

# CSCI 567: Machine Learning

Vatsal Sharan  
Spring 2026

Lecture 7, Feb 27

# Administrivia

- Exam 1 is next week (March 6, 2 hr 20 min, starting at 1pm)
- Students will be split into two rooms, instructions later (DEN students will get separate instructions)
- You can bring one cheat sheet (you can write on both sides), though we will generally provide necessary formulae
- No other books, resources etc.
- Based on materials covered till SVMs (lectures, discussions, HW1 & HW2)



# Multiclass classification

# Setup

Recall the setup:

- input (feature vector):  $\mathbf{x} \in \mathbb{R}^d$
- output (label):  $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping  $f : \mathbb{R}^d \rightarrow [C]$

**Examples:**

- recognizing digits ( $C = 10$ ) or letters ( $C = 26$  or  $52$ )
- predicting weather: sunny, cloudy, rainy, etc
- predicting image category: ImageNet dataset ( $C \approx 20K$ )

# Linear models: Binary to multiclass

Step 1: *What should a linear model look like for multiclass tasks?*

Note: a linear model for binary tasks (switching from  $\{-1, +1\}$  to  $\{1, 2\}$ )

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ 2 & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

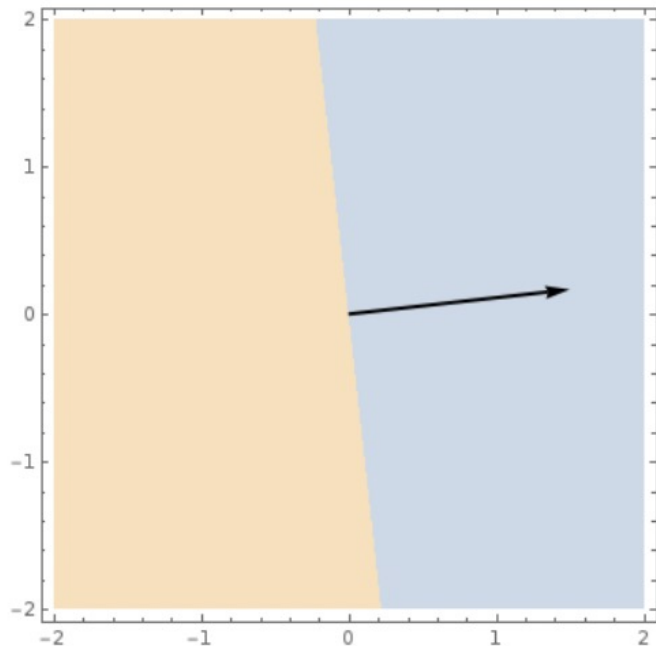
can be written as

$$\begin{aligned} f(\mathbf{x}) &= \begin{cases} 1 & \text{if } \mathbf{w}_1^T \mathbf{x} \geq \mathbf{w}_2^T \mathbf{x} \\ 2 & \text{if } \mathbf{w}_2^T \mathbf{x} > \mathbf{w}_1^T \mathbf{x} \end{cases} \\ &= \operatorname{argmax}_{k \in \{1, 2\}} \mathbf{w}_k^T \mathbf{x} \end{aligned}$$

for any  $\mathbf{w}_1, \mathbf{w}_2$  s.t.  $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$

Think of  $\mathbf{w}_k^T \mathbf{x}$  as **a score for class  $k$** .

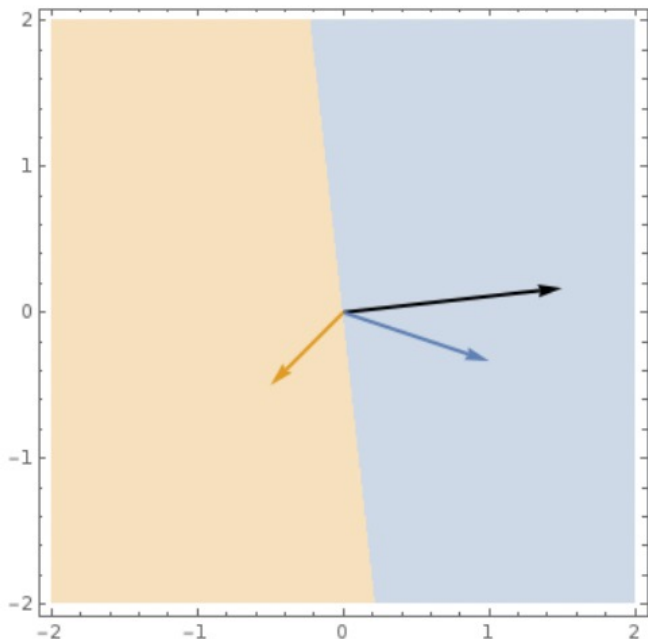
## Linear models: Binary to multiclass



$$w = \left(\frac{3}{2}, \frac{1}{6}\right)$$

- Blue class:  
 $\{x : w^T x \geq 0\}$
- Orange class:  
 $\{x : w^T x < 0\}$

## Linear models: Binary to multiclass



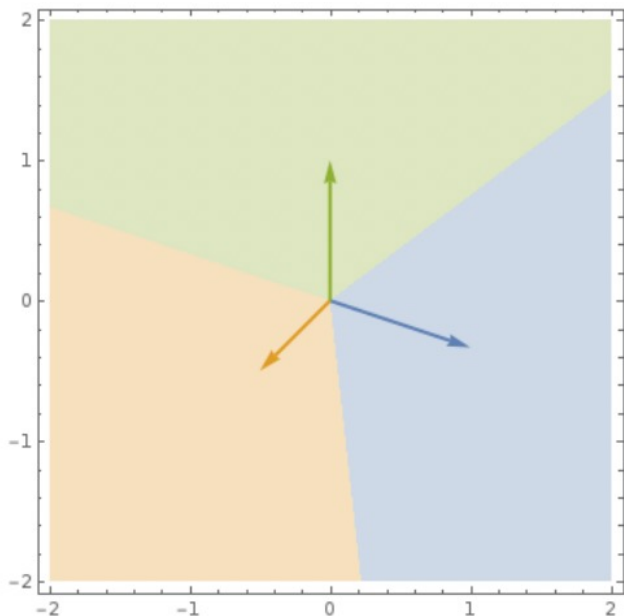
$$\mathbf{w} = \left(\frac{3}{2}, \frac{1}{6}\right) = \mathbf{w}_1 - \mathbf{w}_2$$

$$\mathbf{w}_1 = \left(1, -\frac{1}{3}\right)$$

$$\mathbf{w}_2 = \left(-\frac{1}{2}, -\frac{1}{2}\right)$$

- Blue class:  
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Orange class:  
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$

## Linear models: Binary to multiclass



$$\mathbf{w}_1 = (1, -\frac{1}{3})$$

$$\mathbf{w}_2 = (-\frac{1}{2}, -\frac{1}{2})$$

$$\mathbf{w}_3 = (0, 1)$$

- Blue class:  
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Orange class:  
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Green class:  
 $\{\mathbf{x} : 3 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$



## Function class: Linear models for multiclass classification

$$\begin{aligned}\mathcal{F} &= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^d \right\} \\ &= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} (\mathbf{W} \mathbf{x})_k \mid \mathbf{W} \in \mathbb{R}^{C \times d} \right\}\end{aligned}$$

Next, let's try to generalize the loss functions. Focus on the logistic loss today.

A hand-drawn diagram illustrating the dot product of a weight vector and an input vector. On the left, a square matrix is drawn with horizontal lines representing rows. The label  $\mathbf{W}$  is written in the center, and a vertical ellipsis is to its right. To the right of the matrix, the vectors  $\mathbf{w}_1^T$ ,  $\mathbf{w}_2^T$ , and  $\mathbf{w}_C^T$  are listed vertically, separated by vertical ellipses. In the middle, a vertical rectangle represents the input vector  $\mathbf{x}$ . To its right is an equals sign. Further right is another vertical rectangle, with an arrow pointing to it from the first rectangle, and the label  $\mathbf{w}_k$  next to it, representing the output of the dot product.

# Multinomial logistic regression: a probabilistic view

Observe: for binary logistic regression, with  $w = w_1 - w_2$ :

$$\Pr(y = 1 \mid x; w) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}} = \frac{e^{w_1^T x}}{e^{w_1^T x} + e^{w_2^T x}} \propto e^{w_1^T x}$$

Naturally, for multiclass:

$$\Pr(y = i \mid x; w) = \frac{e^{w_i^T x}}{e^{w_1^T x} + e^{w_2^T x}}$$

$w = w_1 - w_2$

$$\Pr(y = i \mid x; W) = \frac{e^{w_i^T x}}{\sum_{k \in [C]} e^{w_k^T x}} \propto e^{w_i^T x}$$

This is called the *softmax function*.

converts scores  $w_i^T x \rightarrow \Pr(y = i \mid x; w) = \frac{e^{w_i^T x}}{\sum_k e^{w_k^T x}}$

# Let's find the MLE

Maximize probability of seeing labels  $y_1, \dots, y_n$  given  $\mathbf{x}_1, \dots, \mathbf{x}_n$

$$P(\mathbf{W}) = \prod_{i=1}^n \Pr(y_i \mid \mathbf{x}_i; \mathbf{W}) = \prod_{i=1}^n \frac{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_i}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\mathbf{W}) = \sum_{i=1}^n \ln \left( \frac{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_i}}{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i}} \right) = \sum_{i=1}^n \ln \left( 1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i} \right)$$

This is the *multiclass logistic loss*. It is an upper-bound on the 0-1 misclassification loss:

$$\mathbb{I}[f(\mathbf{x}) \neq y] \leq \log_2 \left( 1 + \sum_{k \neq y} e^{(\mathbf{w}_k - \mathbf{w}_y)^T \mathbf{x}} \right) \quad \left( \begin{array}{l} \log_2(1 + e^x) \geq 1 \\ \text{if } x \geq 0 \end{array} \right)$$

When  $C = 2$ , multiclass logistic loss is the same as binary logistic loss (let's verify).

## Relating binary and multiclass logistic loss

$$F(W) = \sum_{i=1}^n \ln \left( 1 + \sum_{k \neq y_i} e^{(w_k - w_{y_i})^T x_i} \right)$$

Consider any  $i \in [n]$ ,

$$\text{For } y_i = 1, \quad \ln \left( 1 + e^{(w_2 - w_1)^T x} \right)$$

$$\text{For } y_i = 2, \quad \ln \left( 1 + e^{(w_1 - w_2)^T x} \right)$$

For  $w = w_1 - w_2$ , and transferring labels from  $\{1, 2\} \rightarrow \{1, -1\}$

$$F(W) = \sum_{i=1}^n \ln \left( 1 + e^{-y_i w^T x_i} \right)$$

## Next, **optimization**

Apply **SGD**: what is the gradient of

$$F(\mathbf{W}) = \ln \left( 1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i} \right) ?$$

It's a  $C \times d$  matrix. Let's focus on the  $k$ -th row:

If  $k \neq y_i$ : *row-vector* *row-vector*

$$\nabla_{\mathbf{w}_k^T} F(\mathbf{W}) = \frac{e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i}}{1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i}} \mathbf{x}_i^T = \frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i} + \sum_{k \neq y_i} e^{\mathbf{w}_k^T \mathbf{x}_i}} \mathbf{x}_i^T = \text{Pr}(y = k \mid \mathbf{x}_i; \mathbf{W}) \mathbf{x}_i^T$$

*row-vector*  $\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_i}$

else:

$$\nabla_{\mathbf{w}_k^T} F(\mathbf{W}) = \frac{- \left( \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i} \right)}{1 + \sum_{k \neq y_i} e^{(\mathbf{w}_k - \mathbf{w}_{y_i})^T \mathbf{x}_i}} \mathbf{x}_i^T = \frac{- \left( \sum_{k \neq y_i} e^{\mathbf{w}_k^T \mathbf{x}_i} \right)}{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i} + \sum_{k \neq y_i} e^{\mathbf{w}_k^T \mathbf{x}_i}} \mathbf{x}_i^T = (\text{Pr}(y = y_i \mid \mathbf{x}_i; \mathbf{W}) - 1) \mathbf{x}_i^T$$

# SGD for multinomial logistic regression

Initialize  $\mathbf{W} = \mathbf{0}$  (or randomly). Repeat:

1. pick  $i \in [n]$  uniformly at random
2. update the parameters

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \begin{pmatrix} \Pr(y = 1 \mid \mathbf{x}_i; \mathbf{W}) \\ \vdots \\ \Pr(y = y_i \mid \mathbf{x}_i; \mathbf{W}) - 1 \\ \vdots \\ \Pr(y = C \mid \mathbf{x}_i; \mathbf{W}) \end{pmatrix} \mathbf{x}_i^T \rightarrow \mathbb{R}^{1+d}$$

$\in \mathbb{R}^{C \times 1}$

Think about why the algorithm makes sense intuitively.

# Probabilities -> Prediction

Having learned  $\mathbf{W}$ , we can either

- make a *deterministic* prediction  $\operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x}$
- make a *randomized* prediction according to  $\Pr(y = k \mid \mathbf{x}; \mathbf{W}) \propto e^{\mathbf{w}_k^T \mathbf{x}}$

To minimize misclassification rate, make deterministic predictions.

## Beyond linear models

Suppose we have any model  $f$  (not necessary linear) which gives some score  $f_k(\mathbf{x})$  for the datapoint  $\mathbf{x}$  having the  $k$ -th label.

For linear model,  $f_k(\mathbf{x}) = \mathbf{w}_k^\top \mathbf{x}$

How can we convert this score to probabilities? Use the *softmax function*!

$$\tilde{f}_k(\mathbf{x}) = \Pr(y = k \mid \mathbf{x}; f) = \frac{e^{f_k(\mathbf{x})}}{\sum_{k' \in [C]} e^{f_{k'}(\mathbf{x})}} \propto e^{f_k(\mathbf{x})}$$

Once we have probability estimates, what is suitable loss function to train the model? Use the *log loss*. Also known as the *cross-entropy loss*.



## Log Loss/Cross-entropy loss: Binary case

Let's start with binary classification again. Consider a model which predicts  $\tilde{f}(x)$  as the probability of label being 1 for labelled datapoint  $(x, y)$ . The log loss is defined as,

$$\begin{aligned} \text{LogLoss} &= \mathbf{1}(y = 1) \ln \left( \frac{1}{\tilde{f}(x)} \right) + \mathbf{1}(y = -1) \ln \left( \frac{1}{1 - \tilde{f}(x)} \right) \\ &= -\mathbf{1}(y = 1) \ln(\tilde{f}(x)) - \mathbf{1}(y = -1) \ln((1 - \tilde{f}(x))). \end{aligned} \quad \left\{ \begin{array}{l} \text{Why?} \\ \text{If } y=1, \text{ want to} \\ \text{max } \ln(\tilde{f}(x)) \\ \Rightarrow \min \ln(1/\tilde{f}(x)) = -\ln(\tilde{f}(x)) \end{array} \right.$$

When the model is linear, this reduces to the logistic regression loss we defined before!

$$\text{Linear model: } \tilde{f}(x) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}, \quad 1 - \sigma(w^T x) = \sigma(-w^T x)$$

$$\begin{aligned} \text{Log loss} &= -\mathbf{1}(y=1) \ln((1 + e^{-w^T x})^{-1}) - \mathbf{1}(y=-1) \ln((1 + e^{w^T x})^{-1}) \\ &= \ln(1 + e^{-y w^T x}) \quad (\text{logistic regression loss}) \end{aligned}$$

## Log Loss/Cross-entropy loss: Multiclass case

This generalizes easily to the multiclass case. For datapoint  $(\mathbf{x}, y)$ , if  $\tilde{f}_k(\mathbf{x})$  is the predicted probability of label  $k$ ,

$$\begin{aligned}\text{LogLoss} &= \sum_{k=1}^C \mathbf{1}(y = k) \ln \left( \frac{1}{\tilde{f}_k(\mathbf{x})} \right) \\ &= - \sum_{k=1}^C \mathbf{1}(y = k) \ln \left( \tilde{f}_k(\mathbf{x}) \right).\end{aligned}$$

When the model is linear, this also reduces to the multiclass logistic regression loss we defined earlier today.

*Exercise!*

## Log Loss/Cross-entropy loss: Multiclass case

By combining the softmax and the log-loss, we have a general loss  $\ell(f(\mathbf{x}), y)$  which we can use to train a multi-class classification model which assigns scores  $f_k(\mathbf{x})$  to the  $k$ -th class. (These scores  $f_k(\mathbf{x})$  are sometimes referred to as logits).

$$\begin{aligned}\ell(f(\mathbf{x}), y) &= - \sum_{k=1}^C \mathbf{1}(y = k) \ln \left( \tilde{f}_k(\mathbf{x}) \right) \\ &= \ln \left( \frac{\sum_{k \in [C]} e^{f_k(\mathbf{x})}}{e^{f_y(\mathbf{x})}} \right) \\ &= \ln \left( 1 + \sum_{k \neq y} e^{f_k(\mathbf{x}) - f_y(\mathbf{x})} \right) .\end{aligned}$$

## Other techniques for multiclass classification

Cross-entropy is the most popular, but there are other *black-box techniques* to convert multiclass classification to binary classification.










- **one-versus-all** (one-versus-rest, one-against-all, etc.)
- **one-versus-one** (all-versus-all, etc.)
- **Error-Correcting Output Codes** (ECOC)
- **tree-based reduction**

# One-versus-all

Idea: train  $C$  binary classifiers to learn “**is class  $k$  or not?**” for each  $k$ .

Training: for each class  $k \in [C]$ ,

- relabel examples with class  $k$  as  $+1$ , and all others as  $-1$
- train a binary classifier  $h_k$  using this new dataset

					
$x_1$ 	$\Rightarrow$	$x_1$ —	$x_1$ +	$x_1$ —	$x_1$ —
$x_2$ 		$x_2$ —	$x_2$ —	$x_2$ +	$x_2$ —
$x_3$ 		$x_3$ —	$x_3$ —	$x_3$ —	$x_3$ +
$x_4$ 		$x_4$ —	$x_4$ +	$x_4$ —	$x_4$ —
$x_5$ 		$x_5$ +	$x_5$ —	$x_5$ —	$x_5$ —
		$\Downarrow$ $h_1$	$\Downarrow$ $h_2$	$\Downarrow$ $h_3$	$\Downarrow$ $h_4$

# One-versus-all

Idea: train  $C$  binary classifiers to learn “**is class  $k$  or not?**” for each  $k$ .

Prediction: for a new example  $\mathbf{x}$

- ask each  $h_k$ : **does this belong to class  $k$ ?** (i.e.  $h_k(\mathbf{x})$ )
- randomly pick among all  $k$ 's s.t.  $h_k(\mathbf{x}) = +1$ .

Issue: can make a mistake *as long as one of  $h_k$  errs*.

# One-versus-one

Idea: train  $\binom{C}{2}$  binary classifiers to learn “**is class  $k$  or  $k'$ ?**”.

Training: for each pair  $(k, k')$ ,

- relabel examples with class  $k$  as  $+1$  and examples with class  $k'$  as  $-1$
- *discard all other examples*
- train a binary classifier  $h_{(k,k')}$  using this new dataset

		■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■
		■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■
$x_1$ ■	⇒	$x_1$ —			$x_1$ —		$x_1$ —
$x_2$ ■			$x_2$ —	$x_2$ +			$x_2$ +
$x_3$ ■				$x_3$ —	$x_3$ +	$x_3$ —	
$x_4$ ■		$x_4$ —			$x_4$ —		$x_4$ —
$x_5$ ■		$x_5$ +	$x_5$ +			$x_5$ +	
		⇓	⇓	⇓	⇓	⇓	⇓
		$h_{(1,2)}$	$h_{(1,3)}$	$h_{(3,4)}$	$h_{(4,2)}$	$h_{(1,4)}$	$h_{(3,2)}$

# One-versus-one

Idea: train  $\binom{C}{2}$  binary classifiers to learn “**is class  $k$  or  $k'$ ?**”.

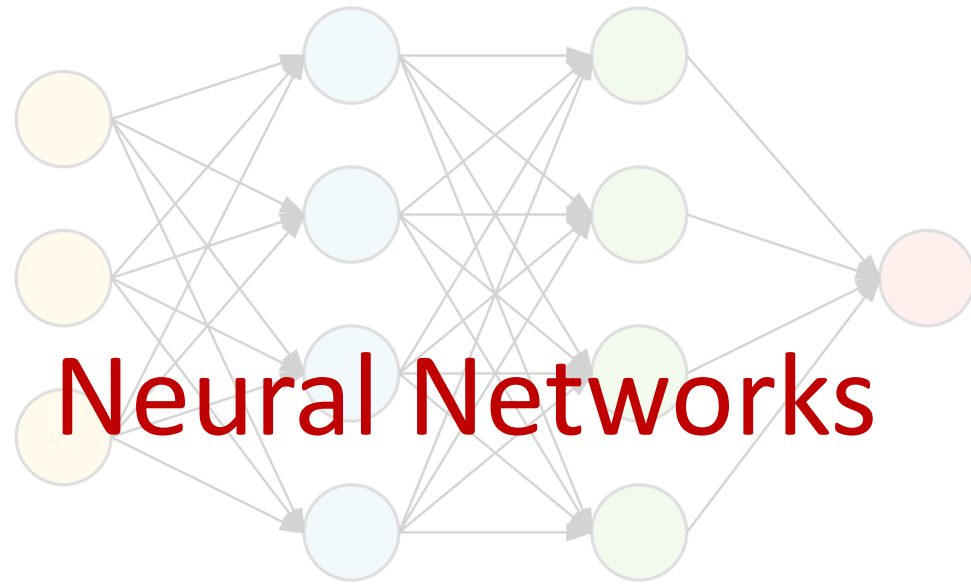
Prediction: for a new example  $x$

- ask each classifier  $h_{(k,k')}$  to **vote for either class  $k$  or  $k'$**
- predict the class with the most votes (break tie in some way)

**More robust** than one-versus-all, but *slower* in prediction.

Other techniques such as tree-based methods and error-correcting codes can achieve intermediate tradeoffs.





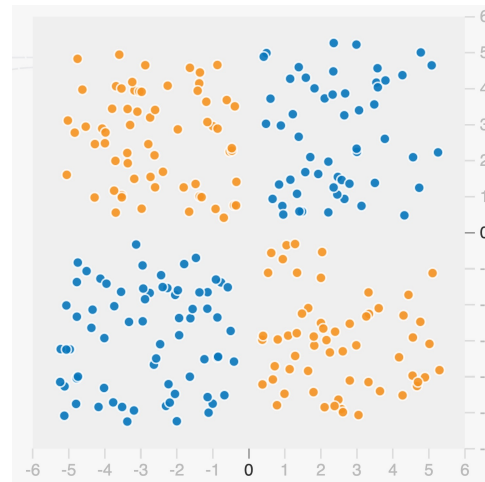
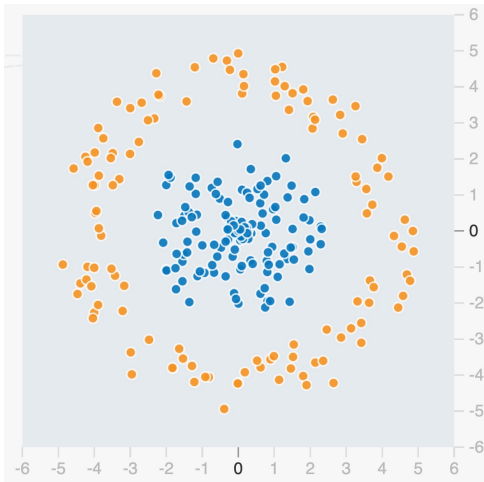
input layer

hidden layer 1

hidden layer 2

output layer

Linear -> Fixed non-linear -> **Learned non-linear map**



Linear models aren't always enough. As we discussed, we can use a nonlinear mapping and learn a linear model in the feature space:

$$\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^d \rightarrow \mathbf{z} \in \mathbb{R}^M$$

But what kind of nonlinear mapping  $\phi$  should be used?

*Can we just learn the nonlinear mapping itself?*

## Supervised learning in one slide

- Loss function:** What is the right loss function for the task?
- Representation:** What class of functions should we use?
- Optimization:** How can we efficiently solve the empirical risk minimization problem?
- Generalization:** Will the predictions of our model transfer gracefully to unseen examples?

*All related! And the fuel which powers everything is **data**.*

# Loss function

For model which makes predictions  $f(\mathbf{x})$  on labelled datapoint  $(\mathbf{x}, y)$ , we can use the following losses.

## Regression:

$$\ell(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2.$$

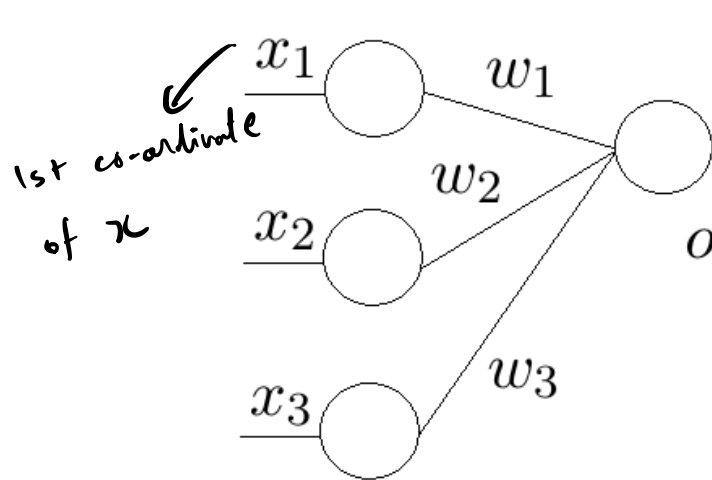
## Classification:

$$\ell(f(\mathbf{x}), y) = \ln \left( \frac{\sum_{k \in [C]} e^{f_k(\mathbf{x})}}{e^{f_y(\mathbf{x})}} \right) = \ln \left( 1 + \sum_{k \neq y} e^{f_k(\mathbf{x}) - f_y(\mathbf{x})} \right).$$

There maybe other, more suitable options for the problem at hand, but these are the most popular for supervised problems.

# Representation: Defining neural networks

Linear model as a one-layer neural network:



$$f(x) = w^T x$$
$$= \sum_{j=1}^3 w_j x_j$$

$$o = h(w^T x)$$

For a linear model,  $h(a) = a$ .

To create non-linearity, can use some nonlinear (differentiable) function:

- Rectified Linear Unit (**ReLU**):  $h(a) = \max\{0, a\}$
- Sigmoid function:  $h(a) = \frac{1}{1+e^{-a}}$
- Tanh:  $h(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$
- many more

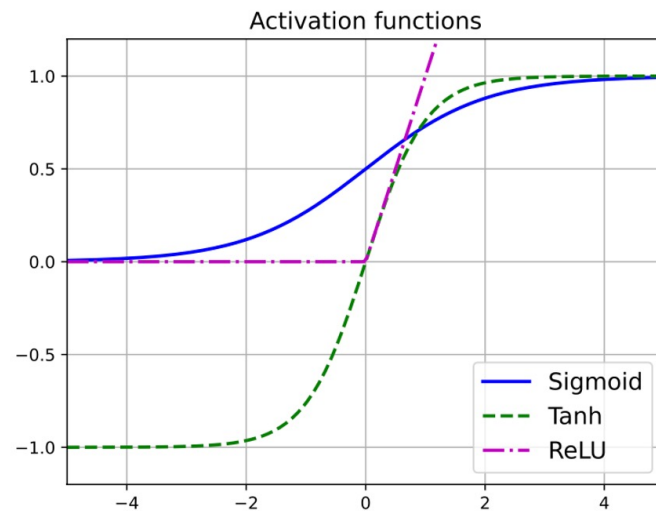
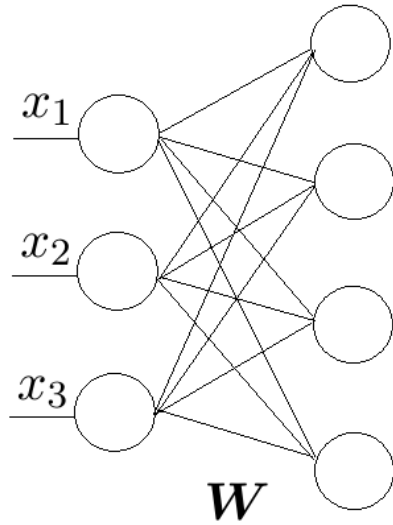
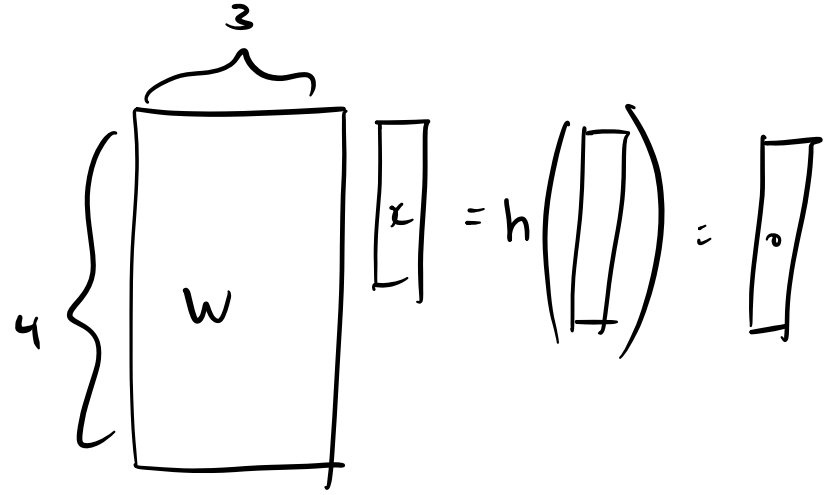


Figure 13.2 from PML

## Adding a layer



$$o = h(Wx)$$



$$W \in \mathbb{R}^{4 \times 3}, \mathbf{h} : \mathbb{R}^4 \rightarrow \mathbb{R}^4 \text{ so } \mathbf{h}(\mathbf{a}) = (h_1(a_1), h_2(a_2), h_3(a_3), h_4(a_4))$$

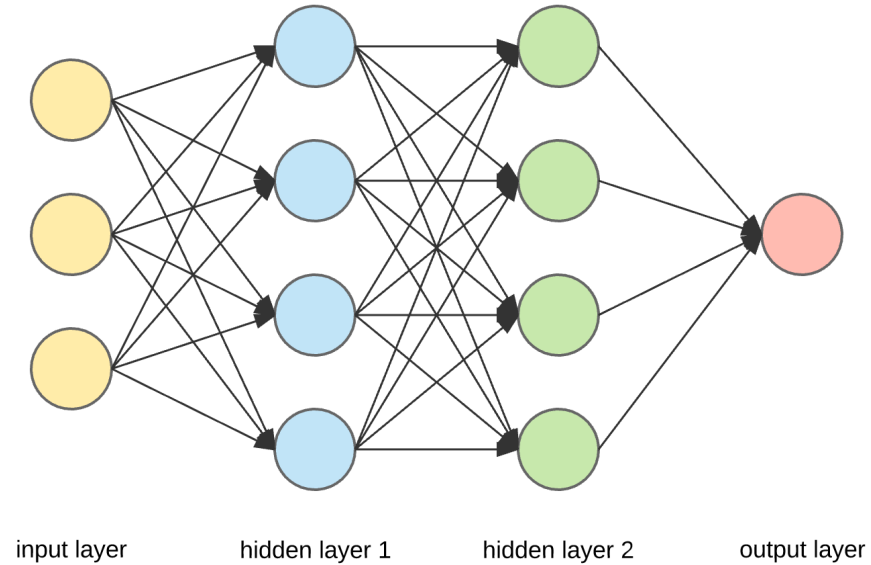
Can think of this as a nonlinear mapping:  $\phi(x) = h(Wx)$

$$\underbrace{h(\mathbf{a})}_{\text{applying function to each entry}} = (h(a_1), h(a_2), h(a_3), h(a_4))$$

# Putting things together: a neural network

We now have a network:

- each node is called a **neuron**
- $h$  is called the **activation function**
  - can use  $h(a) = 1$  for one neuron in each layer to incorporate bias term
  - output neuron can use  $h(a) = a$
- #layers refers to #hidden\_layers (plus 1 or 2 for input/output layers)
- **deep** neural nets can have many layers and *millions* of parameters
- this is a **feedforward, fully connected** neural net, there are many variants (convolutional nets, residual nets, recurrent nets, etc.)



# Neural network: Definition

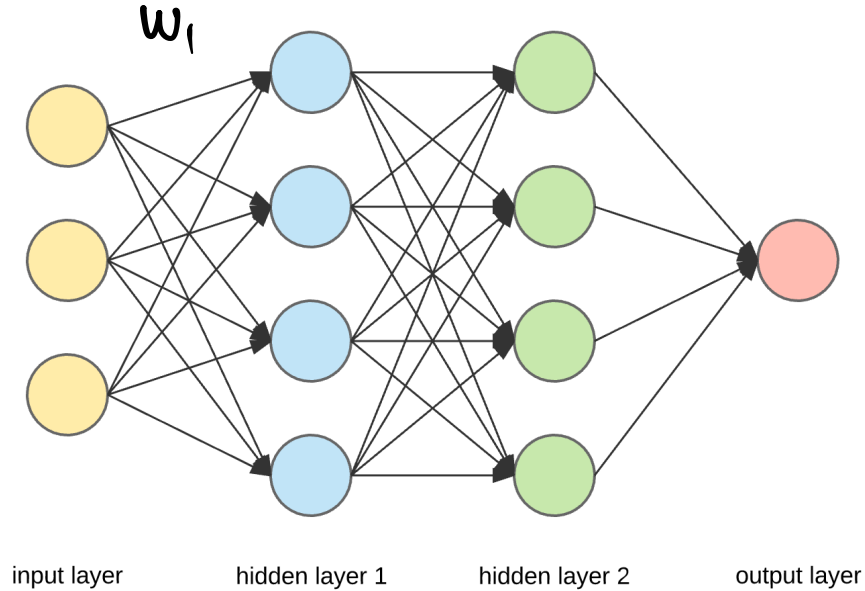
$d_\ell$ : # neurons in layer  $\ell$

An L-layer neural net can be written as

$$f(x) = h_L(W_L h_{L-1}(W_{L-1} \cdots h_1(W_1 x))).$$

Define

- input: output of prev layer  
output: input of next layer
- $W_\ell \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$  is the weights between layer  $\ell - 1$  and  $\ell$
  - $d_0 = d, d_1, \dots, d_L$  are numbers of neurons at each layer
  - $a_\ell \in \mathbb{R}^{d_\ell}$  is input to layer  $\ell$
  - $o_\ell \in \mathbb{R}^{d_\ell}$  is output of layer  $\ell$
  - $h_\ell : \mathbb{R}^{d_\ell} \rightarrow \mathbb{R}^{d_\ell}$  is activation functions at layer  $\ell$



Now, for a given input  $x$ , we have recursive relations:

$$o_0 = x, a_\ell = W_\ell o_{\ell-1}, o_\ell = h_\ell(a_\ell), \quad (\ell = 1, \dots, L).$$



# Optimization

Our optimization problem is to minimize,

$$F(\mathbf{W}_1, \dots, \mathbf{W}_L) = \frac{1}{n} \sum_{i=1}^n F_i(\mathbf{W}_1, \dots, \mathbf{W}_L)$$

where

$$F_i(\mathbf{W}_1, \dots, \mathbf{W}_L) = \begin{cases} \|\mathbf{f}(\mathbf{x}_i) - \mathbf{y}_i\|_2^2 & \text{for regression} \\ \ln \left( 1 + \sum_{k \neq y_i} e^{f_k(\mathbf{x}_i) - f_{y_i}(\mathbf{x}_i)} \right) & \text{for classification} \end{cases}$$

How to solve this? Apply **SGD**!

To compute the gradient efficiently, we use *backpropagation*. More on this soon.

# Generalization

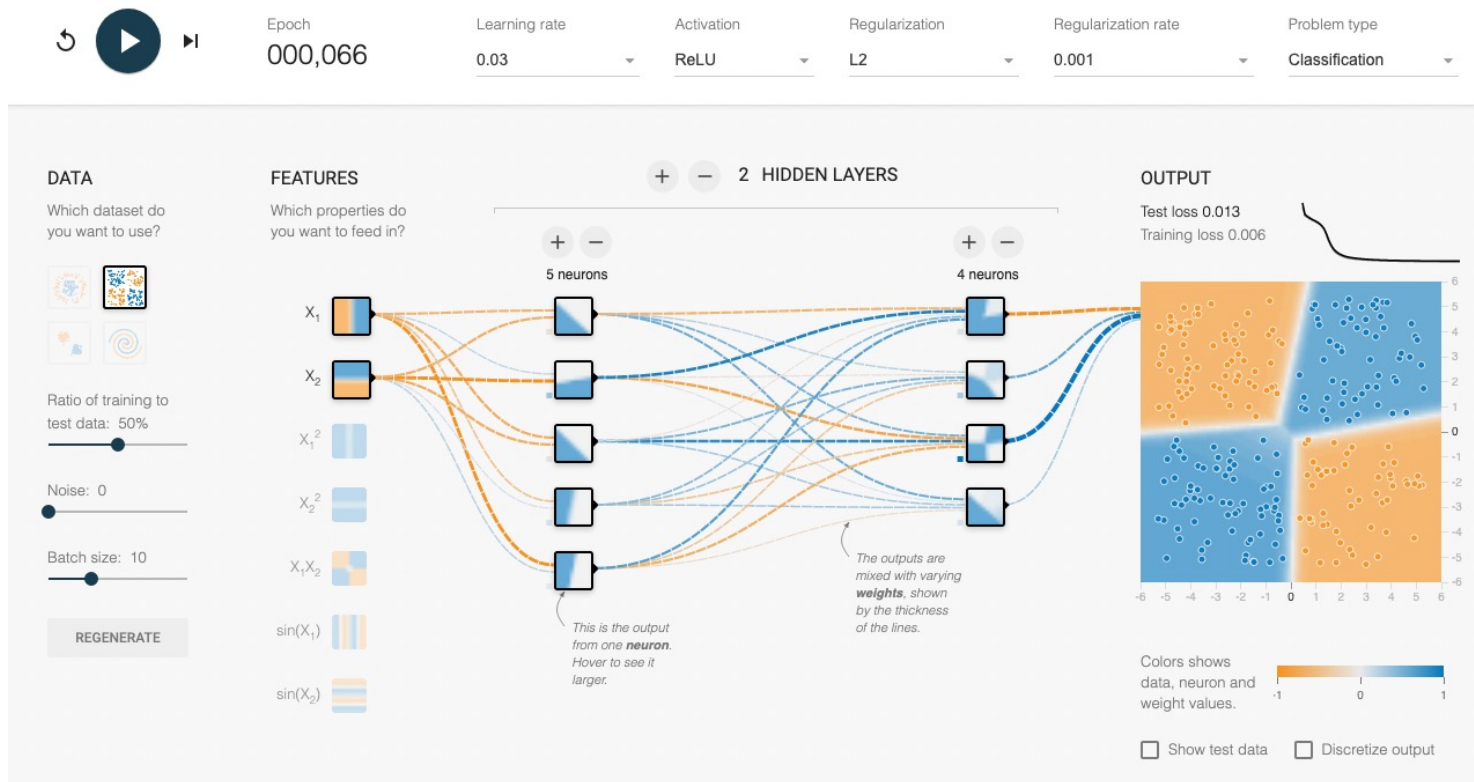
Overfitting is a concern for such a complex model, but there are ways to handle it.

For example, we can add  $\ell_2$  regularization.

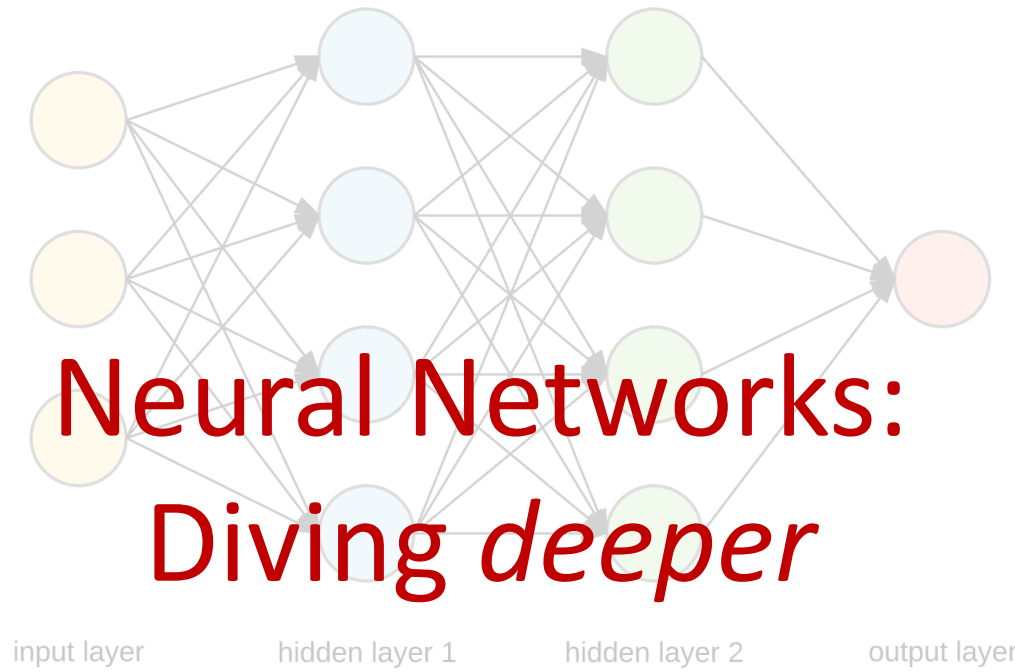
$\ell_2$  regularization: minimize

$$G(\mathbf{W}_1, \dots, \mathbf{W}_L) = F(\mathbf{W}_1, \dots, \mathbf{W}_L) + \lambda \sum_{\substack{\text{all weights } w \\ \text{in network}}} w^2$$

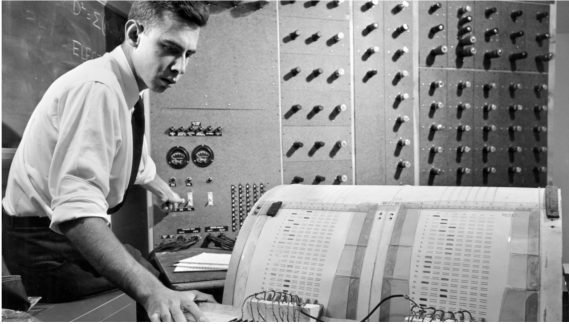
# Demo



<http://playground.tensorflow.org/>

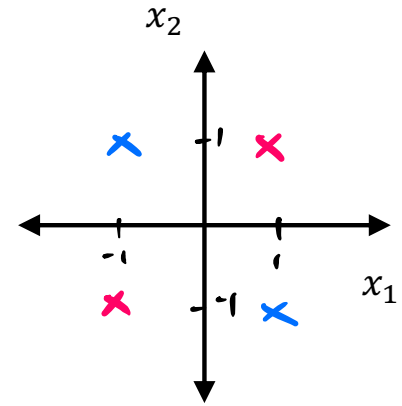
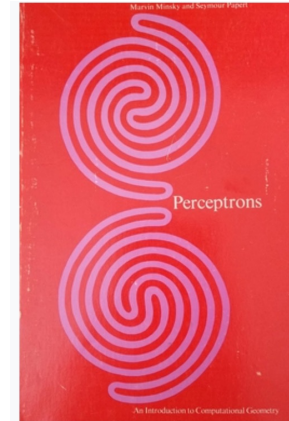


# History: A multilayer perceptron



Frank Rosenblatt invented the perceptron in 1958.

Minsky and Papert wrote a book on it in 1969. The book said how perceptrons were unable to solve the XOR problem, which made the field skeptical of neural networks and drove attention to symbolic AI systems.



Rumelhart, Hinton and Williams (re)discovered backpropagation in 1986 which eventually led to deeper networks.

# Representation: **Very powerful function class!**

**Universal approximation theorem** (Cybenko, 89; Hornik, 91):

*A feedforward neural net with a single hidden layer can approximate any continuous function.*

It might need a huge number of neurons though, and *depth helps!*

Choosing the network architecture is important.

- for feedforward network, need to decide number of hidden layers, number of neurons at each layer, activation functions, etc.

Designing the architecture can be complicated, though various standard choices exist.

# Optimization: Computing gradients efficiently using **Backprop**

To run SGD, need gradients of  $F_i(\mathbf{W}_1, \dots, \mathbf{W}_L)$  with respect to all the weights in all the layers. How do we get the gradient?

Here's a naive way to compute gradients. For some function  $F(w)$  of a univariate parameter  $w$ ,

$$\frac{dF(w)}{dw} = \lim_{\epsilon \rightarrow 0} \frac{F(w + \epsilon) - F(w - \epsilon)}{2\epsilon}$$

If our network has  $m$  weights, this requires  $O(m)$  model evaluation.

(this is still useful for "gradient checking")

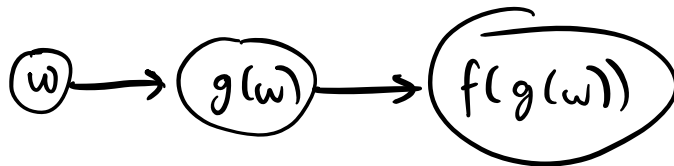
# Backprop

**Backpropagation:** A very efficient way to compute gradients of neural networks using an application of the chain rule (similar to dynamic programming).

*Chain rule:*

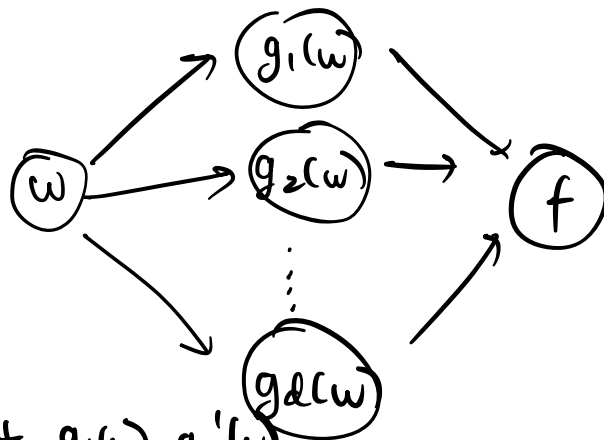
- for a composite function  $f(g(w))$

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$



- for a composite function  $f(g_1(w), \dots, g_d(w))$

$$\frac{\partial f}{\partial w} = \sum_{i=1}^d \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial w}$$

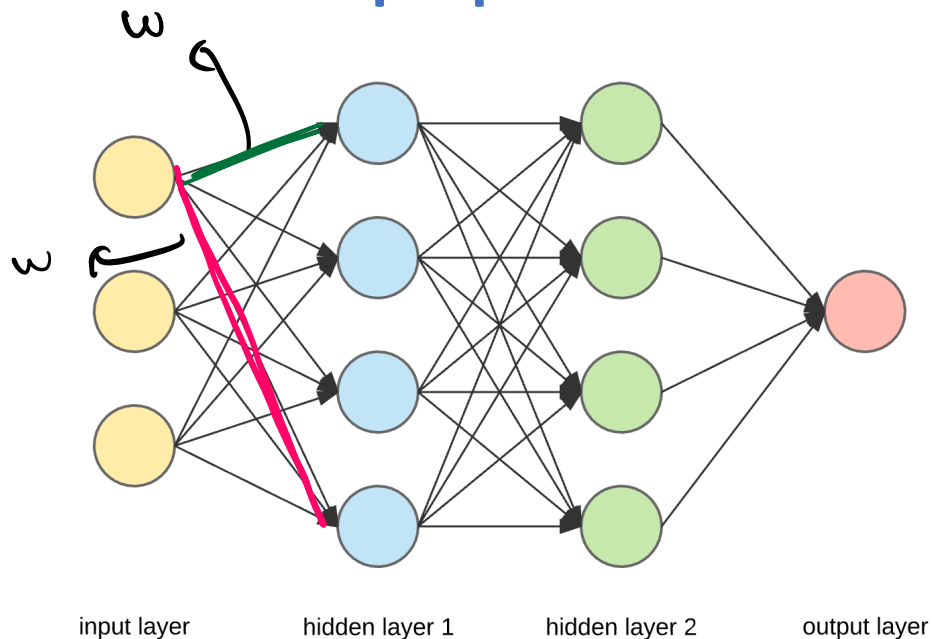


the simplest example  $f(g_1(w), g_2(w)) = g_1(w)g_2(w)$

$$\frac{df}{dw} = \frac{\partial f}{\partial g_1} \cdot \frac{dg_1}{dw} + \frac{\partial f}{\partial g_2} \cdot \frac{dg_2}{dw} = g_2(w) \cdot g_1'(w) + g_1(w) \cdot g_2'(w)$$



## Backprop: Intuition



Naive: apply chain rule for each weight

Backprop: reuse computation by storing gradients wrt  
input to each layer (ae)

# Backprop: Derivation

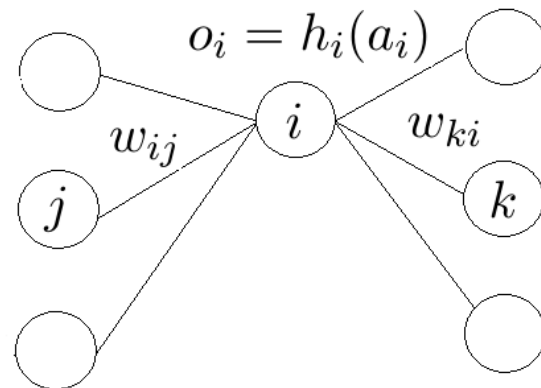
Drop the subscript  $\ell$  for layer for simplicity. For this derivation, refer to the loss function on the  $i$ -th datapoint as  $F$  (instead of  $F_i$ ) for convenience.

Find the **derivative of  $F$  w.r.t. to  $w_{ij}$**

$$\frac{\partial F}{\partial w_{ij}} = \frac{\partial F}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial F}{\partial a_i} \frac{\partial (w_{ij} o_j)}{\partial w_{ij}} = \frac{\partial F}{\partial a_i} o_j$$

$$\left( \frac{\partial F}{\partial a_i} \right) = \frac{\partial F}{\partial o_i} \frac{\partial o_i}{\partial a_i} = \left( \sum_k \frac{\partial F}{\partial a_k} \frac{\partial a_k}{\partial o_i} \right) h'_i(a_i) = \left( \sum_k \frac{\partial F}{\partial a_k} w_{ki} \right) h'_i(a_i)$$

key quantity



$$a_i = \sum_{\text{all neurons } j} w_{ij} o_j$$

# Backprop: Derivation

Adding the subscript for layer:

$$\frac{\partial F}{\partial w_{\ell,ij}} = \frac{\partial F}{\partial a_{\ell,i}} o_{\ell-1,j}$$
$$\frac{\partial F}{\partial a_{\ell,i}} = \left( \sum_k \frac{\partial F}{\partial a_{\ell+1,k}} w_{\ell+1,ki} \right) h'_{\ell,i}(a_{\ell,i})$$

For the last layer, for square loss

$$\frac{\partial F}{\partial a_{L,i}} = \frac{\partial (h_{L,i}(a_{L,i}) - y_m)^2}{\partial a_{L,i}} = 2(h_{L,i}(a_{L,i}) - y_m) h'_{L,i}(a_{L,i})$$

**Exercise:** try to do it for logistic loss yourself.

