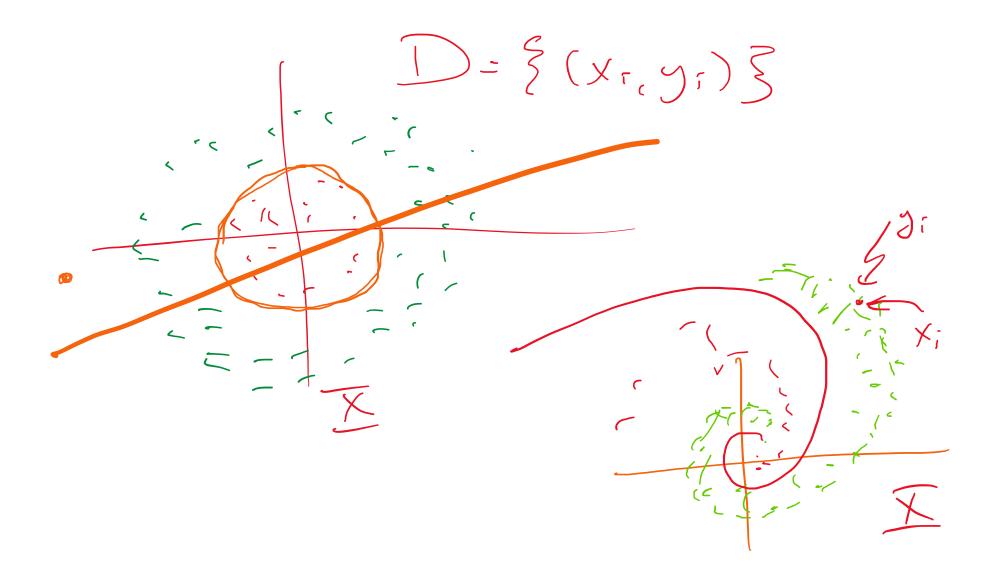


Labelled data Learn a function Pattern Structure

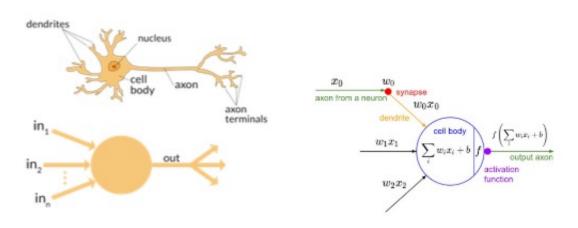
Well...at least one more - re-enforcement learning





Artificial Neural Networks

Neural Networks (NN)



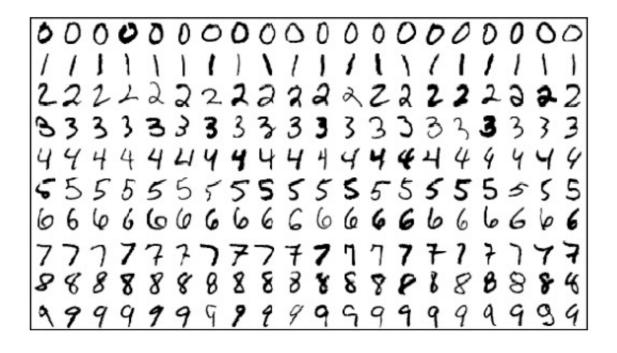
Suddenly ANN is best – why here why now? - I thought it was a crazy idea when I first heard of it.

Advancements of technology

- · Stochastic Gradient Descent.
- The computational GPU (Graphic Processing Unites), and the huge evolution of computations speed.

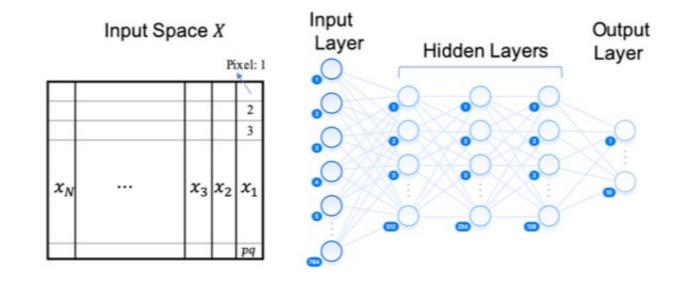
The problem – an example data set USPS Figs –

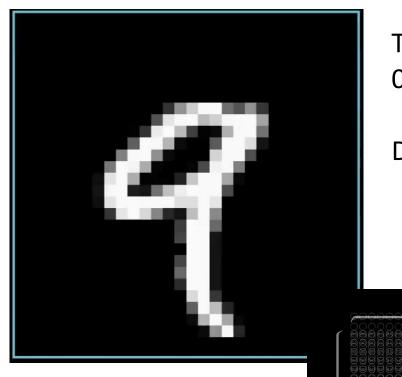
A supervised learning problem



The basis picture of data input into a ANN

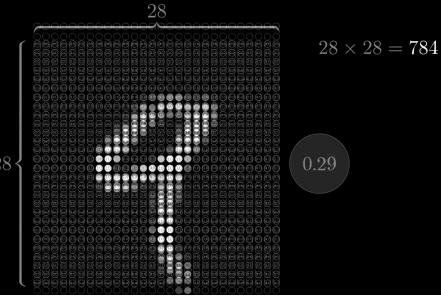
- -input layer
- -hidden inner layers
- -output layer





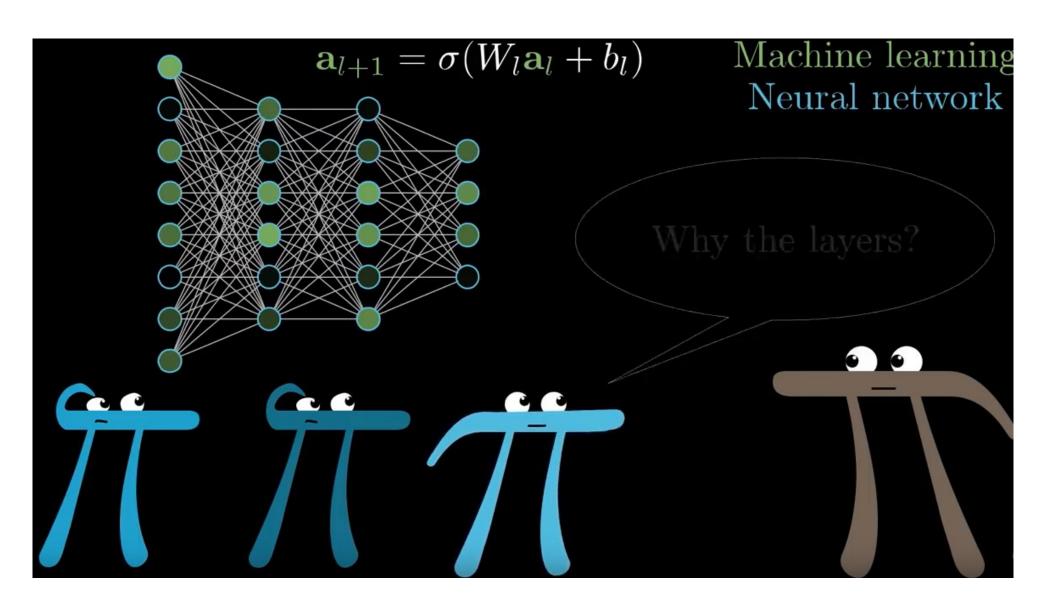
This datum is a 784 numbers between 0 and 1.

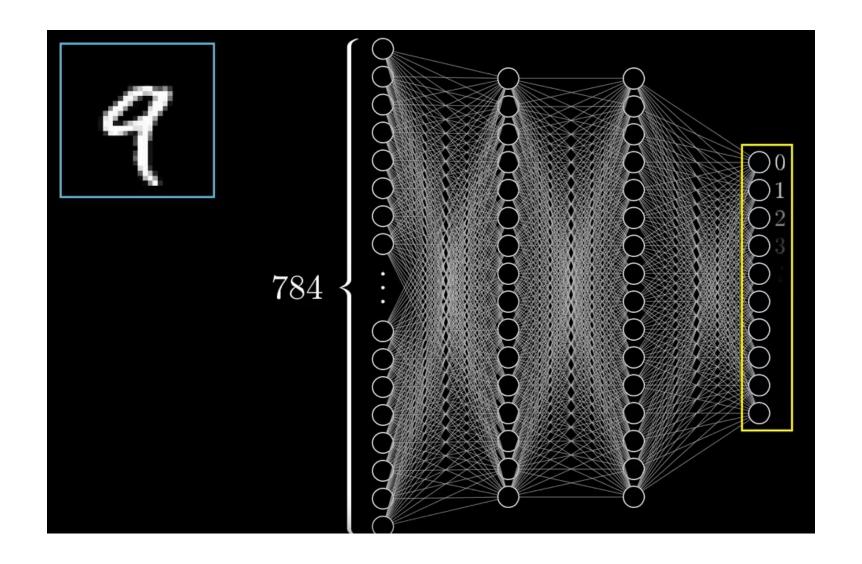
Data {(x_i,y_i)}_i=1..784

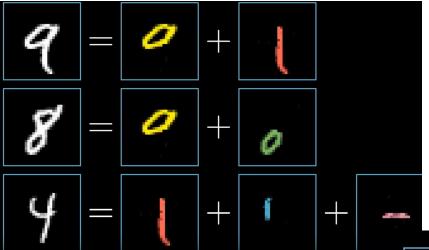


The data is a collection on ordered pairs, vector valued figures together with a symbol

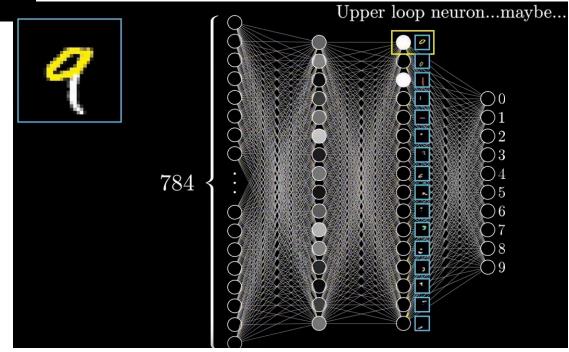
$$\begin{array}{c} (\textbf{0},0)(\textbf{6},6)(\textbf{3},3)(\textbf{6},6)(\textbf{7},7)(\textbf{8},8)(\textbf{0},0)(\textbf{9},9) \\ (\textbf{5},5)(\textbf{4},4)(\textbf{3},3)(\textbf{6},6)(\textbf{5},5)(\textbf{8},8)(\textbf{9},9)(\textbf{5},5) \\ (\textbf{4},4)(\textbf{4},4)(\textbf{7},7)(\textbf{2},2)(\textbf{0},0)(\textbf{3},3)(\textbf{2},2)(\textbf{8},8) \\ (\textbf{9},9)(\textbf{I},1)(\textbf{9},9)(\textbf{2},2)(\textbf{7},7)(\textbf{9},9)(\textbf{I},4) \\ (\textbf{8},8)(\textbf{7},7)(\textbf{4},4)(\textbf{I},1)(\textbf{3},3)(\textbf{I},1)(\textbf{5},5)(\textbf{3},3) \\ (\textbf{2},2)(\textbf{3},3)(\textbf{9},9)(\textbf{0},0)(\textbf{9},9)(\textbf{9},9)(\textbf{I},1)(\textbf{5},5) \\ (\textbf{8},8)(\textbf{4},4)(\textbf{I},7)(\textbf{2},7)(\textbf{4},4)(\textbf{4},4)(\textbf{4},4)(\textbf{4},2) \\ (\textbf{0},0)(\textbf{7},7)(\textbf{2},2)(\textbf{I},4)(\textbf{8},8)(\textbf{2},2)(\textbf{6},6)(\textbf{9},9) \\ (\textbf{1},9)(\textbf{2},2)(\textbf{5},8)(\textbf{2},7)(\textbf{6},6)(\textbf{I},1)(\textbf{I},1)(\textbf{2},2) \\ (\textbf{3},3)(\textbf{9},9)(\textbf{I},1)(\textbf{6},6)(\textbf{5},5)(\textbf{I},1)(\textbf{I},1)(\textbf{0},0) \\ \end{array}$$

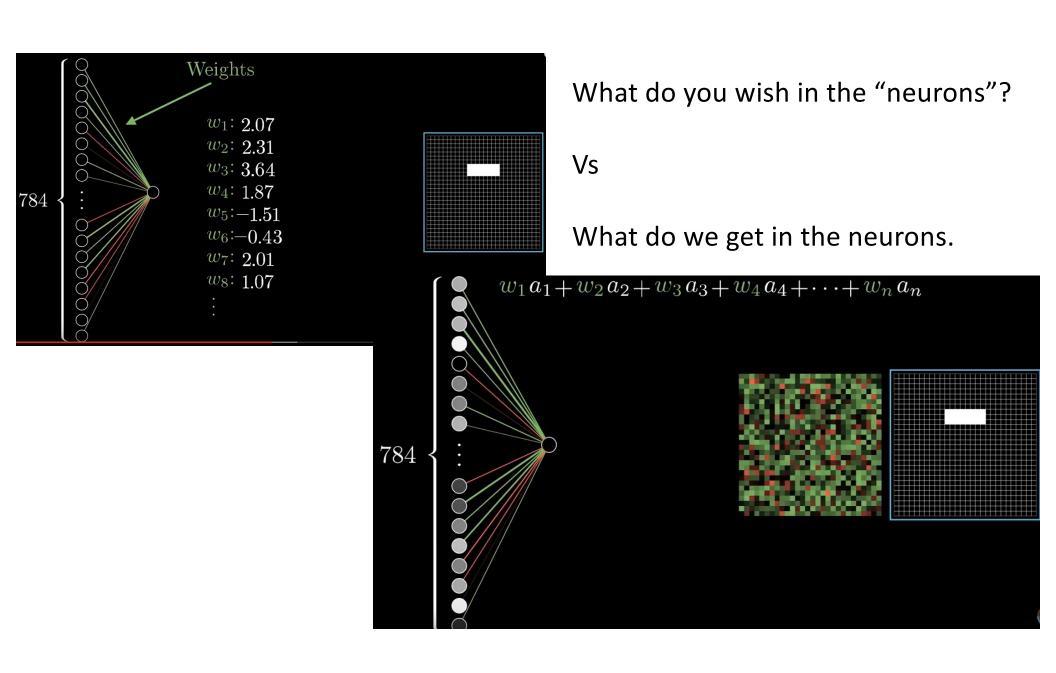


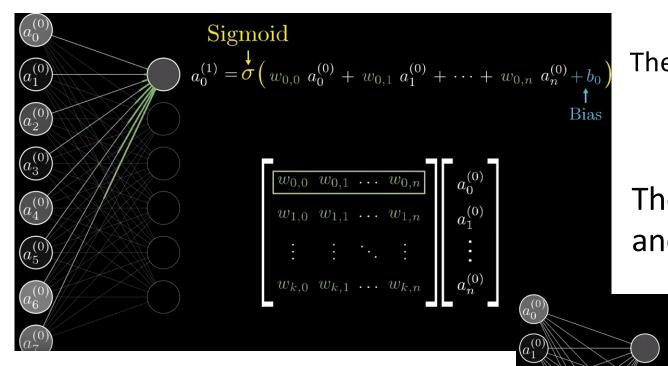




What do you wish in the "neurons"?

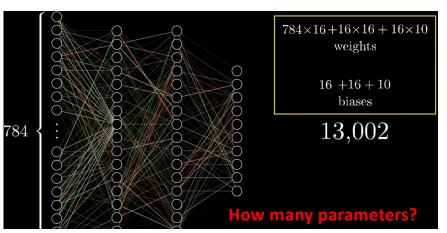






The ANN – Mathematical Process

The matrix form of a layer, and composition



$$\sigmaig(\mathbf{Wa}^{(0)} + \mathbf{b}ig)$$
 $oldsymbol{\sigma}igg(egin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \ w_{1,0} & w_{1,1} & \dots & w_{1,n} \ dots & dots & \ddots & dots \ w_{k,0} & w_{k,1} & \dots & w_{k,n} \ \end{pmatrix} egin{bmatrix} a_0^{(0)} & a_1^{(0)} \ a_n^{(0)} \ \vdots & dots \ a_n^{(0)} \ \end{pmatrix} + egin{bmatrix} b_0 \ b_1 \ dots \ b_n \ \end{pmatrix}$

The basic nodal statement.

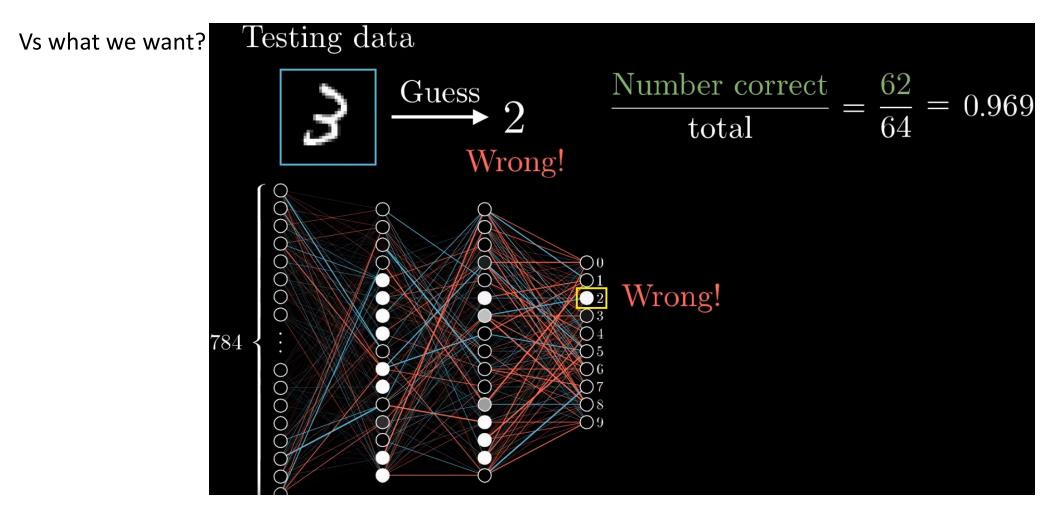
$$y = \sigma \left(\sum_{i=1}^{n} w_i x_i + b \right)$$

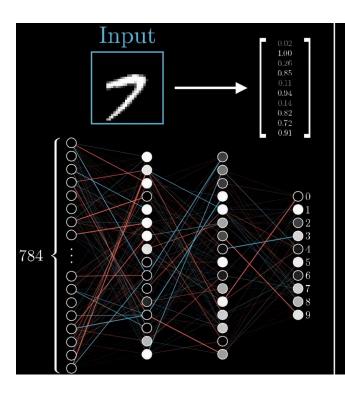
Favorite activation functions

- Linear function: $\sigma(x) = ax + b$, with a and b are constants.
- Sigmoid (Logistic) function: $\sigma(x) = \frac{1}{1+e^{-x}}$.
- Hyperbolic Tangent: $\sigma(x) = \tanh(x)$.
- Rectified Linear Unit (ReLU): σ(x) = max{0, x}.

What's the cost?

Vs what?





Neural network function

Input: 784 numbers (pixels)

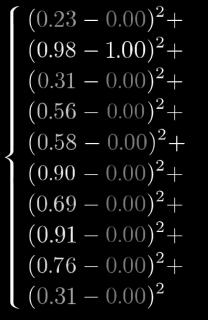
Output: 10 numbers

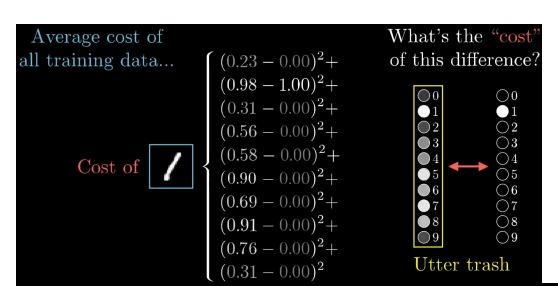
Parameters:

Let's write out the cost function

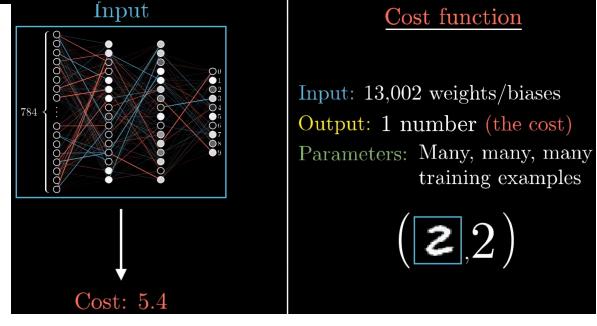
Average cost of all training data...

Cost of /

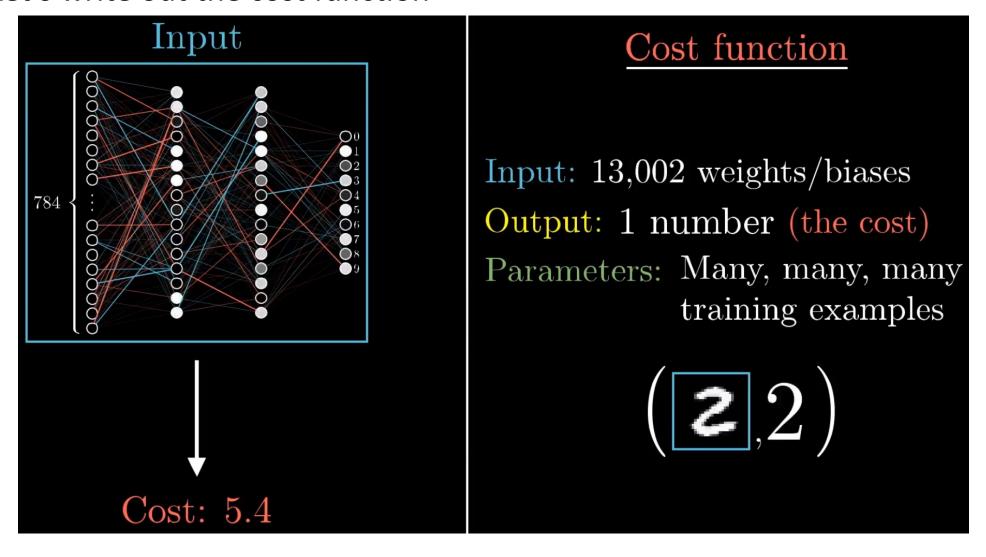


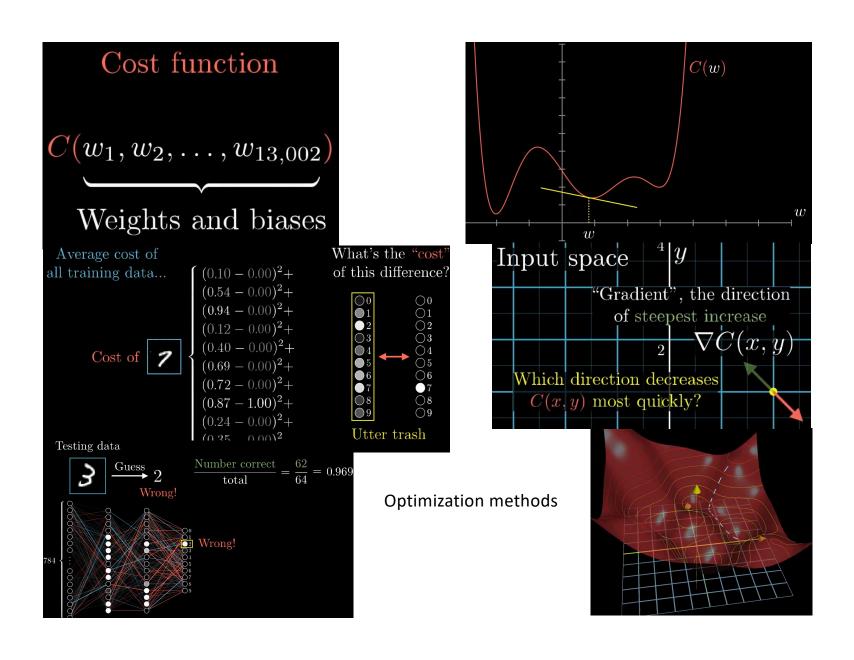


Let's write out the cost function



Let's write out the cost function



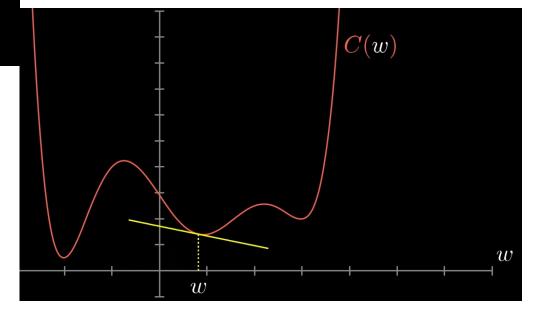


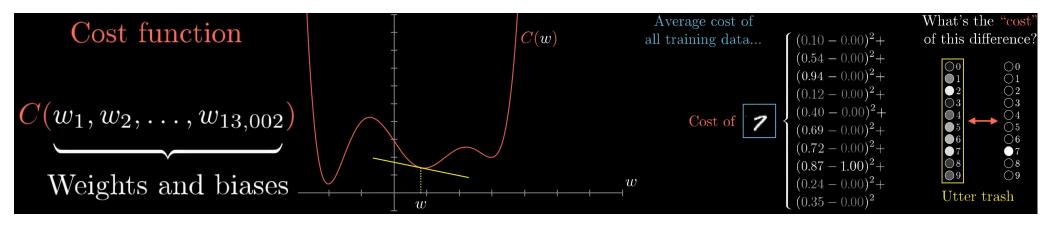
Cost function

An optimization problem

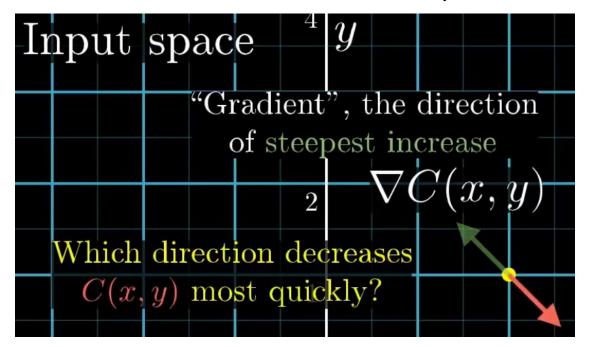
 $C(w_1, w_2, \ldots, w_{13,002})$

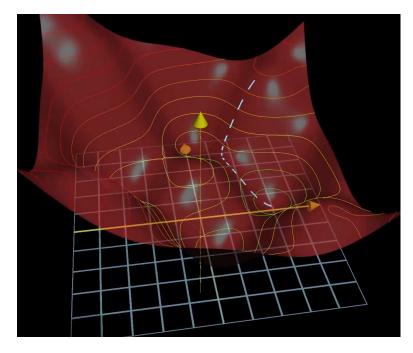
Weights and biases





Optimization methods

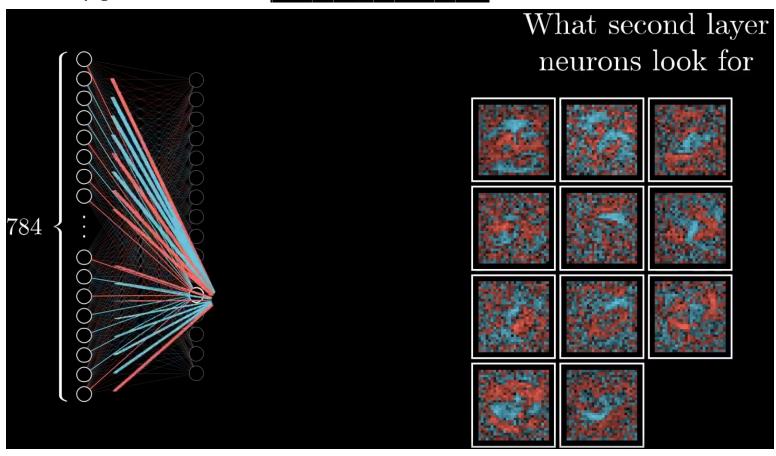




What do you wish in the "neurons"?

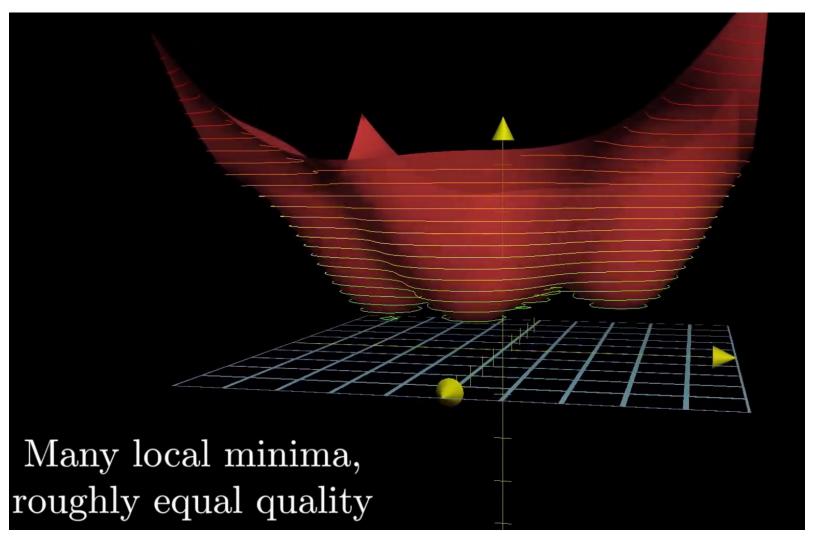
9 = 0 + 1 8 = 0 + 0 4 = 1 + 1 + -

What do we actually get?

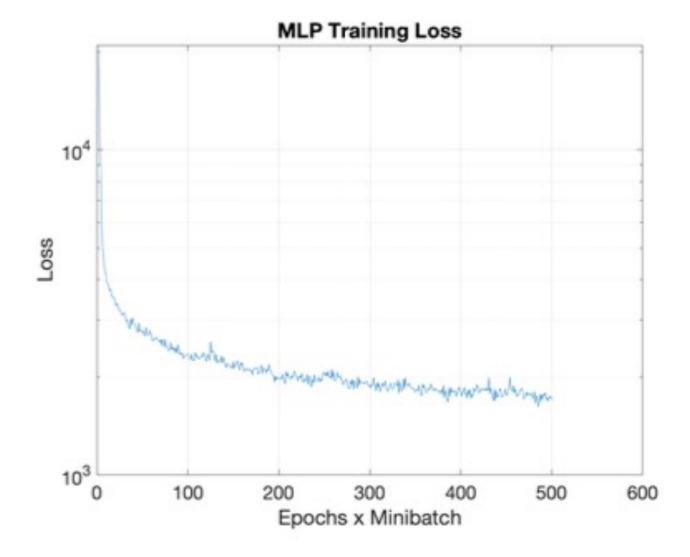


MANY local minimima of the VERY high dimensional so how important is it to find the best

One?

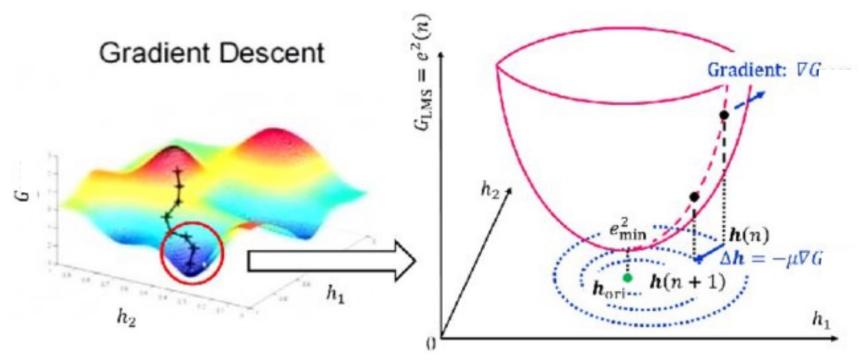




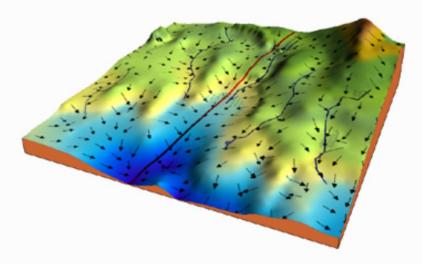


General smooth optimization, G(h):R^d->R

By general gradient descent



More on gradient descent



Given the cost function:

$$f(m, b) = \frac{1}{N} \sum_{i=1}^{n} (y_i - (mx_i + b))^2$$

The gradient can be calculated as:

$$f'(m,b) = \begin{bmatrix} \frac{df}{dm} \\ \frac{df}{db} \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum -2x_i(y_i - (mx_i + b)) \\ \frac{1}{N} \sum -2(y_i - (mx_i + b)) \end{bmatrix}$$

$$\mathbf{w}_{j+1}(\delta) = \mathbf{w}_j - \delta \nabla f_k(\mathbf{w}_j)$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \cdot \nabla_{\theta} L(\theta^{(t)}, \mathbf{y})$$

In this equation:

- $\theta^{(t)}$ is our current estimate of θ^* at the tth iteration
- α is the learning rate
- $\nabla_{\theta} L(\theta^{(t)}, \mathbf{y})$ is the gradient of the loss function
- We compute the next estimate $\theta^{(t+1)}$ by subtracting the product of α and $\nabla_{\theta} L(\theta, \mathbf{y})$ computed at $\theta^{(t)}$

Stochastic Gradient Descent

$$\mathbf{w}_{j+1}(\delta) = \mathbf{w}_j - \delta \nabla f_K(\mathbf{w}_j)$$

Batch (several data points)

$$K \in [k_1, k_2, \cdots k_p]$$

