

Advanced School in Artificial Intelligence

Introduction to Neural Networks

Ing. Zese Riccardo
riccardo.zese@unife.it

Progetto di alta formazione in ambito tecnologico economico e culturale per una regione della conoscenza europea e attrattiva approvato e cofinanziato dalla Regione Emilia-Romagna con deliberazione di Giunta regionale n. 1625/2021



**Università
degli Studi
di Ferrara**

Outline

- Introduction to Python
- Introduction to Neural Networks
- Convolutional NN
- Recurrent NN
- Autoencoders and self supervised learning

Sources

- These slides are taken from:

- Intel Nervana AI Academy Instructional Content

Intel Legal Notices and Disclaimers

This presentation is for informational purposes only. INTEL MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at intel.com.

This sample source code is released under the [Intel Sample Source Code License Agreement](#).

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

*Other names and brands may be claimed as the property of others.

Copyright © 2017, Intel Corporation. All rights reserved.

- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep learning. MIT press, 2016.

https://www.deeplearningbook.org/lecture_slides.html

- Chapter “Neural Networks” of “A Course in Machine Learning” by Hal Daume III.

<http://ciml.info/>

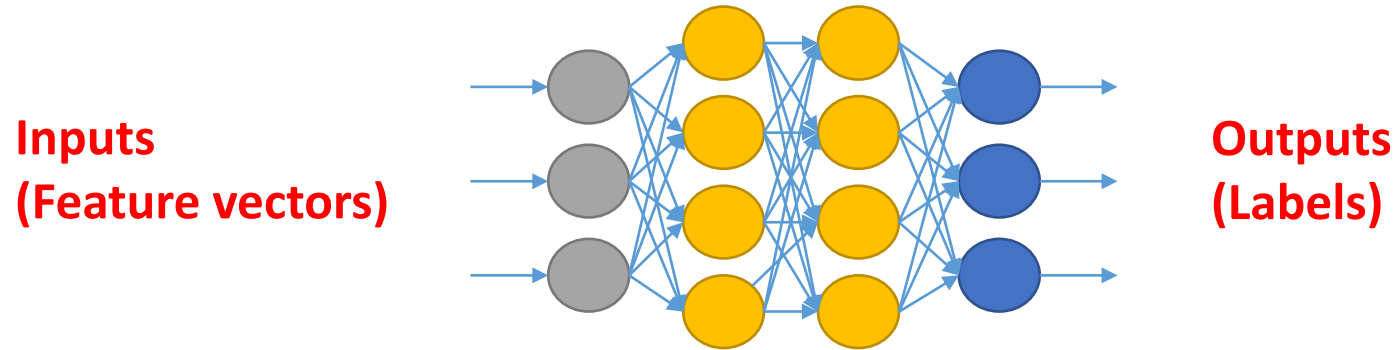
- Some parts from “CS231n: Convolutional Neural Networks for Visual Recognition”, Stanford University


<http://cs231n.stanford.edu/>

Outline

- Introduction to Python
- Introduction to Neural Networks
- Convolutional NN
- Recurrent NN
- Autoencoders and self supervised learning

How to train a Neural Net?



1. Put in Training inputs, get the output: **Forward Propagation**
2. Compare output to correct answers: Look at loss function J
3. Adjust and repeat! 
4. **Backpropagation** (also called **backprop**) tells us how to make a single adjustment using calculus.

Backpropagation

- With the term back-propagation we are not indicating the whole learning algorithm.
- It refers only to the method for computing the gradient, while another algorithm, such as stochastic gradient descent, is used to perform learning using this gradient.

How have we trained before?

- **Gradient Descent.**



For most classical ML algorithms, the training happens here.

1. Make prediction
2. Calculate Loss
3. Calculate gradient of the loss function w.r.t. parameters
4. Update parameters by taking a step in the opposite direction
5. Iterate

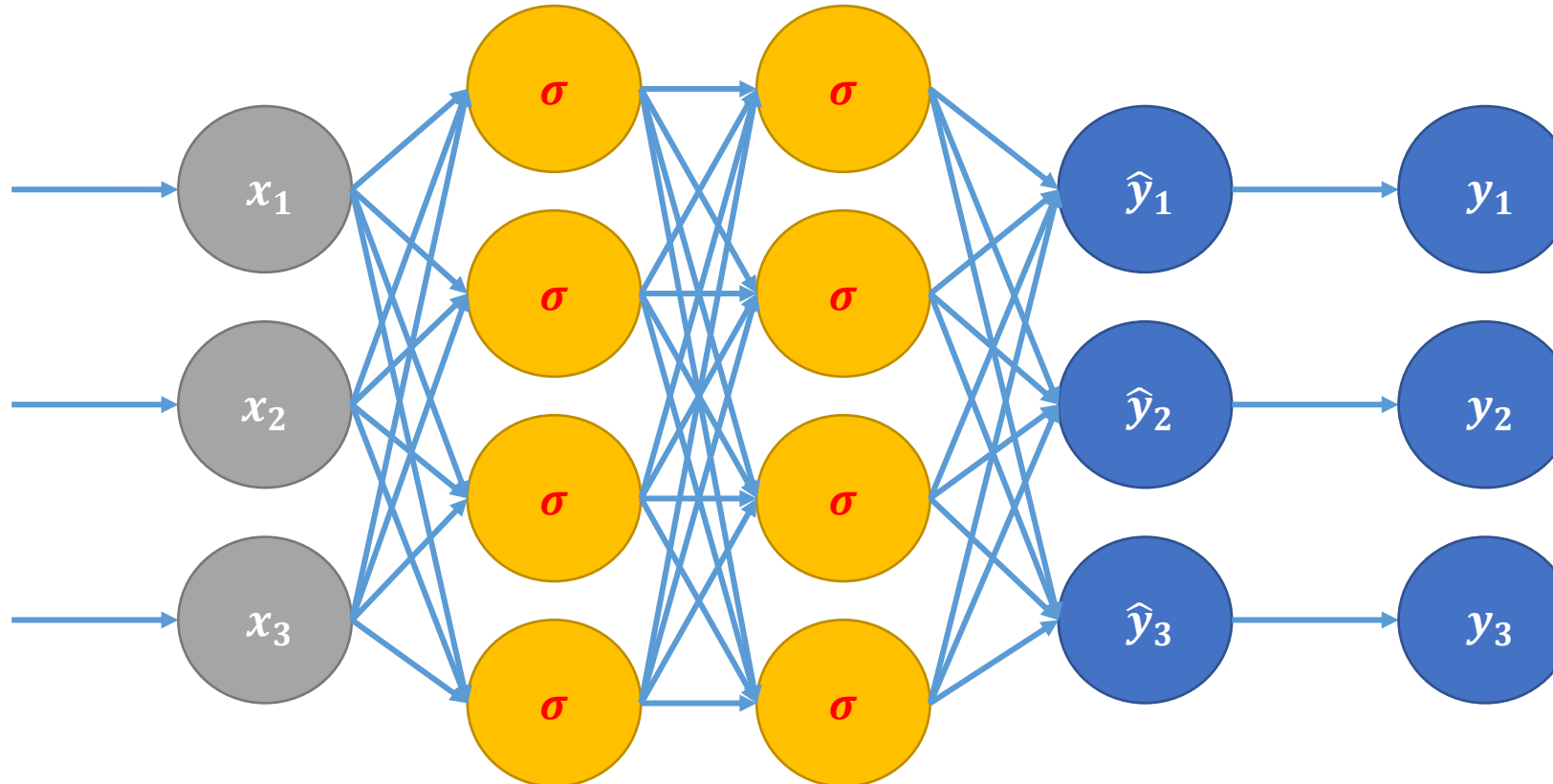
How have we trained before?

- **Gradient Descent.**

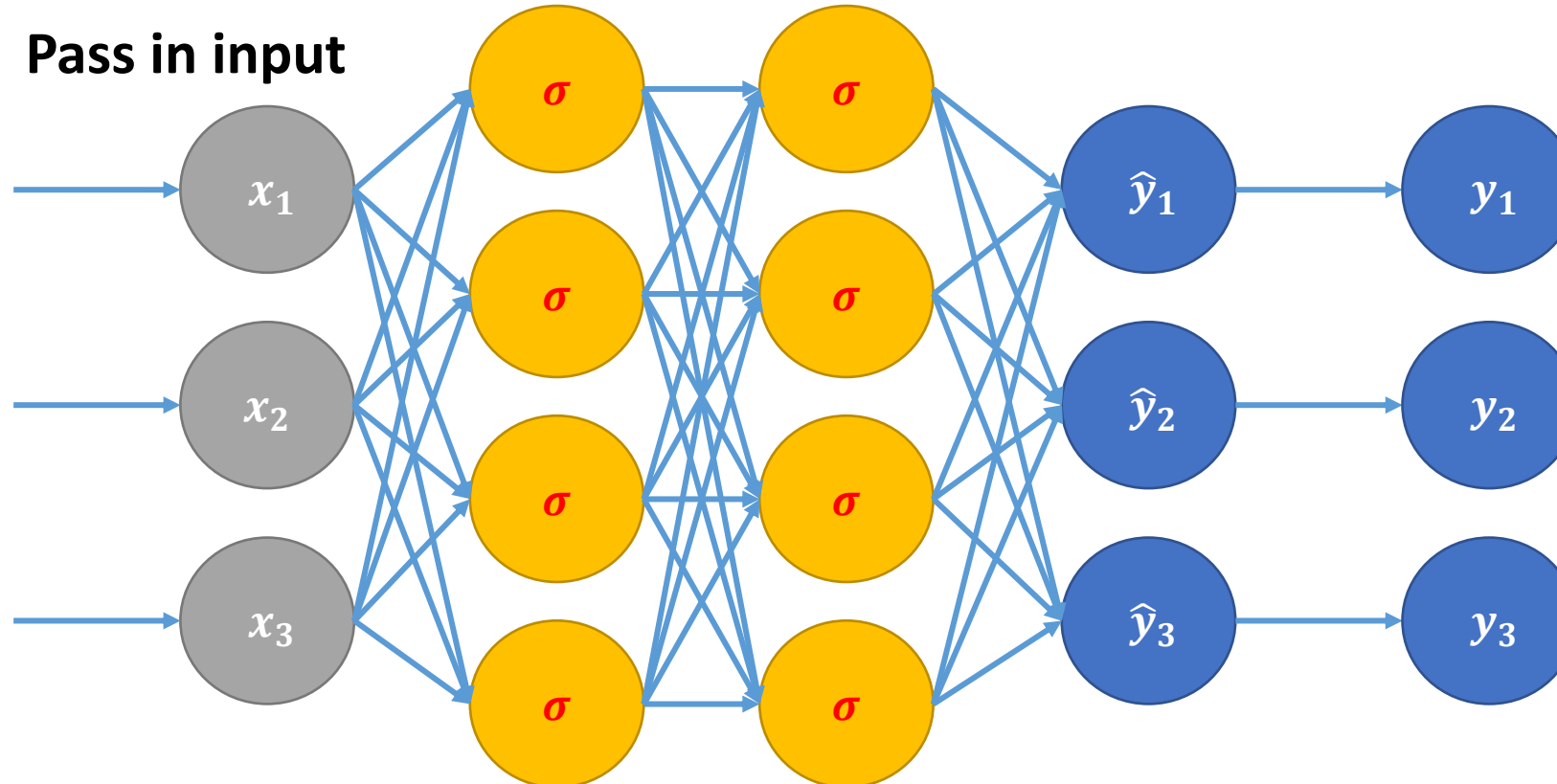
For most classical ML algorithms, the training happens here.

1. Make prediction
2. Calculate Loss
3. Calculate gradient of the loss function w.r.t. parameters 
4. Update parameters by taking a step in the opposite direction
5. Iterate 

Feedforward Neural Network



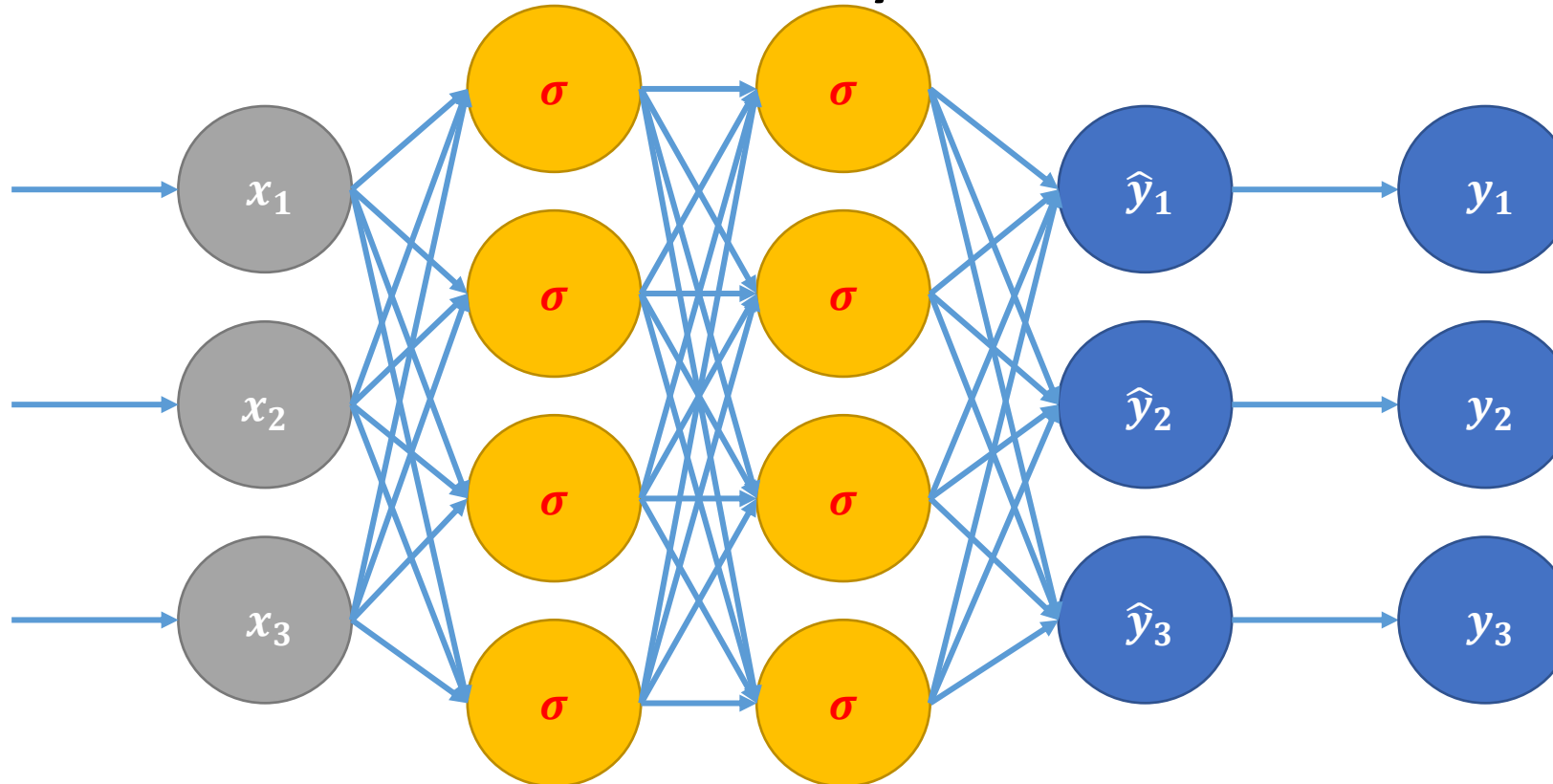
Feedforward Neural Network



Inputs are rows of our training data

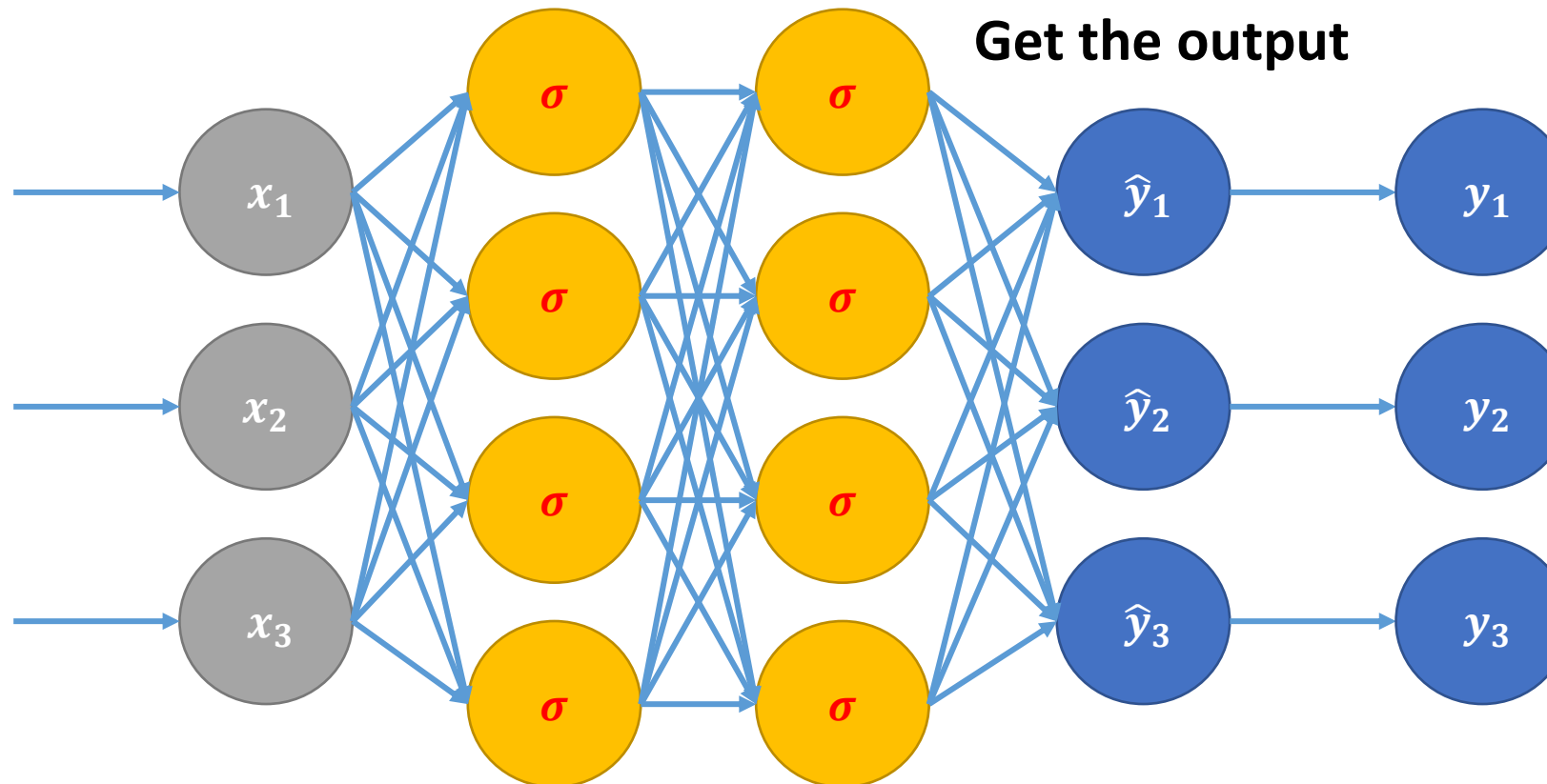
Feedforward Neural Network

Calculate each layer



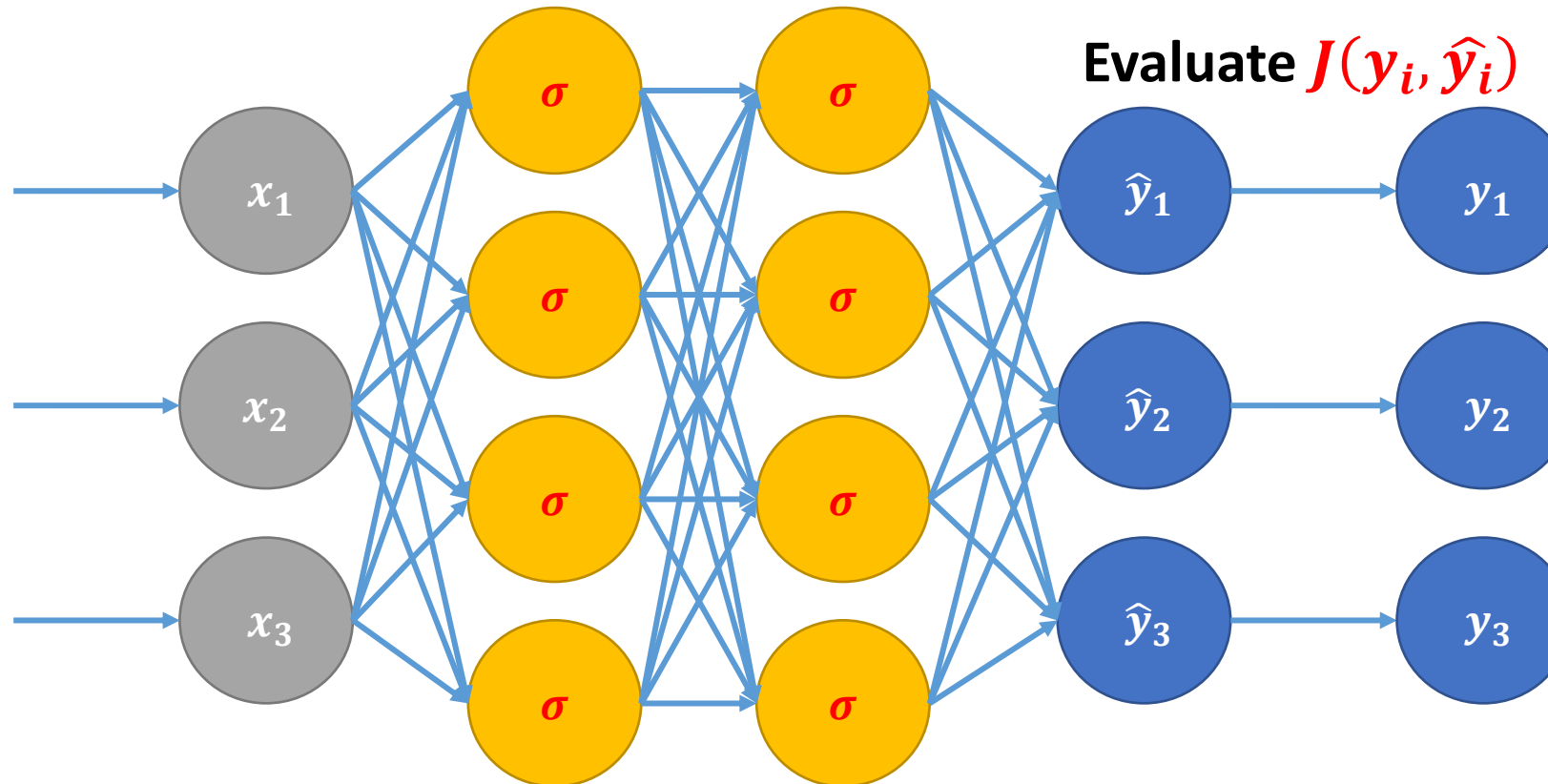
Perform the matrix multiplications and activation functions in order to calculate each layer.

Feedforward Neural Network



Output is the prediction made by the net.

Feedforward Neural Network



Compare the predictions to the known ground truths.
Specifically, calculate the loss function $J(y_i, \hat{y}_i)$.

How have we trained before?

- **Gradient Descent.**

For most classical ML algorithms, the training happens here.

1. Make prediction
2. Calculate Loss
3. Calculate gradient of the loss function w.r.t. parameters
4. Update parameters by taking a step in the opposite direction
5. Iterate

How to calculate gradient

- **Chain rule:** used to compute the derivatives of functions formed by composing other functions whose derivatives are known.
- Back-propagation is an algorithm that computes the chain rule, with a specific order of operations that is highly efficient.

Backpropagation

- You can summarize backpropagation as:

Backpropagation = gradient descent + chain rule

Backpropagation

- Example: consider a two layers network whose output is

$$\hat{y} = \sum_i w_{2,i} f(w_{1,i} x)$$

where W_1 is a matrix containing weights of the first layer and W_2 a vector containing the weights of the second layer. f is the activation function of the first layer, which is usually non-linear (ReLU, tanh,...).

Backpropagation

- The overall objective is to minimize the loss function by choosing values for the weights:

$$\min_{W_1, W_2} \sum_n \frac{1}{2} \left(y_n - \sum_i w_{2,i} f(w_{1,i} x_n) \right)^2$$

- The easy case is to differentiate it with respect to the weights of the output unit W_2

$$-y_n + \sum_i w_{2,i} f(w_{1,i} x_n)$$

- From this perspective, it is just a linear model, attempting to minimize squared error \rightarrow the input here is $f(W_1 x)$, not x .

Backpropagation

- Therefore, the gradient is

$$\nabla_{W_2} = - \sum_n \left(-y_n + \sum_i w_{2,i} f(w_{1,i} x_n) \right) f(W_1 x_n)$$

- This is exactly like the linear case.
- It shows how the output weights have to change to make the prediction better.
 - It is easy to measure how their changes affect the output.

Backpropagation

- When we move on the weights of the first layer the problem becomes more complicated.
- These weights usually are not trying to produce certain values (e.g., 0 or 1) but they are trying to produce values able to activate output units to return certain values.
 - So the change they want to make depends crucially on how the output layer interprets them.

How to train a NN?

- How could we change the weights to make our Loss Function smaller?
- Think of the neural net as a function $f: x \rightarrow y$
- f is a complex computation involving many weights W_k
- Given the structure, the weights “define” the function f (and therefore define our model)
- Loss Function is $J(y, f(x))$

How to train a NN?


- How could we change the weights to make our Loss Function smaller?
- Think of the neural net as a function $f: x \rightarrow y$
- f is a complex computation involving many weights W_k
- Given the structure, the weights “define” the model (our model)
- Loss Function is $J(y, f(x))$

The goal is to change the weights to make the loss function smaller.

How to train a NN?

- Get $\frac{\partial J}{\partial W_k}$ for every weight in the network.
- This tells us what direction to adjust each W_k if we want to lower our loss function.
- Make an adjustment and repeat!

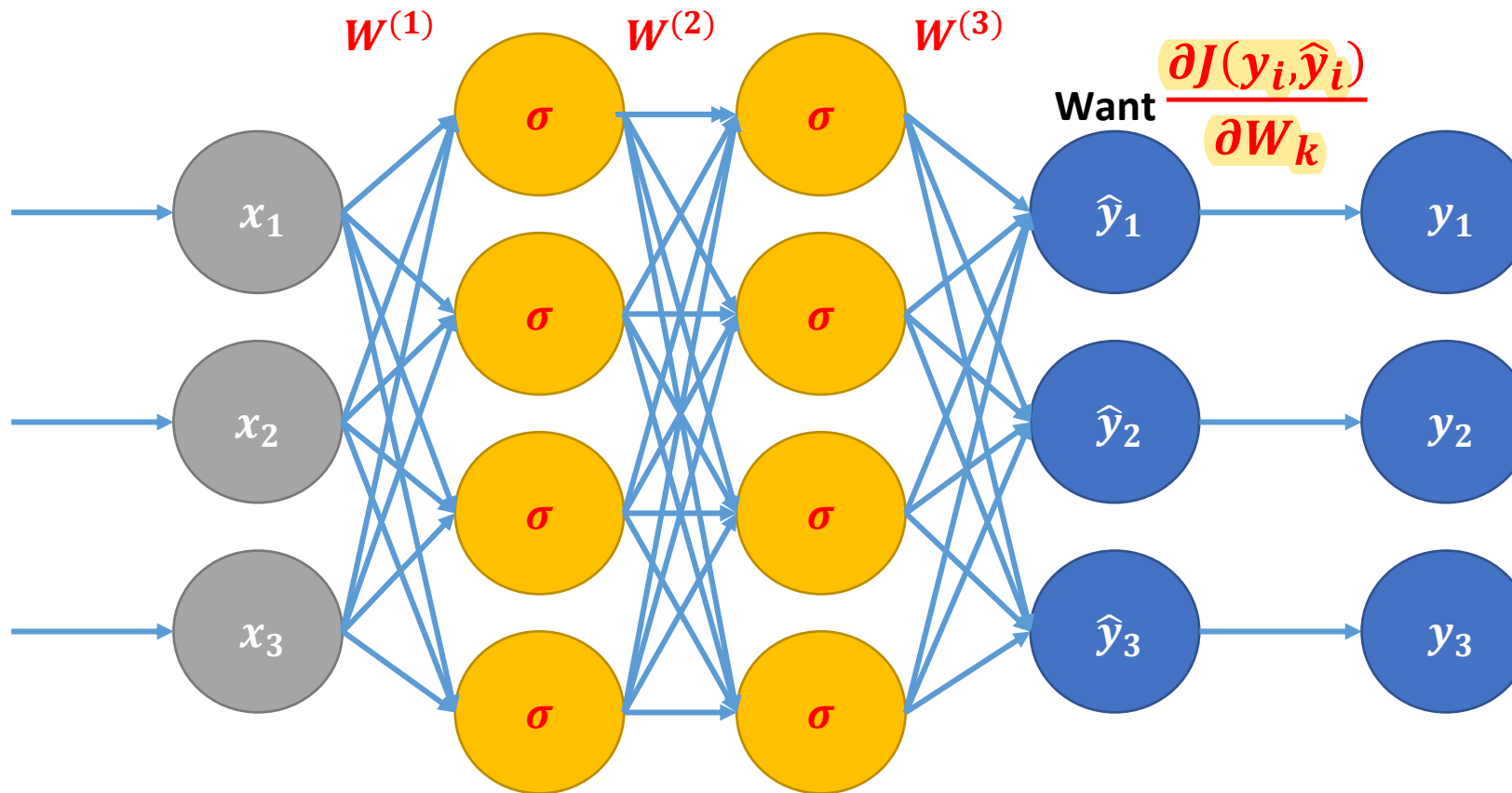
How to train a NN?

- Get $\frac{\partial J}{\partial W_k}$ for every weight in the network. 
- This tells us what direction to adjust each W_k if we want to lower our loss function.
- Make an adjustment and repeat!

With this setting, from an abstract mathematical standpoint, there is no difference between this and classical ML gradient descent.

It's just that the function involved is much more complicated, and computation of gradients is mathematically and computationally more challenging.

Feedforward Neural Network



So we want to be able to compute the partial derivative of the loss function w.r.t. the weights W_k

An Example...

Consider Cross-Entropy (Log-loss)

$$\frac{\partial J}{\partial W^{(3)}} = (\hat{y} - y) \cdot a^{(3)} \quad \text{💬}$$

$$\frac{\partial J}{\partial W^{(2)}} = (\hat{y} - y) \cdot W^{(3)} \cdot \sigma'(z^{(3)}) \cdot a^{(2)}$$

$$\frac{\partial J}{\partial W^{(1)}} = (\hat{y} - y) \cdot W^{(3)} \cdot \sigma'(z^{(3)}) \cdot W^{(2)} \cdot \sigma'(z^{(2)}) \cdot x$$

- Recall that: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Though they appear complex, above are easy to compute!

An Example...

$$\frac{\partial J}{\partial W^{(3)}} = (\hat{y} - y) \cdot a^{(3)}$$

$$\frac{\partial J}{\partial W^{(2)}} = (\hat{y} - y) \cdot W^{(3)} \cdot \sigma'(z^{(3)}) \cdot a^{(2)}$$

$$\frac{\partial J}{\partial W^{(1)}} = (\hat{y} - y) \cdot W^{(3)} \cdot \sigma'(z^{(3)}) \cdot W^{(2)} \cdot \sigma'(z^{(2)}) \cdot x$$

- Recall that: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Though they appear complex, above are easy to compute!

Early layers reuse computation from later layers → **BACK** propagation.
E.g., the gradient of $W^{(1)}$ uses the gradient of $W^{(2)}$



An Example...

$$\frac{\partial J}{\partial W^{(3)}} = (\hat{y} - y) \cdot a^{(3)}$$

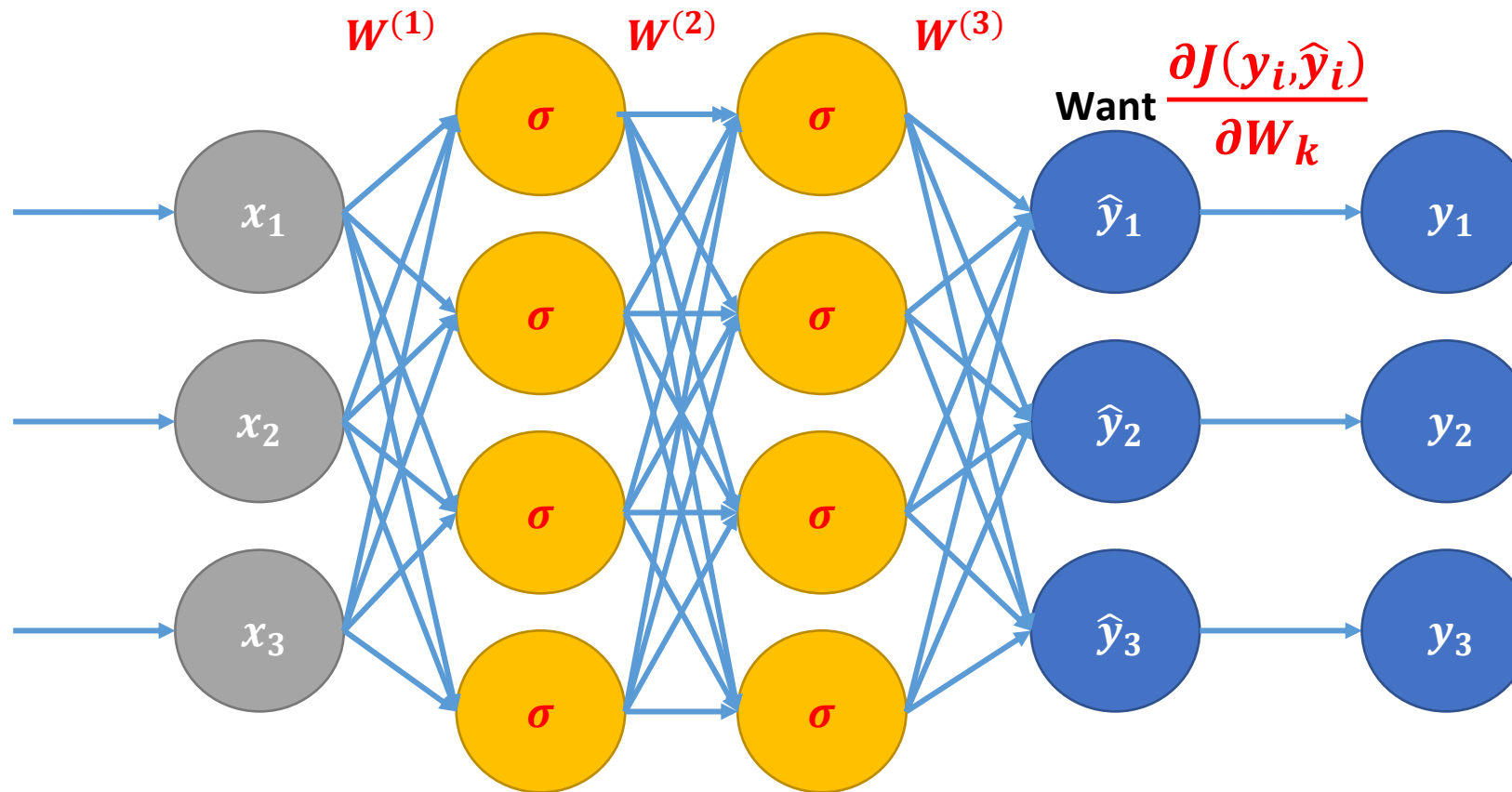
$$\frac{\partial J}{\partial W^{(2)}} = (\hat{y} - y) \cdot W^{(3)} \cdot \sigma'(z^{(3)}) \cdot a^{(2)}$$

$$\frac{\partial J}{\partial W^{(1)}} = (\hat{y} - y) \cdot W^{(3)} \cdot \sigma'(z^{(3)}) \cdot W^{(2)} \cdot \sigma'(z^{(2)}) \cdot x$$

Early layers have more terms \rightarrow smaller numbers \rightarrow vanishing gradient

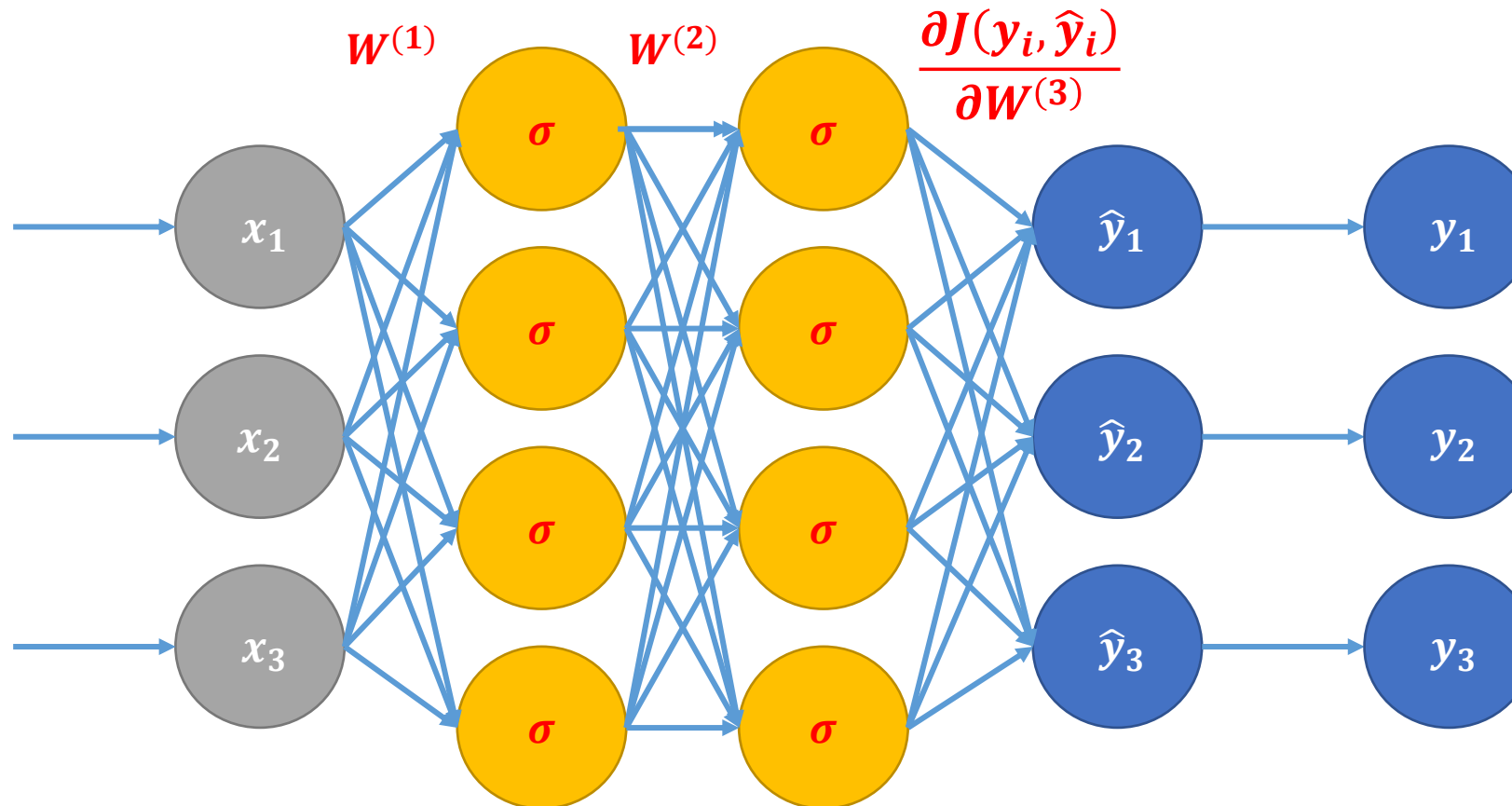
- Recall that: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Though they appear complex, above are easy to compute!

Backpropagation



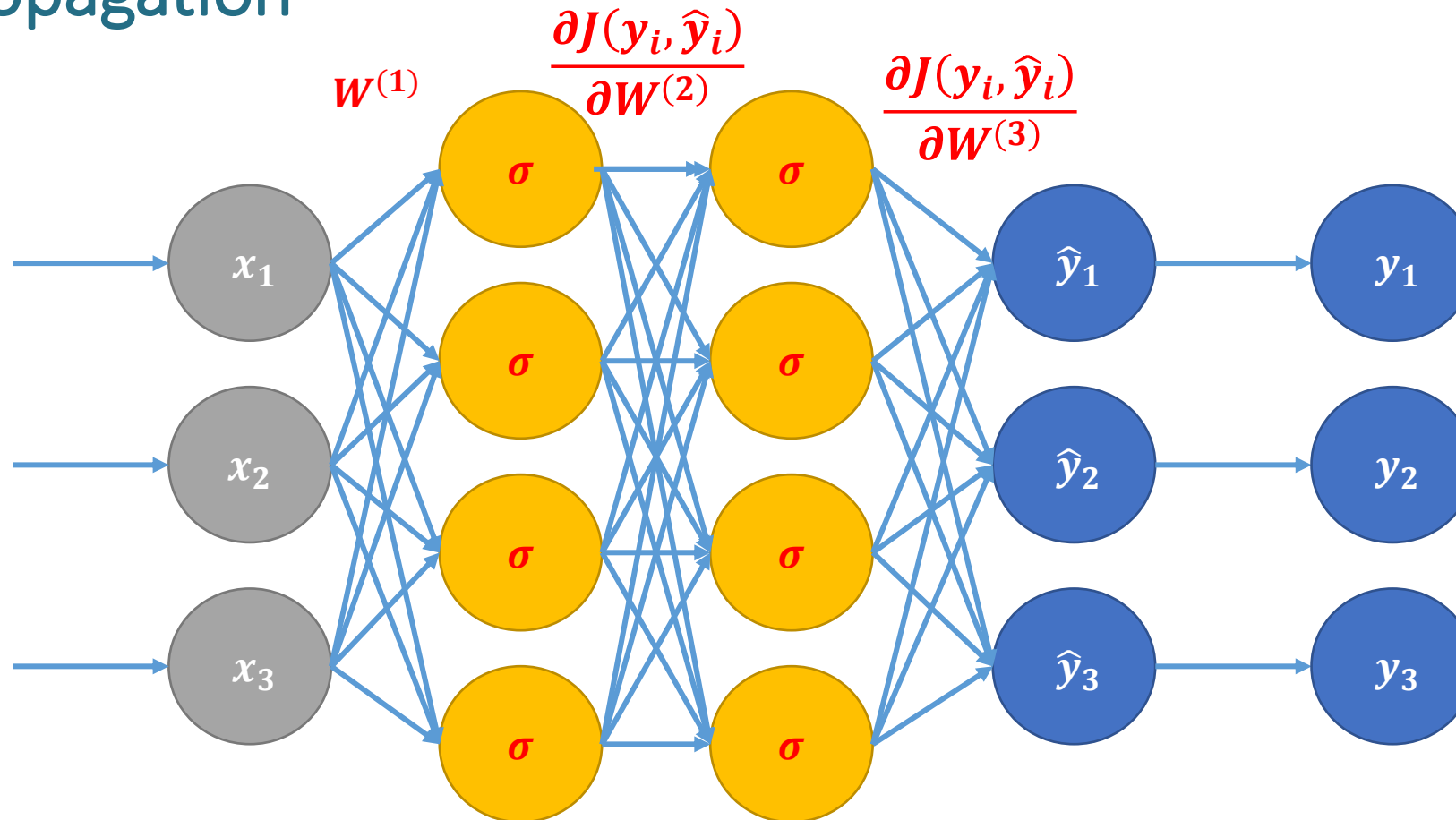
So we want to be able to compute the partial derivative of the loss function w.r.t. the weights W_k

Backpropagation



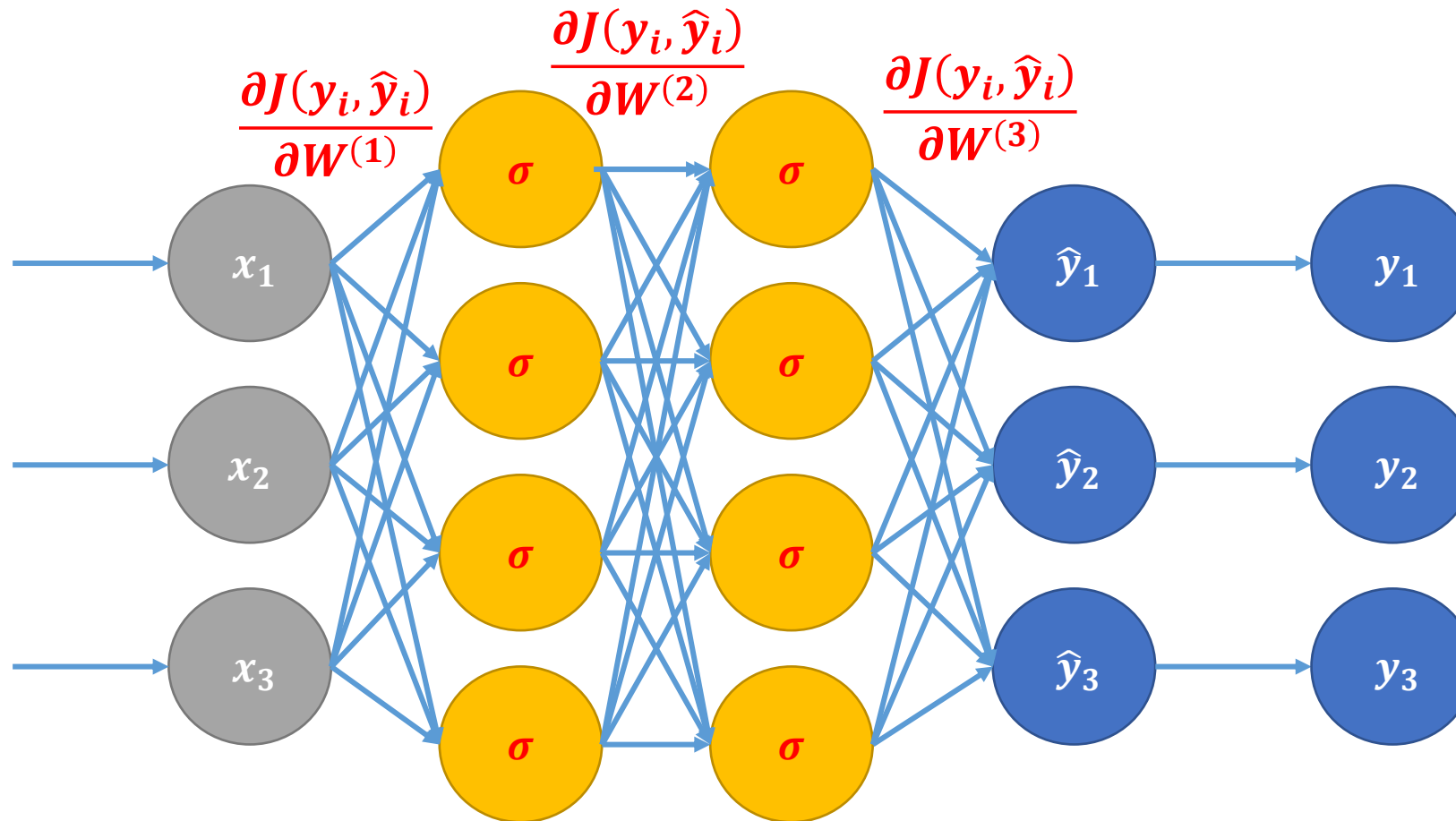
We first compute this, math turns out to be simpler for the last layer.

Backpropagation



Using that value, we compute the partial derivative for the previous layer → **Back propagation!**

Backpropagation



How have we trained before?

- **Gradient Descent.**

For most classical ML algorithms, the training happens here.

1. Make prediction
2. Calculate Loss
3. Calculate gradient of the loss function w.r.t. parameters
4. Update parameters by taking a step in the opposite direction
5. Iterate

Gradient descent

Update parameters by taking a step in the opposite direction

- Once we have the gradient, we just take a step in the opposite direction.
- Specifically; $w = w - \text{learning_rate} * \text{gradient}$
- This part is the same as in ML.