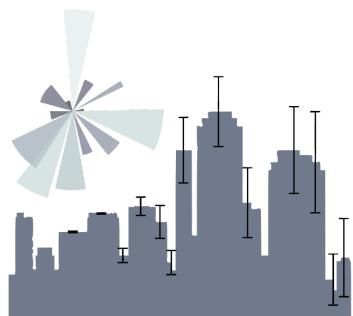


principles of Urban Science 14



Neural Networks: autoencoders

dr.federica bianco

fbb.space



fedhere



fedhere

this slide deck:

http://slides.com/federicabianco/pus2020_14

NN are a vast topics and we only have 2 weeks!

Some FREE references!

Neural Networks and Deep Learning

Neural Networks and Deep Learning is a free online book. The book will teach you about:

- Neural networks, a beautiful biologically-inspired programming paradigm which enables a computer to learn from observational data
- Deep learning, a powerful set of techniques for learning in neural networks

Neural networks and deep learning currently provide the best solutions to many problems in image recognition, speech recognition, and natural language processing. This book will teach you many of the core concepts behind neural networks and deep learning.

For more details about the approach taken in the book, [see here](#). Or you can jump directly to [Chapter 1](#) and get started.

[Neural Networks and Deep Learning](#)
[What this book is about](#)
[On the exercises and problems](#)
► [Using neural nets to recognize handwritten digits](#)
► [How the backpropagation algorithm works](#)
► [Improving the way neural networks learn](#)
► [A visual proof that neural nets can compute any function](#)
► [Why are deep neural networks hard to train?](#)
► [Deep learning](#)
[Appendix: Is there a simple algorithm for intelligence?](#)
[Acknowledgements](#)
[Frequently Asked Questions](#)

If you benefit from the book, please make a small donation. I suggest \$5, but you can choose the amount.



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

michael nielsen

better pedagogical approach, more basic, more clear

[Deep Learning](#)

An MIT Press book in preparation

Ian Goodfellow, Yoshua Bengio and Aaron Courville

[Book](#) [Exercises](#) [External Links](#)

Lectures

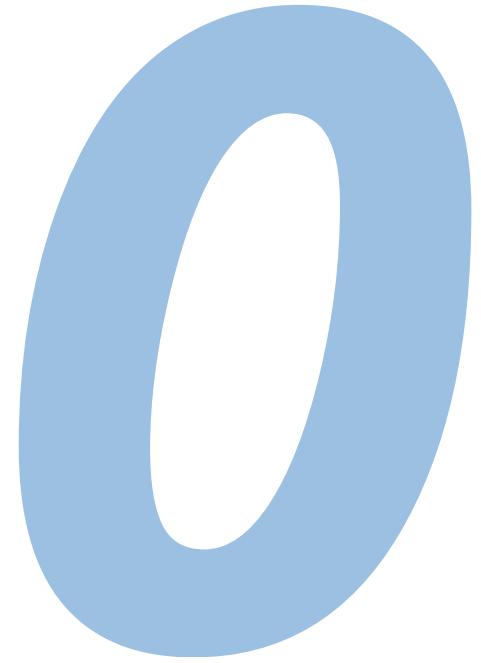
We plan to offer lecture slides accompanying all chapters of this book. We currently offer slides for only some chapters. If you are a course instructor and have your own lecture slides that are relevant, feel free to contact us if you would like to have your slides linked or mirrored from this site.

1. [Introduction](#)
 - Presentation of Chapter 1, based on figures from the book [[key](#)] [[pdf](#)]
 - [Video](#) of lecture by Ian and discussion of Chapter 1 at a reading group in San Francisco organized by Alena Kruchkova
2. [Linear Algebra](#) [[key](#)] [[pdf](#)]
3. [Probability and Information Theory](#) [[key](#)] [[pdf](#)]
4. [Numerical Computation](#) [[key](#)] [[pdf](#)] [[youtube](#)]
5. [Machine Learning Basics](#) [[key](#)] [[pdf](#)]
6. [Deep Feedforward Networks](#) [[key](#)] [[pdf](#)]
 - [Video](#) (.flv) of a presentation by Ian and a group discussion at a reading group at Google organized by Chintan Kaur.

<https://www.deeplearningbook.org/>

ian goodfellow

mathematical approach, more advanced, unfinished



Recap

Machine Learning

Data driven models for exploration of structure, prediction that learn parameters from data.

Machine Learning

Data driven models for exploration of structure, prediction that learn parameters from data.

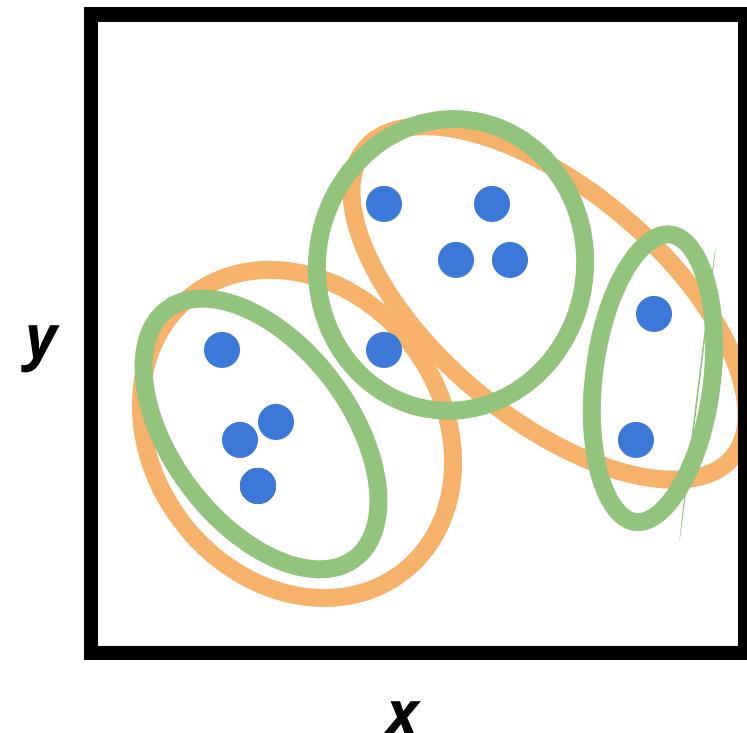
unupervised

set up: All features known for all observations

Goal: explore structure in the data

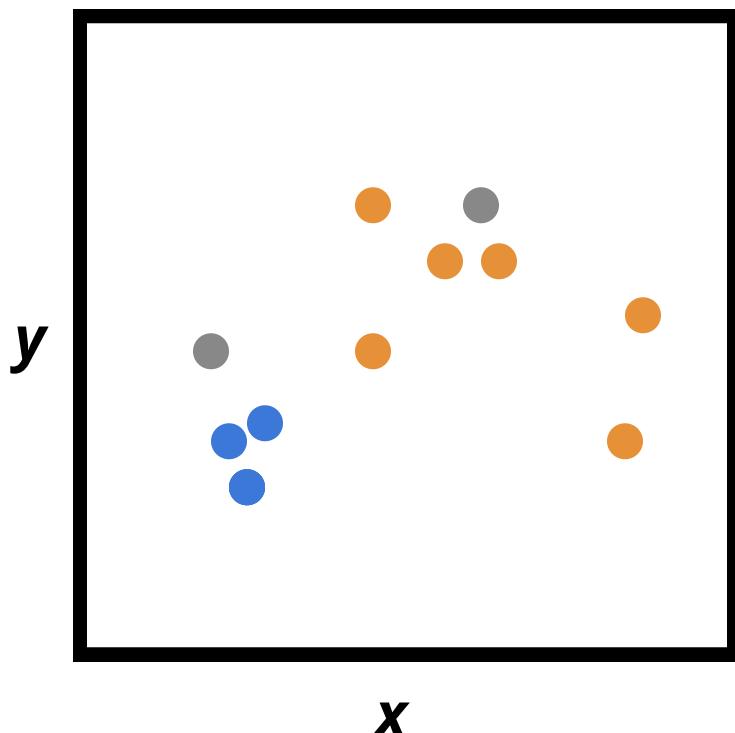
- data compression
- understanding structure

Algorithms: ***Clustering, (...)***



Machine Learning

Data driven models for exploration of structure, prediction that learn parameters from data.



supervised

set up: All features known for a subset of the data; one feature cannot be observed for the rest of the data

Goal: predicting missing feature

- classification
- regression

Algorithms: *regression, SVM, tree methods, k-nearest neighbors, neural networks, (...)*

Machine Learning

unupervised

set up: All features known for all observations

Goal: explore structure in the data

- data compression
- understanding structure

Algorithms: *k-means clustering,*
agglomerative clustering,
density based clustering, (...)

supervised

set up: All features known for a sunbset
of the data; one feature cannot be
observed for the rest of the data

Goal: predicting missing feature

- classification
- regression

Algorithms: *regression, SVM, tree*
methods, k-nearest neighbors,
neural networks, (...)

Machine Learning

Learning relies on the definition of a ***loss function***

model parameters are learned by calculating a loss function for different parameter sets and trying to minimize loss
(or a target function and trying to maximize)

e.g.

$$L1 = |target - prediction|$$

Machine Learning

Learning relies on the definition of a ***loss function***

learning type	loss / target
unsupervised	intra-cluster variance / inter cluster distance
supervised	distance between prediction and truth

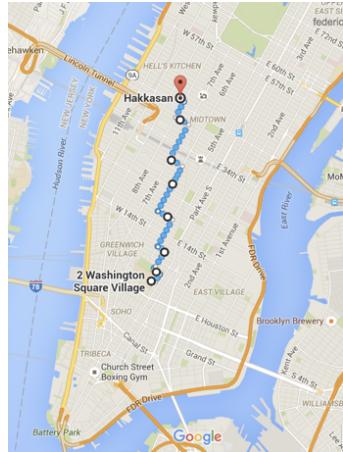
Machine Learning

The definition of a loss function requires the definition of *distance* or *similarity*

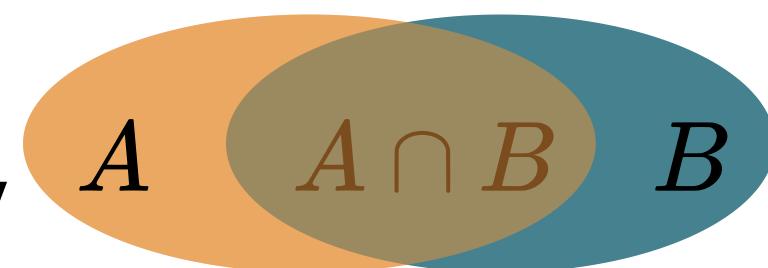
Machine Learning

The definition of a loss function requires the definition of *distance* or *similarity*

Minkowski distance



Jaccard similarity



Great circle distance



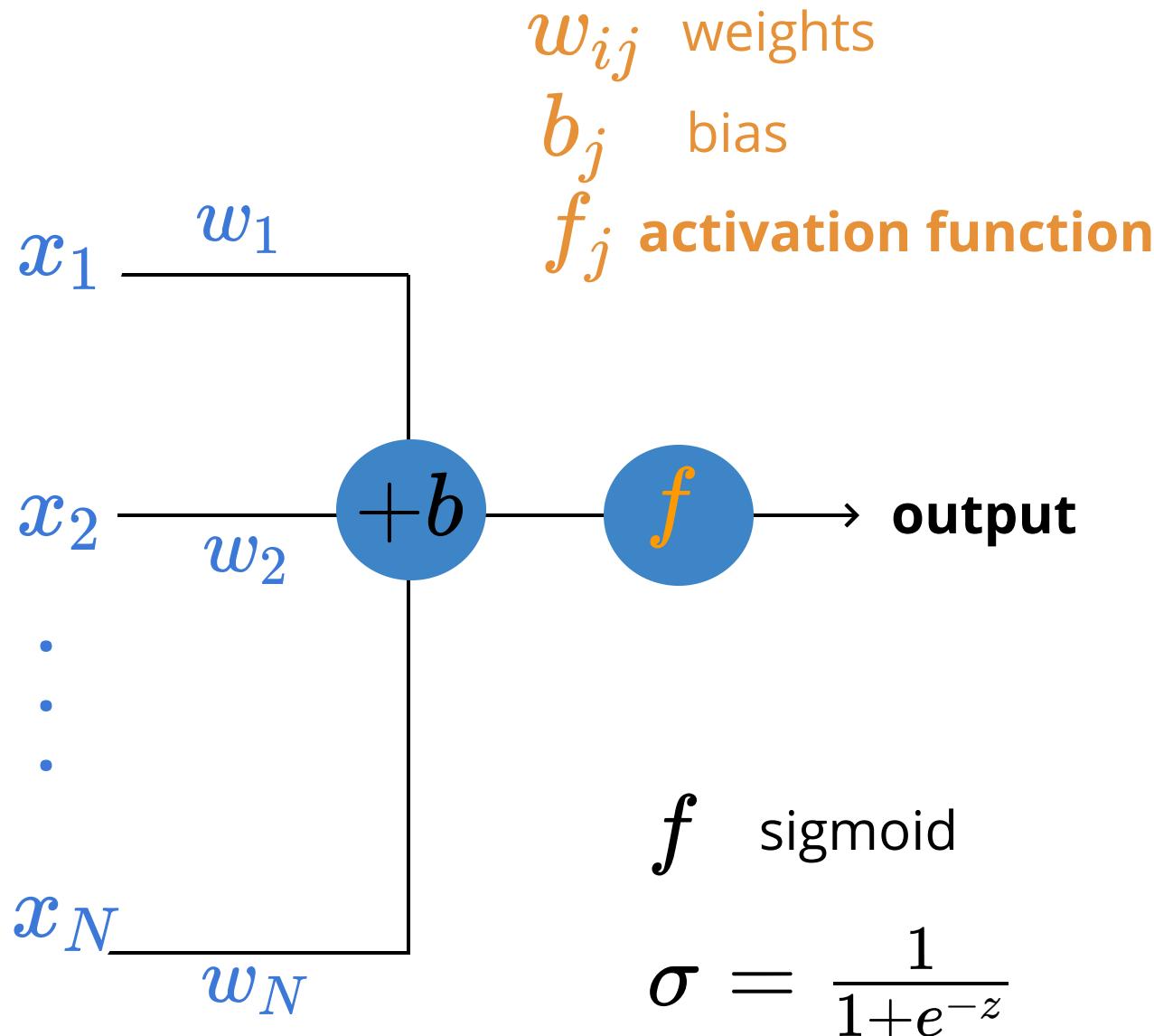
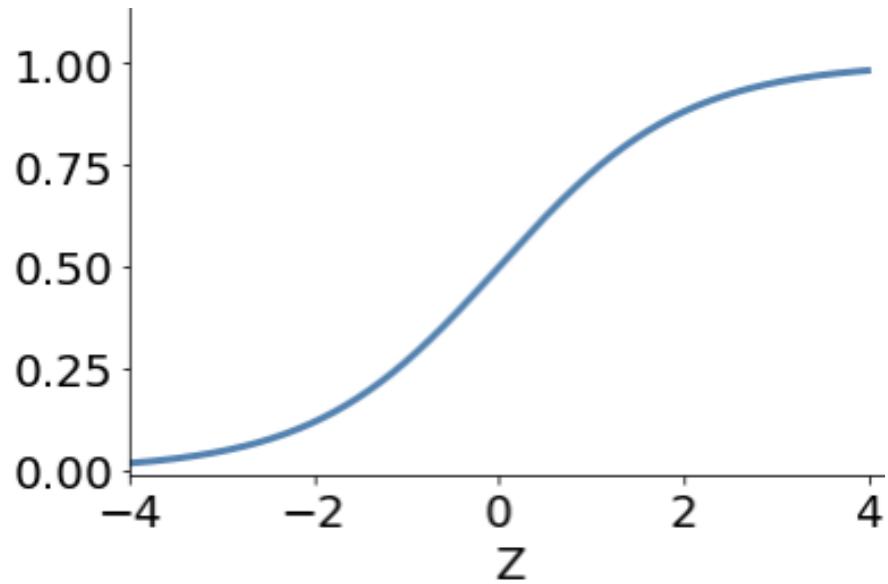
NN:
1
Neural Networks

Perceptrons are *linear classifiers*:

makes predictions based on a linear predictor function

combining a set of weights (=parameters) with the feature vector.

$$y = f(\sum_i w_i x_i + b)$$

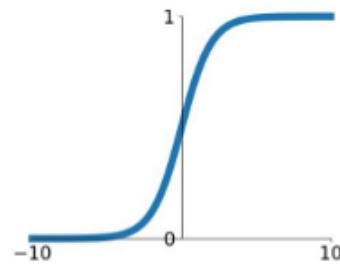


f

activation function

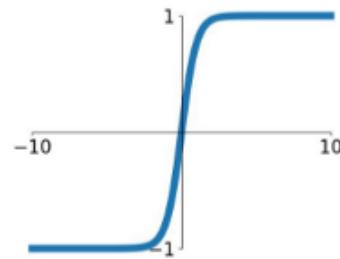
Sigmoid

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



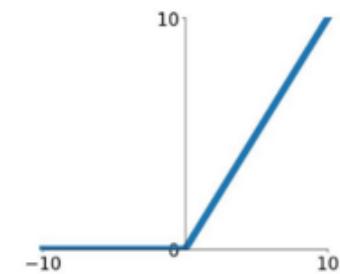
tanh

$$\tanh(x)$$



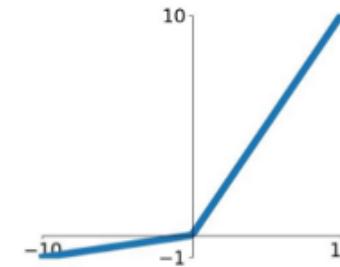
ReLU

$$\max(0, x)$$



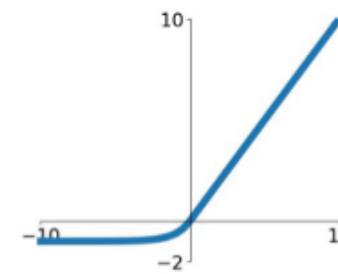
Leaky ReLU

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



Maxout

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

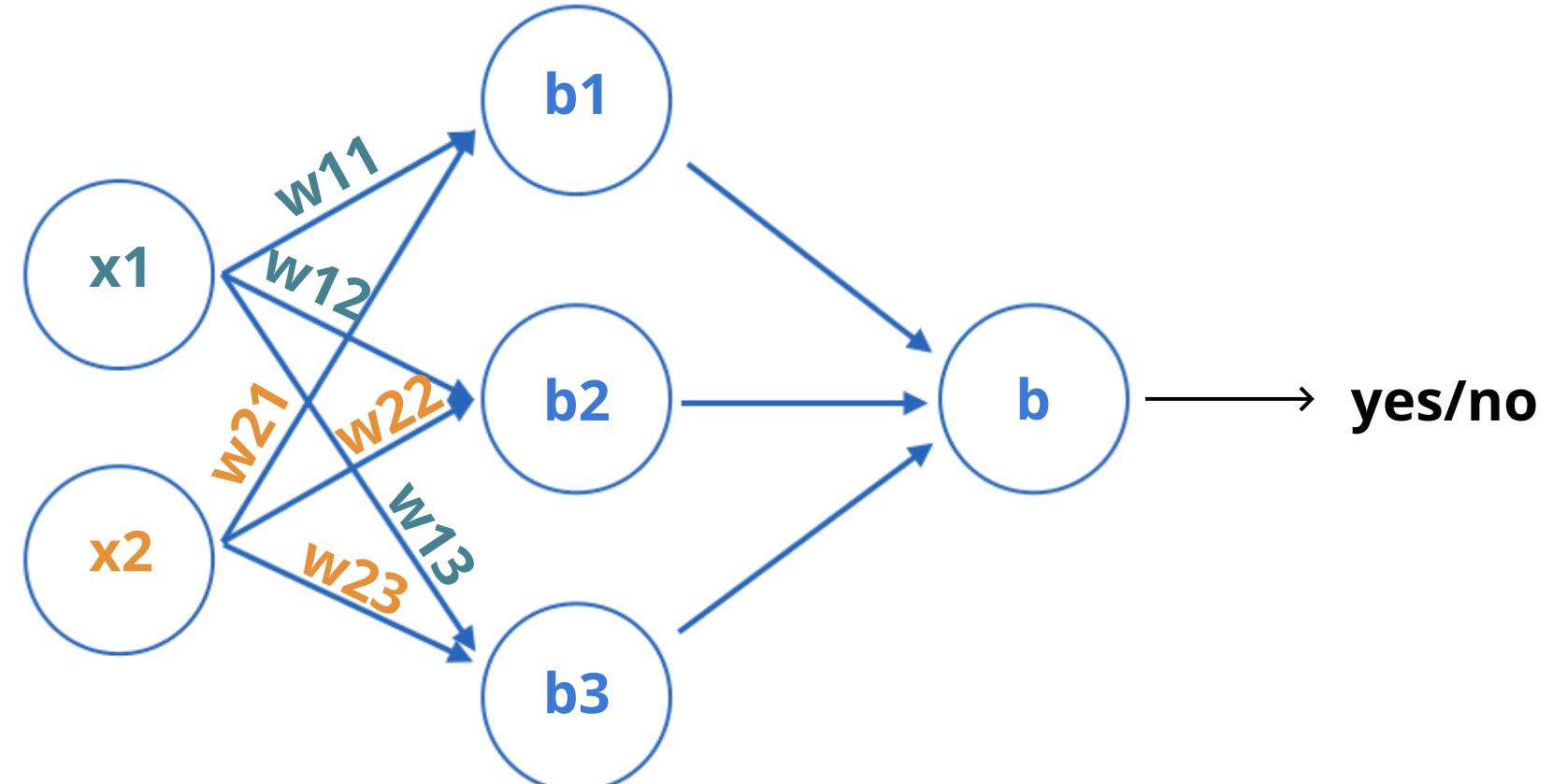


ELU

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

Turn a linear prediction into a binary or probabilistic classification

multilayer perceptron



w: weight

sets the sensitivity of a neuron

b: bias:

up-down weights a neuron

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

multilayer perceptron

connected: all nodes go to all nodes of the next layer.

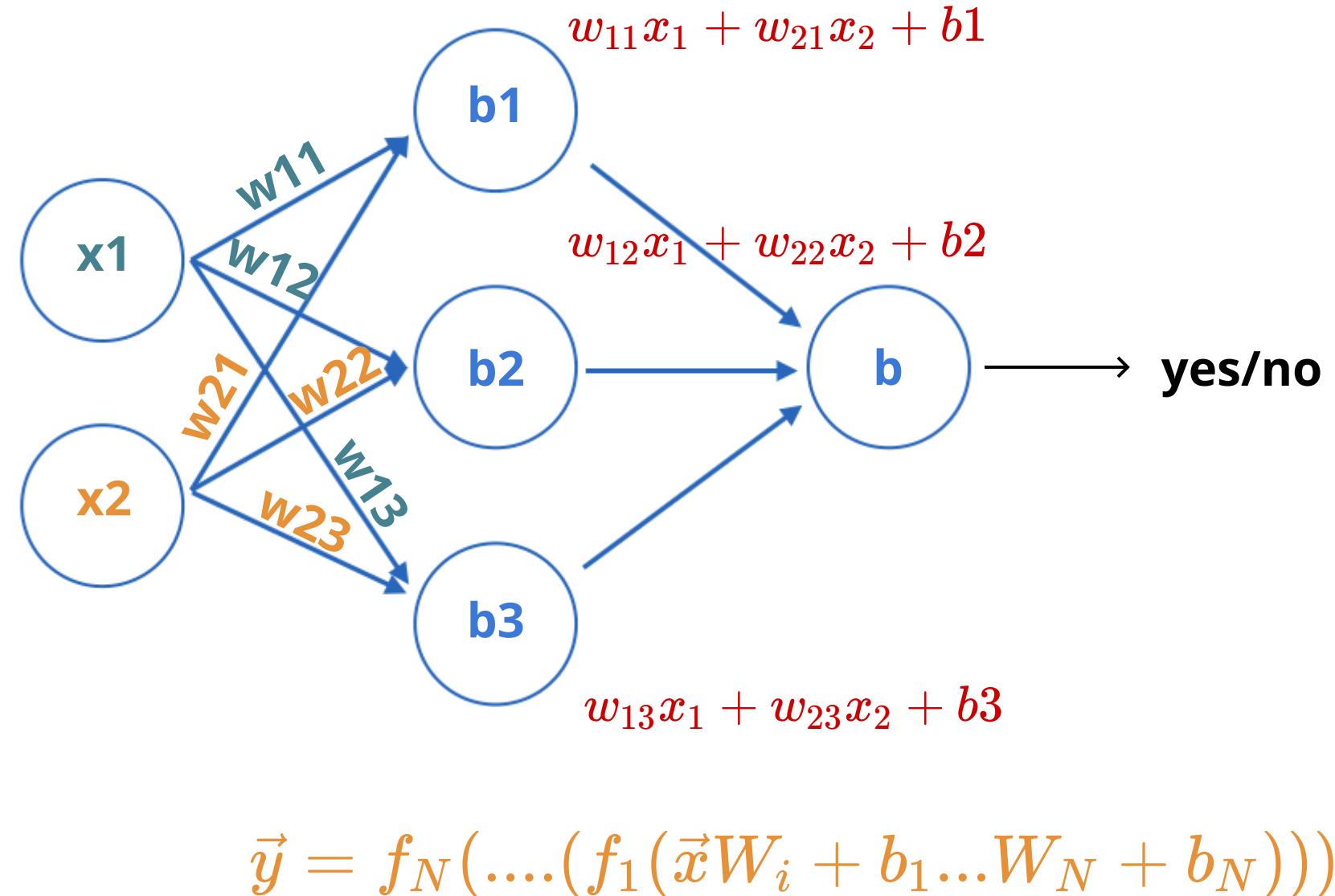
activation function

w: weight

sets the sensitivity of a neuron

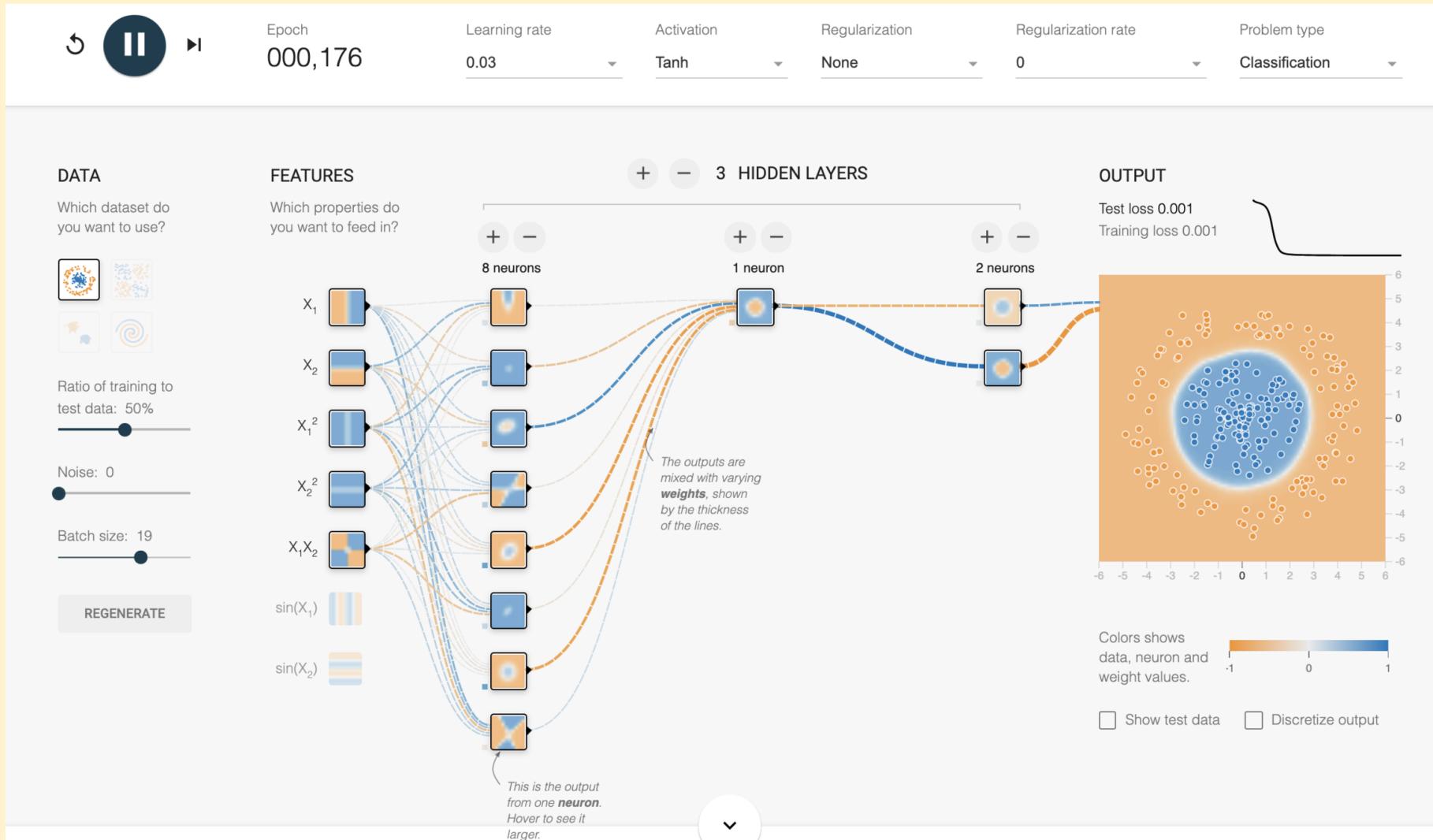
b: bias:

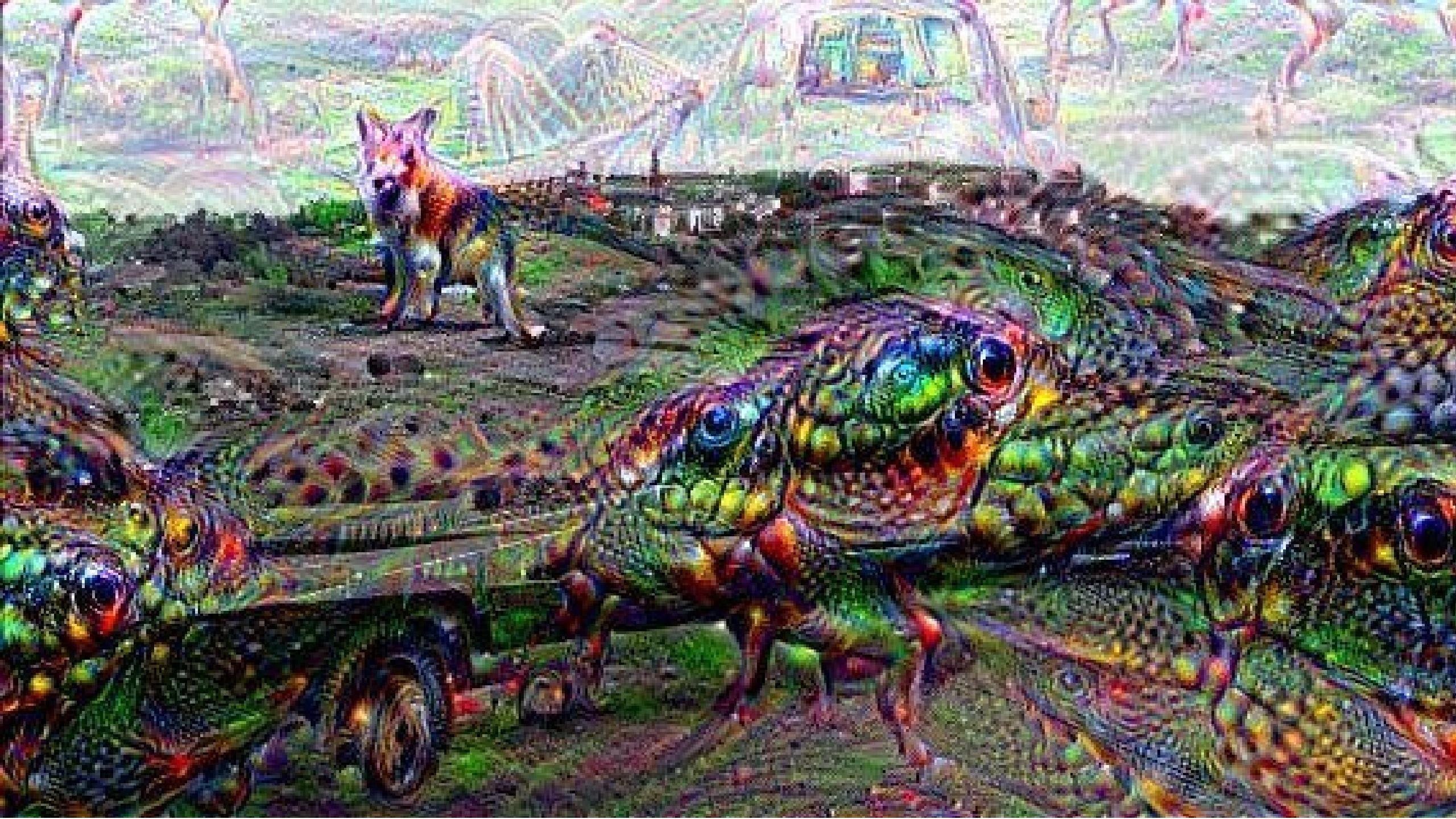
up-down weights a neuron



EXERCISE

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



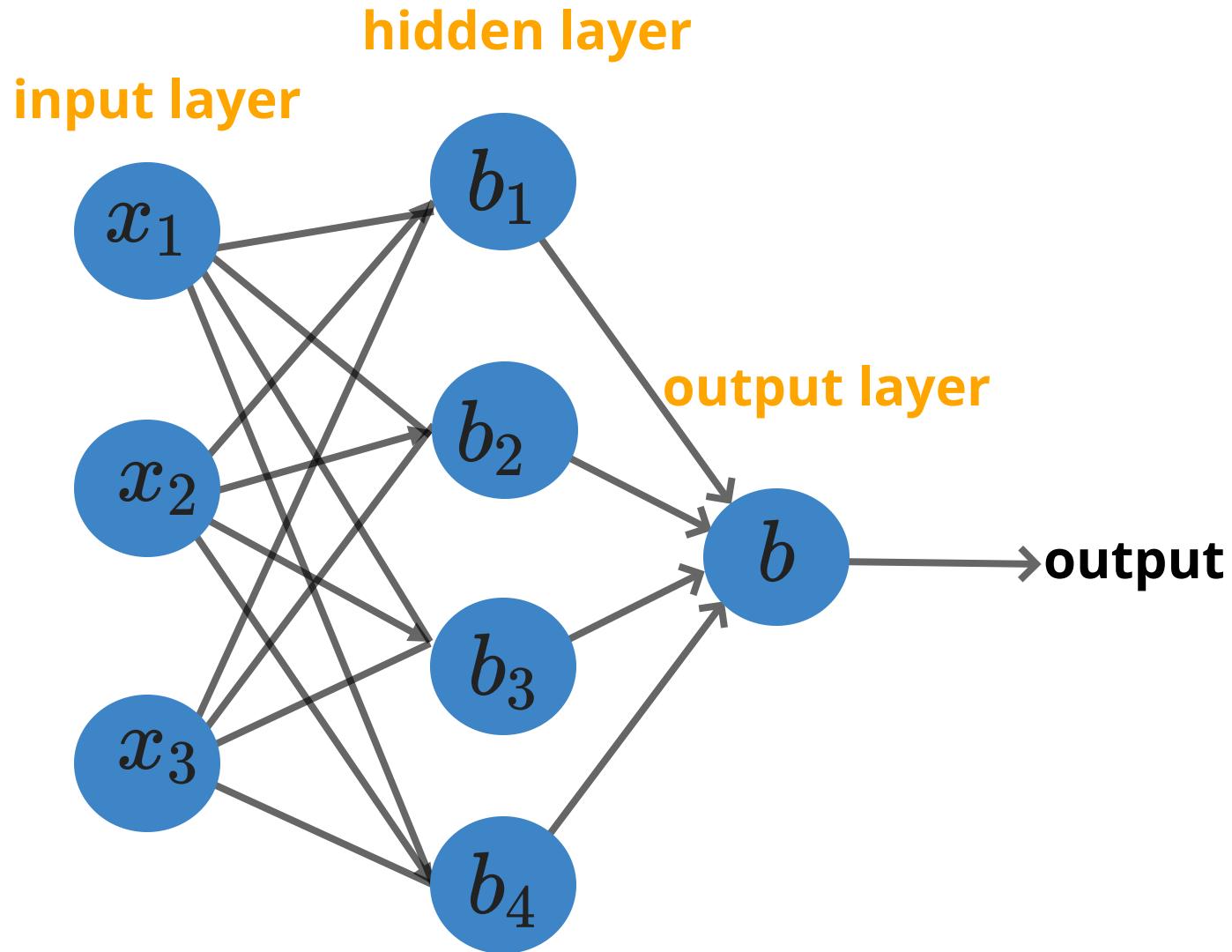


A large, light blue graphic of a question mark is positioned above the text. It has a thick, rounded top curve and a sharp, triangular point at the bottom, partially obscured by the word "DNN".

DNN: *Deep Neural Networks*

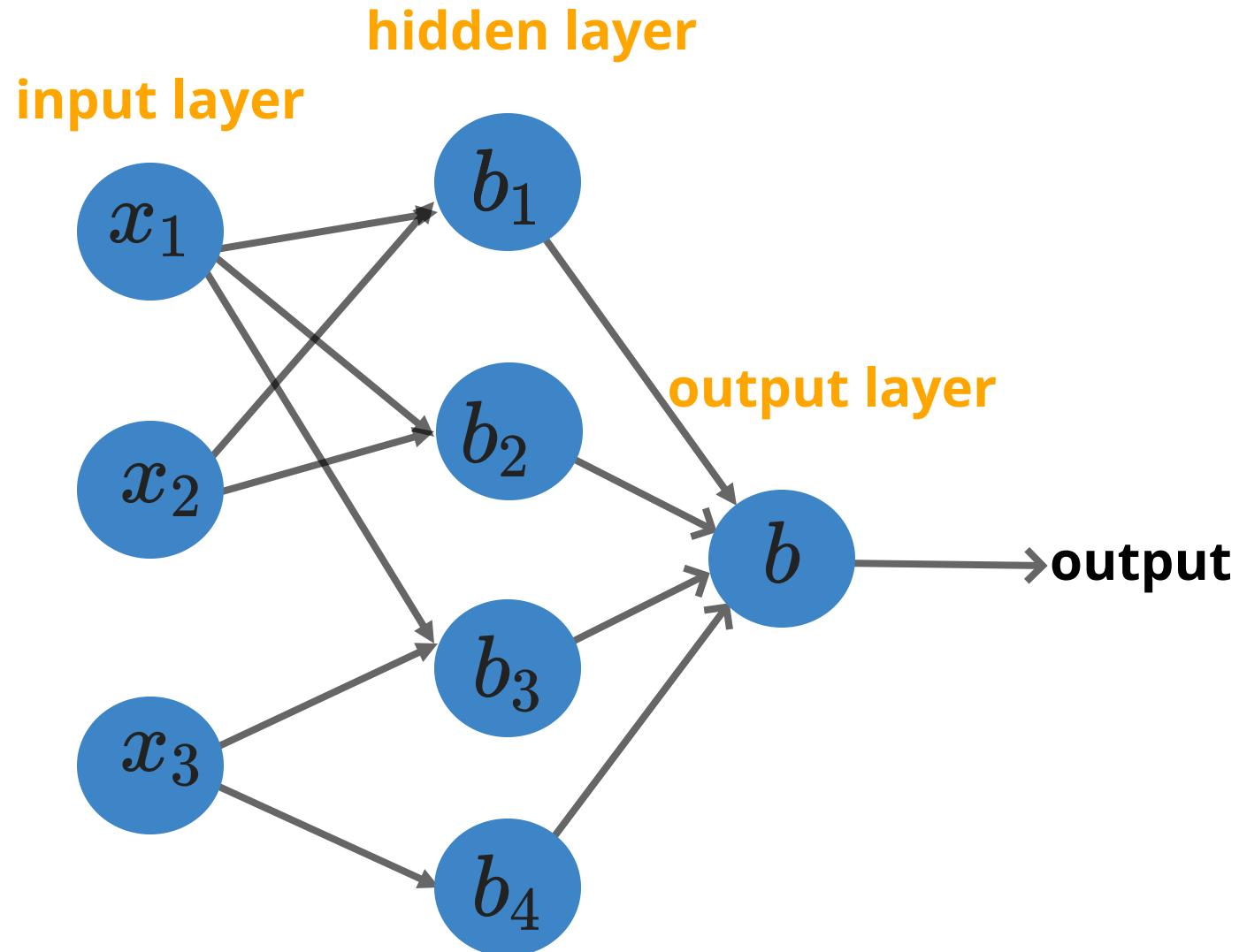
layer connectivity

Fully connected: all nodes go to all nodes of the next layer.



layer connectivity

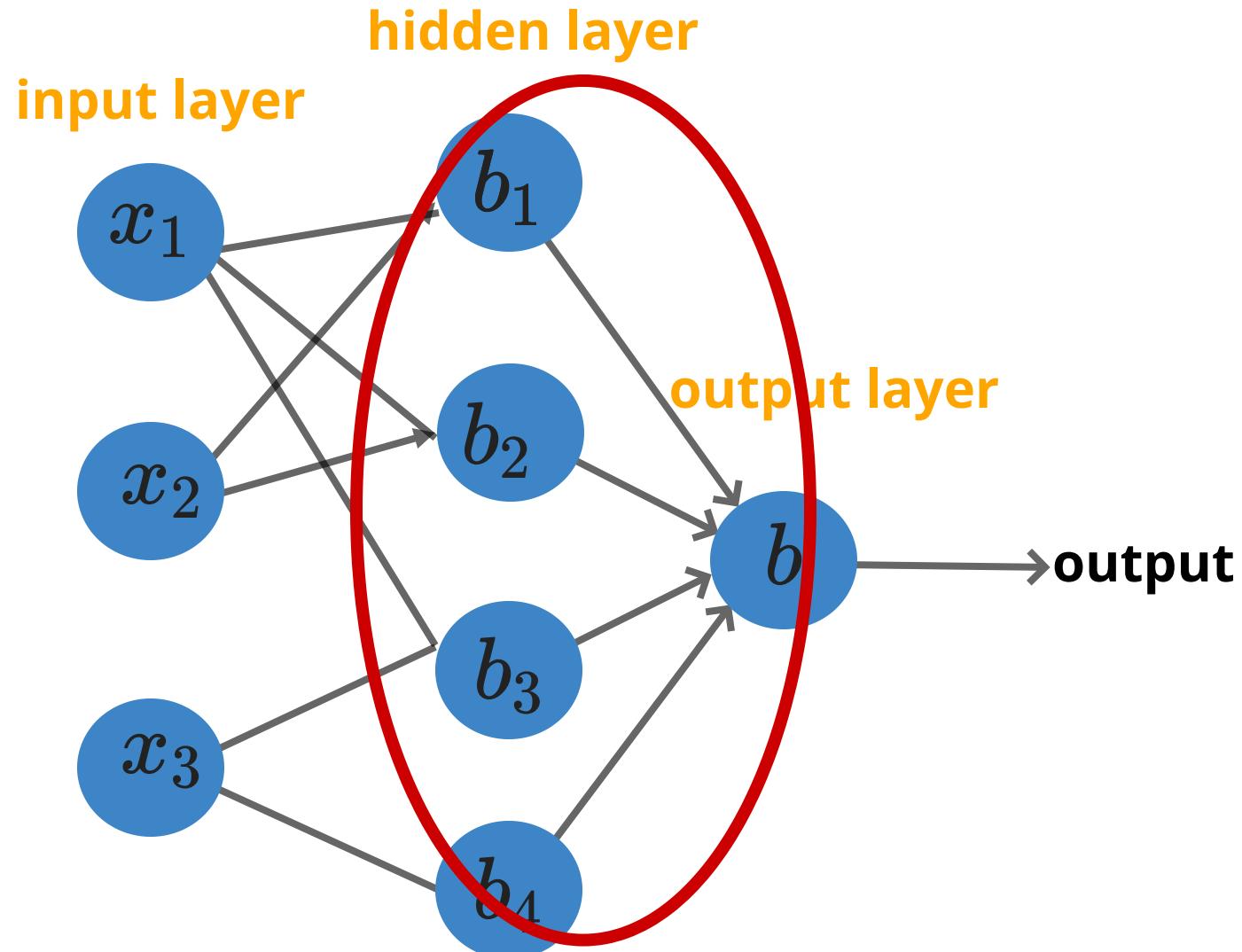
Sparcely connected: all nodes go to all nodes of the next layer.



layer connectivity

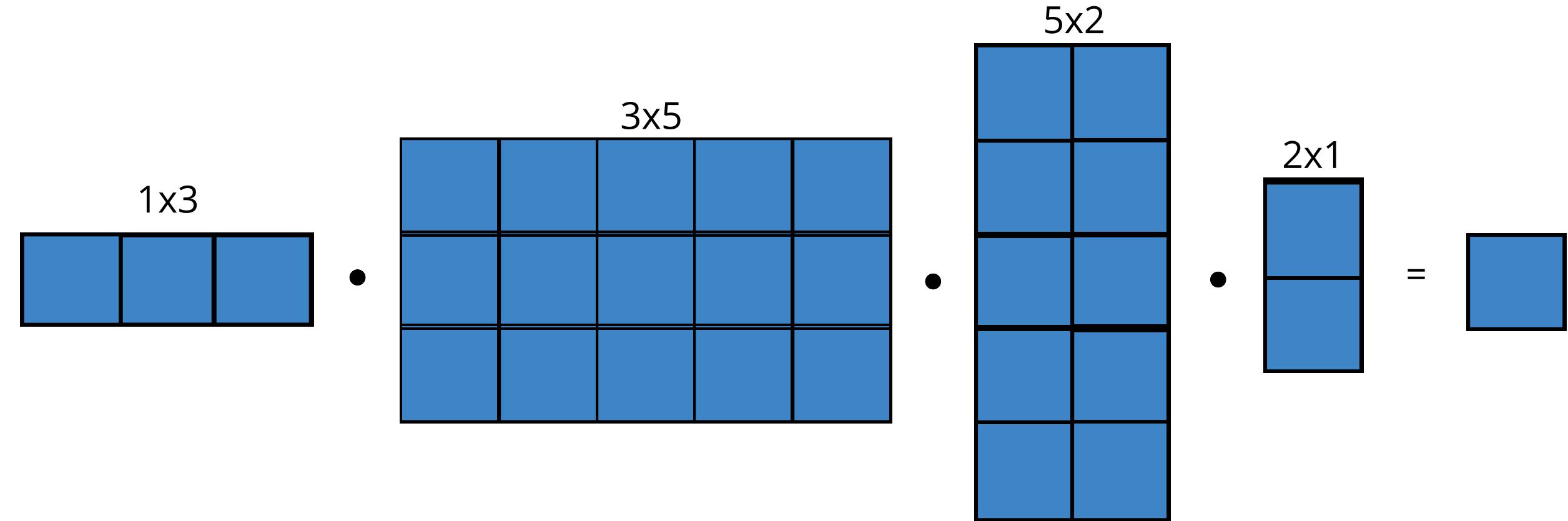
Sparcely connected: all nodes go to all nodes of the next layer.

The last layer is always connected



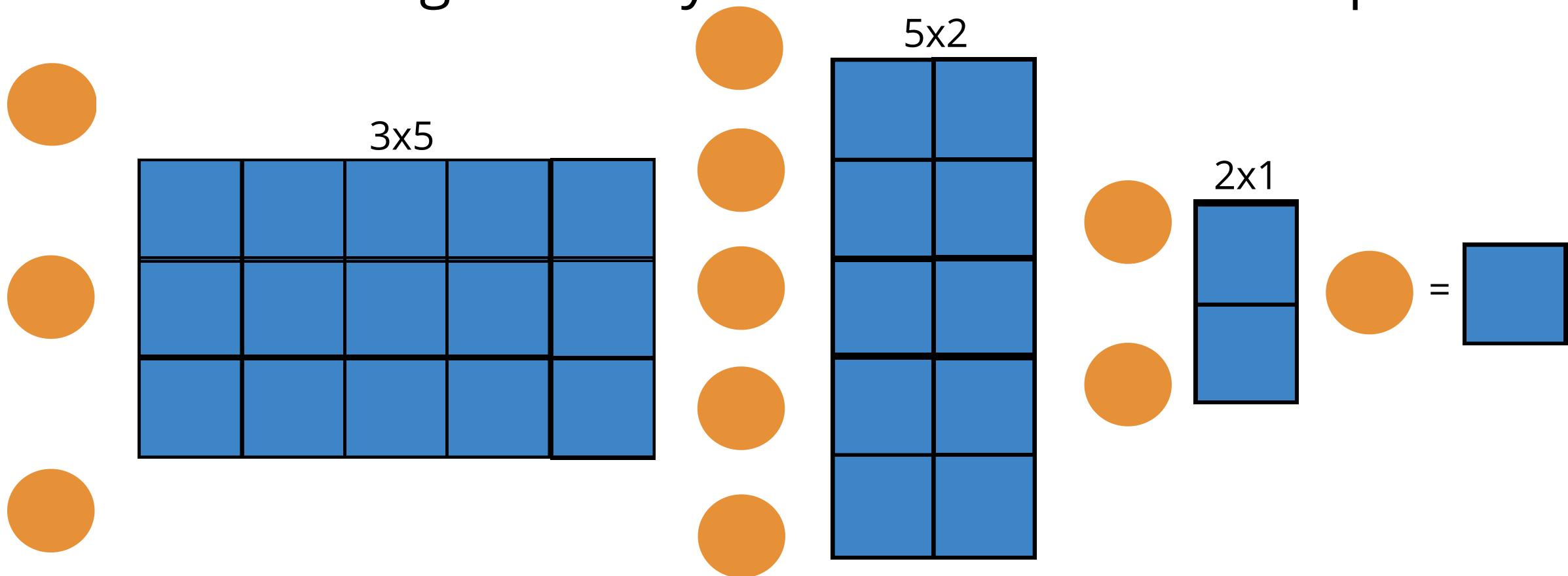
Deep Neural Network

what we are doing is just a series of matrix multiplications.



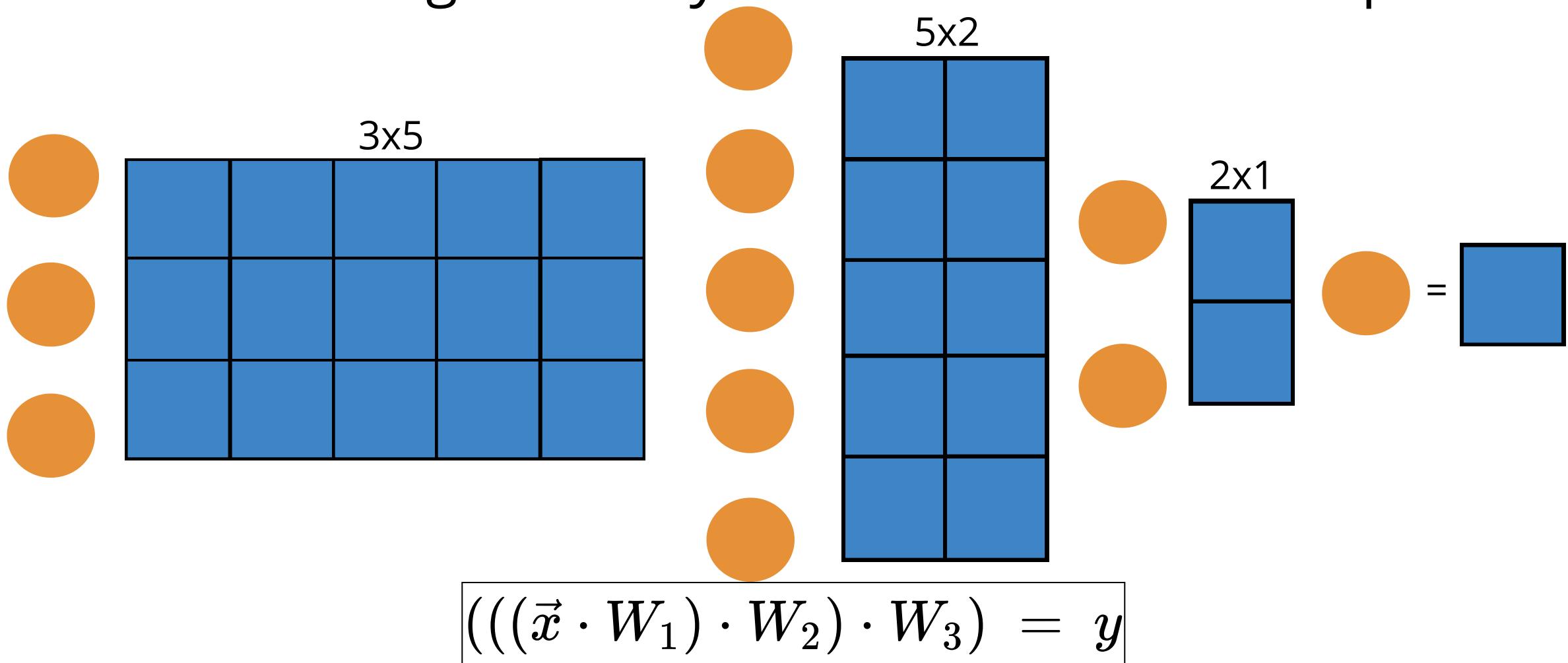
Deep Neural Network

what we are doing is exactly a series of matrix multiplications.



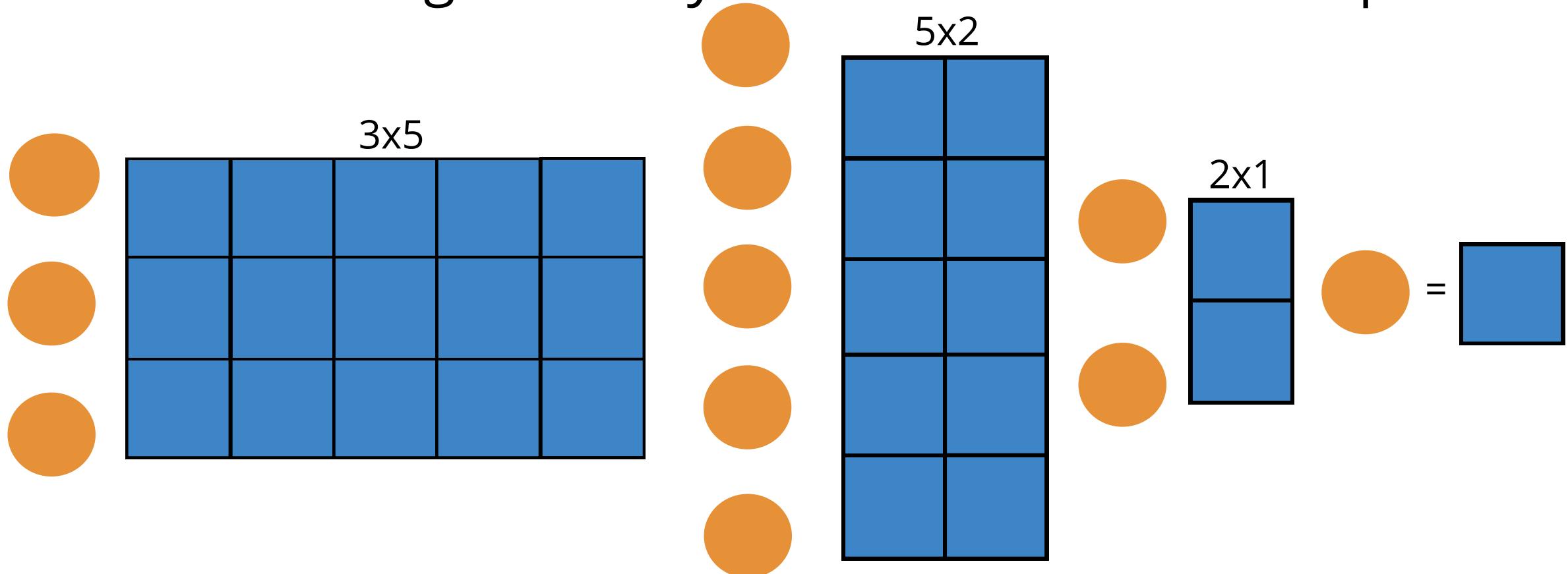
Deep Neural Network

what we are doing is exactly a series of matrix multiplications.



Deep Neural Network

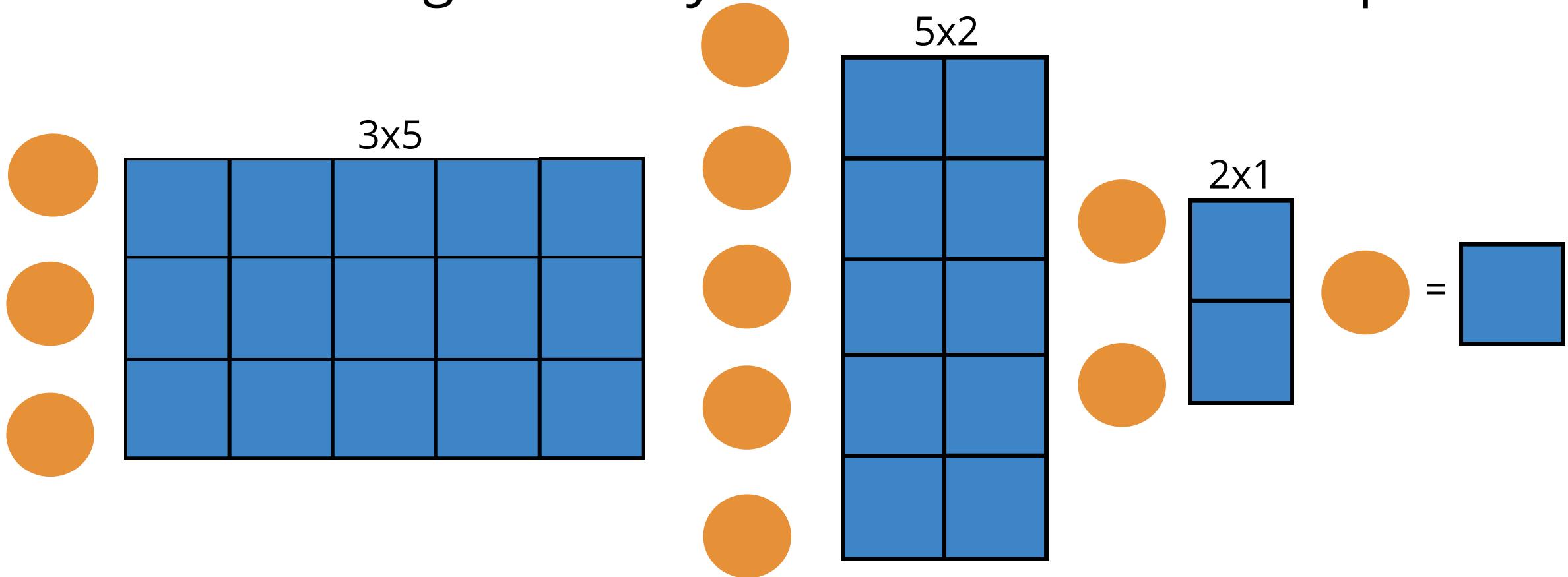
what we are doing is exactly a series of matrix multiplications.



$$(((\vec{x} \cdot W_1 + \vec{b}_1) \cdot W_2 + \vec{b}_2) \cdot W_3 + \vec{b}_3) = y$$

Deep Neural Network

what we are doing is exactly a series of matrix multiplications.



Deep Neural Network

what we are doing is exactly a series of matrix multiplications.

The purpose is to
approximate a function φ

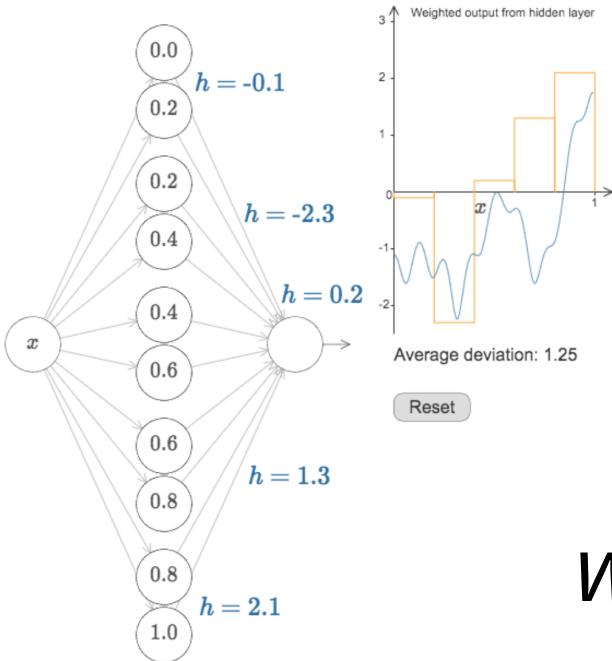
$$\mathbf{y} = \varphi(\mathbf{x})$$

*which (in general) is not linear
with linear operations*

$$\phi(\vec{x}) \sim f^{(3)}(f^{(2)}(f^{(1)}(\vec{x} \cdot W_1 + \vec{b}_1) \cdot W_2 + \vec{b}_2) \cdot W_3 + \vec{b}_3) = y$$

Deep Neural Network

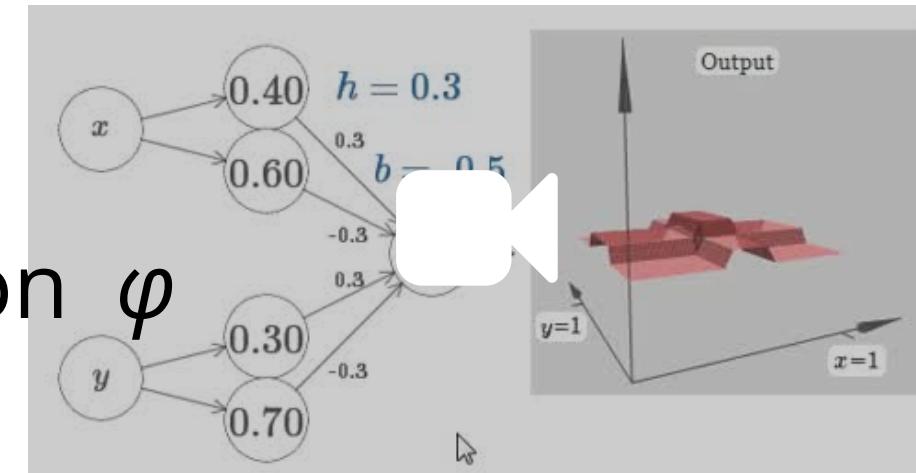
<http://neuralnetworksanddeeplearning.com/chap4.html>



The purpose is to
approximate a function φ

$$\mathbf{y} = \varphi(\mathbf{x})$$

*which (in general) is not linear
with linear operations*



$$\phi(\vec{x}) \sim f^{(3)}(f^{(2)}(f^{(1)}(\vec{x} \cdot W_1 + \vec{b}_1) \cdot W_2 + \vec{b}_2) \cdot W_3 + \vec{b}_3) = y$$

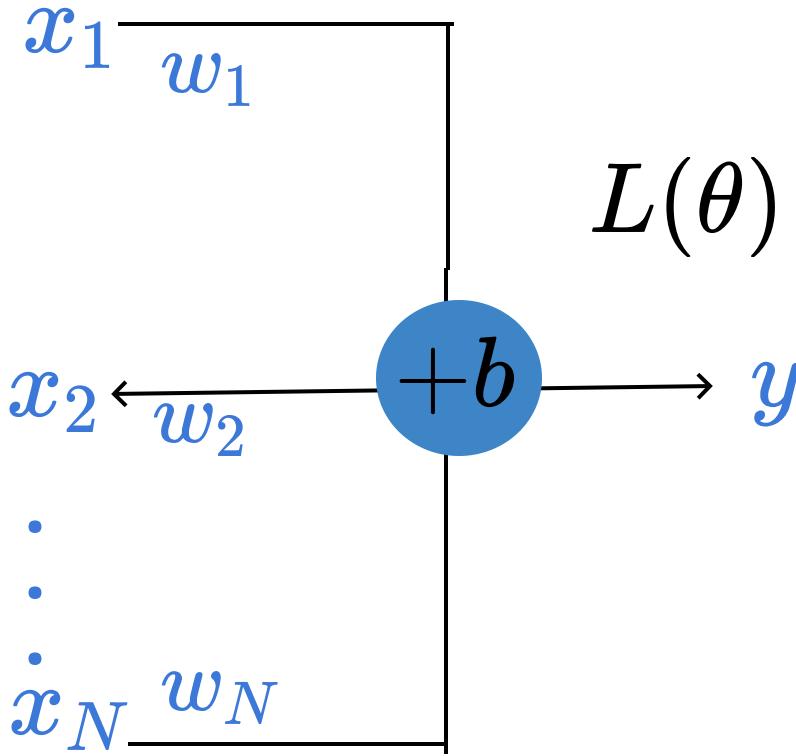
3 *training DNN*

https://colab.research.google.com/drive/13c9uJ_fPGjszgsyEuYWafR2F4_n-lXeZ

Gradient Descent

y : prediction

Any linear model:



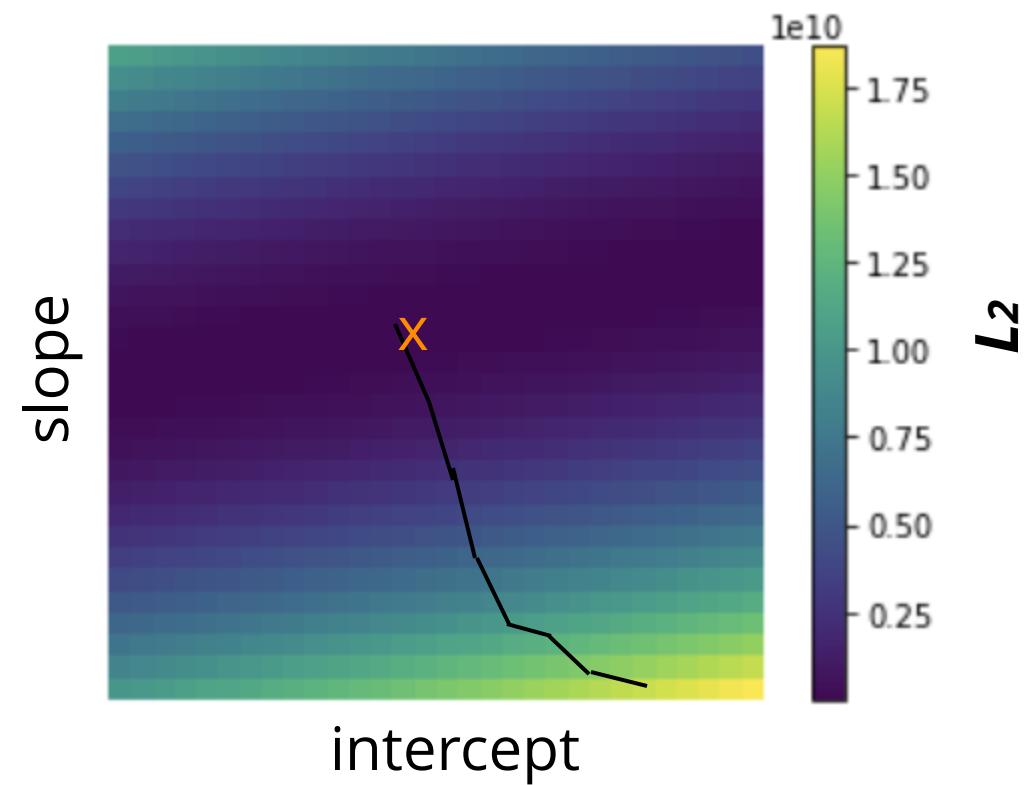
Error: e.g.

$$L(\theta) = |y - y_{\text{model}}|^2$$

$$\vec{y} = \vec{x}W + b$$

Find the best parameters by
finding the minimum of the L2
hyperplane

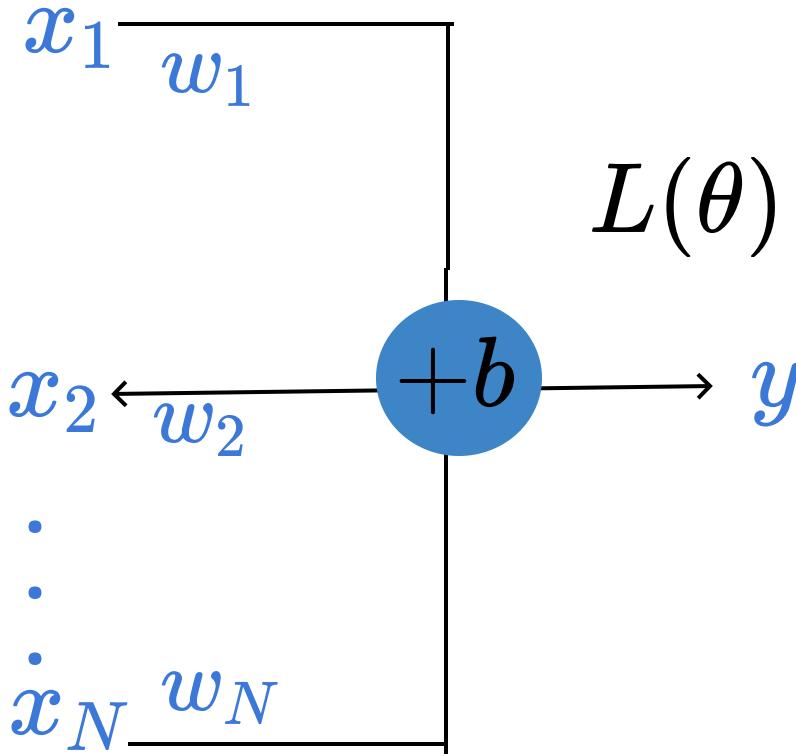
at every step look around and
choose the best direction



Gradient Descent

y : prediction

Any linear model:



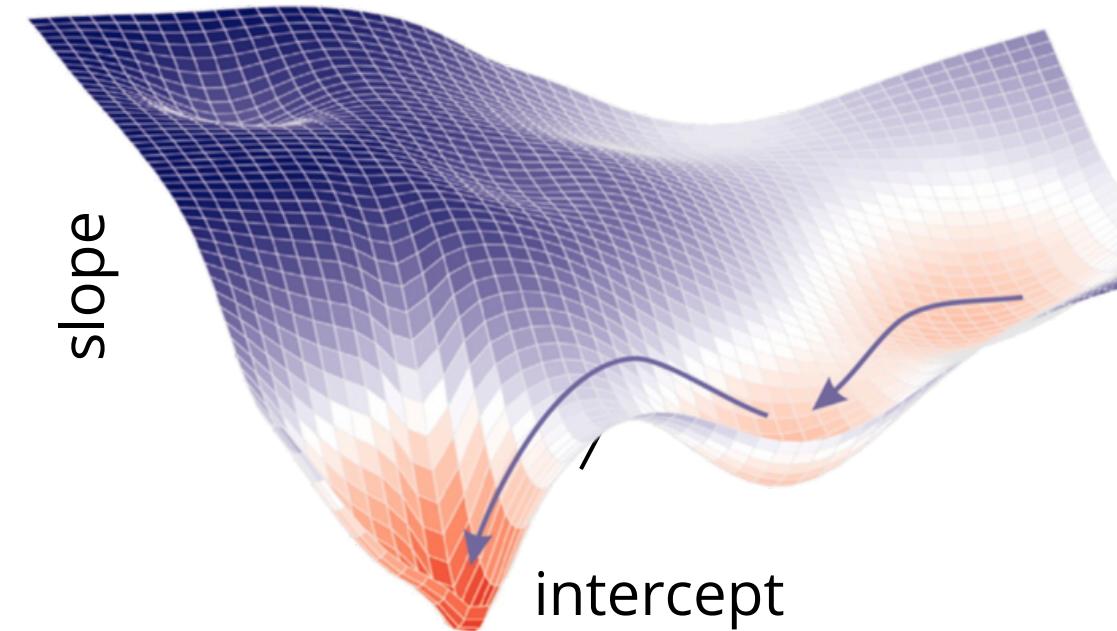
Error: e.g.

$$L(\theta) = |y - y_{\text{model}}|^2$$

$$\vec{y} = \vec{x}W + b$$

Find the best parameters by finding the minimum of the L2 hyperplane

at every step look around and choose the best direction

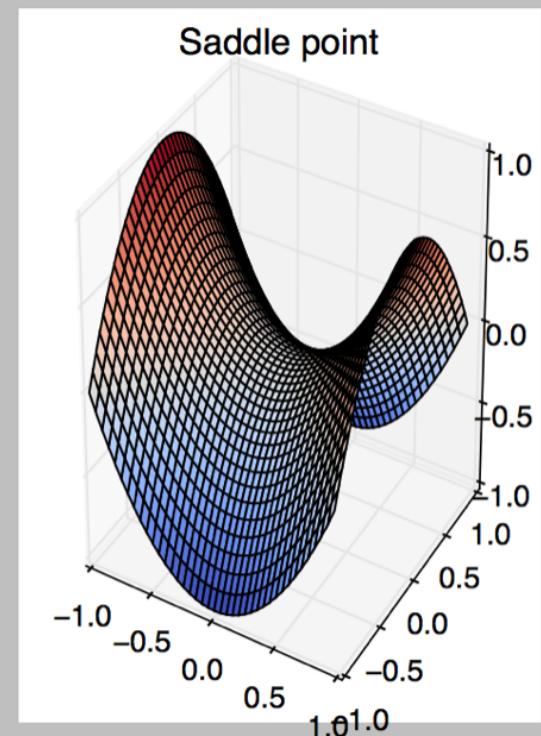
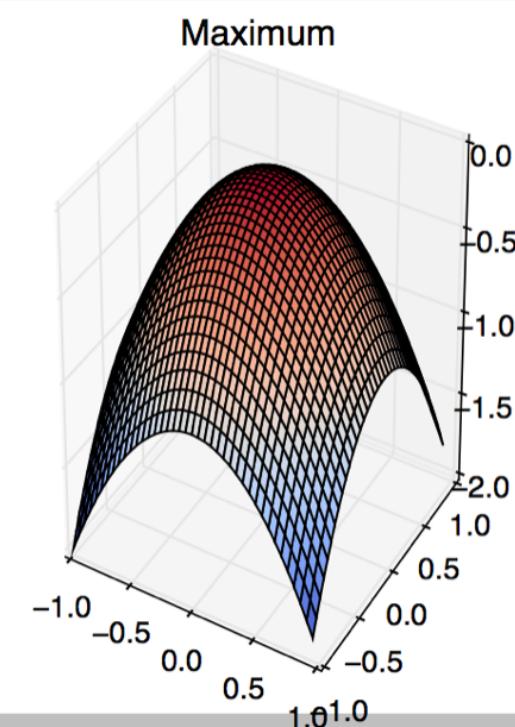
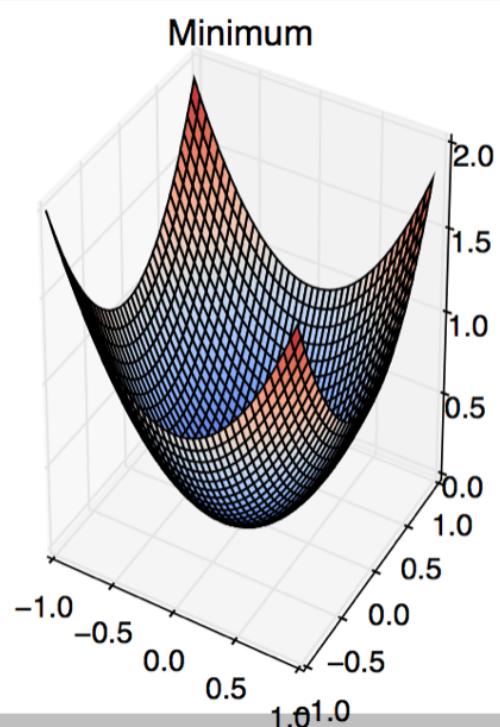


Gradient Descent

Critical points

at every step look around and choose the best direction

Zero gradient, and Hessian with...



All positive eigenvalues

Some positive
and some negative

Training a feed-forward DNN

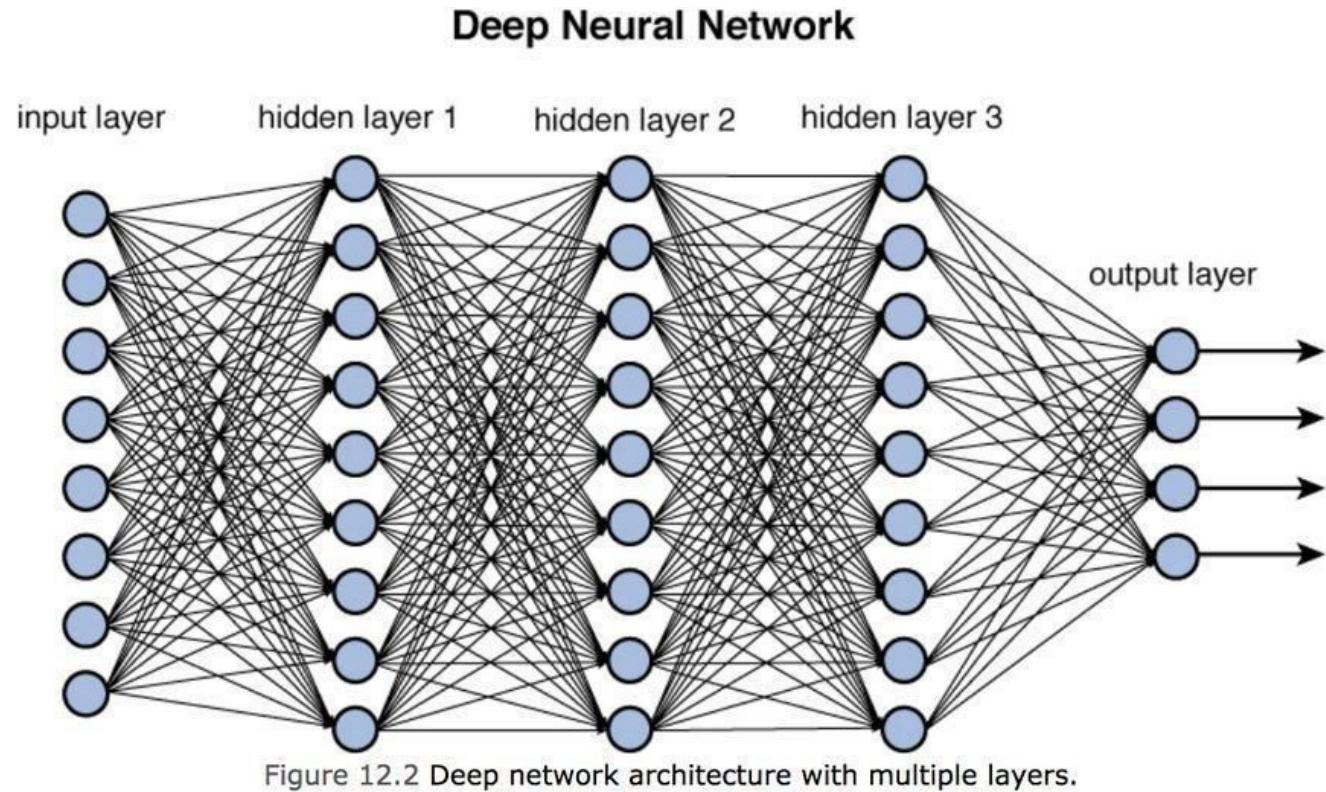
Training models with this many parameters requires a lot of care:

- . defining the metric
- . optimization schemes
- . training/validation/testing sets

But just like our simple linear regression case, the fact that small changes in the parameters leads to small changes in the output for the right activation functions.

define a cost function, e.g.

$$C = \frac{1}{2} |y - a^L|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$



feed data forward through network and calculate cost metric

for each layer, calculate effect of small changes on next layer

$$\vec{y} = f^{(N)}(\dots(f^{(1)}(\vec{x}W_i + b_1 \dots W_N + b_N)))$$

Training a feed-forward DNN

Training models with this many parameters requires a lot of care:

- . defining the metric
- . optimization schemes
- . training/validation/testing sets

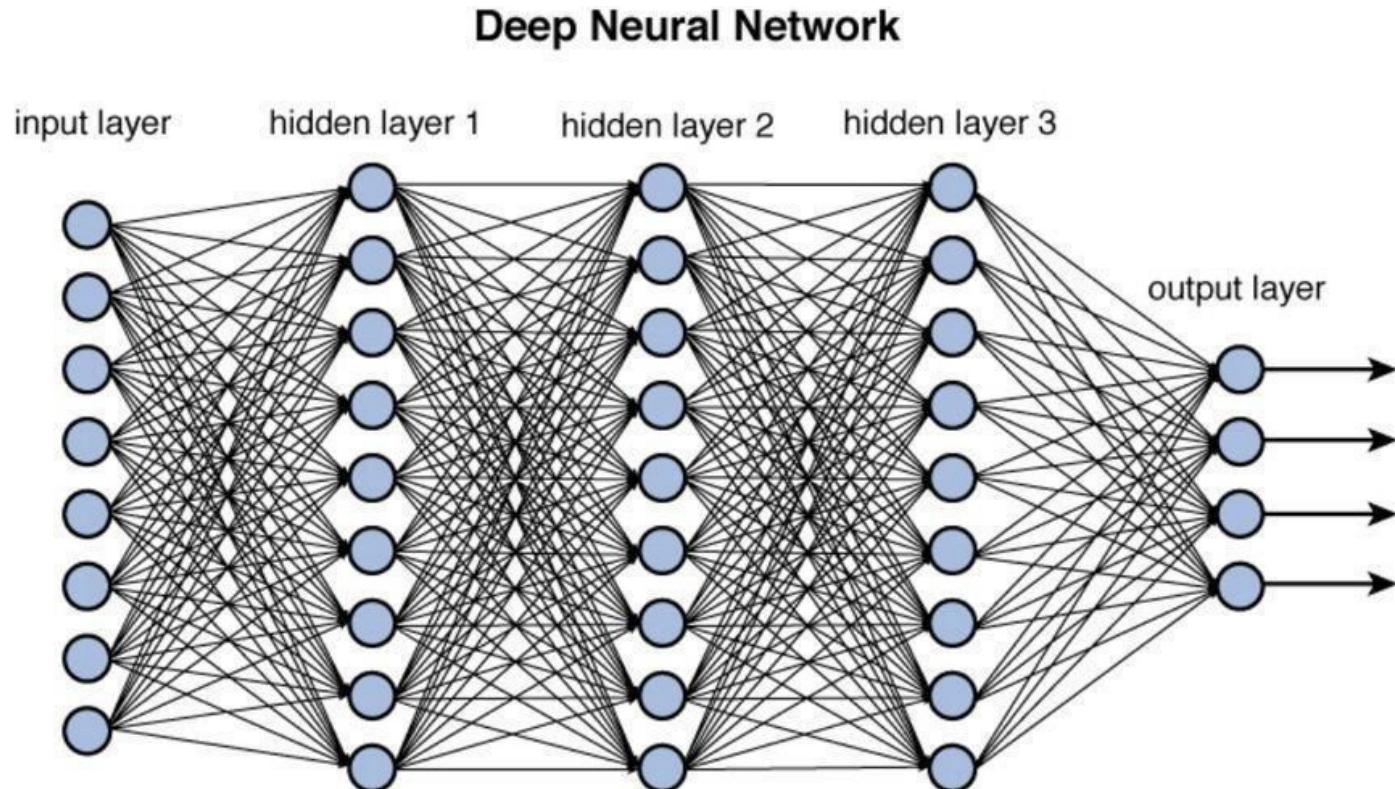


Figure 12.2 Deep network architecture with multiple layers.

earlier layers learn more slowly

$$\vec{y} = f^{(N)}(\dots(f^{(1)}(\vec{x}W_i + b_1\dots W_N + b_N)))$$

Gradient Descent

Loss functions: with NN you often encounter this
loss function

$$L(\theta) = -E_{x,y \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y|x)$$

negative loglikelihood or cross entropy

$$p_{\text{model}}(y|x) = N(y; f(x; \theta), I)$$

$$L(\theta) = \frac{1}{2} E_{x,y \sim \hat{p}_{\text{data}}} \|y - f(x; \theta)\|^2 + c \sim 2MSE$$

Gradient Descent

Loss functions: with NN you often encounter this
loss function

$$L(\theta) = -E_{x,y \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(y|x)$$

negative loglikelihood or cross entropy

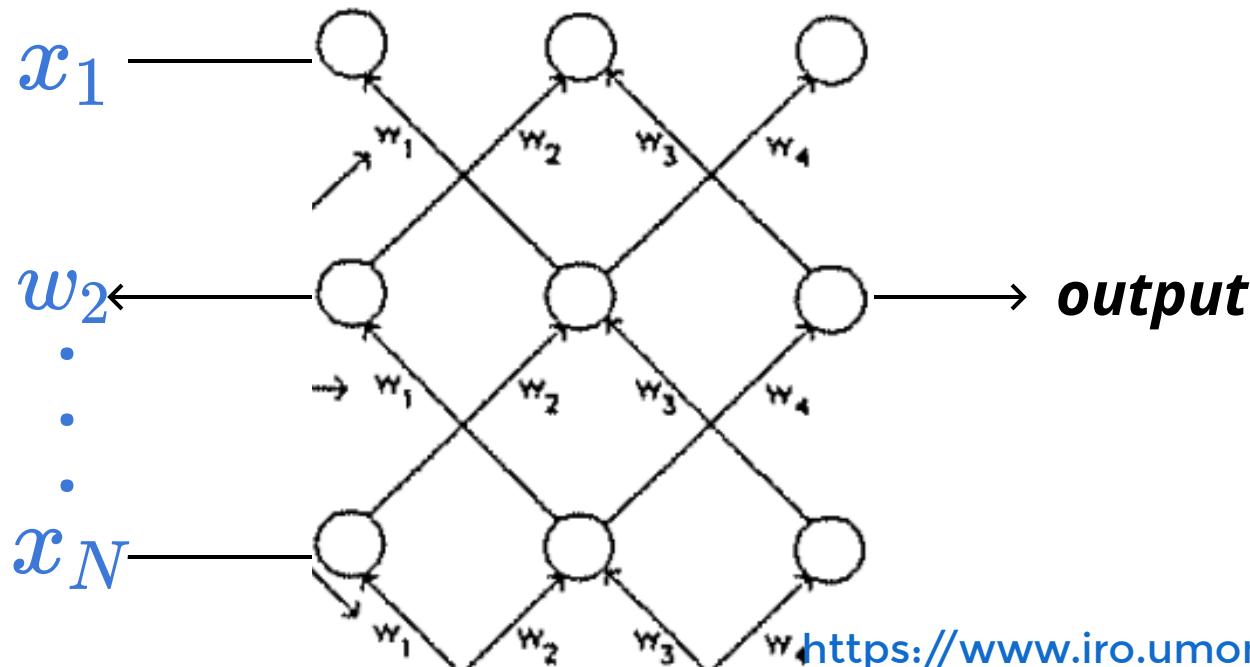
if $p_{\text{model}}(y|x) = N(y; f(x; \theta), I)$

$$L(\theta) = \frac{1}{2} E_{x,y \sim \hat{p}_{\text{data}}} \|y - f(x; \theta)\|^2 + c \sim 2MSE$$

Training a feed-forward DNN back-propagation

how does linear descent look when you have a whole network structure with hundreds of weights and biases to optimize??

$$x_j = \sum_i y_i w_{ji} \quad y_j = \frac{1}{1+e^{-x_j}}$$



nature

Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

Nature 323, 533–536(1986) | Cite this article

22k Accesses | 7872 Citations | 167 Altmetric | Metrics

Abstract

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

Training a feed-forward DNN back-propagation

how does linear descent look when you have a whole network structure with hundreds of weights and biases to optimize??

- we want to get the gradient to use it in downhill optimization



Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

Nature 323, 533–536(1986) | Cite this article

22k Accesses | 7872 Citations | 167 Altmetric | Metrics

Abstract

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal ‘hidden’ units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

Training a feed-forward DNN back-propagation

backprop is a *dynamic programming* algorithm that calculates all gradients than looks them up

- we want to get the gradient to use it in downhill optimization

- chain rule

$$0 \frac{\partial C}{\partial x} = \frac{\partial C}{\partial y} \frac{\partial y}{\partial x}$$



Learning representations by back-propagating errors

David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams

Nature 323, 533–536(1986) | [Cite this article](#)

22k Accesses | 7872 Citations | 167 Altmetric | [Metrics](#)

Abstract

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal ‘hidden’ units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

Training a feed-forward DNN

This is the simplest deep NN: one neuron per layer

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



$$a_1 = \sigma(z_1) = \sigma(w_1 a_0 + b_1)$$

$$\Delta a_1 = \frac{\partial \sigma(w_1 a_0 + b_1)}{\partial b_1} \Delta b_1 = \sigma'(z_1) \Delta b_1$$

$$\vec{y} = f^{(N)}(\dots(f^{(1)}(\vec{x}W_i + b_1 \dots W_N + b_N)))$$

Training a feed-forward DNN

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



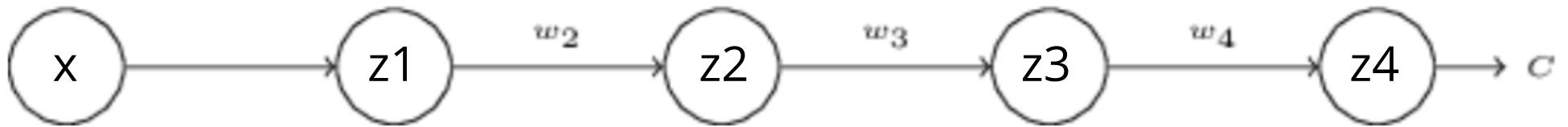
$$a_1 = \sigma(z_1) = \sigma(w_1 a_0 + b_1)$$

$\frac{\partial C}{\partial w}, \frac{\partial C}{\partial b}$ these are the changes on the last layer w respect to w and b

$$\vec{y} = f^{(N)}(\dots(f^{(1)}(\vec{x}W_i + b_1\dots W_N + b_N)))$$

Training a feed-forward DNN

This is the simplest deep NN: one neuron per layer



$$\begin{aligned}\frac{\partial C}{\partial x} &= \frac{\partial C}{\partial z4} \frac{\partial C}{\partial z3} \frac{\partial C}{\partial z2} \frac{\partial C}{\partial z1} \frac{\partial C}{\partial x} = \\ &= f'(f(f(f(z))))f'(f(f(z)))f'(f(z))f'(z) =\end{aligned}$$

$$\vec{y} = f^{(N)}(\dots(f^{(1)}(\vec{x}W_i + b_1\dots W_N + b_N)))$$

Training a DNN back-propagation

how does linear descent look when you have a whole network structure with hundreds of weights and biases to optimize??

think of applying just gradient to a function of a function of a function... use:

- 1) partial derivatives, 2) chain rule

define a cost function, e.g. $C = \frac{1}{2}|y - a^L|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$

An equation for the error in the output layer, δ^L : The components of δ^L are given by

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L). \quad (\text{BP1})$$

matrix-based form, as

$$\delta^L = \nabla_a C \odot \sigma'(z^L). \quad (\text{BP1a})$$

Here, $\nabla_a C$ is defined to be a vector whose components are the partial derivatives $\partial C / \partial a_j^L$. You can think of $\nabla_a C$ as expressing the rate of change of C with respect to the output activations

The backpropagation equations provide us with a way of computing the gradient of the cost function. Let's explicitly write this out in the form of an algorithm:

1. **Input x :** Set the corresponding activation a^1 for the input layer.
2. **Feedforward:** For each $l = 2, 3, \dots, L$ compute $z^l = w^l a^{l-1} + b^l$ and $a^l = \sigma(z^l)$.
3. **Output error δ^L :** Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$.
4. **Backpropagate the error:** For each $l = L-1, L-2, \dots, 2$ compute $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$.
5. **Output:** The gradient of the cost function is given by

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j^l} = \delta_j^l.$$

Examining the algorithm you can see why it's called *backpropagation*. We compute the error vectors δ^l backward, starting from the final layer. It may seem peculiar that we're going

<http://neuralnetworksanddeeplearning.com/chap2.html>

$$\vec{y} = f_N(\dots(f_1(\vec{x}W_i + b_1\dots W_N + b_N)))$$

Training a DNN

Summary: the equations of backpropagation

$$\delta^L = \nabla_a C \odot \sigma'(z^L)$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

Output: The gradient of the cost function is given by

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j^l} = \delta_j^l.$$

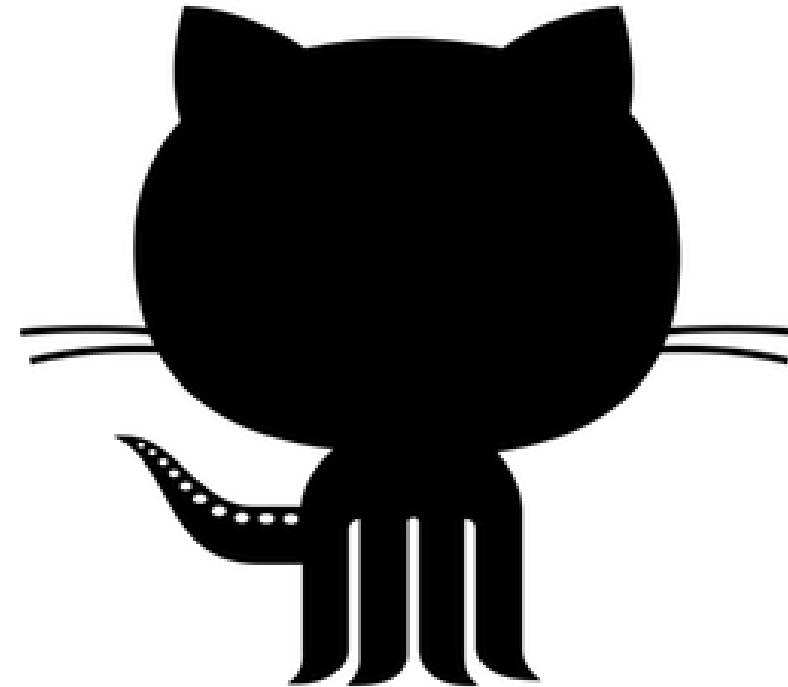
Output error δ^L : Compute the vector $\delta^L = \nabla_a C \odot \sigma'(z^L)$.

Gradient descent: For each $l = L, L - 1, \dots, 2$ update the weights according to the rule $w^l \rightarrow w^l - \frac{\eta}{m} \sum_x \delta^{x,l} (a^{x,l-1})^T$, and the biases according to the rule $b^l \rightarrow b^l - \frac{\eta}{m} \sum_x \delta^{x,l}$.

<http://neuralnetworksanddeeplearning.com/chap2.html#backpropsummary>

$$\vec{y} = f^{(N)}(\dots(f^{(1)}(\vec{x}W_i + b_1 \dots W_N + b_N)))$$

https://colab.research.google.com/drive/13c9uJ_fPGjszgsyEuYWafR2F4_n-IXeZ



build a DNN from scratch using numpy

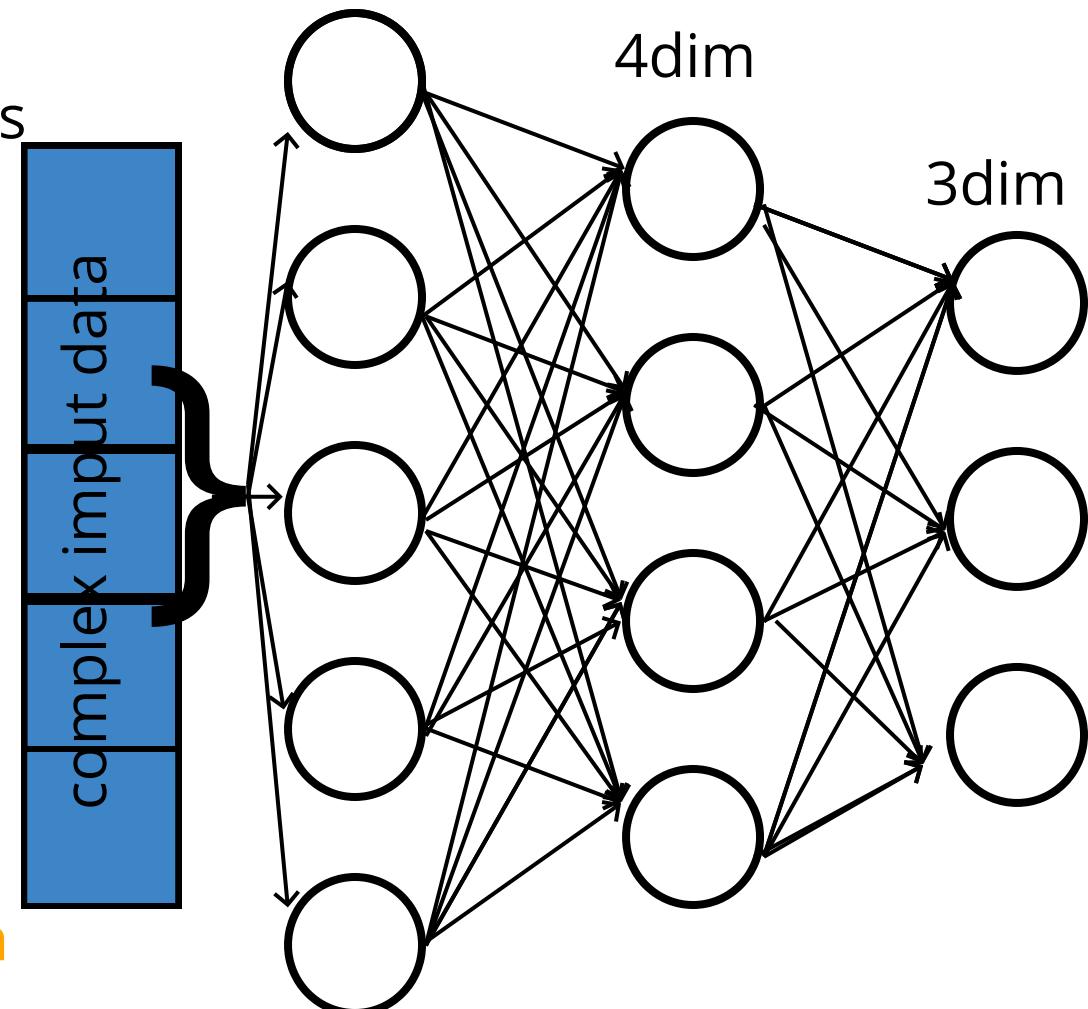
4

Autoencoders

Unsupervised learning with Neural Networks

What do NN do? approximate complex functions with series of linear functions

5dim representation



.... so if my layers are smaller what I have is a compact representation of the data

Unsupervised learning with Neural Networks

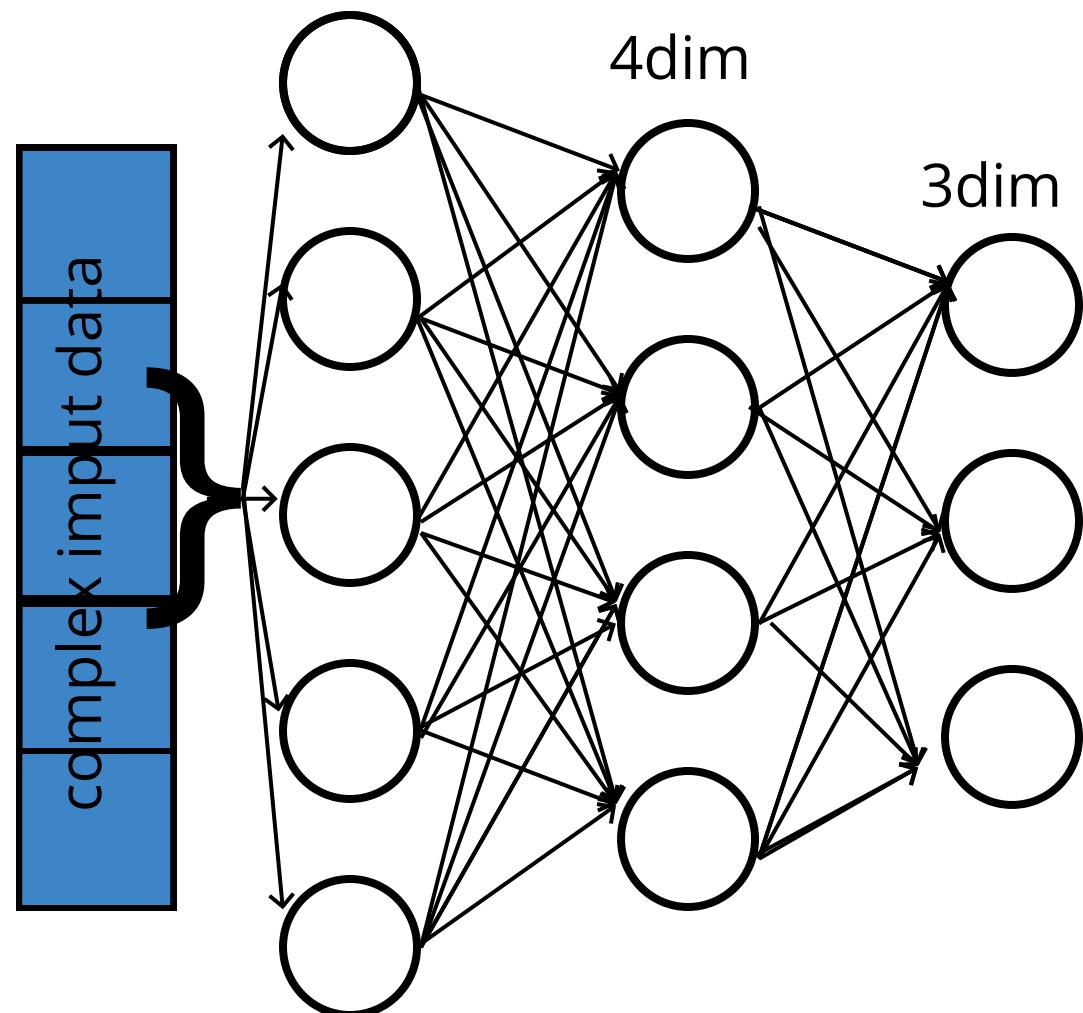
What do NN do? approximate complex functions with series of linear functions

To do that they extract information from the data

Each layer of the DNN produces a representation of the data a "latent representation".

.... so if my layers are smaller what I have is a compact representation of the data

5dim representation



Unsupervised learning with Neural Networks

What do NN do? approximate complex functions with series of linear functions

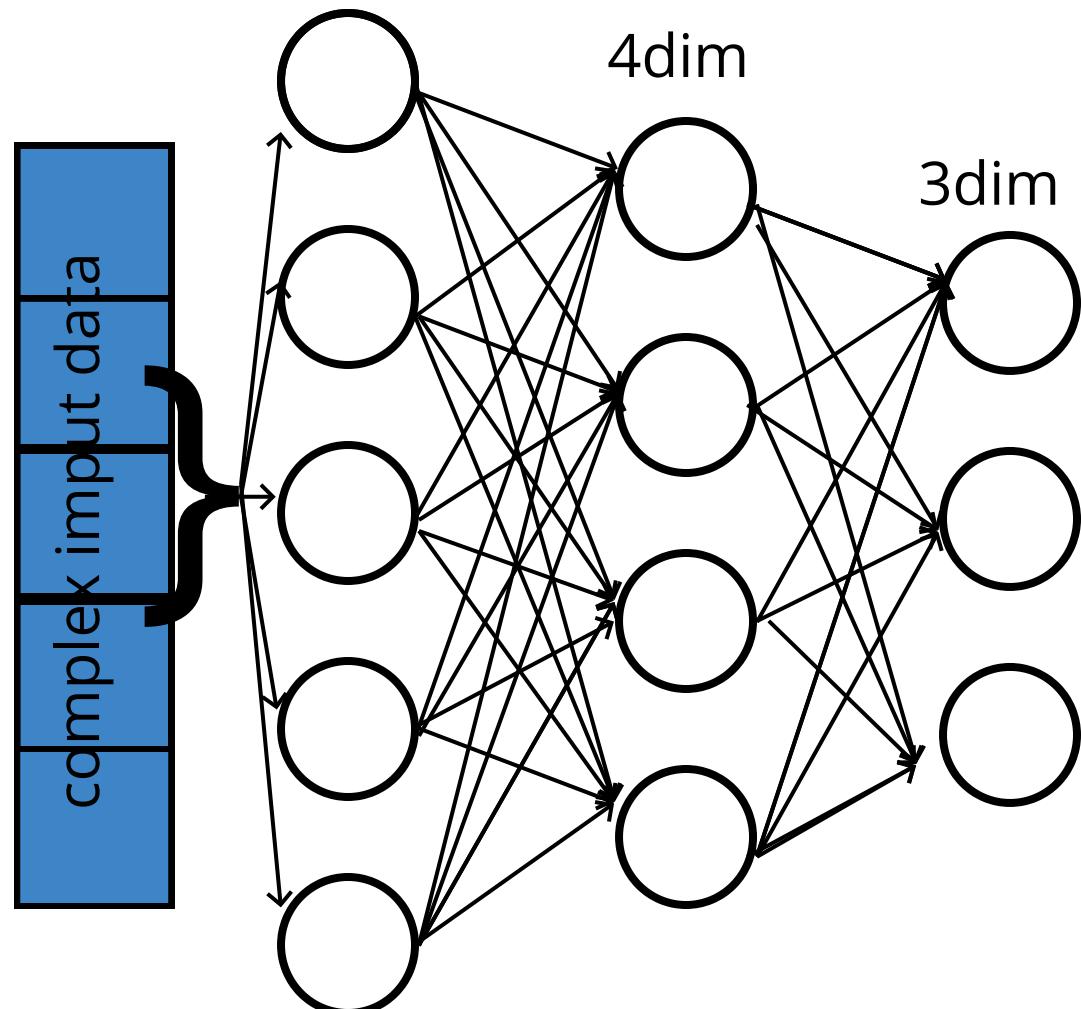
To do that they extract information from the data

Each layer of the DNN produces a representation of the data a "latent representation".

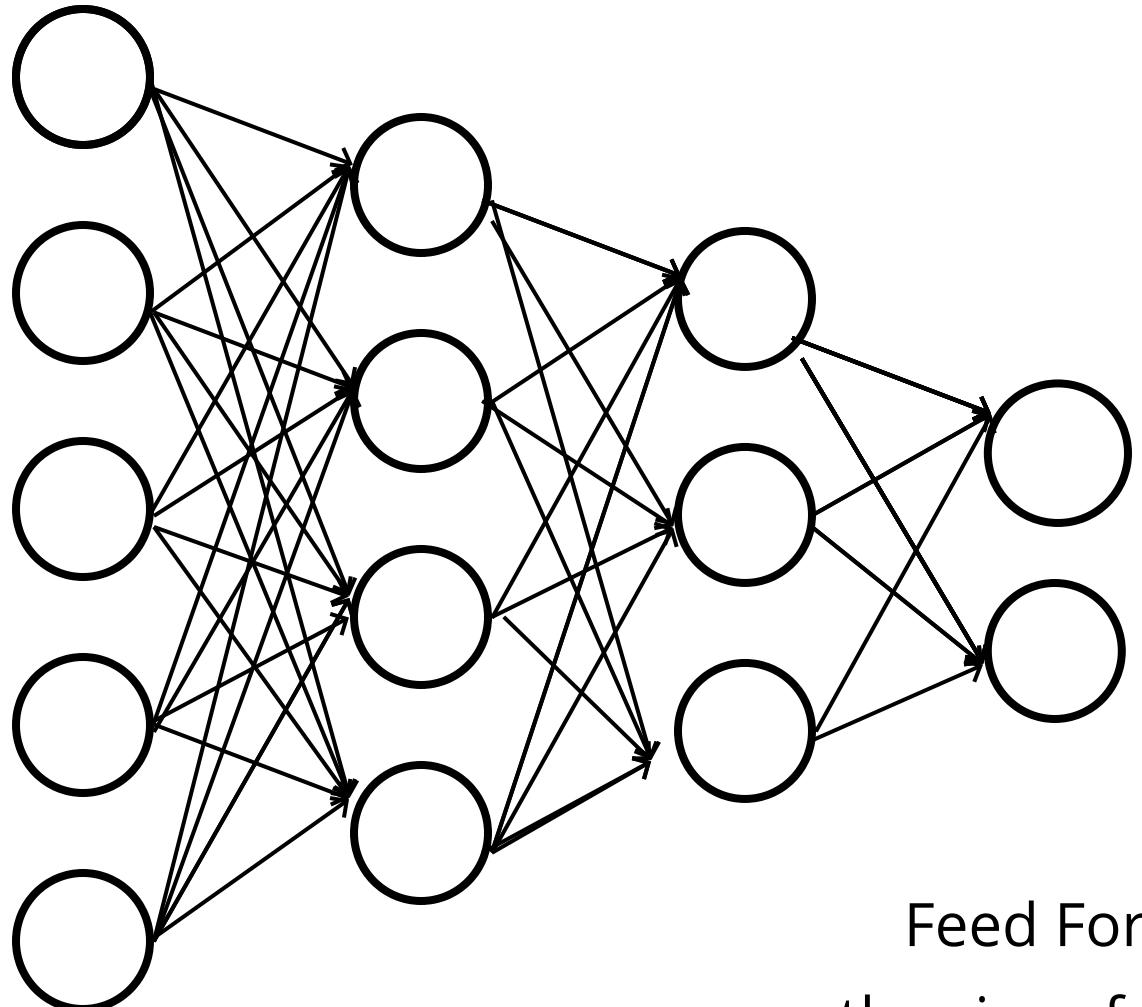
The dimensionality of that latent representation is determined by the size of the layer (and its connectivity, but we will ignore this bit for now)

.... so if my layers are smaller what I have is a compact representation of the data

5dim representation

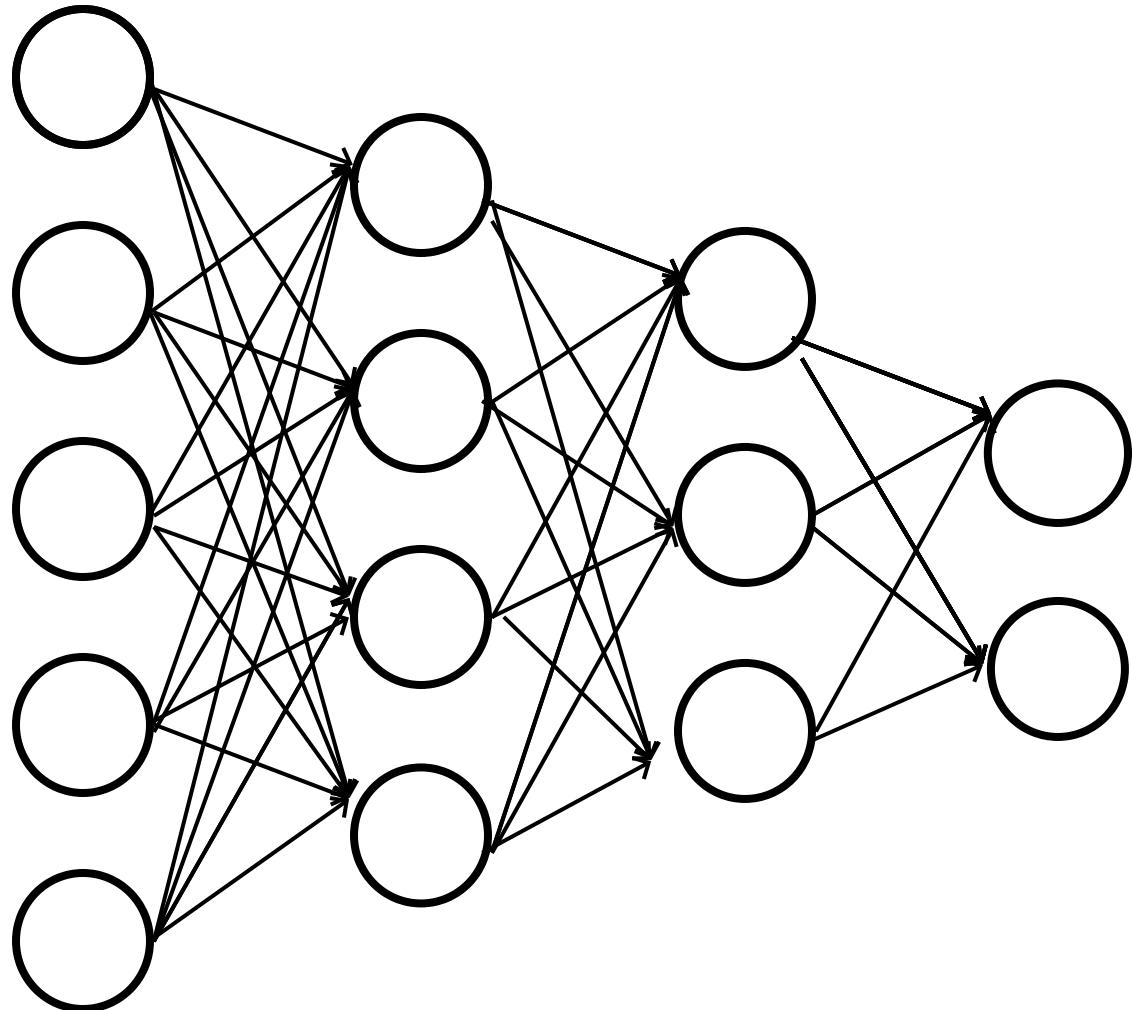


Autoencoder Architecture



Feed Forward DNN:
the size of the input is 5,
the size of the last layer is 2

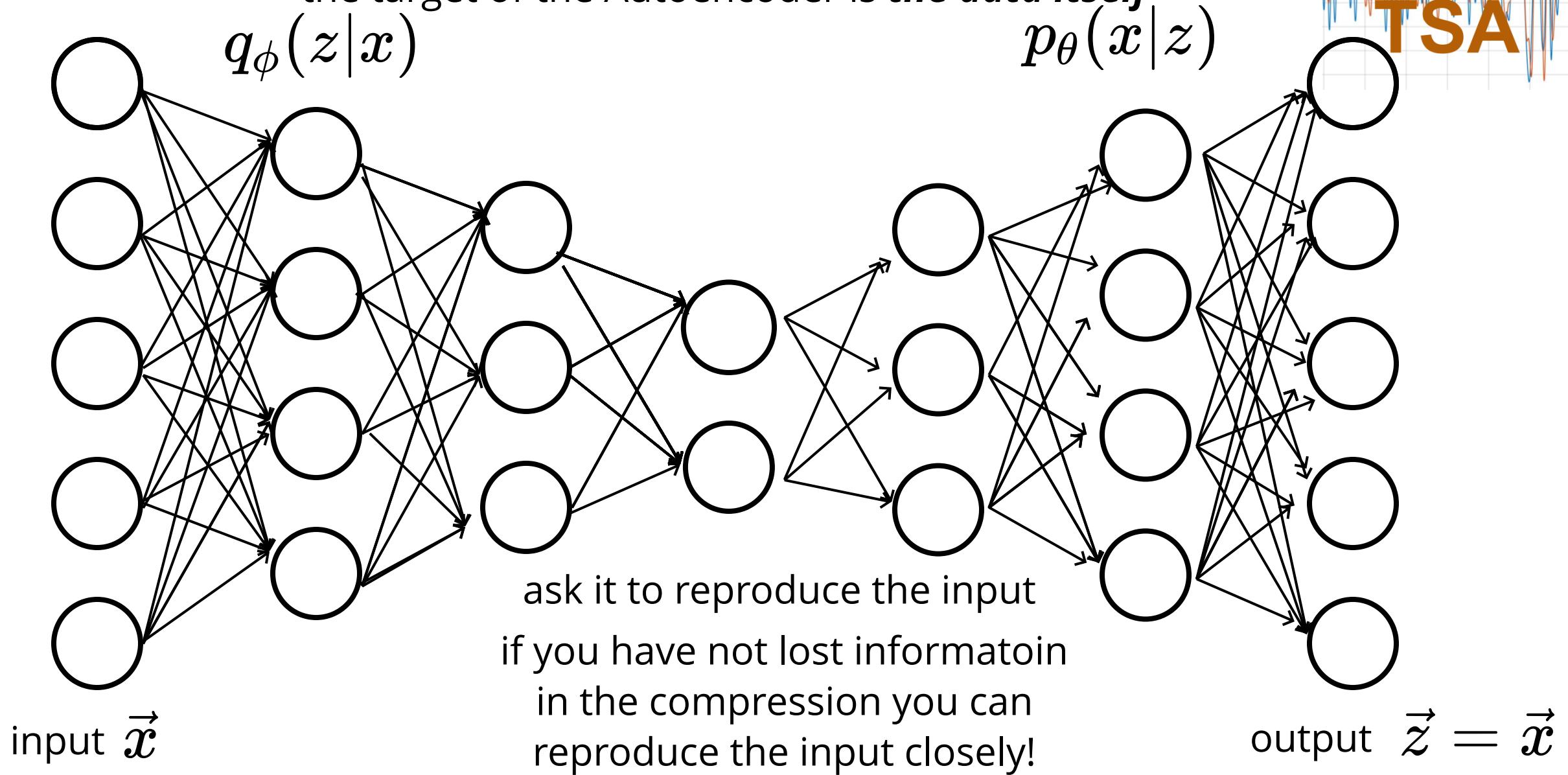
Autoencoder Architecture



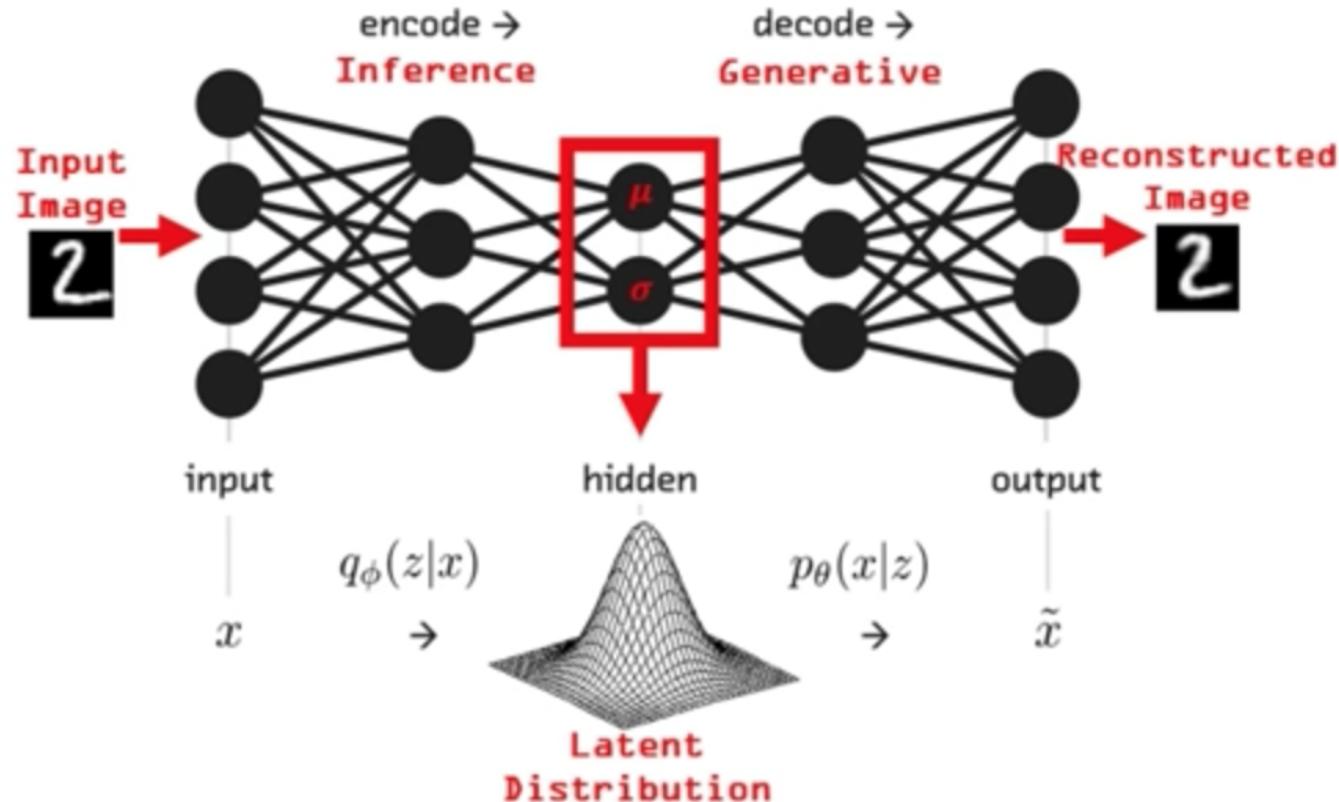
replicat the same structure backwards

Autoencoder Architecture

the target of the Autoencoder is ***the data itself***

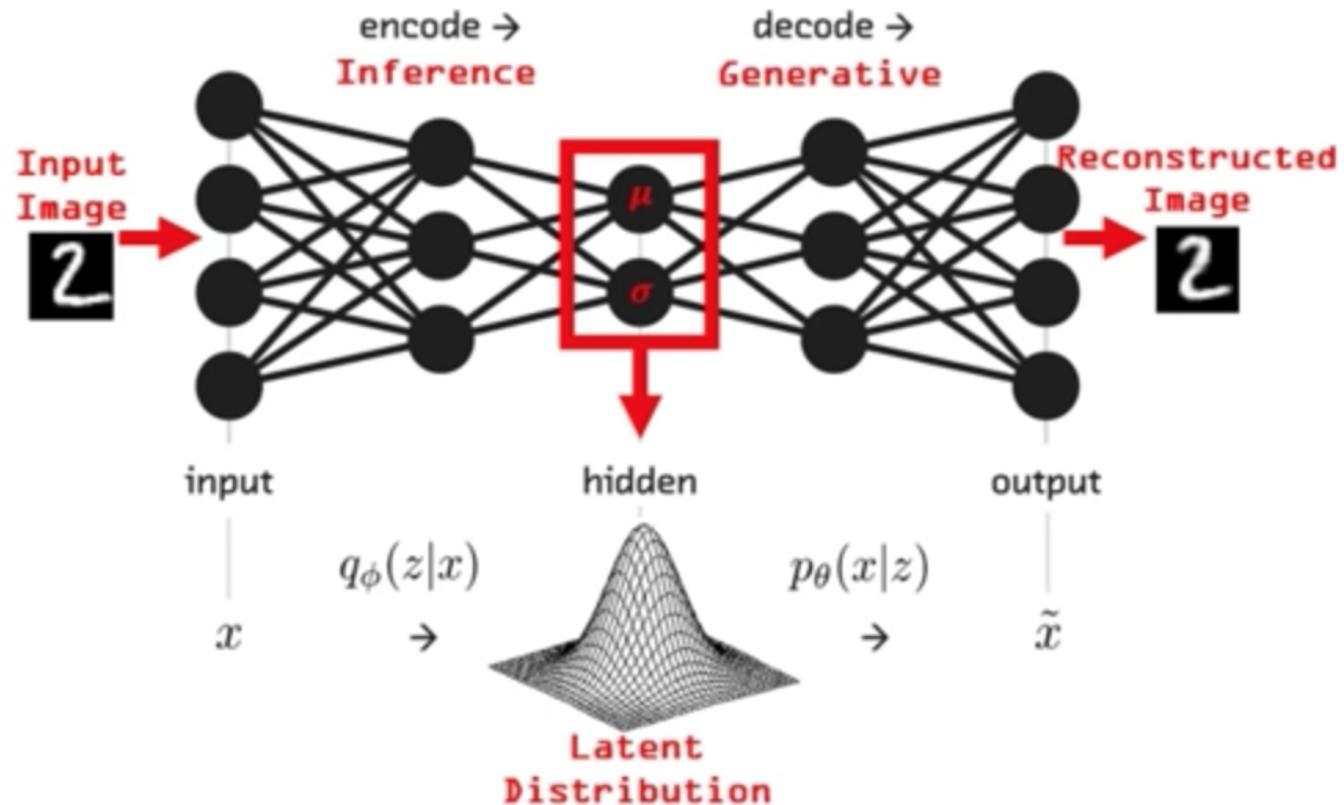


Autoencoder Architecture



- **Encoder:** outputs a lower dimensional representation z of the data x (similar to PCA, tSNE...)
- **Decoder:** Learns how to reconstruct x given z : learns $p(x|z)$

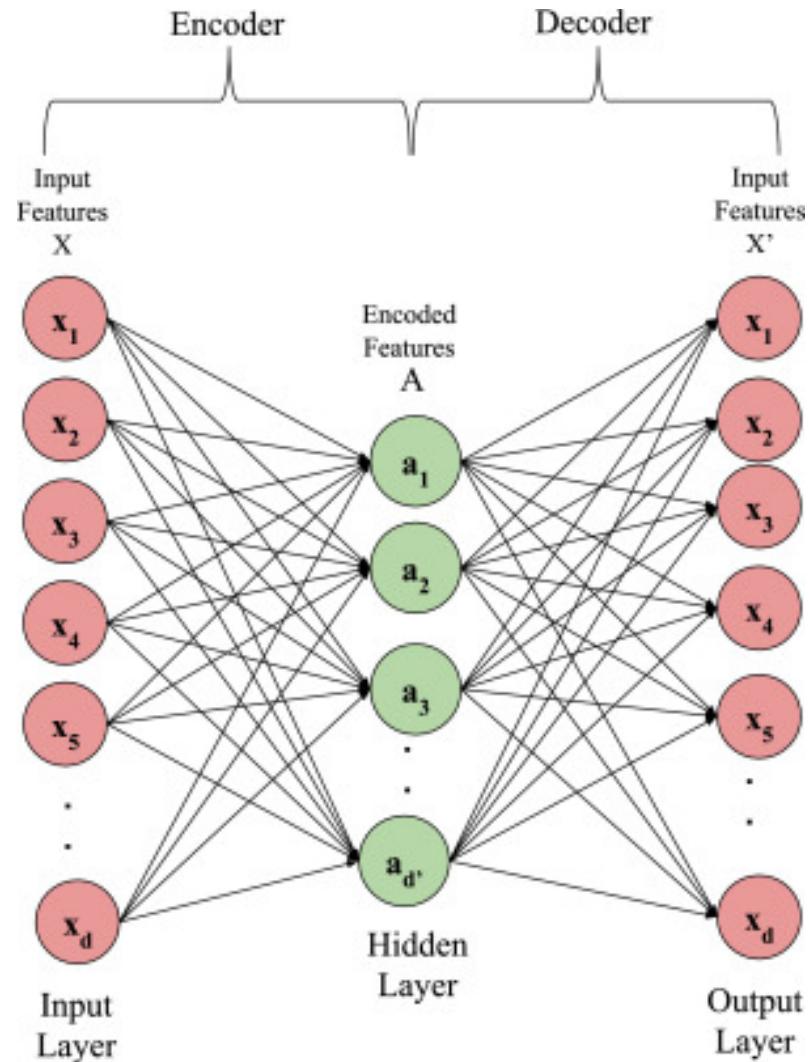
Autoencoder Architecture



```
1 from keras.layers import Dense, Flatten, Reshape
2 from keras.models import Sequential, Model
3
4 def build_autoencoder(image_shape, bn_size):
5     # Encoder
6     encoder = Sequential()
7     encoder.add(InputLayer(img_shape))
8     encoder.add(Flatten())
9     encoder.add(Dense(bn_size))
10
11    # Decoder
12    decoder = Sequential()
13    decoder.add(InputLayer((bn_size,)))
14    decoder.add(Dense(np.prod(image_shape)))
15    decoder.add(Reshape(image_shape))
16
```

Autoencoder Architecture

https://link.springer.com/chapter/10.1007/978-981-13-6661-1_3



Building a DNN with keras and tensorflow



Trivial to build, but the devil is in the details!

Building a DNN with keras and tensorflow



Trivial to build, but the devil is in the details!

```
1 from keras.models import Sequential
2 #can upload pretrained models from keras.models
3 from keras.layers import Dense, Conv2D, MaxPooling2D
4 #create model
5 model = Sequential()
6
7
8 #create the model architecture by adding model layers
9 model.add(Dense(10, activation='relu', input_shape=(n_cols,)))
10 model.add(Dense(10, activation='relu'))
11 model.add(Dense(1))
12
13 #need to choose the loss function, metric, optimization scheme
14 model.compile(optimizer='adam', loss='mean_squared_error')
15
16 #need to learn what to look for - always plot the loss function!
17 model.fit(x_train, y_train, validation_data=(x_test, y_test),
18             epochs=20, batch_size=100, verbose=1)
19 #note that the model allows to give a validation test,
20 #this is for a 3fold cross valiation: train-validate-test
21 #model.evaluate
```

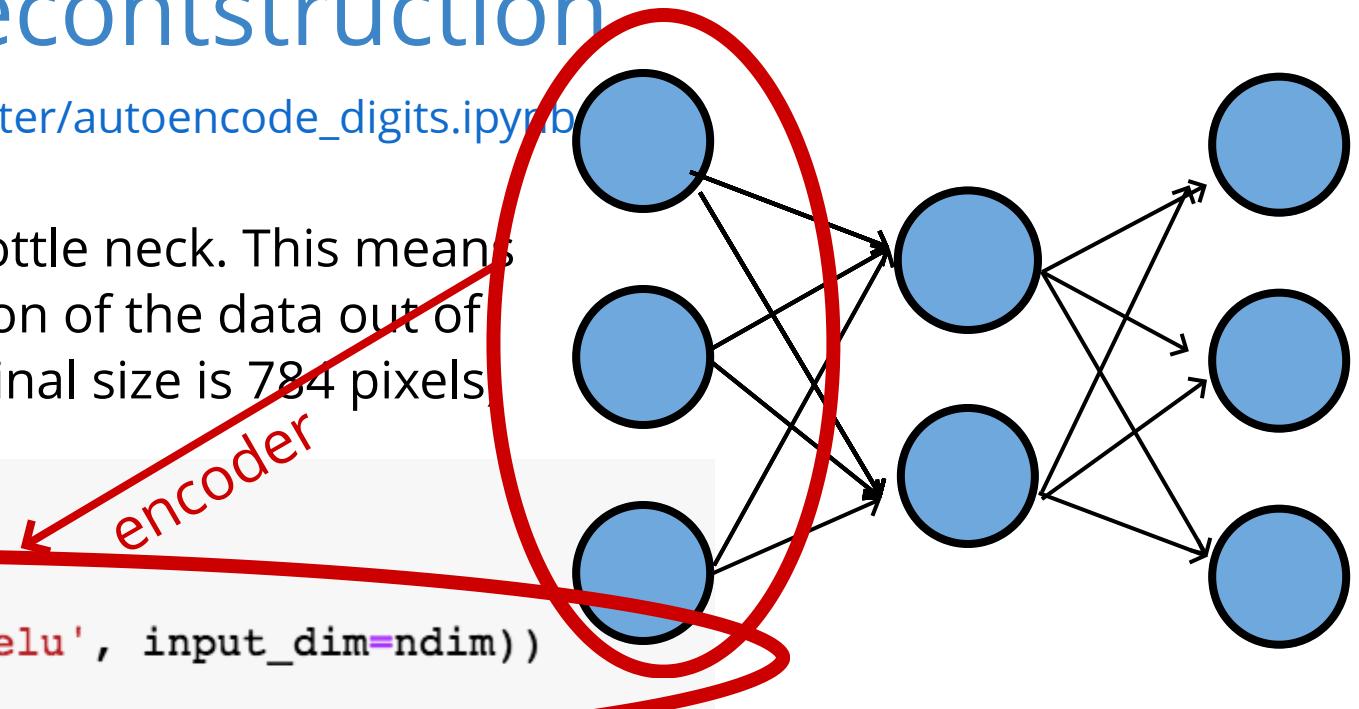


Building a DNN with keras and tensorflow autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

This autoencoder model has a 64-neuron bottle neck. This means it will generate a compressed representation of the data out of that layer which is 16-dimensional (the original size is 784 pixels).

```
model_digits64 = Sequential()
## encoder
# input layer and the output size
model_digits64.add(Dense(128, activation='relu', input_dim=ndim))
#compression layer
model_digits64.add(Dense(64, activation='relu'))
## deencoder
#decompression layer, same size as in the encoder
model_digits64.add(Dense(128, activation='relu'))
#output layer, same size as input
model_digits64.add(Dense(ndim, activation='linear'))
```



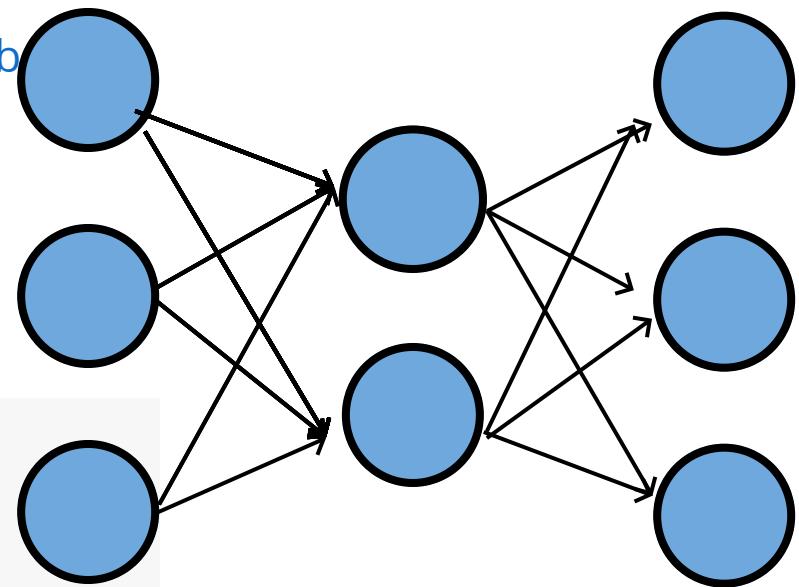
Building a DNN with keras and tensorflow autoencoder for image reconstruction



https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

This autoencoder model has a 64-neuron bottle neck. This means it will generate a compressed representation of the data out of that layer which is 16-dimensional (the original size is 784 pixels)

```
model_digits64 = Sequential()
## encoder
# input layer and the output size
model_digits64.add(Dense(128, activation='relu', input_dim=ndim))
#compression layer
model_digits64.add(Dense(64, activation='relu'))
## deencoder
#decompression layer, same size as in the encoder
model_digits64.add(Dense(128, activation='relu'))
#output layer, same size as input
model_digits64.add(Dense(ndim, activation='linear'))
```



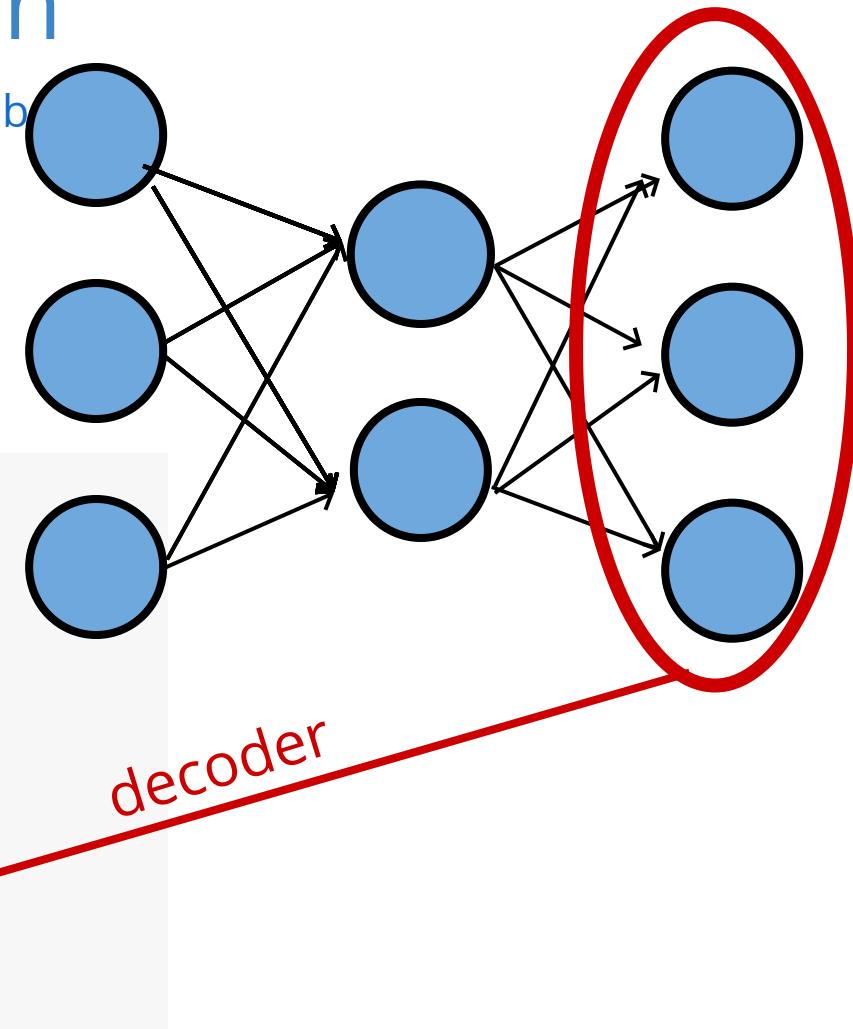
Building a DNN with keras and tensorflow autoencoder for image reconstruction



https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

This autoencoder model has a 64-neuron bottle neck. This means it will generate a compressed representation of the data out of that layer which is 16-dimensional (the original size is 784 pixels)

```
model_digits64 = Sequential()
## encoder
# input layer and the output size
model_digits64.add(Dense(128, activation='relu', input_dim=ndim))
#compression layer
model_digits64.add(Dense(64, activation='relu'))
## deencoder
#decompression layer, same size as in the encoder
model_digits64.add(Dense(128, activation='relu'))
#output layer, same size as input
model_digits64.add(Dense(ndim, activation='linear'))
```



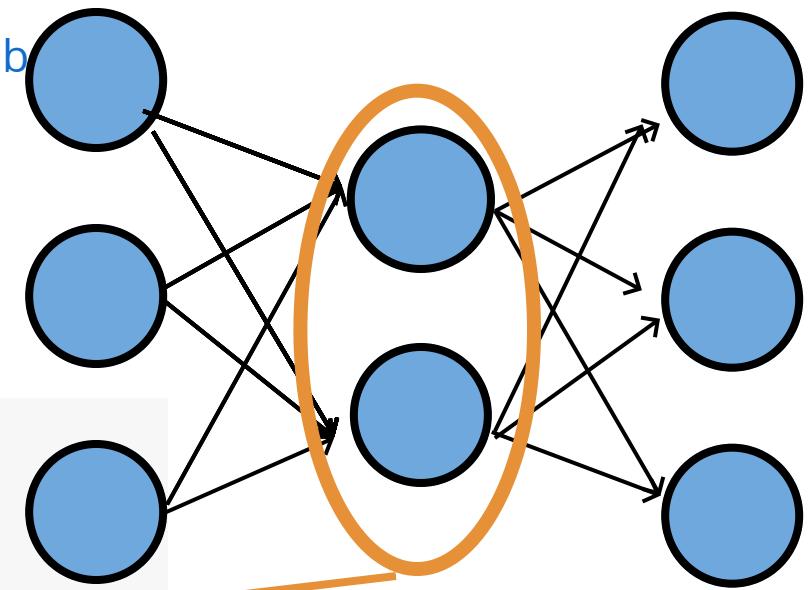


Building a DNN with keras and tensorflow autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

This autoencoder model has a 64-neuron bottle neck. This means it will generate a compressed representation of the data out of that layer which is 16-dimensional (the original size is 784 pixels)

```
model_digits64 = Sequential()
## encoder
# input layer and the output size
model_digits64.add(Dense(128, activation='relu', input_dim=ndim))
#compression layer
model_digits64.add(Dense(64, activation='relu')) ← bottle neck
## deencoder
#decompression layer, same size as in the encoder
model_digits64.add(Dense(128, activation='relu'))
#output layer, same size as input
model_digits64.add(Dense(ndim, activation='linear'))
```





Building a DNN with keras and tensorflow autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

```
# choose the optimizer and loss appropriately!
model_digits64.compile(optimizer="adadelta", loss="mean_squared_error")
```

```
print(model_digits64.summary())
```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 128)	100480
dense_2 (Dense)	(None, 64)	8256
dense_3 (Dense)	(None, 128)	8320
dense_4 (Dense)	(None, 784)	101136

Total params: 218,192

Trainable params: 218,192

Non-trainable params: 0

None

This simple odel has 200000 parameters!

My original choice is to train it with "adadelta" with a mean squared loss function, all activation functions are relu, appropriate for a linear regression



Building a DNN with keras and tensorflow autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

3.0.1 regression

- loss='mean_squared_error' L2: default loss to use for regression problems. => linear activation funct in output layer, one node out

alternatives: loss='mean_squared_logarithmic_error', 'mean_absolute_error' (which is L1 instead of L2)

3.0.2 binary classification

- loss='binary_crossentropy' => sigmoid activation function in output layer, one node out

alternatives: 'hinge'

3.0.3 multiclass classification

categorical encoded as numerical

- loss='categorical_crossentropy' => softmax n nodes out

onehot encoded categorical

- 'parse_categorical_crossentropy' => softmax n nodes out
- 'kullback Leibler Divergence Loss' => probabilistic categorical classification; $\log(P/Q)$

What should I choose for the loss function and how does that relate to the activation function and optimization?

Building a DNN with keras and tensorflow autoencoder for image reconstruction



https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

What should I choose for the loss function and how does that relate to the activation function and optimization?

loss	good for	activation last layer	size last layer
mean_squared_error	regression	linear	one node
mean_absolute_error	regression	linear	one node
mean_squared_logarithmit_error	regression	linear	one node
binary_crossentropy	binary classification	sigmoid	one node
categorical_crossentropy	multiclass classification	sigmoid	N nodes
Kullback_Divergence	multiclass classification, probabilistic interpretation	sigmoid	N nodes

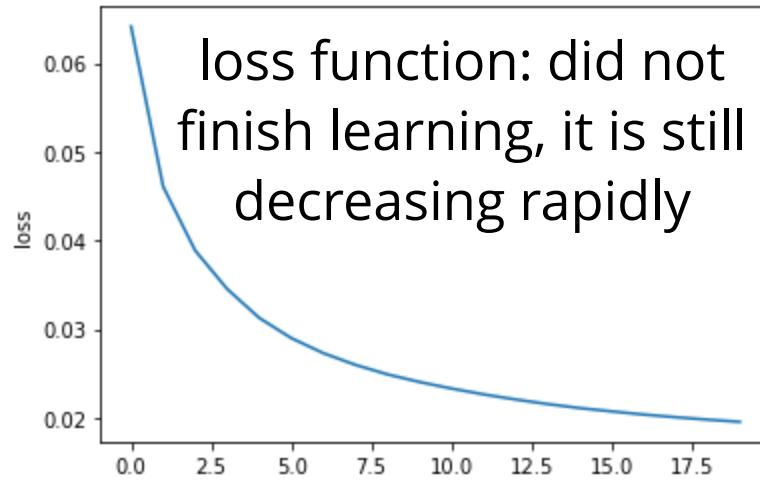
autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

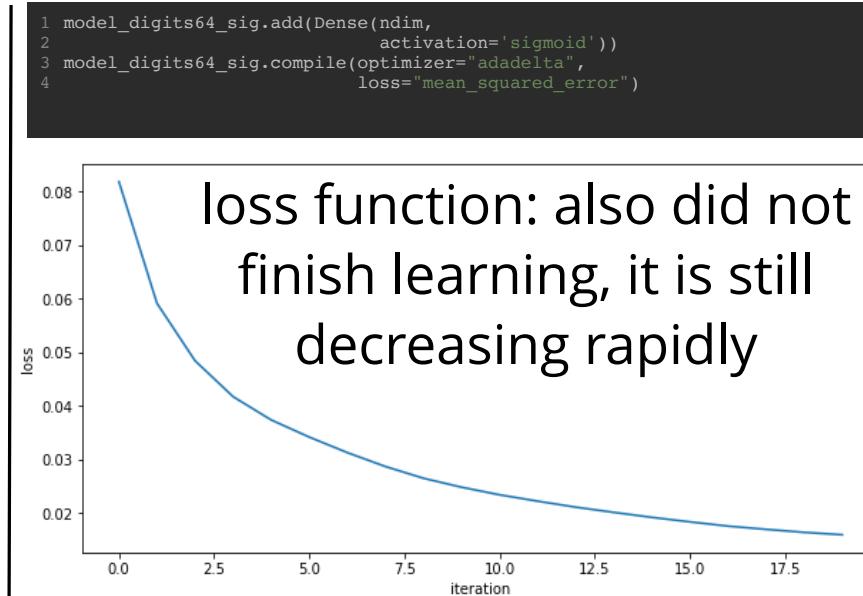
```
1 model_digits64.add(Dense(ndim,  
2                         activation='linear'))  
3 model_digits64_sig.compile(optimizer="adadelta",  
4                           loss="mean_squared_error")
```

```
1 model_digits64_sig.add(Dense(ndim,  
2                             activation='sigmoid'))  
3 model_digits64_sig.compile(optimizer="adadelta",  
4                           loss="mean_squared_error")
```

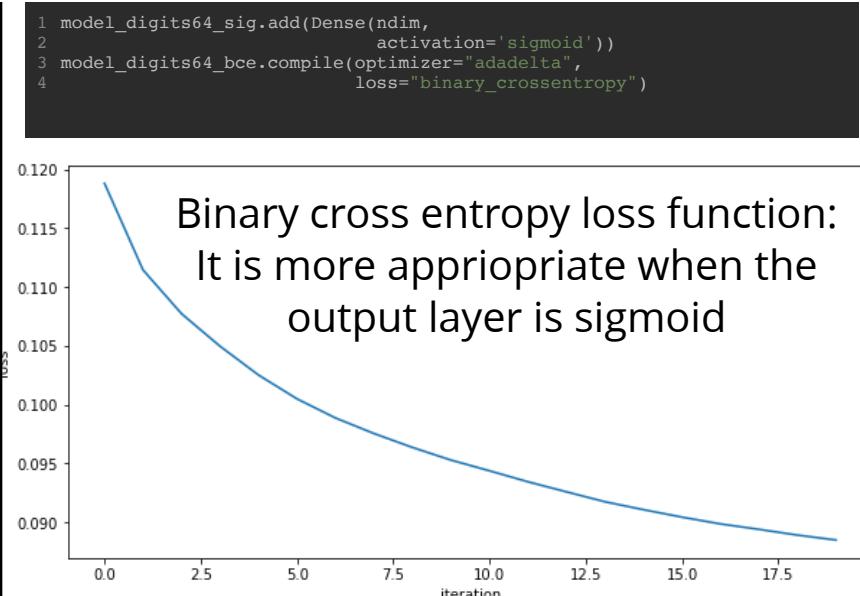
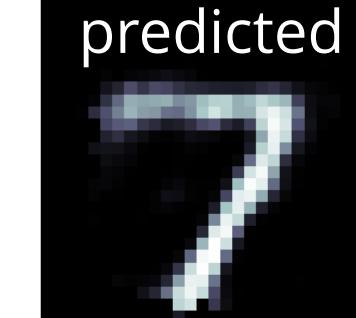
```
1 model_digits64_sig.add(Dense(ndim,  
2                             activation='sigmoid'))  
3 model_digits64_bce.compile(optimizer="adadelta",  
4                           loss="binary_crossentropy")
```



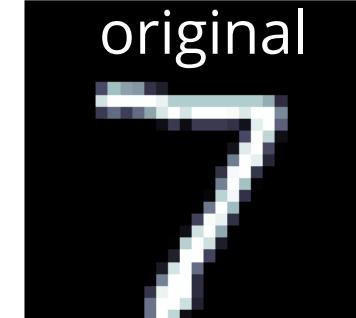
The predictions are far too detailed. While the input is not binary, it does not have a lot of details. Maybe approaching it as a binary problem (with a sigmoid and a binary cross entropy loss) will give better results



A sigmoid gives activation gives a much better result!



Even better results!



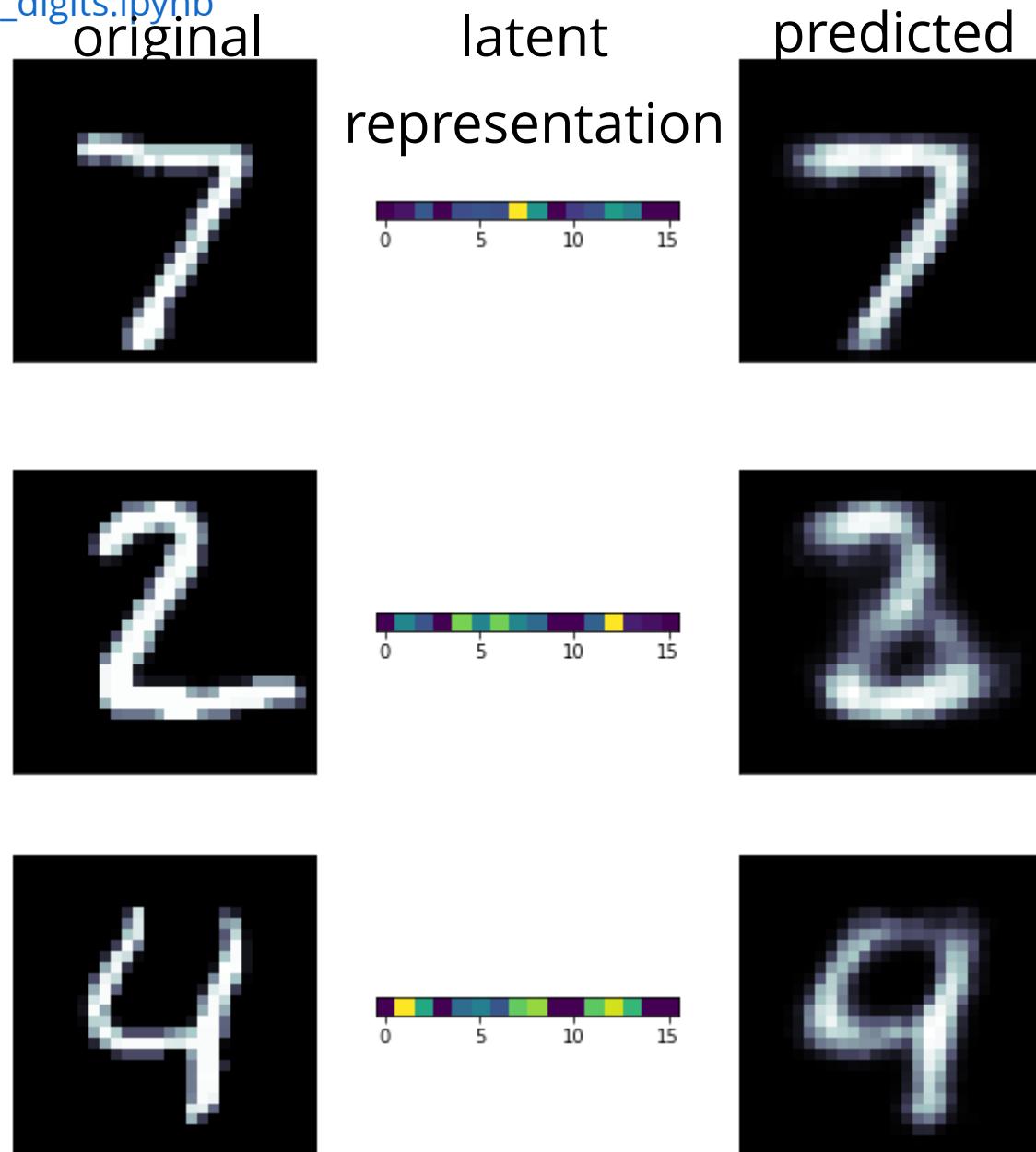
autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

A more ambitious model has a 16 neurons bottle neck: we are trying to extract 16 numbers to reconstruct the entire image! its pretty remarkable! those 16 number are **extracted features** from the data

Model: "sequential_2"

Layer (type)	Output Shape	Param #
dense_9 (Dense)	(None, 128)	100480
dense_10 (Dense)	(None, 64)	8256
dense_11 (Dense)	(None, 32)	2080
dense_12 (Dense)	(None, 16)	528
dense_13 (Dense)	(None, 32)	544
dense_14 (Dense)	(None, 64)	2112
dense_15 (Dense)	(None, 128)	8320
dense_16 (Dense)	(None, 784)	101136
<hr/>		
Total params: 223,456		
Trainable params: 223,456		
Non-trainable params: 0		



autoencoder for image reconstruction

https://github.com/fedhere/MLTSA_FBianco/blob/master/autoencode_digits.ipynb

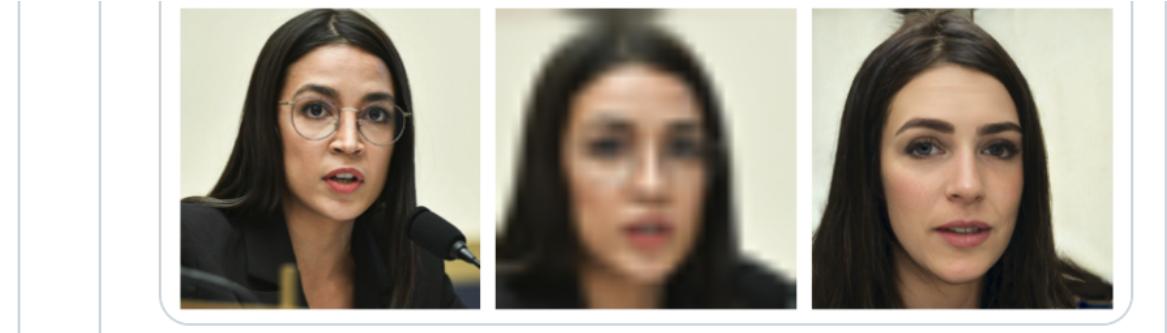


🔥👤Robert Osazuwa Ness👤🔥 @osazuwa · Jun 20, 2020 🐦

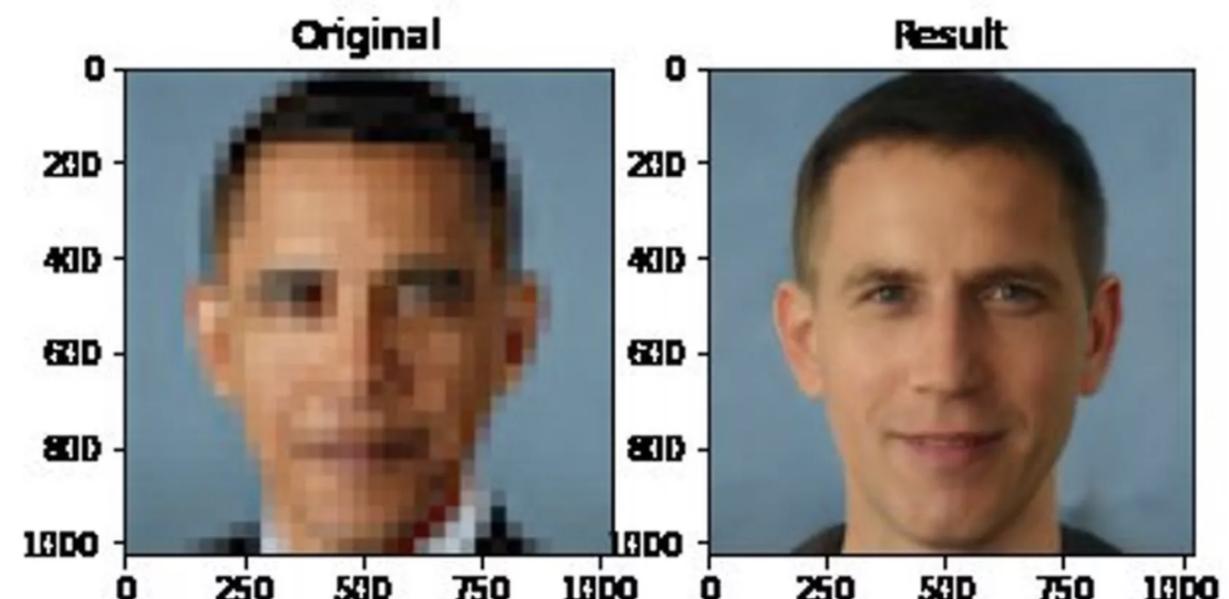
An image of [@BarackObama](#) getting upsampled into a white guy is floating around because it illustrates racial bias in [#MachineLearning](#). Just in case you think it isn't real, it is, I got the code working locally. Here is me, and here is [@AOC](#).



🔥👤Robert Osazuwa Ness👤🔥
@osazuwa



🔥👤Robert Osazuwa Ness👤🔥
@osazuwa



~~models are neutral, the bias is in the data~~

The bias is in the **data**

The bias is in the **models** and
the decision we make

The bias is in **how we choose to**
optimize our model

**The bias is society that
provides the framework to
validate our biased models**

Should AI reflect
who we are
(and *enforce and grow our bias*)
or should it reflect who we
aspire to be?
(*and who decides what that is?*)

none of this is new

~~models are neutral,~~
~~the bias is in the data~~

The bias is in the **data**

The bias is in the **models** and
the decision we make

The bias is in **how we choose to**
optimize our model

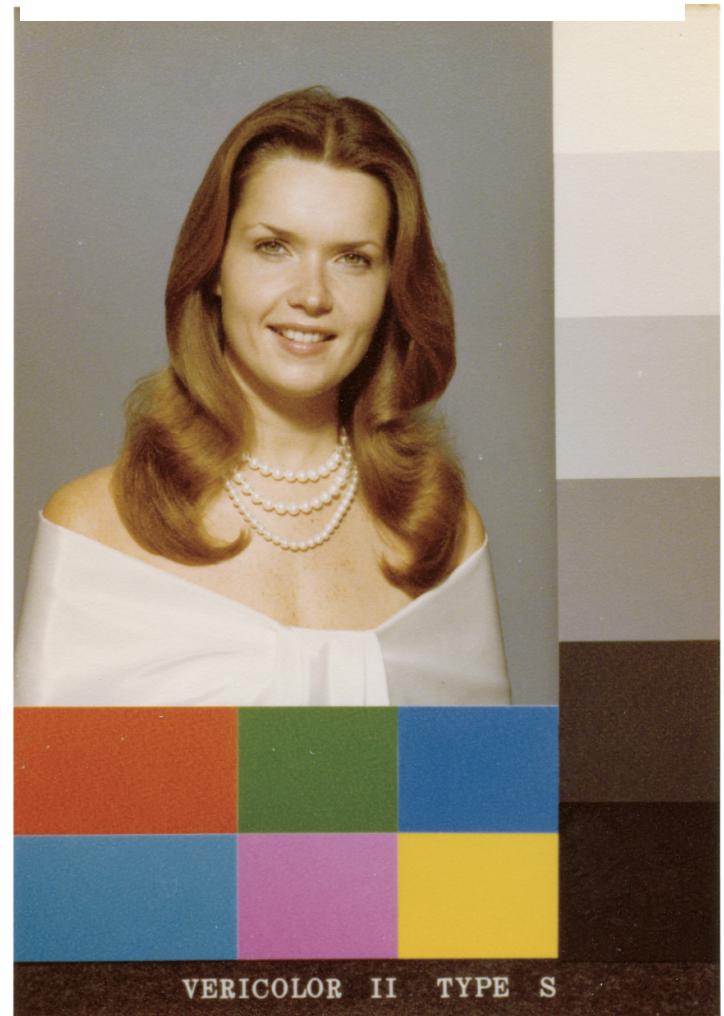
**The bias is society that
provides the framework to
validate our biased models**

Lens

LENS

The Racial Bias Built Into
Photography

Sarah Lewis explores the relationship between racism and the camera.



Shirley Card, 1978. Courtesy of Hermann Zschiegner

[https://www.nytimes.com/2019/04/25/lens/sar ah-lewis-racial-bias-photography.html](https://www.nytimes.com/2019/04/25/lens/sarah-lewis-racial-bias-photography.html)

Neural Network and Deep Learning

an excellent and free book on NN and DL

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

Deep Learning An MIT Press book in preparation

Ian Goodfellow, Yoshua Bengio and Aaron Courville

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

Resources

History of NN
<https://cs.stanford.edu/people/eroberts/courses/soco/projects/neural-networks/History/history2.html>

Gradient Descent

https://ml-cheatsheet.readthedocs.io/en/latest/gradient_descent.html

resources