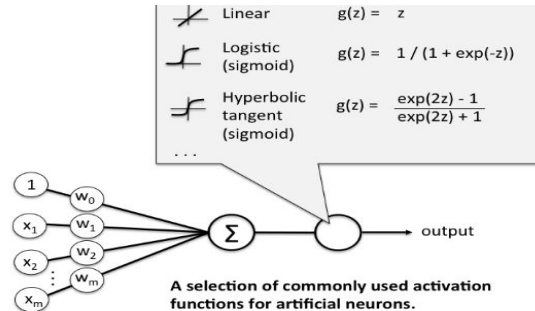


# Activation Functions:

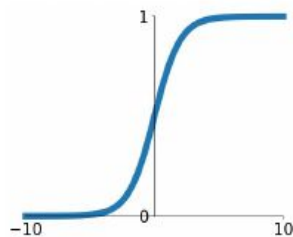
- They introduce non-linear properties to our Network. **Their main purpose is to convert a input signal of a node in a A-NN to an output signal.**
- That output signal now is used as a input in the next layer in the stack.
- If we do not apply a Activation function then the output signal would simply be a simple **linear function.**
- Now, a linear equation is easy to solve but they are limited in their complexity and have less power to learn complex functional mappings from data.



# Activation Functions

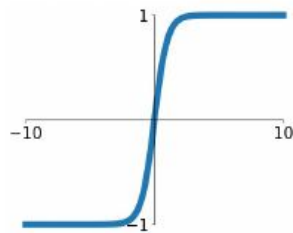
## Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



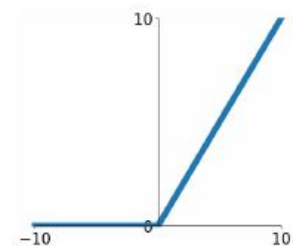
## tanh

$$\tanh(x)$$



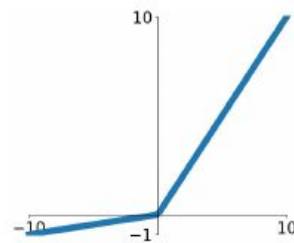
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

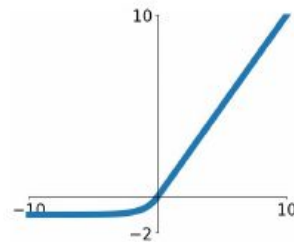


## Maxout

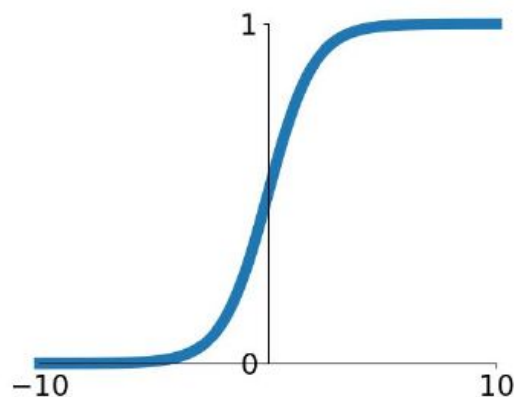
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Activation Functions



**Sigmoid**

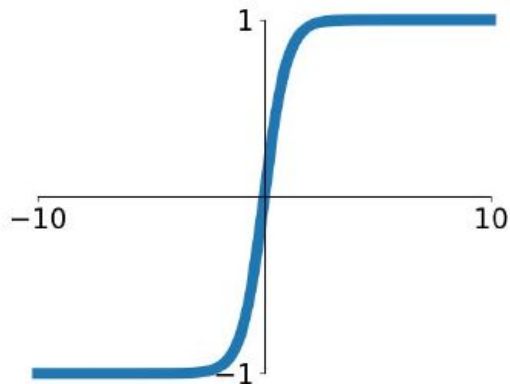
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3.  $\exp()$  is a bit compute expensive

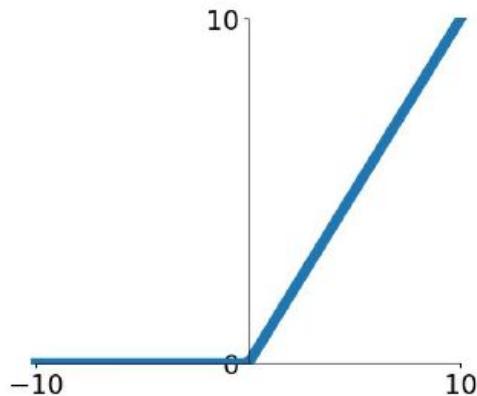
# Activation Functions



**$\tanh(x)$**

- Squashes numbers to range  $[-1,1]$
- zero centered (nice)
- still kills gradients when saturated :(

# Activation Functions



**ReLU**  
(Rectified Linear Unit)

- Computes  $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Not zero-centered output
- An annoyance:

hint: what is the gradient when  $x < 0$ ?

LOGITS  
SCORES



SOFTMAX

PROBABILITIES

$y$   $\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$

→

$$S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

→  $p = 0.7$

→  $p = 0.2$

→  $p = 0.1$

$X$  ← INPUT

LINEAR  
MODEL

$$WX + b$$

LOGIT

$Y$

2.0

1.0

0.1

SOFTMAX

$S(Y)$

0.7

0.2

0.1

CROSS-  
ENTROPY

$$D(S, L)$$

1-HOT  
LABELS  
 $L$

1.0

0.0

0.0

# Cross Entropy Loss :

In binary classification, where the number of classes  $M$  equals 2, cross-entropy can be calculated as:

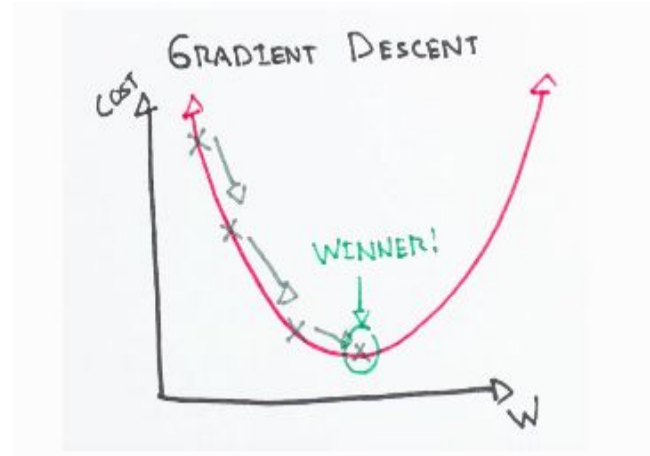
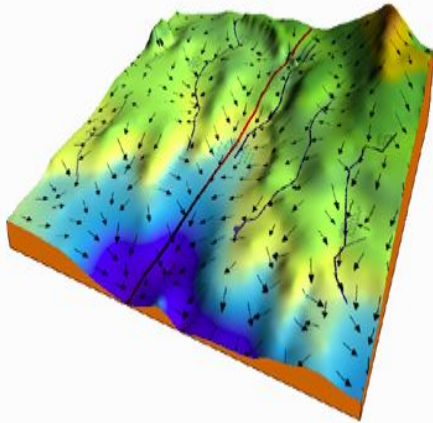
$$-(y \log(p) + (1 - y) \log(1 - p))$$

If  $M > 2$  (i.e. multiclass classification), we calculate a separate loss for each class label per observation and sum the result.

$$-\sum_{c=1}^M y_{o,c} \log(p_{o,c})$$



# Gradient Descent



# Stochastic Gradient Descent

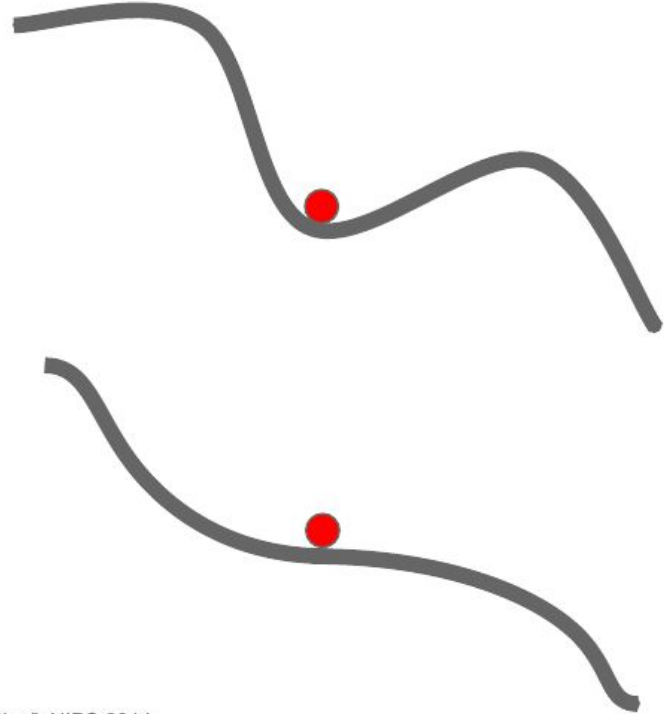
- In typical Gradient Descent optimization, like Batch Gradient Descent, the batch is taken to be the whole dataset.
- if you use a typical Gradient Descent optimization technique, you will have to use all of the one million samples for completing one iteration while performing the Gradient Descent, and it has to be done for every iteration until the minima is reached. Hence, it becomes computationally very expensive to perform.
- This problem is solved by Stochastic Gradient Descent. In SGD, it uses only a single sample, i.e., a batch size of one, to perform each iteration.

$$\theta = \theta - \alpha \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$$

# Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point**?

Saddle points much more common in high dimension



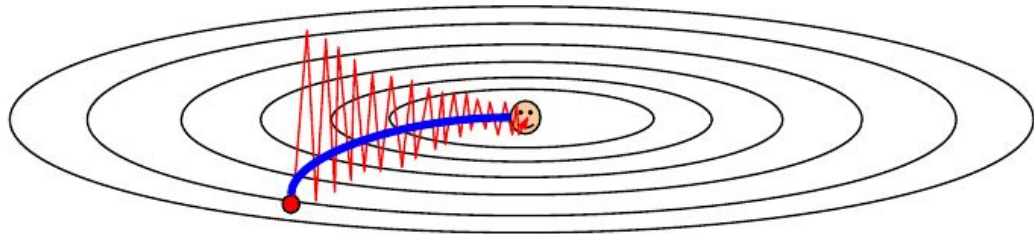
# SGD + Momentum

Local Minima

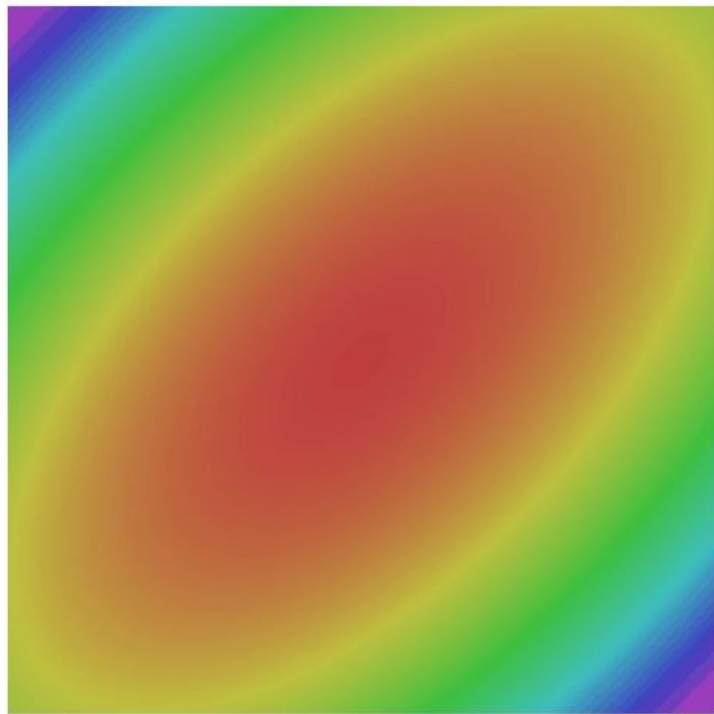
Saddle points



Poor Conditioning



Gradient Noise



— SGD

— SGD+Momentum

# SGD + Momentum

## SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:
    dx = compute_gradient(x)
    x -= learning_rate * dx
```

## SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x -= learning_rate * vx
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99