

# Perceptron vs Gradient Descent

## GRADIENT DESCENT ALGORITHM:

Change  $w_i$  to  $w_i + \alpha(y - \hat{y})x_i$

## PERCEPTRON ALGORITHM:

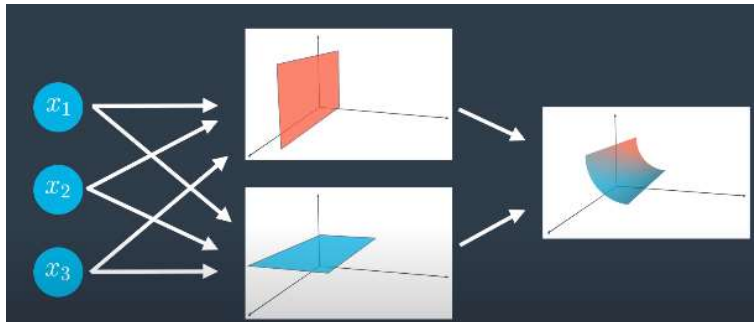
If  $x$  is misclassified:

Change  $w_i$  to  $\begin{cases} w_i + \alpha x_i & \text{if positive} \\ w_i - \alpha x_i & \text{if negative} \end{cases}$

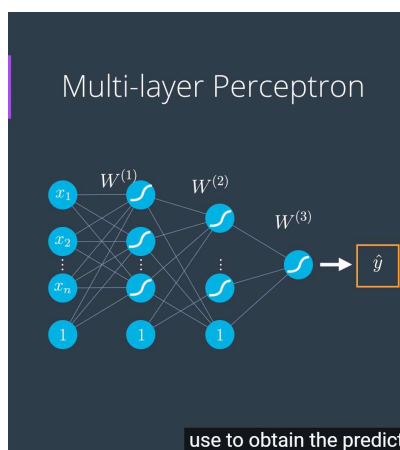
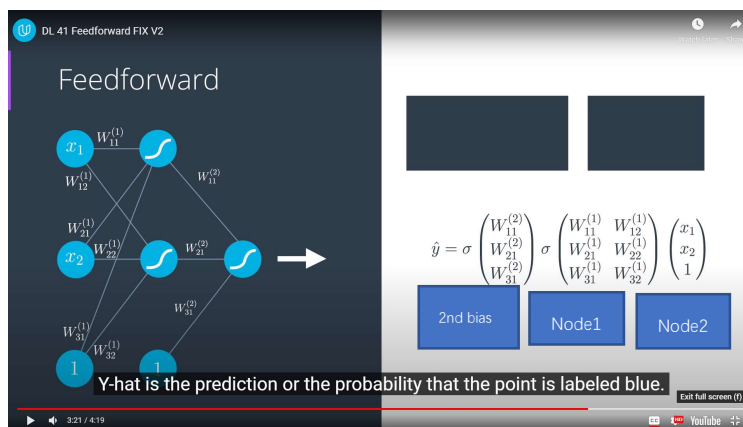
If correctly classified:  $y - \hat{y} = 0$

If misclassified:  $\begin{cases} y - \hat{y} = 1 & \text{if positive} \\ y - \hat{y} = -1 & \text{if negative} \end{cases}$

## Multi layer



## Feedforward



## PREDICTION

$$\hat{y} = \sigma \circ W^{(3)} \circ \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

DL 42 Neural Network Error Function (1)

## Multi-layer Perceptron

### PREDICTION

$$\hat{y} = \sigma \circ W^{(3)} \circ \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

### ERROR FUNCTION

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

And still, this function will tell us how badly a point gets misclassified.

## Perceptron

### PREDICTION

$$\hat{y} = \sigma(Wx + b)$$

### ERROR FUNCTION

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

### GRADIENT OF THE ERROR FUNCTION

$$\nabla E = \left( \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n}, \frac{\partial E}{\partial b} \right)$$

They correspond to these edges over here,

Calculating The Gradient 1

## Backpropagation

$$\hat{y} = \sigma W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

$$W^{(1)} = \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \quad W^{(2)} = \begin{pmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{pmatrix}$$

$$\nabla E = \begin{pmatrix} \frac{\partial E}{\partial W_{11}^{(1)}} & \frac{\partial E}{\partial W_{12}^{(1)}} & \frac{\partial E}{\partial W_{11}^{(2)}} \\ \frac{\partial E}{\partial W_{21}^{(1)}} & \frac{\partial E}{\partial W_{22}^{(1)}} & \frac{\partial E}{\partial W_{21}^{(2)}} \\ \frac{\partial E}{\partial W_{31}^{(1)}} & \frac{\partial E}{\partial W_{32}^{(1)}} & \frac{\partial E}{\partial W_{31}^{(2)}} \end{pmatrix}$$

the learning rate times the partial derivative of E with respect to that same weight.

Regra da cadeia

## Chain Rule

$$\frac{\partial B}{\partial x} = \frac{\partial B}{\partial A} \frac{\partial A}{\partial x}$$

when composing functions, that derivatives just multiply,

Watch on youtube.com

DL 46 Calculating The Gradient 2 V2 (2)

## Feedforward

$$h_1 = W_{11}^{(1)} x_1 + W_{21}^{(1)} x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{32}^{(1)}$$

$$h = W_{11}^{(2)} \sigma(h_1) + W_{21}^{(2)} \sigma(h_2) + W_{31}^{(2)}$$

$$\hat{y} = \sigma(h)$$

$$\hat{y} = \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

the sigmoid of W superscript 1 applied to the input x and that is feedforward.

Watch on youtube.com

DL 46 Calculating The Gradient 2 V2 (2)

## Backpropagation

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

$$E(W) = E(W_{11}^{(1)}, W_{12}^{(1)}, \dots, W_{31}^{(2)})$$

$$\nabla E = \left( \frac{\partial E}{\partial W_{11}^{(1)}}, \dots, \frac{\partial E}{\partial W_{31}^{(2)}} \right)$$

$$\frac{\partial E}{\partial W_{11}^{(1)}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial h_1} \frac{\partial h_1}{\partial W_{11}^{(1)}}$$

So, let us calculate the other ones.

Watch on youtube.com

DL 46 Calculating The Gradient 2 V2 (2)

## Backpropagation

$$h = W_{11}^{(2)} \sigma(h_1) + W_{21}^{(2)} \sigma(h_2) + W_{31}^{(2)}$$

$$\frac{\partial h}{\partial h_1} = W_{11}^{(2)} \sigma(h_1) [1 - \sigma(h_1)]$$

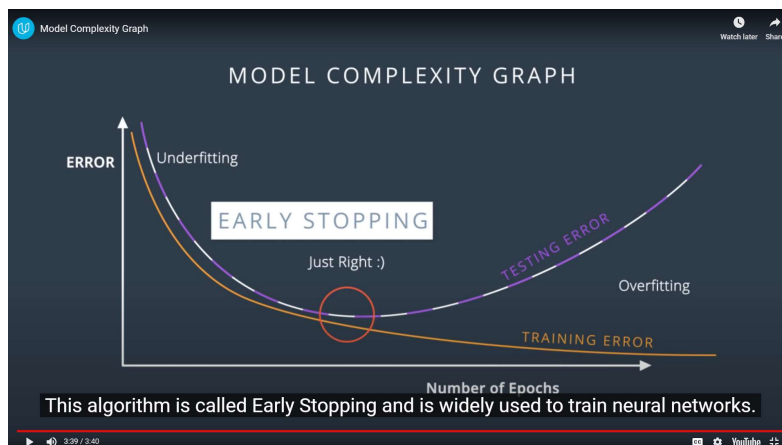
precisely sigmoid of h times 1 minus sigmoid of h. Again,

### Calculation of the derivative of the sigmoid function

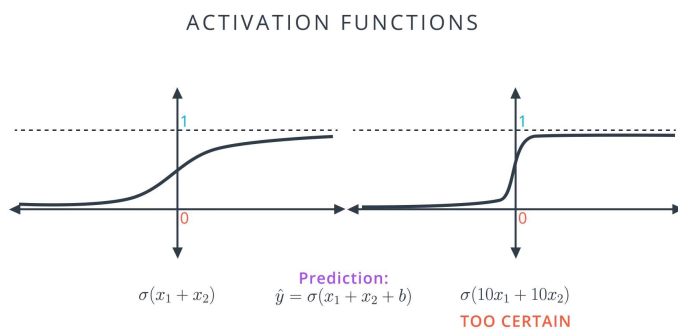
Recall that the sigmoid function has a beautiful derivative, which we can see in the following calculation. This will make our backpropagation step much cleaner.

$$\begin{aligned} \sigma'(x) &= \frac{\partial}{\partial x} \frac{1}{1+e^{-x}} \\ &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} \\ &= \sigma(x)(1 - \sigma(x)) \end{aligned}$$

### Early Stopping



### Regularization:



will generate large errors and it will be hard to tune the model to correct them.

Regularization

## Solution: Regularization

LARGE COEFFICIENTS → OVERFITTING

PENALIZE LARGE WEIGHTS  
 $(w_1, \dots, w_n)$

ERROR FUNCTION =  $-\frac{1}{m} \sum_{i=1}^m (1 - y_i) \ln(1 - \hat{y}_i) + y_i \ln(\hat{y}_i) + \lambda(|w_1| + \dots + |w_n|)$  **L1**

ERROR FUNCTION =  $-\frac{1}{m} \sum_{i=1}^m (1 - y_i) \ln(1 - \hat{y}_i) + y_i \ln(\hat{y}_i) + \lambda(w_1^2 + \dots + w_n^2)$  **L2**

As you can see, these two are large if the weights are large.

Regularization

## L1 vs L2 Regularization

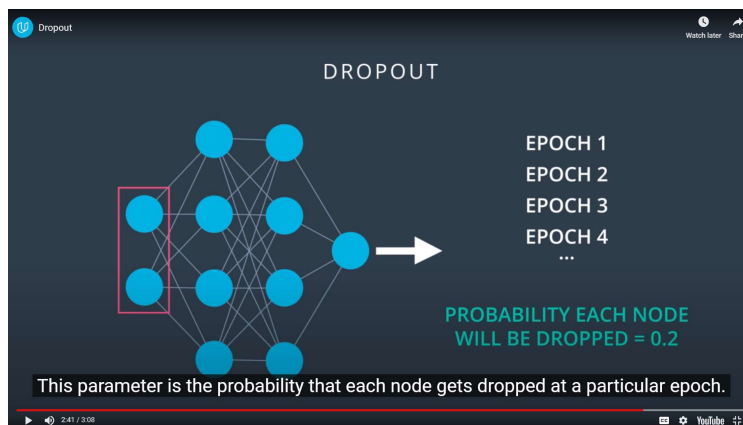
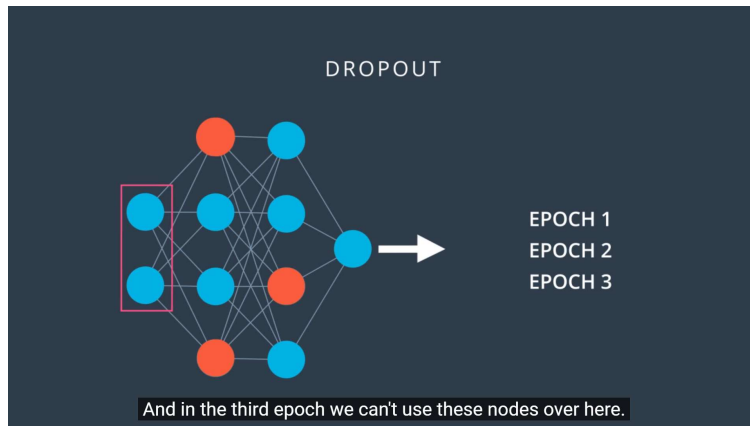
L1	L2
SPARSITY: (1, 0, 0, 1, 0)	SPARSITY: (0.5, 0.3, -0.2, 0.4, 0.1)
GOOD FOR FEATURE SELECTION	NORMALLY BETTER FOR TRAINING MODELS
	$(1, 0) \rightarrow (0.5, 0.5)$ $1^2 + 0^2 = 1 \quad 0.5^2 + 0.5^2 = 0.5$
	since this one produces a smaller sum of squares.

L1 tends to end up with sparse vectors, That means, a small weights will tend to go to zero. So if we want to reduce the number of weights and end up with a small set, we can use L1. This is also good for feature selections and sometimes we have a problem with hundreds of features, and L1 regularization will help us select which ones are important and it will turn the rest into zeros.

L2 tries to maintain all the weights homogeneously small, This one normally gives better results for training models. So why would L1 regularization produce vectors with sparse weights and L2 regularization will produce vectors with small homogeneous weights?

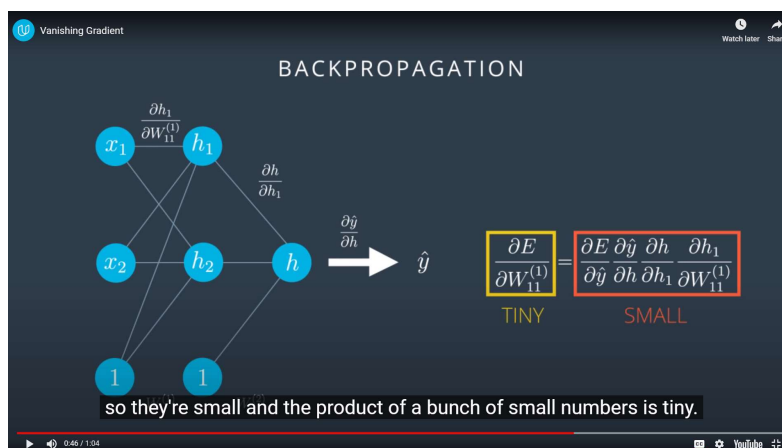
L2 regularization will prefer the vector point(0.5,0.5) over the vector(1,0) since this one produces a smaller sum of squares.

## Dropout

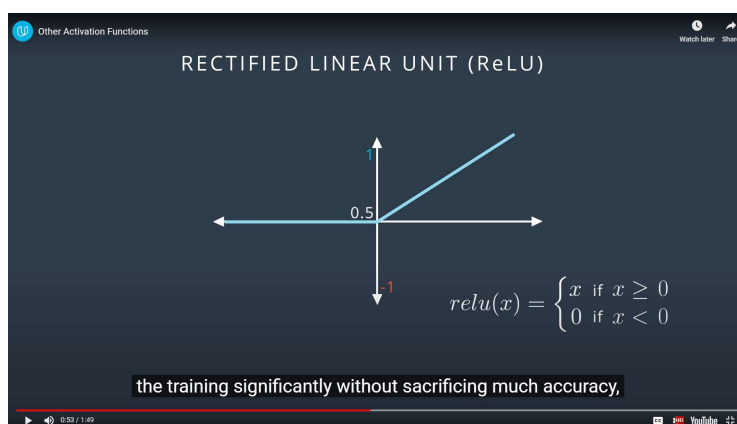
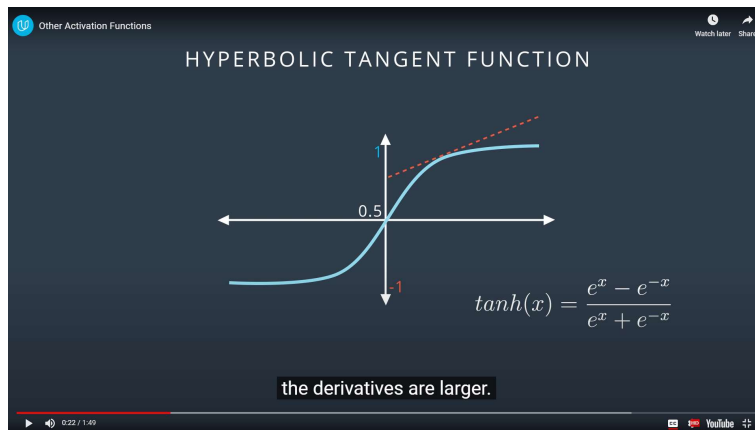


## Vanishing Gradient

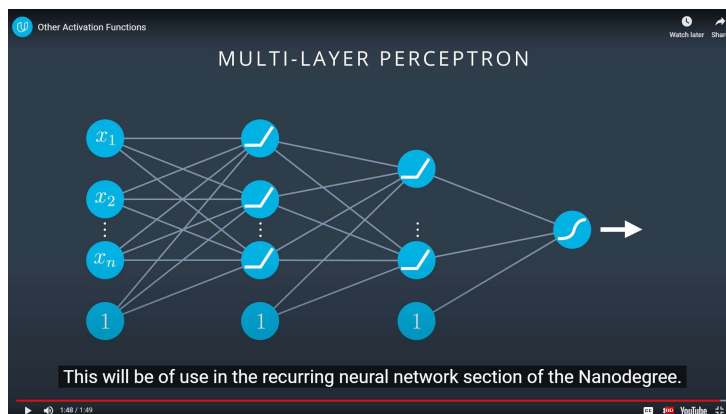
Sigmoid has very small derivatives when it closes to 1 or -1.  
This makes weight update very difficult.



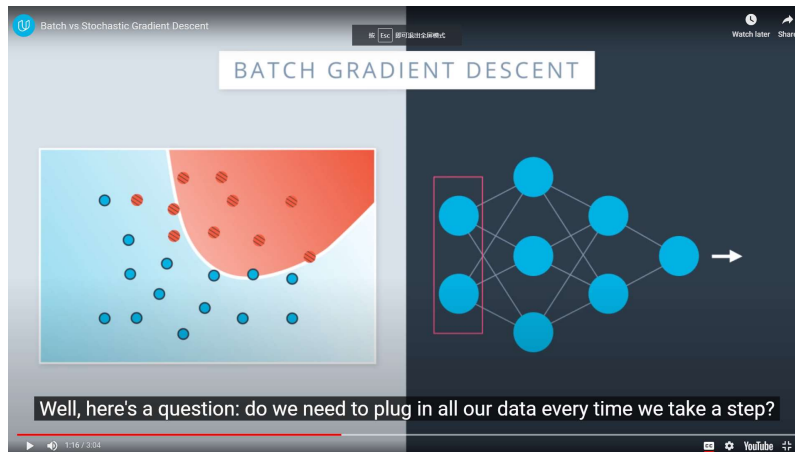
## Other activation function



The relu can improve the training significantly without sacrificing much accuracy.

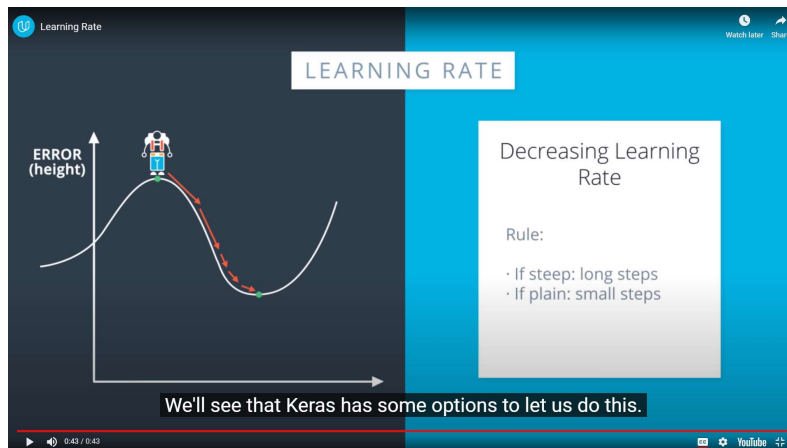


Stochastic Gradient

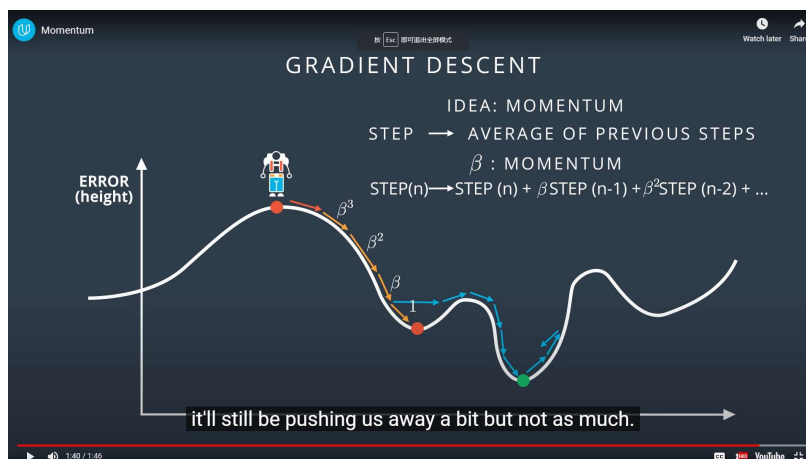


The idea behind stochastic gradient descent is simply that we take small subsets of data, run them through the nn, calculate the gradient of error function based on those points and move one step in that direction. Split the data into several batches.

## Learning rate



To solve the local minima to have a random restart or momentum



The previous step matters a lot and the steps before that matter less and less



## Neural Network For regression

This is classification:

