Linear Perceptron





Multilayer Perceptrons (MLP)



Multilayer Perceptrons (MLP)

$$J(\mathcal{D}, W, b) = \sum_{x, y \in \mathcal{D}} \log p(y|x) \quad \text{Log-Likelihood}$$

$$\log p(y|x) = \cancel{x^T W + b + c(x)}$$

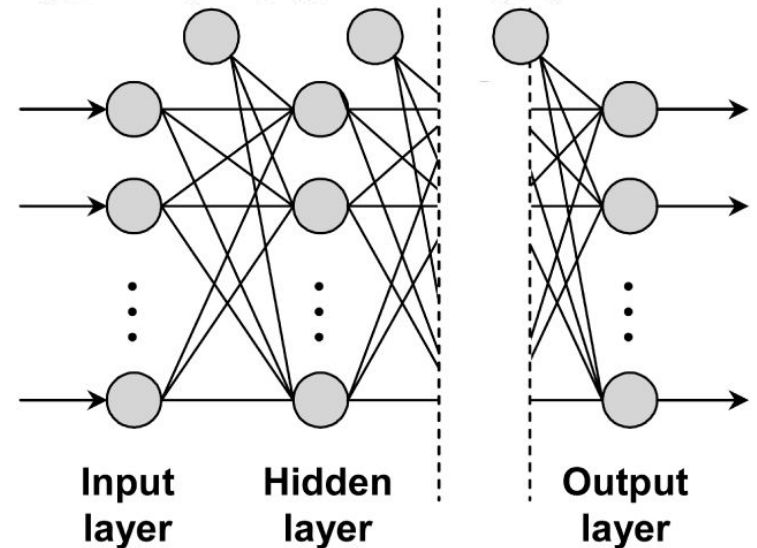
Single layer softmax



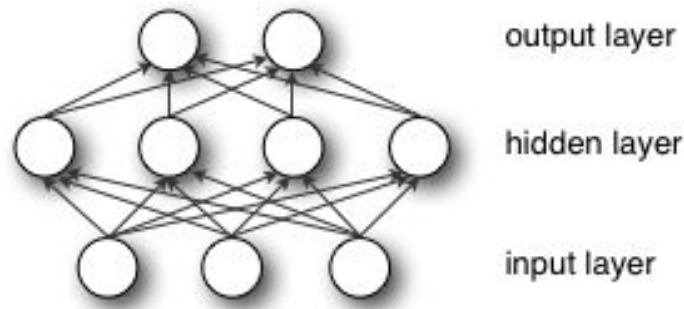
MLP output

$$f(x) = g_L(g_{L-1}(\dots g_2(g_1(x)) \dots))$$

$$g_\ell(x) = \sigma(x^T W^{(\ell)} + b^{(\ell)})$$

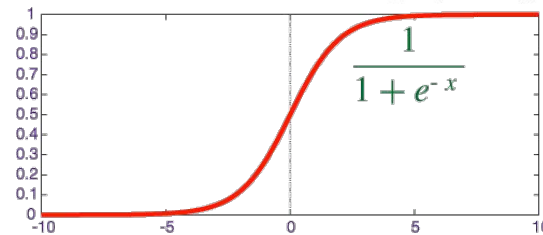


Multilayer Perceptrons (MLP)



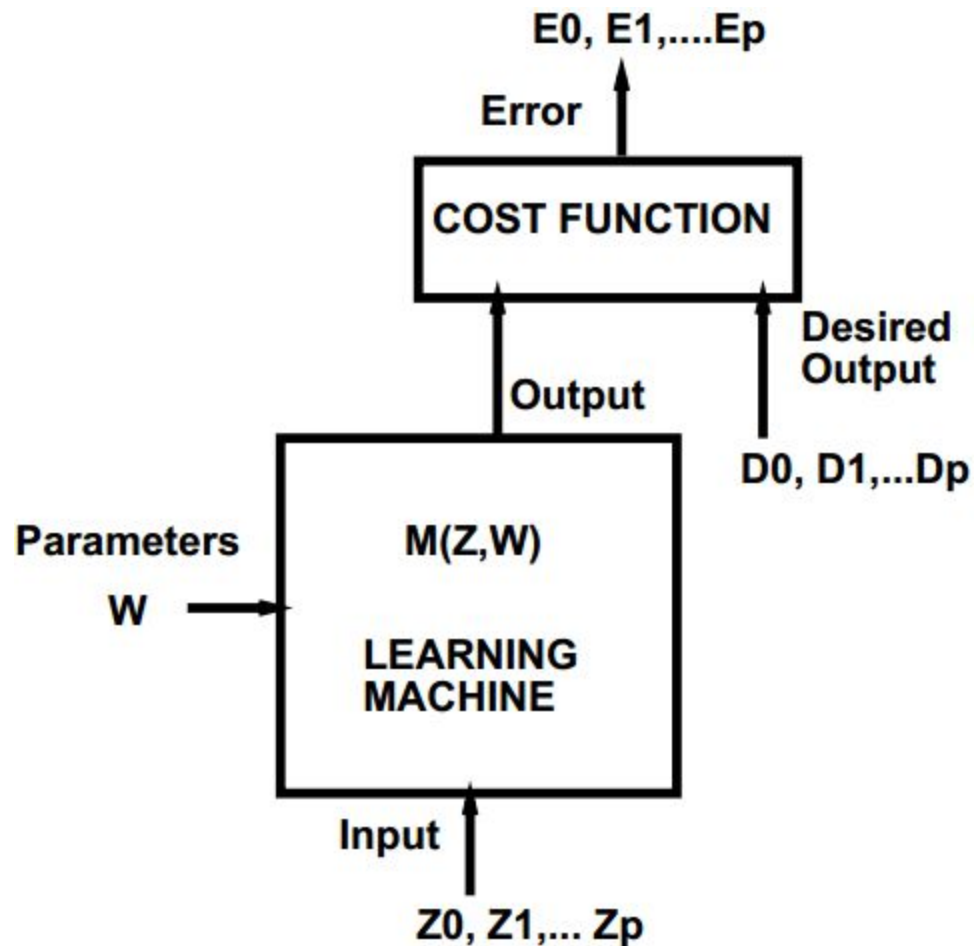
$$g_2(g_1) = \text{softmax}(g_1^T W^{(2)} + b^{(2)})$$

$$g_1(x) = \sigma(x^T W^{(1)} + b^{(1)}) \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$



$$f(x) = \text{softmax} \left(\left(\sigma(x^T W^{(1)} + b^{(1)}) \right)^T W^{(2)} + b^{(2)} \right)$$

Gradient Based Learning Machine



$$E^p = \frac{1}{2}(D^p - M(Z^p, W))^2$$

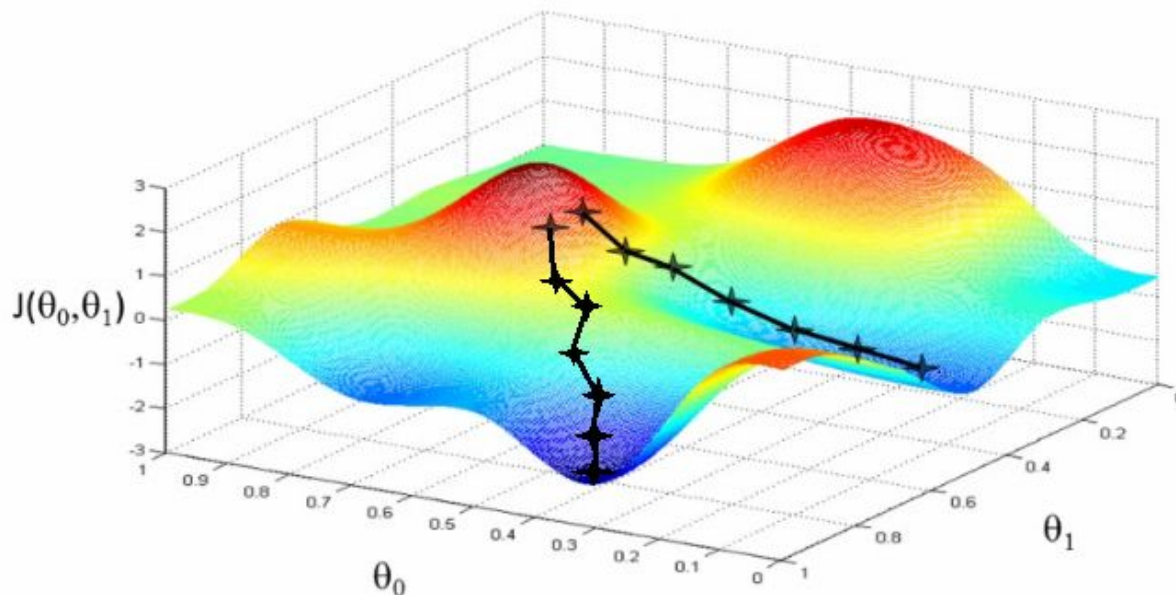
$$E_{train} = \frac{1}{P} \sum_{p=1} E^p$$

$$W_k = W_{k-1} - \epsilon \frac{\partial E(W)}{\partial W}$$

Gradient Descent

Algorithm 1 GRADIENT DESCENT

```
while True do
  loss =  $f(\text{params})$ 
  d_loss_wrt_params = ... ▷ compute gradient
  params -= learning_rate * d_loss_wrt_params
  if stopping condition is met then return params
  end if
end while
```



Stochastic Gradient Descent

Algorithm 1 GRADIENT DESCENT

```
while True do
    loss =  $f(\text{params})$ 
    d_loss_wrt_params = ... ▷ compute gradient
    params -= learning_rate * d_loss_wrt_params
    if stopping condition is met then return params
end if
end while
```

Algorithm 2 STOCHASTIC GRADIENT DESCENT

```
1: for  $(x_i, y_i) \in \mathcal{D}_{train}$  do
2:     ▷ imagine an infinite generator that may repeat
3:     ▷ examples (if there is only a finite training set)
4:     loss =  $f(\text{params}, x_i, y_i)$ 
5:     d_loss_wrt_params = ... ▷ compute gradient
6:     params - = learning_rate * d_loss_wrt_params
7:     if stopping condition is met then return params
8:     end if
9: end for
```

Mini-batch Gradient Descent

Algorithm 2 STOCHASTIC GRADIENT DESCENT

```
1: for  $(x_i, y_i) \in \mathcal{D}_{train}$  do
2:                                     ▷ imagine an infinite generator that may repeat
3:                                     ▷ examples (if there is only a finite training set)
4:   loss =  $f(\text{params}, x_i, y_i)$ 
5:   d_loss_wrt_params = ...                                     ▷ compute gradient
6:   params − = learning_rate * d_loss_wrt_params
7:   if stopping condition is met then return params
8:   end if
9: end for
```


Algorithm 3 MINIBATCH SGD

```
1: for (x_batch, y_batch) ∈ train_batches do
2:                                     ▷ imagine an infinite generator
3:                                     ▷ that may repeat examples
4:   loss =  $f(\text{params}, x\_batch, y\_batch)$ 
5:   d_loss_wrt_params = ...                                     ▷ compute gradient
6:   params − = learning_rate * d_loss_wrt_params
7:   if stopping condition is met then return params
8:   end if
9: end for
```

Hyperparameters: Learning Rate

$$W_i \leftarrow W_i - \eta \frac{\partial E(W_i)}{\partial W_i}$$

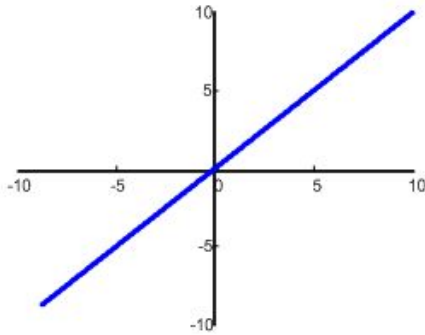
learning rate



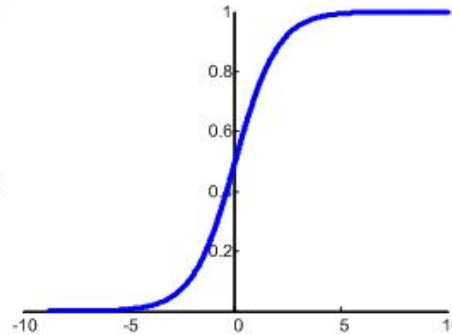
- constant learning rate (simplest solution)
- logarithmic grid search (10^{-1} , 10^{-2} , ...)
- decreasing learning rate over time:

$$\eta_t = \frac{\eta_0}{1 + at}$$

Hyperparameters: Nonlinearity

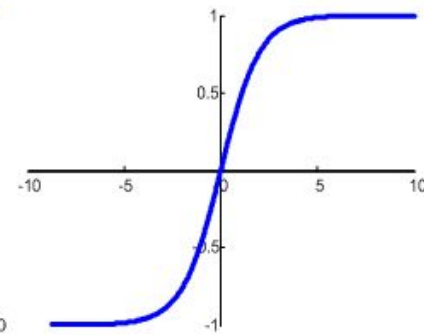


linear



logistic

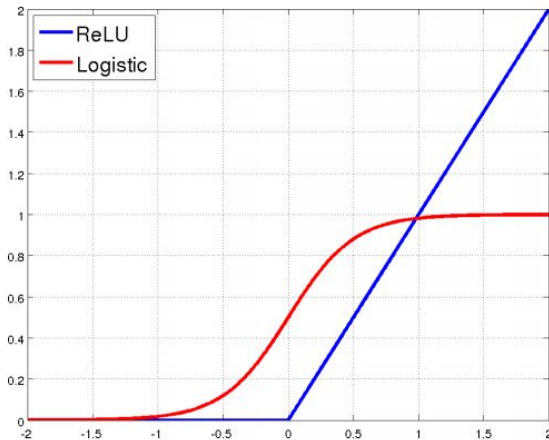
$$f(u) = 1/(1 + \exp(-u))$$



tanh

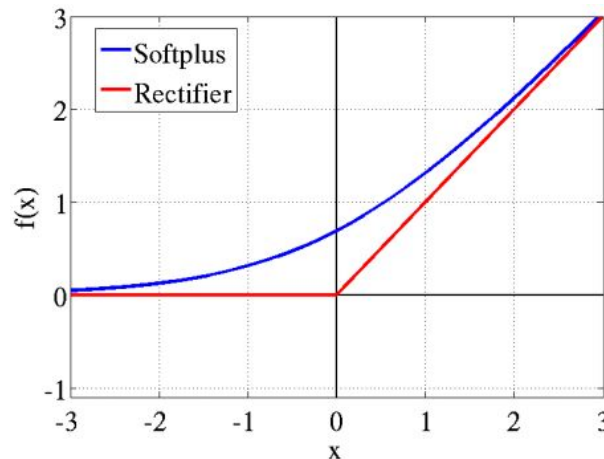
ReLU

$$f(u) = \max(0, u)$$

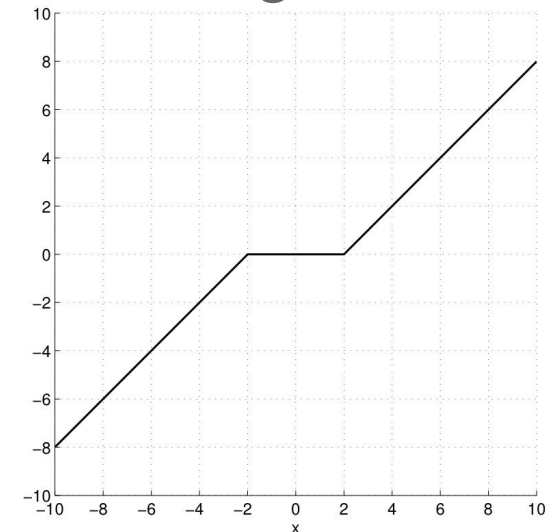


Softplus

$$\text{softplus}(x) = \log(1 + e^x)$$



Shrinkage



Hyperparameters: Weight Initialization

$$W \sim U\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right] \quad \text{for tanh activations}$$

$$W \sim U\left[-\frac{4\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{4\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right] \quad \text{for logistic activations}$$

Glorot, Xavier, and Yoshua Bengio. "Understanding the difficulty of training deep feedforward neural networks." *International Conference on Artificial Intelligence and Statistics* 2010: 249-256.

Hyperparameters: momentum

$$\Delta\theta_i(t) = v_i(t) = \alpha v_i(t-1) - \epsilon \frac{dE}{d\theta_i}(t)$$

momentum learning rate

Hinton, Geoffrey E. "A practical guide to training restricted boltzmann machines." *Neural Networks: Tricks of the Trade* (2012): 599-619.