

# Lecture 13. Feed-forward Neural Networks and Backpropagation

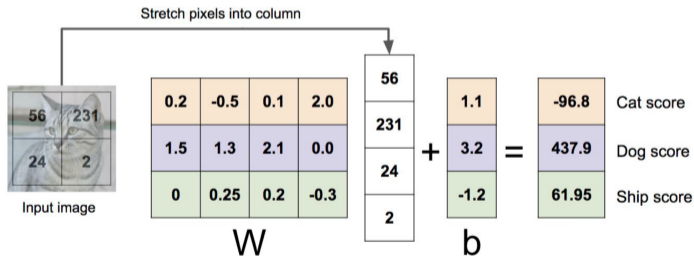
Bao Wang  
Department of Mathematics  
Scientific Computing and Imaging Institute  
University of Utah  
Math 5750/6880, Fall 2021

# Multi-class Logistic Regression

## Linear classifier

$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}\mathbf{x} + \mathbf{b}.$$

Consider an image classification task with 3 classes (cat/dog/frog), where each image has 4 pixels. In this case,  $\mathbf{W} \in \mathbb{R}^{3 \times 4}$  and  $\mathbf{b} \in \mathbb{R}^3$ .



Is the output good? How to find optimal  $\mathbf{W}$  and  $\mathbf{b}$ ?

## How to find $\mathbf{W}$ and $\mathbf{b}$ ?

- Define a **loss function** that quantifies our unhappiness with the scores across the training data.
- **Optimization:** Find the parameters  $\mathbf{W}$  and  $\mathbf{b}$  that minimize the loss function. Note that we can absorb  $\mathbf{b}$  into  $\mathbf{W}$  by introduce an additional dimension (all 1s) to the feature.

## Loss function

A loss function tells how good our current classifier is.

- Given a dataset of examples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , where  $\mathbf{x}_i$  is image and  $y_i$  is (integer) label.
- Loss over the dataset is the average of loss over examples

$$L = \frac{1}{N} \sum_i L_i(f(\mathbf{x}_i, \mathbf{W}), y_i).$$

## Regularization

$$L(\mathbf{W}) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(\mathbf{x}_i, \mathbf{W}), y_i)}_{\text{data loss}} + \underbrace{\lambda R(\mathbf{W})}_{\text{regularization}},$$

where  $\lambda$  is the regularization strength.

- $L_2$  regularization:  $R(\mathbf{W}) = \sum_k \sum_l W_{k,l}^2$ .
- $L_1$  regularization:  $R(\mathbf{W}) = \sum_k \sum_l |W_{k,l}|$ .
- Elastic net:  $R(\mathbf{W}) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$ .

## Softmax classifier – Multinomial logistic regression

Interpret raw classifier scores as probabilities.

Given

$$\mathbf{s} = f(\mathbf{x}_i; \mathbf{W}),$$

how to convert  $\mathbf{s}$  into probabilities?

## Softmax classifier – Multinomial logistic regression

Interpret raw classifier scores as probabilities.

$$\mathbf{s} = f(\mathbf{x}_i; \mathbf{W}) \Rightarrow P(Y = k | X = \mathbf{x}_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad (\text{softmax function}).$$

We need to maximize probability of correct class, i.e., minimize

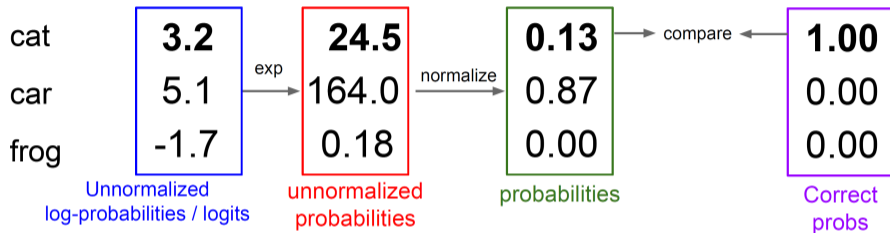
$$L_i = -\log P(Y = y_i | X = \mathbf{x}_i),$$

and therefore the loss contributed by  $\mathbf{x}_i$  is

$$L_i = -\log \left( \frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right).$$

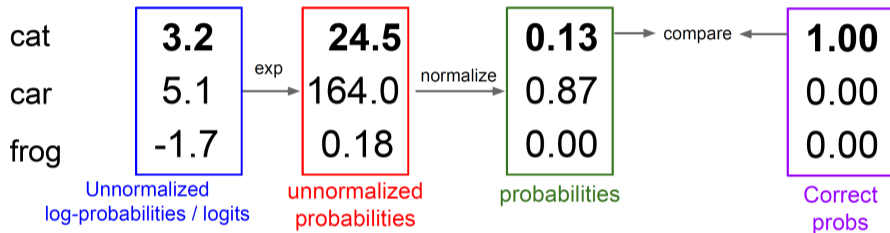
When  $L_i$  will be minimized?





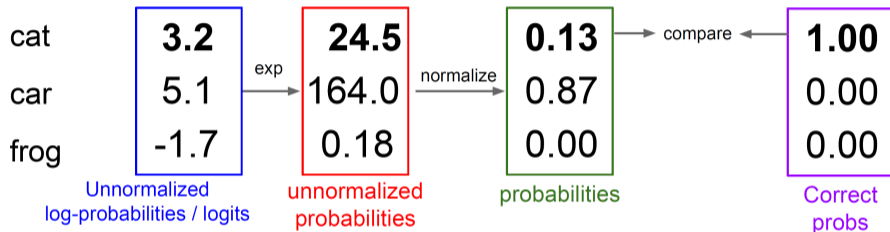
$$L_i = -\log P(Y = y_i | X = x_i) = -\log \left( \frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right) = -\log(0.13) = 2.04.$$

## Kullback-Leibler (KL) divergence



$$D_{KL}(P||Q) = \sum_y P(y) \log \frac{P(y)}{Q(y)}.$$

## Cross entropy



$$H(P, Q) = H(P) + D_{KL}(P||Q),$$

where  $H(P) = -\sum_y P(y) \log P(y)$  is the entropy of the distribution  $P$ .

## Gradient descent (GD)

So far, we have the following loss function

$$L(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N L_i(\mathbf{x}_i, y_i, \mathbf{W}) + \lambda R(\mathbf{W}),$$

and the gradient is

$$\nabla_{\mathbf{W}} L(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{W}} L_i(\mathbf{x}_i, y_i, \mathbf{W}) + \lambda \nabla_{\mathbf{W}} R(\mathbf{W}),$$

then perform the following gradient descent

$$\mathbf{W}^{k+1} = \mathbf{W}^k - s \nabla_{\mathbf{W}} L(\mathbf{W}^k).$$

Compute  $\nabla_{\mathbf{W}} L(\mathbf{W})$  can be very expensive when  $N$  is large.

## Stochastic gradient descent (SGD)

Replace the true gradient above with the following mini-batch gradient

$$\frac{1}{n} \sum_{j=1}^n \nabla_{\mathbf{W}} L_{i_j}(\mathbf{x}_{i_j}, v_{i_j}, \mathbf{W}) + \lambda \nabla_{\mathbf{W}} R(\mathbf{W}), \quad n \ll N.$$

# Neural Networks

## Linear function

- Linear function:  $f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x}$ .

## Neural network

- Linear function:  $f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x}$ .
- 2-layer neural network

$$f(\mathbf{x}; \mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_2 \max(0, \mathbf{W}_1 \mathbf{x}),$$

where  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{W}_q \in \mathbb{R}^{H \times d}$ , and  $\mathbf{W}_2 \in \mathbb{R}^{C \times H}$ .

$\max(0, *)$  is the ReLU activation.



## Neural network

- Linear function:  $f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x}$ .
- 2-layer neural network

$$f(\mathbf{x}; \mathbf{W}_1, \mathbf{W}_2) = \mathbf{W}_2 \max(0, \mathbf{W}_1 \mathbf{x}),$$

where  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{W}_q \in \mathbb{R}^{H \times d}$ , and  $\mathbf{W}_2 \in \mathbb{R}^{C \times H}$ .

- 3-layer neural network

$$f(\mathbf{x}; \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3) = \mathbf{W}_3 \max(0, \mathbf{W}_2 \max(0, \mathbf{W}_1 \mathbf{x})),$$

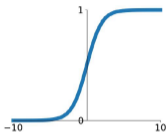
where  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{W}_1 \in \mathbb{R}^{H_1 \times d}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{H_2 \times H_1}$ , and  $\mathbf{W}_3 \in \mathbb{R}^{C \times H_2}$ .

$\max(0, *)$  is the ReLU activation.

## Activation functions

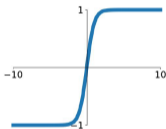
### Sigmoid

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



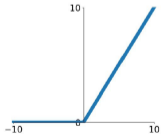
### tanh

$$\tanh(x)$$



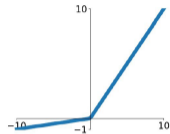
### ReLU

$$\max(0, x)$$



### Leaky ReLU

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

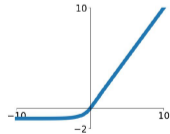


### Maxout

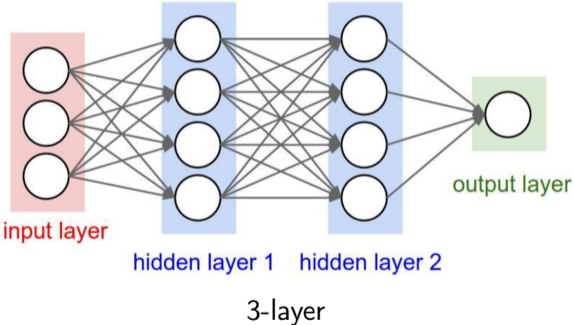
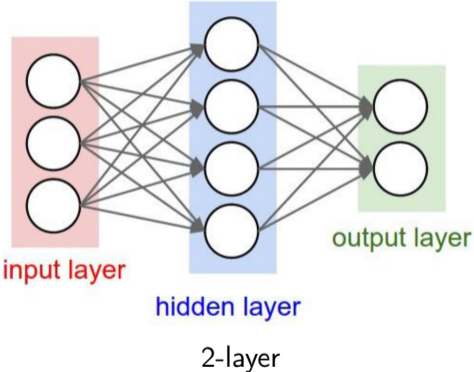
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

### ELU

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



# Architectures



Backpropagation

## A simple example

Consider the following function

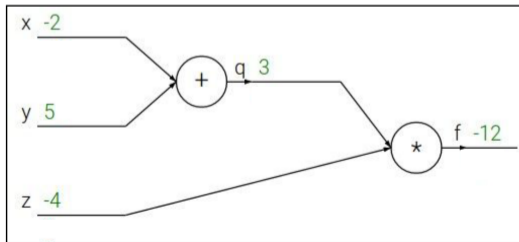
$$f(x, y, z) = (x + y)z,$$

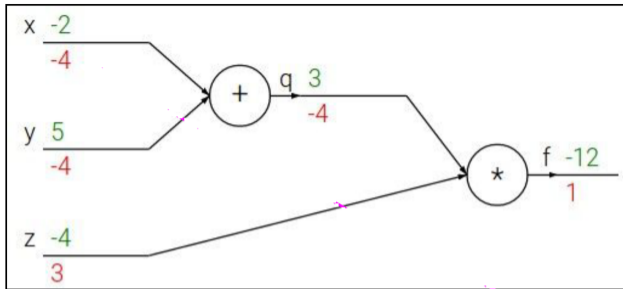
e.g.,  $x = -2, y = 5, z = -4$ . Note that  $f(x, y, z)$  can be decomposed as follows

$$q = x + y, \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1; \quad f = qz, \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q.$$

We need to compute  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$ .

Note the above function can be represented as the following computational graph





Note that

$$\frac{\partial f}{\partial f} = 1,$$

then

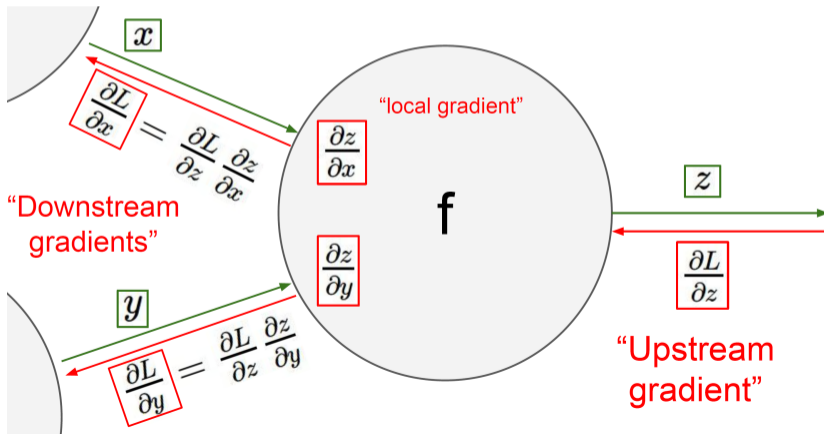
$$\frac{\partial f}{\partial z} = q = x + y = 3, \quad \frac{\partial f}{\partial q} = z = -4,$$

thus

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = -4 * 1 = -4; \quad \frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = -4 * 1 = -4.$$

## Chain rule

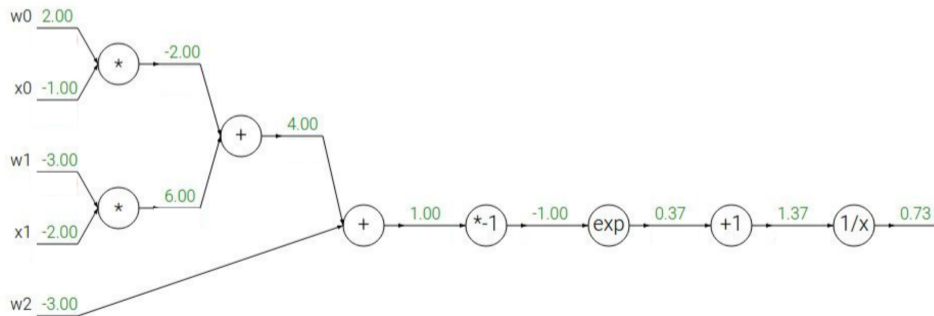
$$\underbrace{\frac{\partial f}{\partial x}}_{\text{Downstream gradient}} = \underbrace{\frac{\partial f}{\partial q}}_{\text{Upstream gradient}} \times \underbrace{\frac{\partial q}{\partial x}}_{\text{Local gradient}}$$



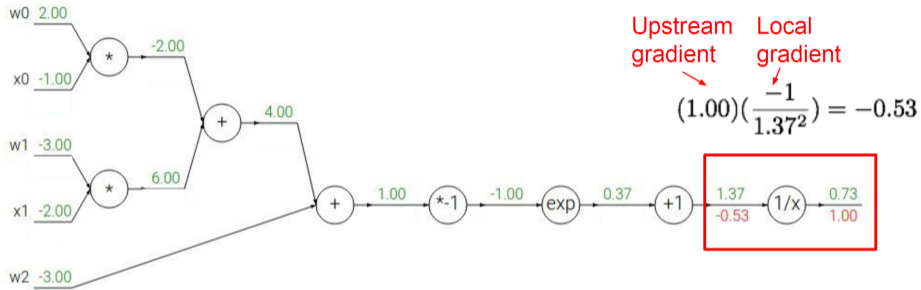
## Another example

$$f(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}},$$

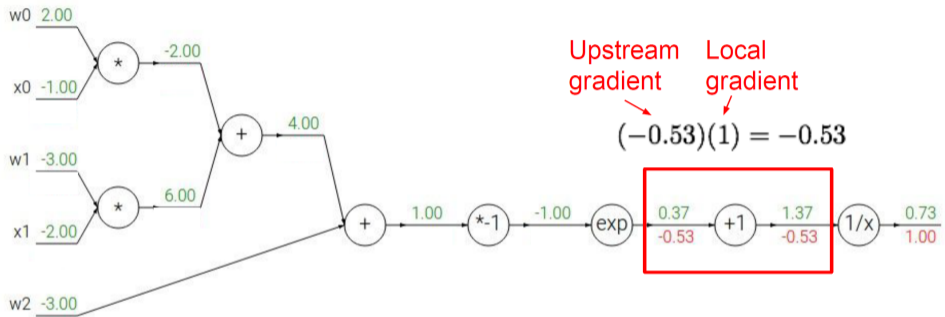
where  $w_0 = 2$ ,  $w_1 = -3$ ,  $w_2 = -3$ ,  $x_0 = -1$ , and  $x_1 = -2$ . The function can be represented by the following computational graph



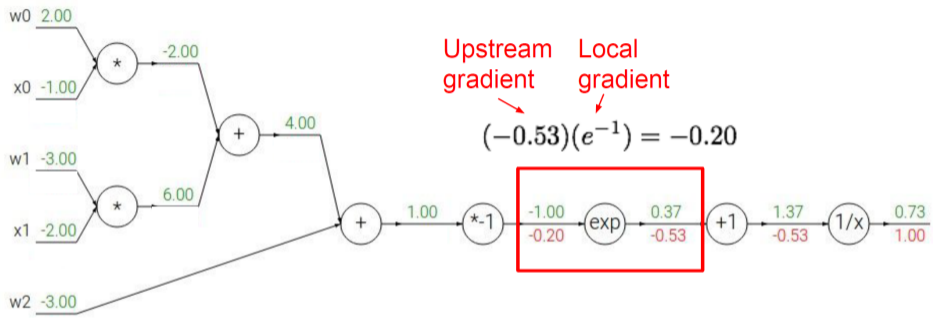




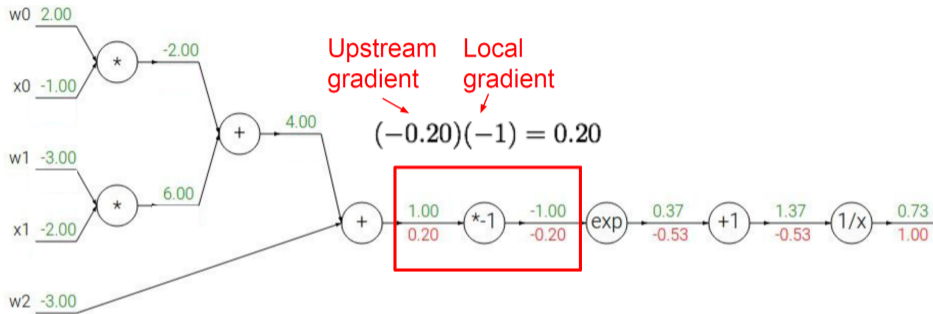
$$f(x) = \frac{1}{x} \Rightarrow \frac{df}{dx} = -\frac{1}{x^2}.$$



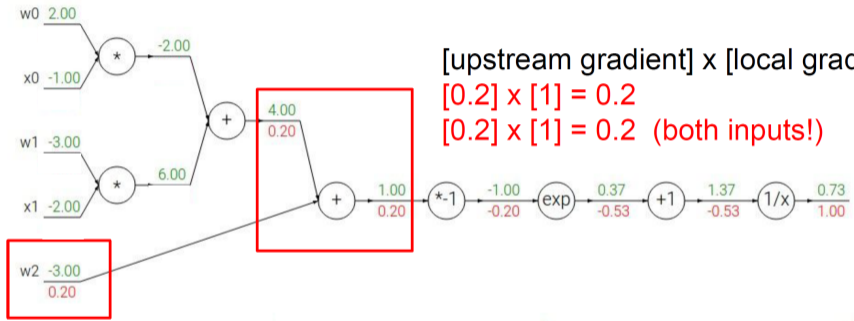
$$f_c(x) = c + x \Rightarrow \frac{df_c}{dx} = 1.$$

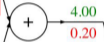
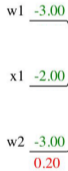
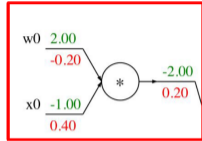


$$f(x) = e^x \Rightarrow \frac{df}{dx} = e^x.$$



$$f_a(x) = ax \Rightarrow \frac{df}{dx} = a.$$





[upstream gradient] x [local gradient]

$w_0: [0.2] \times [-1] = -0.2$

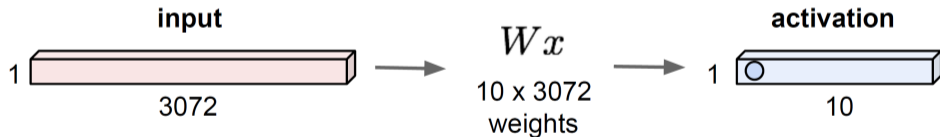
$x_0: [0.2] \times [2] = 0.4$

# Convolutional Neural Networks

## How linear model classify images?

**Image:** height  $\times$  width  $\times$  channels, e.g.  $32 \times 32 \times 3$ .

- First stretch the image to a vector  $\mathbf{x} \in \mathbb{R}^{3072 \times 1}$ .
- Feed  $\mathbf{x}$  into the model  $\mathbf{y} = \mathbf{W}\mathbf{x}$  (fully-connected layer).

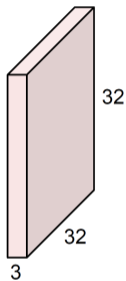




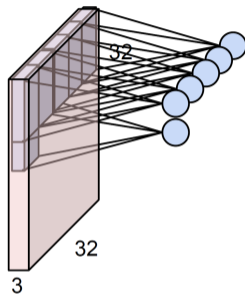
## Convolution layer

Preserve the spatial structure of the input image.

32x32x3 image



5x5x3 filter



Convolve the filter with the image, i.e. “slide over the image spatially over height and width/or over height, width, and channel, computing dot products”.

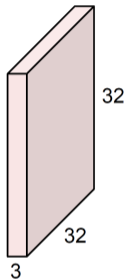
Recall convolution of two functions

$$(f * g)(y) = \int_{-\infty}^{\infty} f(x) \cdot g(y - x) dx.$$

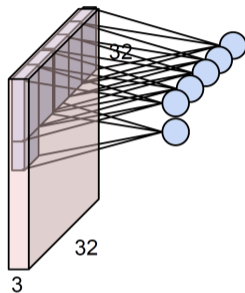
## Convolution layer

Preserve the spatial structure of the input image.

32x32x3 image



5x5x3 filter

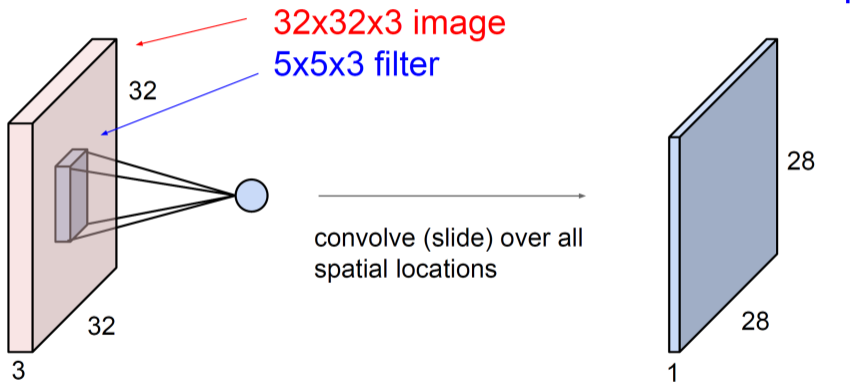


Convolve the filter with the image, i.e. “slide over the image spatially over height and width/over height, width, and channel, computing dot products”.

Recall convolution of two signals

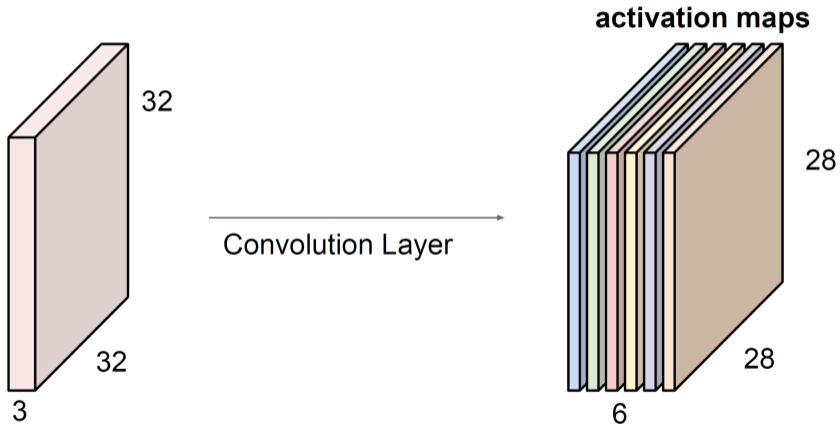
$$f(x, y) * g(x, y) = \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f(n_1, n_2) \cdot g(x - n_1, y - n_2).$$

# Activation map



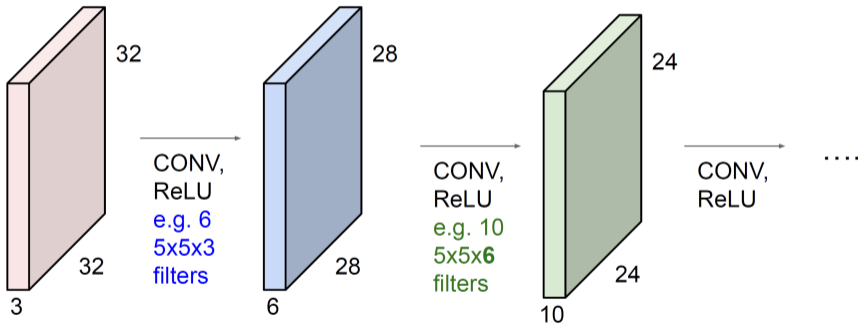
## Multiple activation maps

We can use multiple filters to get multiple activation maps.

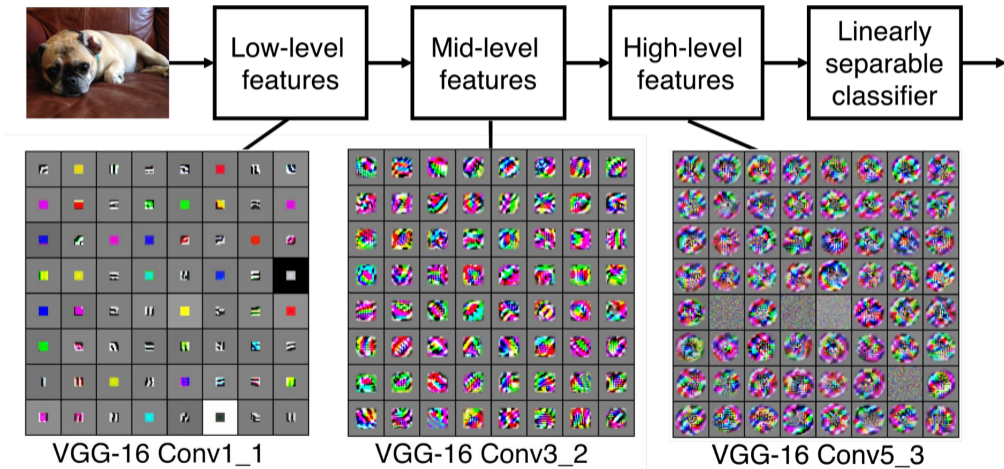


## Convolutional neural networks (CNNs)

CNN is a sequence of convolution layers, interspersed with activation functions.

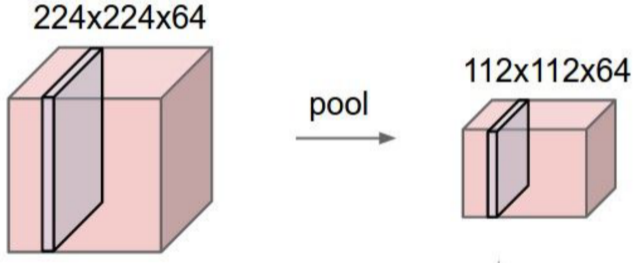


# CNN for representation learning

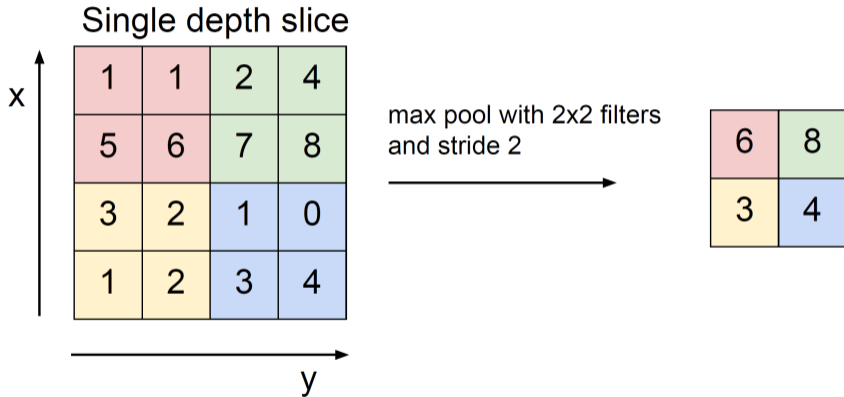


## Pooling layer

Pooling layer is used to make the representations smaller and more manageable, which operates over each activation map independently.



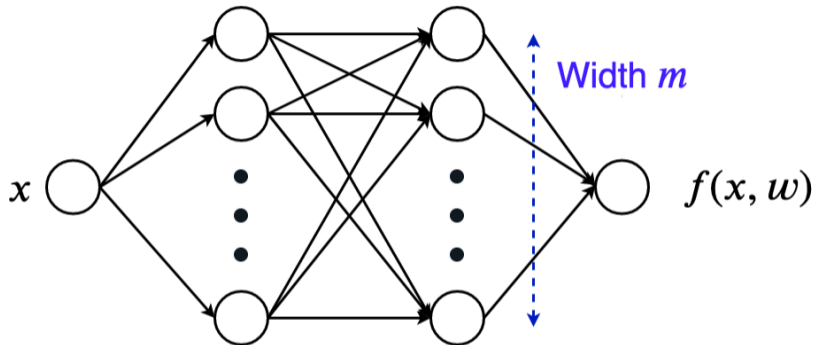
# Max pooling





## Neural Networks vs. Kernel Methods

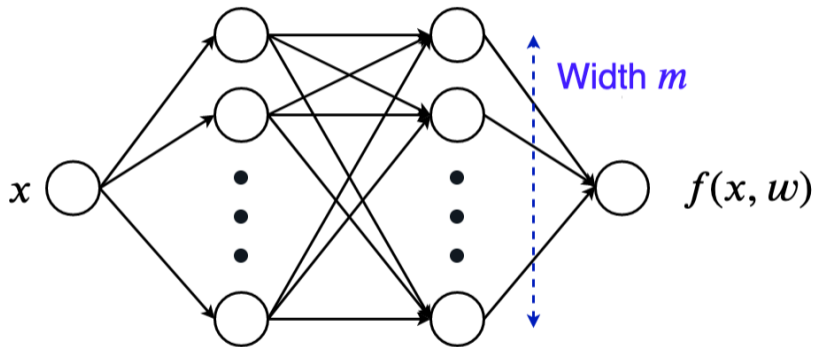
How does neural network relate to kernel method?



Neural networks in the infinite-width regime simplify to linear models with a kernel called the neural tangent kernel. Under this regime, we can show that gradient descent will converge to 0 training loss.

## A 1D example

Let us start with a very simple example with 1D input and 1D output. In particular, a simple 2-hidden layer ReLU network with width  $m$ .



- Denote the neural network function as  $f(x, \mathbf{w})$  where  $x$  is the input and  $\mathbf{w}$  is the combined vector of weights (say of size  $p$ ).
- Assume the training dataset is  $\{\bar{x}_i, \bar{y}_i\}_{i=1}^N$ .

Let us consider performing full-batch gradient descent on the following least squares loss

$$L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \left( f(\bar{x}_i, \mathbf{w}) - \bar{y}_i \right)^2.$$

- Stack all the output dataset value  $\bar{y}_i$  into a single vector of size  $N$ , and call it  $\bar{\mathbf{y}}$ .
- Stack all the model outputs for each input,  $f(\bar{x}_i, \mathbf{w})$  into a single prediction vector  $\mathbf{y}(\mathbf{w}) \in \mathbb{R}^N$ , i.e.,  $\mathbf{y}(\mathbf{w})_i = f(\bar{x}_i, \mathbf{w})$ .
- Our loss simplifies to

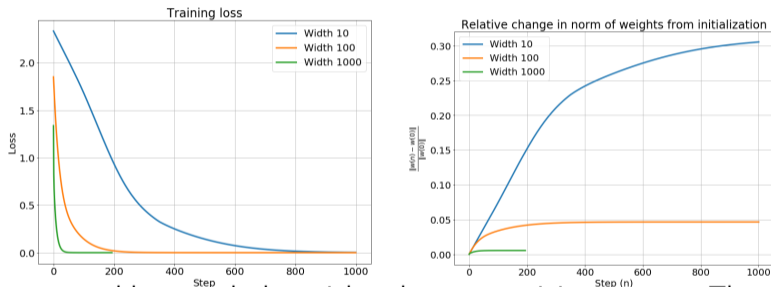
$$L(\mathbf{w}) = \frac{1}{N} \cdot \frac{1}{2} \|\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}\|_2^2.$$

We can rescale it and consider  $L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}\|_2^2$ .

How does the weight change under gradient descent?

Consider the *relative change* in the norm of the weight vector from initialization:

$$\frac{\|\mathbf{w}(n) - \mathbf{w}_0\|_2}{\|\mathbf{w}_0\|_2}$$



**Figure:** Loss curves and how much the weights change as training progress. The weight norms are calculated using all parameters in the model stacked into a single vector. Also, all other hyperparameters (like learning rate) are kept constant.

The weights don't change much at all for larger hidden widths. – What can we do to analyze  $f(x, \mathbf{w})$ ?

## Linearize neural networks

Taylor expand the network function **w.r.t. the weights** around its **initialization**. (Note in training, the weights do not change much.)

$$f(x, \mathbf{w}) \approx f(x, \mathbf{w}_0) + \nabla_{\mathbf{w}} f(x, \mathbf{w}_0)^\top (\mathbf{w} - \mathbf{w}_0).$$

We've turned the non-linear neural network function into a linear function of the weights. Using our more concise vector notation for the model outputs on a specific dataset we can rewrite as:

$$\mathbf{y}(\mathbf{w}) \approx \mathbf{y}(\mathbf{w}_0) + \nabla_{\mathbf{w}} \mathbf{y}(\mathbf{w}_0)^\top (\mathbf{w} - \mathbf{w}_0),$$

where  $\mathbf{y}(\mathbf{w}) \in \mathbb{R}^{n \times 1}$ ,  $\nabla_{\mathbf{w}} \mathbf{y}(\mathbf{w}_0)^\top \in \mathbb{R}^{n \times p}$ , and  $\mathbf{w} \in \mathbb{R}^{p \times 1}$  with  $n$  and  $p$  are the number of datapoints and parameters, respectively.

Note that the initial output  $\mathbf{y}(\mathbf{w}_0)$  and the model Jacobian  $\nabla_{\mathbf{y}}(\mathbf{w}_0)$  are just constants. Thus

$$\mathbf{y}(\mathbf{w}) \approx \mathbf{y}(\mathbf{w}_0) + \nabla_{\mathbf{w}}\mathbf{y}(\mathbf{w}_0)^\top(\mathbf{w} - \mathbf{w}_0),$$

is just a **linear model in the weights**, so minimizing the least squares loss reduces to just doing **linear regression!**

But, notice that the model function is still **non-linear in the input**, because finding the gradient of the model is definitely not a linear operation. In fact, this is just a linear model using a **feature map**  $\phi(\mathbf{x})$  which is the gradient vector at initialization:

$$\phi(\mathbf{x}) = \nabla_{\mathbf{w}}f(\mathbf{x}, \mathbf{w}_0).$$

This feature map induces a kernel on the input; called the **neural tangent kernel**.

## Gradient descent dynamics

We've considered linearization of neural networks. Let us now get into their training dynamics under gradient descent.

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_k).$$

Rewriting this equation a bit, we get:

$$\frac{\mathbf{w}_{k+1} - \mathbf{w}_k}{\eta} = -\nabla_{\mathbf{w}} L(\mathbf{w}_k).$$

Let  $\eta \rightarrow 0$ , we have

$$\frac{d\mathbf{w}(t)}{dt} = -\nabla_{\mathbf{w}} L(\mathbf{w}(t)).$$

This is called a **gradient flow**. To simplify notation, we denote time derivatives with a dot, the gradient flow is:

$$\dot{\mathbf{w}}(t) = -\nabla L(\mathbf{w}(t)).$$



## Gradient flow

Dropping the time variable, and substituting for the loss and taking the gradient, we get (note  $L(\mathbf{w}(t)) = \frac{1}{2}\|\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}\|_2^2$ .)

$$\dot{\mathbf{w}} = -\nabla \mathbf{y}(\mathbf{w})(\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}). \quad (\nabla \mathbf{y}(\mathbf{w}) = \nabla_{\mathbf{w}} \mathbf{y}(\mathbf{w}))$$

We can now derive the dynamics of the model outputs  $\mathbf{y}(\mathbf{w})$  (this is basically the **dynamics in function space**) induced by this gradient flow using the chain rule:

$$\dot{\mathbf{y}}(\mathbf{w}) = \nabla \mathbf{y}(\mathbf{w})^\top \dot{\mathbf{w}} = -\nabla \mathbf{y}(\mathbf{w})^\top \nabla \mathbf{y}(\mathbf{w}) (\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}).$$

The quantity  $\mathbf{H}(\mathbf{w}) := \nabla \mathbf{y}(\mathbf{w})^\top \nabla \mathbf{y}(\mathbf{w})$  is called the **neural tangent kernel** (NTK).

## Neural tangent kernel

Recall that the linearized model has a feature map  $\phi(x) = \nabla_{\mathbf{w}} f(x, \mathbf{w}_0)$ . The kernel matrix corresponding to this feature map is obtained by taking pairwise inner products between the feature maps of all the data points. This is exactly  $\mathbf{H}(\mathbf{w}_0)$ .

$$\begin{aligned} \begin{array}{c} \updownarrow n \\ \left[ \begin{array}{c} \mathbf{H}(\mathbf{w}_0) \end{array} \right] \\ \leftarrow n \end{array} &= \begin{array}{c} \left[ \begin{array}{c} \nabla_{\mathbf{w}} \mathbf{y}(\mathbf{w}_0)^T \end{array} \right] \\ \leftarrow p \end{array} \begin{array}{c} \left[ \begin{array}{c} \nabla_{\mathbf{w}} \mathbf{y}(\mathbf{w}_0) \end{array} \right] \\ \updownarrow p \\ \leftarrow n \end{array} \\ \text{NTK} & \\ &= \begin{array}{c} \left[ \begin{array}{c} \text{--- } \phi(\bar{x}_1)^T \text{ ---} \\ \vdots \\ \text{--- } \phi(\bar{x}_n)^T \text{ ---} \end{array} \right] \begin{array}{c} \left[ \begin{array}{c} | \\ \phi(\bar{x}_1) \cdots \phi(\bar{x}_n) \\ | \end{array} \right] \end{array} \end{array} \end{aligned}$$

## Gradient flow

If the model is close to its linear approximation, the Jacobian of the model outputs **does not change as training progresses**, i.e.,

$$\nabla \mathbf{y}(\mathbf{w}(t)) \approx \nabla \mathbf{y}(\mathbf{w}_0) \Rightarrow \mathbf{H}(\mathbf{w}(t)) \approx \mathbf{H}(\mathbf{w}_0).$$

This is referred to as the **kernel regime**, because the tangent kernel stays constant during training. The training dynamics now reduces to a very simple **linear ordinary differential equation**:

$$\dot{\mathbf{y}}(\mathbf{w}) = -\mathbf{H}(\mathbf{w}_0)(\mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}).$$

## Gradient flow

Clearly,  $\mathbf{y}(\mathbf{w}) = \bar{\mathbf{y}}$  is an equilibrium of this ODE, and it corresponds to a train loss of 0. Let  $\mathbf{u} = \mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}$ , the flow simplifies to

$$\dot{\mathbf{u}} = -\mathbf{H}(\mathbf{w}_0)\mathbf{u}.$$

The solution of this ODE is given by a matrix exponential

$$\mathbf{u}(t) = \mathbf{u}(0)e^{-\mathbf{H}(\mathbf{w}_0)t}.$$

When over-parameterized ( $p > n$ ), the NTK  $\nabla \mathbf{y}(\mathbf{w}_0)^\top \nabla \mathbf{y}(\mathbf{w}_0)$  is positive definite (ignoring any degeneracy in the dataset that would cause  $\nabla \mathbf{y}(\mathbf{w}_0)$  to not have full column rank).

$$\mathbf{u}(t) = \mathbf{u}(0)e^{-\mathbf{H}t},$$

where  $\mathbf{H} := \mathbf{H}(\mathbf{w}_0)$  is positive definite. Let  $\mathbf{H} = \mathbf{U}\mathbf{D}\mathbf{U}^{-1}$  be the spectral decomposition of  $\mathbf{H}$ , where  $\mathbf{U}$  is orthogonal and  $\mathbf{D}$  is diagonal with diagonal entries be positive. Therefore, using the fact that

$$e^{-\mathbf{U}\mathbf{D}\mathbf{U}^{-1}t} = \sum_{k=0}^{\infty} \frac{1}{k!} \left( \mathbf{U}(-\mathbf{D}t)\mathbf{U}^{-1} \right)^k = \mathbf{U} \left( \sum_{k=0}^{\infty} \frac{1}{k!} (-\mathbf{D}t)^k \right) \mathbf{U}^{-1} = \mathbf{U}e^{-\mathbf{D}t}\mathbf{U}^{-1},$$

we have

$$\mathbf{u}(t) = \mathbf{u}(0)\mathbf{U}e^{-\mathbf{D}t}\mathbf{U}^{-1} \rightarrow 0,$$

where  $(e^{-\mathbf{D}t})_{ii} = e^{-d_{ii}t}$ .

# Universal Approximation Theorem

## Capacity of the Perceptron

**What kind of function can a neural network represent?** Let us start with the simplest neural network, the Perceptron, a NN with a single hidden layer having 1 hidden unit with an activation function  $\sigma$ . Both the input layer and the weights are a  $1 \times n$  vector.

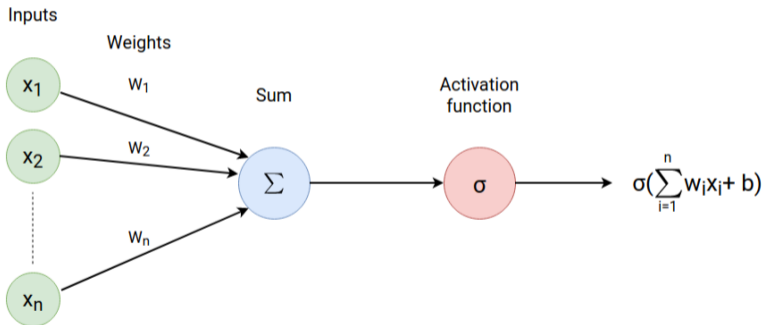


Figure: The perceptron.

Perceptron:

$$\sigma\left(\sum_{i=1}^n w_i x_i + b\right) = \sigma(\mathbf{w}^\top \mathbf{x} + b),$$

where  $\sigma$  can be sigmoid:  $\frac{1}{1+e^{-x}}$ , tanh:  $\frac{e^{2x}-1}{e^{2x}+1}$ , ReLU:  $\max(0, x)$  ...

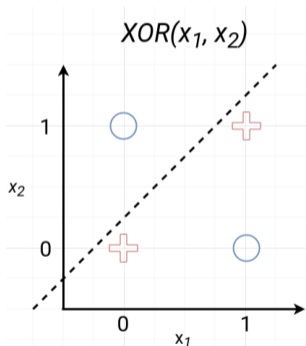


Figure: xOR is not linearly separable.

Perceptron is a linear classifier, and as such it can't model an XOR.



## Capacity of multiple neurons

By allowing ourselves more than 1 neuron in the hidden layer, we can model a XOR and in fact, we get the simplest **universal approximator**.

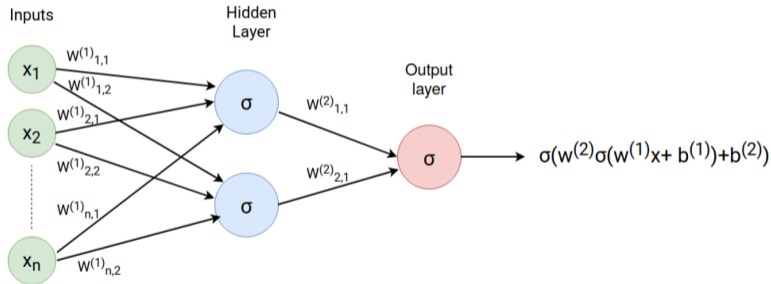


Figure: A NN with 2 hidden units.

## Universal approximation theorem

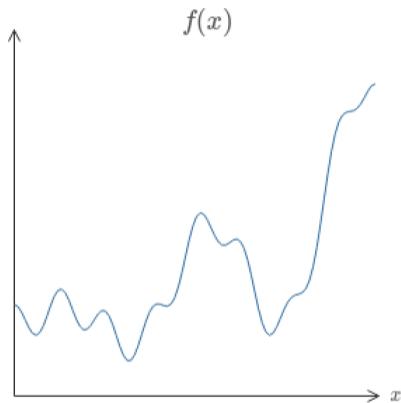
**Universal approximation theorem:** any continuous function  $f : [0, 1]^n \rightarrow [0, 1]$  can be approximated arbitrarily well by a neural network with at least one hidden layer with a finite number of weights.

A visual proof is available at

<http://neuralnetworksanddeeplearning.com/chap4.html>.

## Visual proof of universal approximation

Suppose we want to approximate a function with 1 input and 1 output like:



**Figure:** A continuous function.

We will first consider a simple NN with 2 hidden neurons that have a sigmoid activation function, and for now the output neuron will just be linear.

## Visual proof of universal approximation – Step 1

Make a step function with 1 of the neuron.

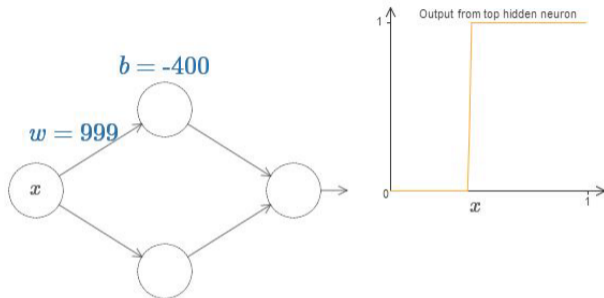


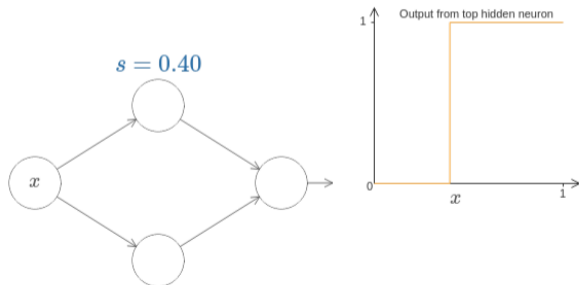
Figure: Making a step function with the top neuron.

Let us focus on the top hidden neuron first, by using a big weight on the top neuron we can approximate the step function with a sigmoid arbitrarily well, and by adjusting the bias we can place it anywhere. (The same argument can be made for the tanh activation, but not for ReLU.) Note that  $\sigma(wx)$  becomes steeper as  $w$  increases.

## Visual proof of universal approximation – Step 1

In the toy example above, we won't be interested in changing the weights of the first layer, they just have to be high enough, so we will just consider them to be constant. Additionally, to make the plots clearer, we will display the position of the step instead of of the bias, which is easily computed with  $s = -b/w$ .

With these changes, the plot above becomes:



**Figure:** Making a step function with the top neuron.

## Visual proof of universal approximation – Step 2

Make a “bin” with an opposite step function.

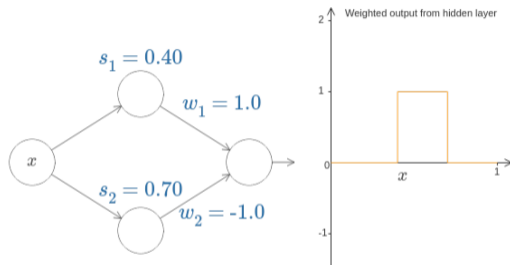
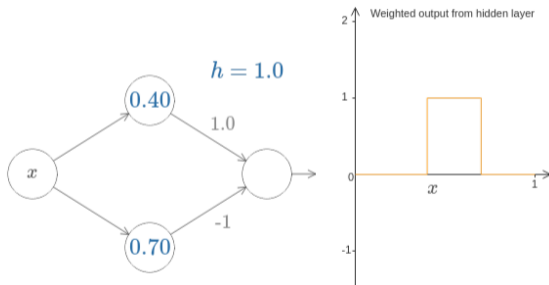


Figure: Making a bin with 2 opposite step functions.

As illustrated above, by using the other neuron to make a step function, and setting opposing weights in the second layer, we can effectively approximate a bin and control its position, size and height.

## Visual proof of universal approximation – Step 2

Now you can probably see where this is going, to make things even clearer, we will just use 1 value for both  $w_1$  and  $-w_2$ , called  $h$ , representing the height of the “bin”.



**Figure:** Making a bin with 2 opposite step functions.

## Visual proof of universal approximation – Step 3

### Discretize the function

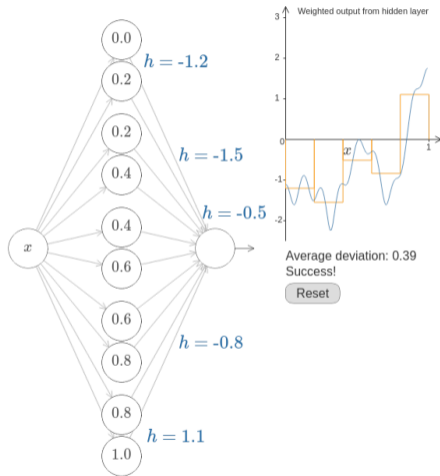


Figure: Approximating  $f$  with an histogram.



In this final step, we combine several “bias” to make an histogram approximating the function. Illustrated above is a very rough approximation using only 5 bins (10 hidden units), but we can obviously make it sharper by simply adding more bins.

**Question:** If we want to use this technique to approximate an  $L$ -Lipschitz function  $f : \mathbb{R} \rightarrow \mathbb{R}$  on the interval  $[0, 1]$  with an error at most  $\epsilon$  at any point, how many bins would we need?

**What if we don't want linear neurons in the output layer?** The above network's output layer is linear, giving us the histogram approximating  $f$ , if we add a sigmoid activation function on the output we just have to approximate  $\sigma^{-1} \circ f$  instead of  $f$ , which we can do with the same method.

## Cybenko approximation by superposition of sigmoid function

We first define a *sigmoidal function*  $\sigma$  as:

$$\sigma(x) \rightarrow \begin{cases} 1 & \text{as } x \rightarrow +\infty \\ 0 & \text{as } x \rightarrow -\infty \end{cases}$$

While sigmoidal functions are usually assumed to be monotonic increasing, this assumption is not necessary for this result.

**Theorem.** Let  $C([0, 1]^n)$  denote the set of all continuous function  $[0, 1]^n \rightarrow \mathbb{R}$ , let  $\sigma$  be any sigmoid activation function then the finite sum of the form  $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i \sigma(\mathbf{w}_i^\top \mathbf{x} + b_i)$  is dense in  $C([0, 1]^n)$ .

Informally, the above theorem says that for any  $g \in C([0, 1]^n)$  for any  $\epsilon > 0$ , there exists  $f : \mathbf{x} \rightarrow \sum_{i=1}^N \alpha_i \sigma(\mathbf{w}_i^\top \mathbf{x} + b)$  such that  $|f(\mathbf{x}) - g(\mathbf{x})| < \epsilon$  for all  $\mathbf{x} \in [0, 1]^n$ .

