

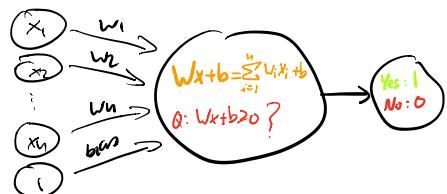
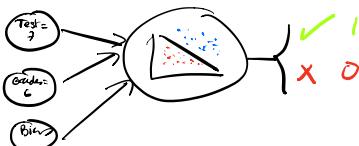
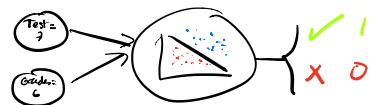
Equation of the hyperplane : $w_1x_1 + w_2x_2 + \dots + w_nx_n + b = 0$

$$[W] [x] + [b] = 0$$

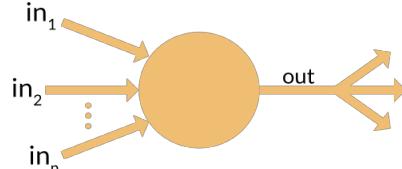
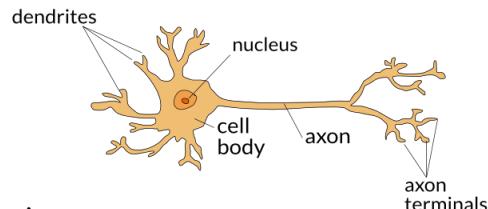
dim: $(1 \times n)$ $(n \times 1)$ (1×1) ↗ single output

Perceptrons:

Perception: Building block for our equation



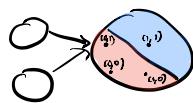
A perception is similar to a neuron in the brain.



Perceptions as logical operators

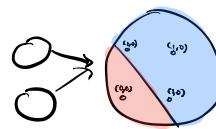
And operator

in	in	out
✓	✓	✓
✓	✗	✗
✗	✗	✗
✓	✗	✗



Or Operator

in	in	out
✓	✓	✓
✓	✗	✓
✗	✓	✓
✗	✗	✗



Perception Trick

Starting at random and gradually improving

Perception Algorithm

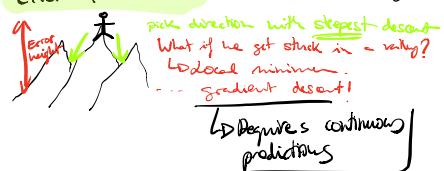
- Start with random weights: w_1, \dots, w_n, b
 - For every misclassified point (x_1, \dots, x_n) :

2.1 If prediction = 0:
 For $i = 1 \dots n$ learning Rule
 change $W_i + \alpha X_i$

2.2 If prediction = 1:
 For $i = 1 \dots n$ learning Rate
 Change $w_i = w_i + \alpha x_i$
 Change b to $b + \alpha$

3 Repeat until small error or # iterations is reached

Error Function Tells us how far away we are from the solution



How do we tell the computer how far it is from a perfect solution?
→ count mistakes for ~~dist~~

perfect solution?
→ count mistakes for distance to live due to small steps

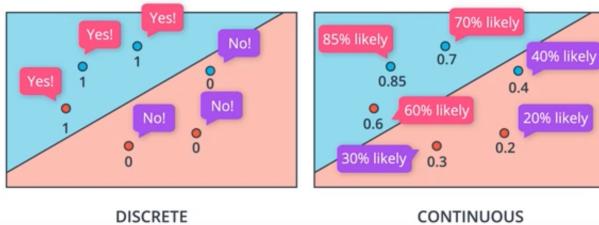
Discrete vs Continuous

In order to do gradient descent, we need to have a continuous output from the perception.

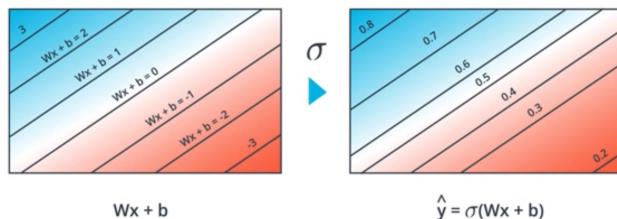
In order to do gradient descent, we need to use a different activation function (like \tilde{J}) in order to transform the output to continuous.

We then need to use a sigmoid function to map the raw scores to probabilities since all outputs range from [0, 1] as shown in Figure 1.

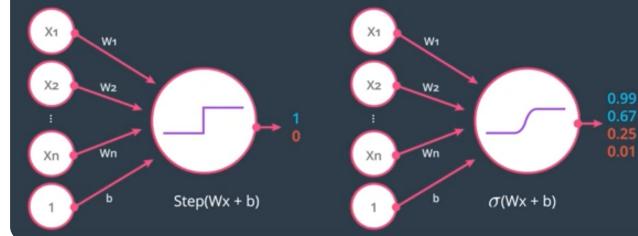
Predictions



Predictions



Perceptron



Softmax

Classification Problem

$P(\text{gift}) = 0.8$
 $P(\text{no gift}) = 0.2$
 $\text{Score}(\text{gift}) = \text{linear function}$
 $P(\text{gift}) = \text{f}(\text{Score})$

But what if we had more options? ex Is it a cat?
 N 11 Dog?
 C 11 Dog?

$$\begin{aligned} P(\text{duck}) &= 0.67 \\ P(\text{beaver}) &= 0.24 \\ P(\text{chicken}) &= 0.09 \end{aligned}$$

They all should add to 1.

$$\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \rightarrow \text{Normalize} \quad \frac{2}{2+1+0} = \frac{2}{3} \quad \frac{1}{2+1+0} = \frac{1}{3} \quad \frac{0}{2+1+0} = 0$$

Problem: What happens if scores are negative? Doesn't work!

We first need to turn all scores into positive! $\rightarrow e^x$

Softmax definition

$$P(\text{class } i) = \frac{e^{z_i}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

$$\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \rightarrow \text{Softmax} \quad \begin{aligned} \frac{e^2}{e^2 + e^1 + e^0} &= 0.67 \\ \frac{e^1}{e^2 + e^1 + e^0} &= 0.24 \\ \frac{e^0}{e^2 + e^1 + e^0} &= 0.09 \end{aligned} \quad \text{adds to 1.}$$

One-hot encoding for categorical data

Animal	Duck?	Beaver?	Walrus?
Duck	1	0	0
Beaver	0	1	0
Walrus	0	0	1

Maximum likelihood and maximizing probability
 Which model is more accurate? The one that gives the highest probabilities to the events

Maximizing the probability \equiv minimizing error We calculate the error

If we had many numbers product scales really bad. Sums are better in changing over time

We can then use logarithms since $\log(a \cdot b) = \log(a) + \log(b)$

Cross Entropy

$0.6 \cdot 0.2 \cdot 0.1 \cdot 0.7 \rightarrow \ln(0.6) + \ln(0.2) + \ln(0.1) + \ln(0.7)$ they are all negative numbers. We then need to take the negatives

A good model will give us a low cross entropy (high probabilities will give smaller numbers)

The misclassified points will be the biggest indicators of high entropy

\approx errors at each point.

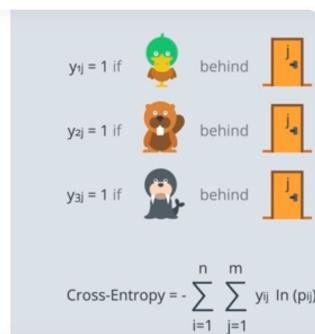
Maximizing probability \equiv minimizing cross entropy

Muticlass Cross-Entropy

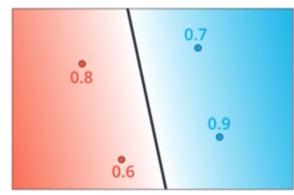
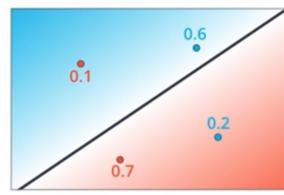
Multi-Class Cross-Entropy

Essentially a generalization of the cross entropy of several classes beyond Yes (p_j) no ($1-p_j$)

ANIMAL	DOOR 1	DOOR 2	DOOR 3
duck	p_{11}	p_{12}	p_{13}
beaver	p_{21}	p_{22}	p_{23}
walrus	p_{31}	p_{32}	p_{33}



Maximum Likelihood



$$0.7 \cdot 0.9 \cdot 0.8 \cdot 0.6 = 0.3024$$

Cross-Entropy

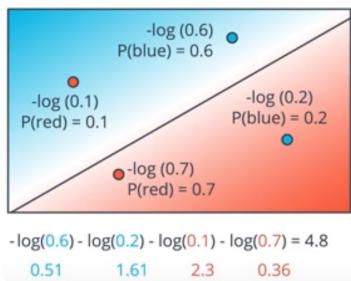
$p_1 = 0.8$	$p_2 = 0.7$	$p_3 = 0.1$	Cross-Entropy
			$-\ln(0.8) - \ln(0.7) - \ln(0.9)$
p_1	p_2	$1 - p_3$	
$y_1 = 1$	$y_2 = 1$	$y_3 = 0$	

$$\text{Cross-Entropy} = - \sum_{i=1}^m y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)$$

$$\text{CE}([(1,1,0), (0.8,0.7,0.1)]) = 0.69$$

Error Function

Error Function



If $y = 1$
 $P(\text{blue}) = \hat{y}$
 $\text{Error} = -\ln(\hat{y})$

If $y = 0$
 $P(\text{red}) = 1 - P(\text{blue}) = 1 - \hat{y}$
 $\text{Error} = -\ln(1 - \hat{y})$

Error = $-(1-y)\ln(1-\hat{y}) - y\ln(\hat{y})$

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

Error Function

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

$$E(W, b) = -\frac{1}{m} \sum_{i=1}^m (1-y_i)\ln(1-\sigma(Wx^{(i)}+b)) + y_i\ln(\sigma(Wx^{(i)}+b))$$

$\underbrace{\hat{y}_i}_{\text{ans}} \text{ aus } \underbrace{\sigma(x)}_{f(x) = \sigma(Wx^{(i)}+b)}$

GOAL
Minimize Error Function

Error Function γ/λ

Multi class error



ERROR FUNCTION:

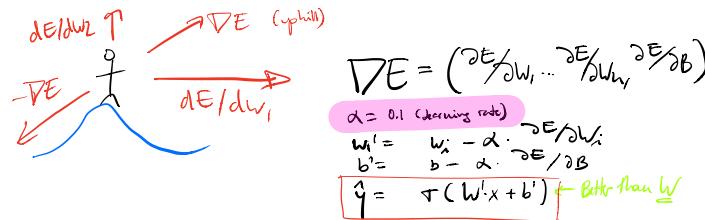
$$-\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1-y_i) \ln(1-\hat{y}_i)$$



ERROR FUNCTION:

$$-\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n y_{ij} \ln(\hat{y}_{ij})$$

Gradient Descent



Logistic Regression Algorithm

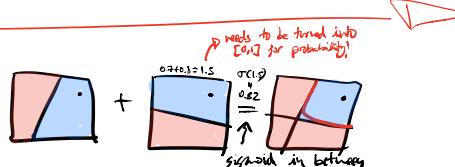
- Start with random weights:
 w_1, \dots, w_n, b
- For every point (x_1, \dots, x_n) :
 - For $i=1 \dots n$
 - Update $w_i^1 \leftarrow w_i - \alpha (y - \hat{y})x_i$
 - Update $b^1 \leftarrow b - \alpha (y - \hat{y})$
- Repeat until error is small
or a fixed # of times (epoch)
Very similar to the perceptron!

Non Linear Models

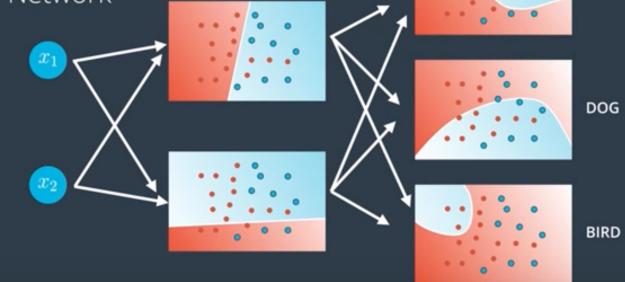
This can be achieved via Neural Networks:

Neural Network architecture

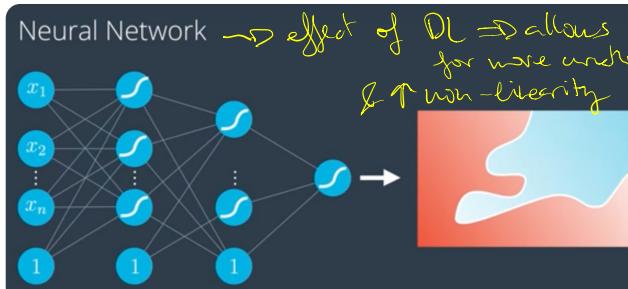
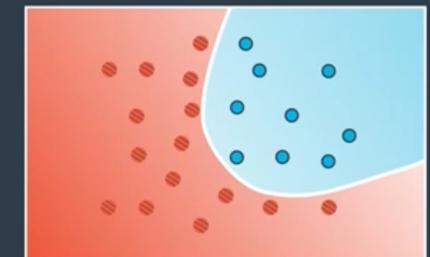
non-linear models can be achieved via the combination of two such linear models



Deep Neural Network

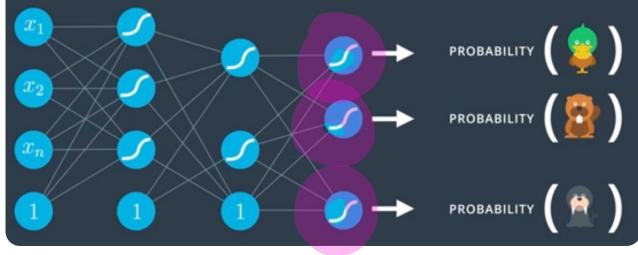


Non-Linear Regions



Multi Class Classification

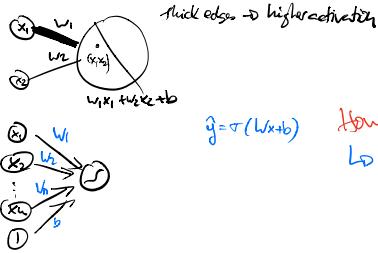
Neural Network



No need to create \leftarrow separate W for each a node, we just need to add more nodes in the output layer

Feed Forward

Process in which neural networks turn an input into an output



How to do it for multilayer perceptron?

Is same error function, with more complicated boundary

Backpropagation

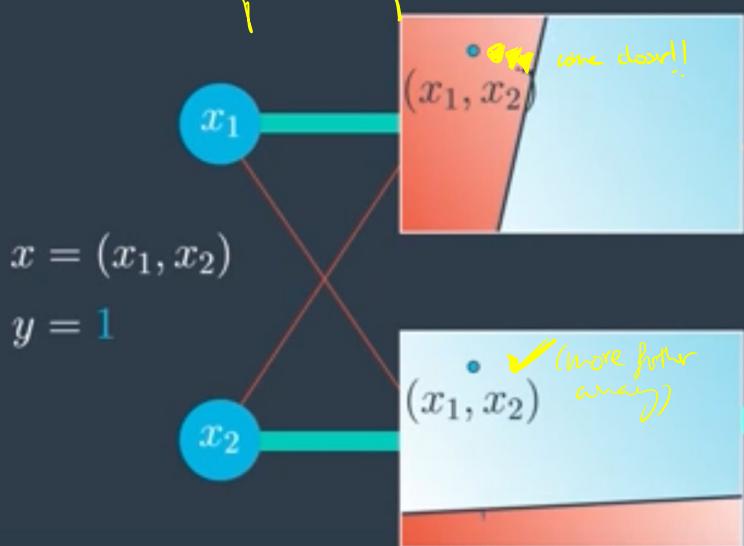
→ Feed forward: Perceptron plots \leftarrow point and returns the probability of it being the

→ Back propagation asks the point: how would you like to change the model? | Closer for any

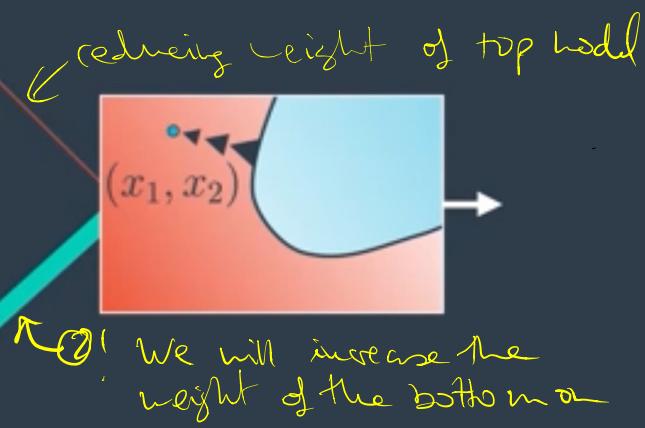
How can we do it for Multilayer Perceptron?

We obtain new weights They get closer to the point

what can this models do to classify you? ↗
Backpropagation

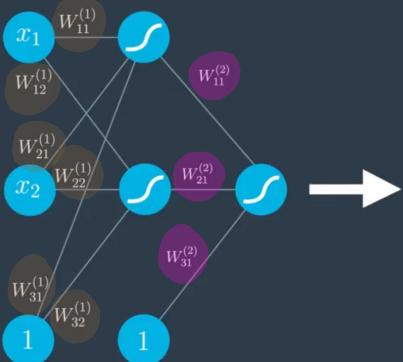


Which model is doing better? This model is not good! (missclass)



Backpropagation Math

Backpropagation



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

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

Note: written for convenience
Bias

$$\nabla E = \left(\begin{array}{c} \frac{\partial E}{\partial W_{11}^{(1)}} \quad \frac{\partial E}{\partial W_{12}^{(1)}} \quad \frac{\partial E}{\partial W_{11}^{(2)}} \\ \frac{\partial E}{\partial W_{21}^{(1)}} \quad \frac{\partial E}{\partial W_{22}^{(1)}} \quad \frac{\partial E}{\partial W_{21}^{(2)}} \\ \frac{\partial E}{\partial W_{31}^{(1)}} \quad \frac{\partial E}{\partial W_{32}^{(1)}} \quad \frac{\partial E}{\partial W_{31}^{(2)}} \end{array} \right)$$

$W'_{ij}^{(k)} \leftarrow W_{ij}^{(k)} - \alpha \frac{\partial E}{\partial W_{ij}^{(k)}}$

Updating each weight

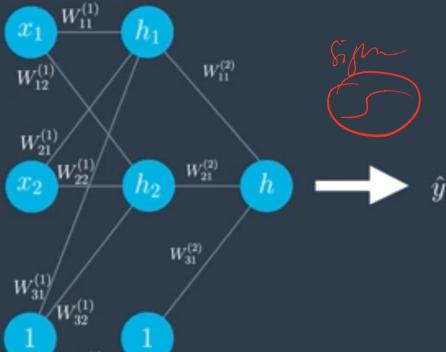
so it will give us new updated weight $w_{i,j}$ super k prime.

Chain Rule

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

When composing functions
(ex: just forward), derivatives just multiply

Feedforward



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

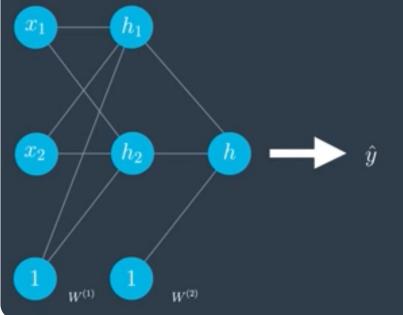
$$h_2 = W_{12}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{32}^{(1)}$$

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

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

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

Backpropagation



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

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

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

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

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

Backpropagation



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

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

derivative of σ w.r.t h_1 ($\frac{\partial \sigma}{\partial h_1}$)

Training / Optimization

Underfitting: killing too little w/ a poor

Overfitting: killing a lot w/ a loose

high variance. Does not generalize well

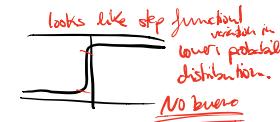
Train / Test \approx **Study / Exam**

Early Stopping
→ When testing error starts increasing.
Early quitting on epochs!

Subtle Overfitting

$10x_1 + 0x_2$ gives less error than $x_1 + x_2$
due to $\Gamma(20) = 0.99 \dots \rightarrow \Gamma(2) = 0.8$
 $\Gamma(-20) = 0. \dots \circ 21 \rightarrow \Gamma(2) = 0.2$

But is this better?



Solution: Regularization

LARGE COEFFICIENTS → OVERFITTING

PENALIZE LARGE WEIGHTS
(w_1, \dots, w_n)

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

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

$$0.5^2 + 0.5^2 = 0.5 \quad \text{vs} \quad 1^2 + 0^2 = 1$$

L2 prefers smaller coeff.

L1 good for feature selection

Sparsity (1, 0, 0, 1, 0)

Sparsity (0.5, 0.3, -0.2, 0.4, 0.1)

↳ L2 usually better for training models

Dropout

Sometimes one part of the network has very large weights and outshadows the other part, which doesn't play any role at all...
In order to solve this we will randomly turn off some of the nodes in the epoch. ⇒ Better generalization.

Vanishing gradient

Occurs due to the shape of the sigmoid

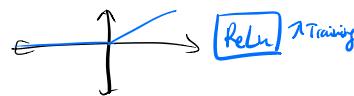
Batch vs Stochastic Gradient Descent

Epoch = Steps

Stochastic gradient descent → smaller sub-steps
In order to do this, we split the data into several batches

Adaptive Learning Rate

Better if it decreases as we get closer to solution



Momentum

A way to get over the humps (local minima)

Constant $\beta \in [0, 1]$

$$\text{Step}(u) + \beta \text{Step}(u-1) + \beta^2 \text{Step}(u-2) + \dots$$

