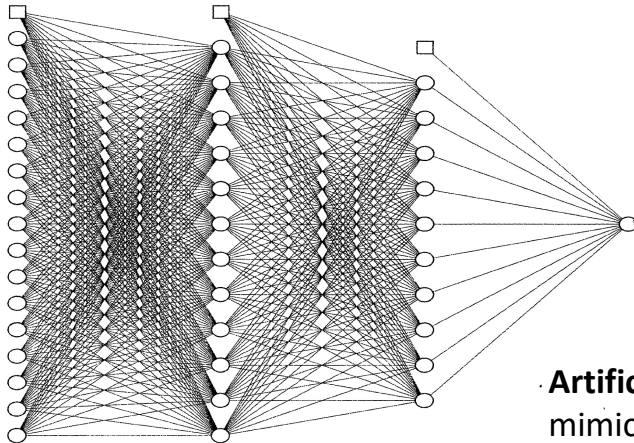
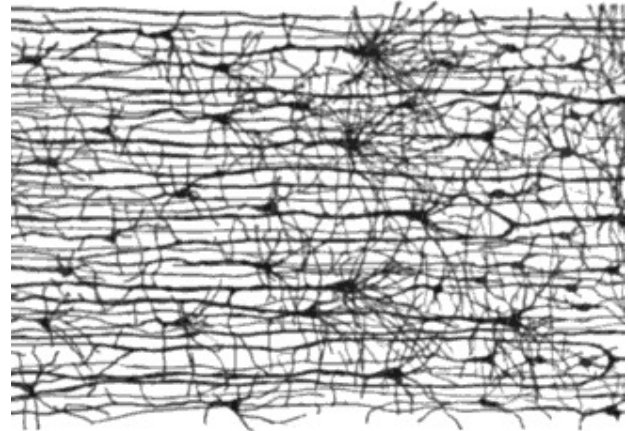
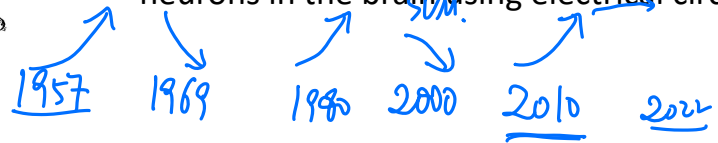


Section 13 Neural Network

- Perceptron
- Neural Network
- Backpropagation



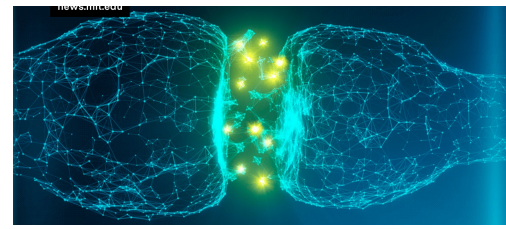
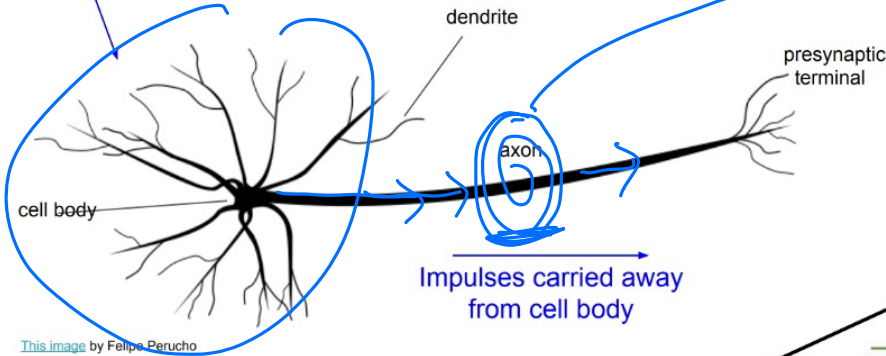
Human Neural Networks was introduced in 1943 by neurophysiologist Warren McCulloch and mathematician Walter Pitts to model neurons in the brain using electrical circuits.



Artificial Neural networks are a series of algorithms that mimic the operations of a human brain to recognize relationships between vast amounts of data. It's a very broad term that encompasses any form of Deep Learning model.

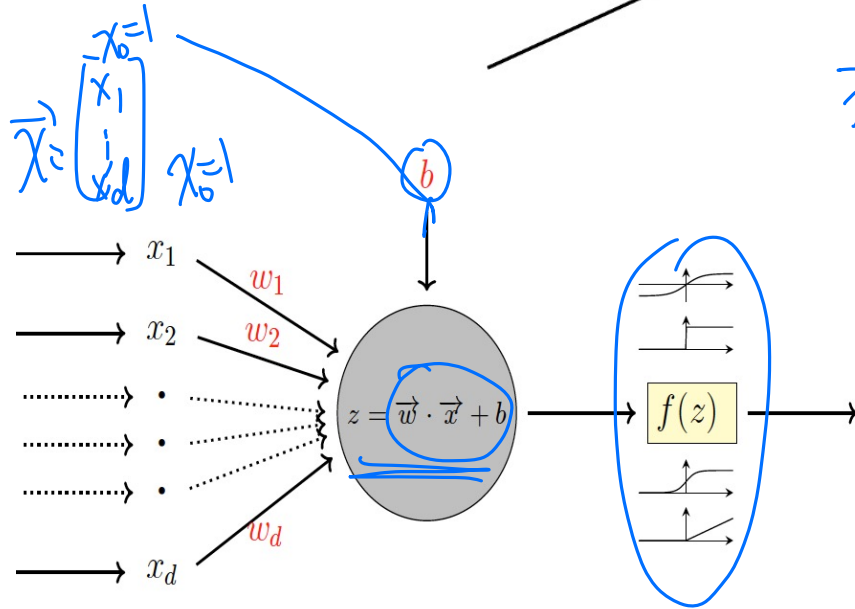
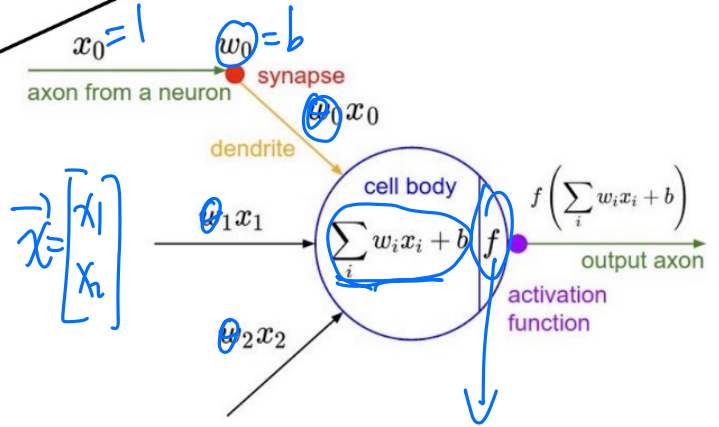
Impulses carried toward cell body

A biological Neuron



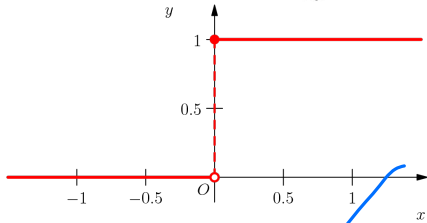
This image by Felipe Perucho is licensed under CC-BY 3.0

An artificial Neuron



Activation function example: Heaviside step function

$$f(x) = \text{step}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



➤ Other activation functions:

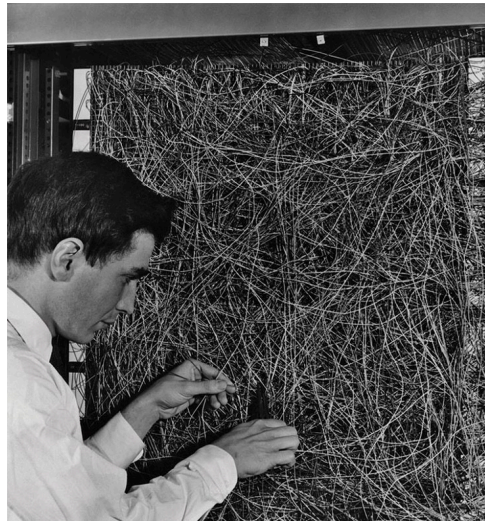
Handwritten notes:
 $z < 0$
 $z > 0$

Sigmoid	Tanh	ReLU	Leaky ReLU
$g(z) = \frac{1}{1 + e^{-z}}$	$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	$g(z) = \max(0, z)$	$g(z) = \max(\epsilon z, z)$ with $\epsilon \ll 1$

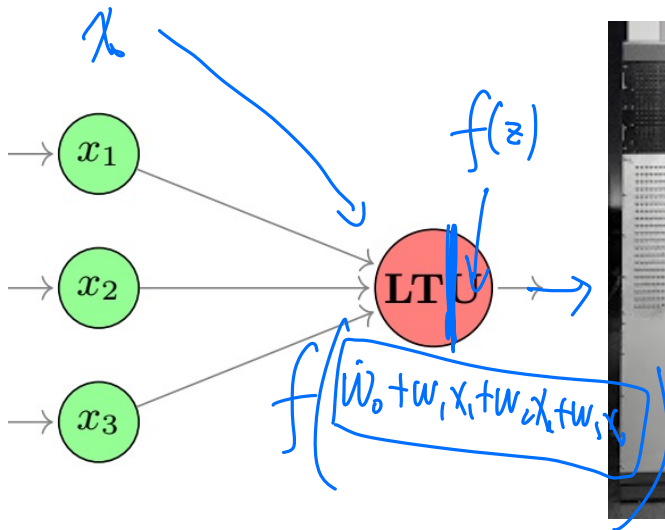
Handwritten notes:
 $z > 0$
 $\epsilon z < z$
 $\epsilon = 0.01$

➤ Perceptron

The perceptron is based around a *linear threshold unit* (LTU).



In 1957 Frank Rosenblatt designed and invented the perceptron which is a type of neural network. A neural network acts like your brain; the brain contains billions of cells called neurons that are connected together in a network. The perceptron connects a web of points where simple decisions are made, which come together in the larger program to solve more complex problems.



The New Yorker, December 6, 1958 P. 44

Talk story about the perceptron, a new electronic brain which hasn't been built, but which has been successfully simulated on the I.B.M. 704. Talk with Dr. Frank Rosenblatt, of the Cornell Aeronautical Laboratory, who is one of the two men who developed the prodigy; the other man is Dr. Marshall C. Yovits, of the Office of Naval Research, in Washington. Dr. Rosenblatt defined the perceptron as the first non-biological object which will achieve an organization o its external environment in a meaningful way. It interacts with its environment, forming concepts that have not been made ready for it by a human agent. If a triangle is held up, the perceptron's eye picks up the image & conveys it along a random succession of lines to the response units, where the image is registered. It can tell the difference betw. a cat and a dog, although it wouldn't be able to tell whether the dog was to theleft or right of the cat. Right now it is of no practical use, Dr. Rosenblatt conceded, but he said that one day it might be useful to send one into outer space to take in impressions for us.

<https://www.newyorker.com/magazine/1958/12/06/rival-2>

<https://news.cornell.edu/stories/2019/09/professors-perceptron-paved-way-ai-60-years-too-soon>

➤ Perceptron

Training Data: $\mathcal{D} = (\vec{x}^{(i)}, y^{(i)})$ for $i = 1 \dots n$.

Assumptions:

- Binary classification (i.e. $y^{(i)} \in \{-1, +1\}$)
- Data is linearly separable, i.e., there exists a hyperplane that separates all the sample points in class A from classes B.

Classifier:

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} = \begin{bmatrix} -b \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

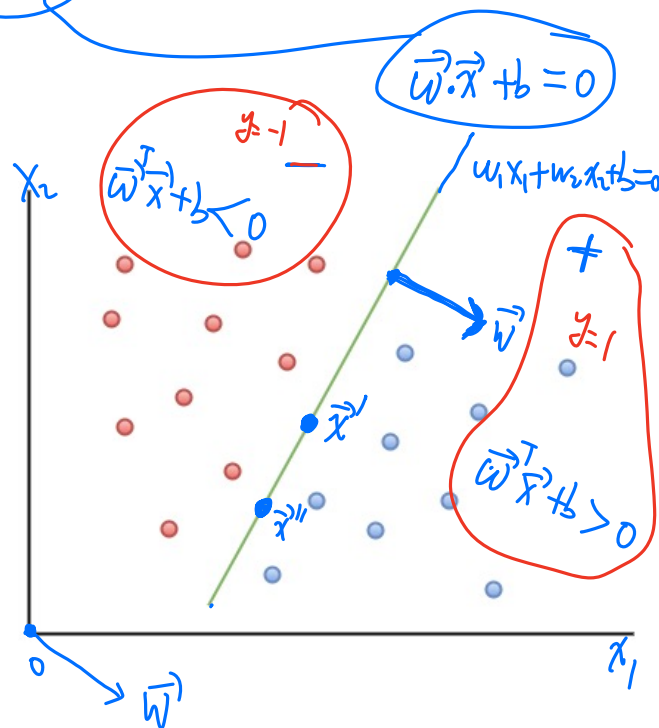
$$h(\vec{x}) = \text{sign}(\vec{\theta}^T \vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$$

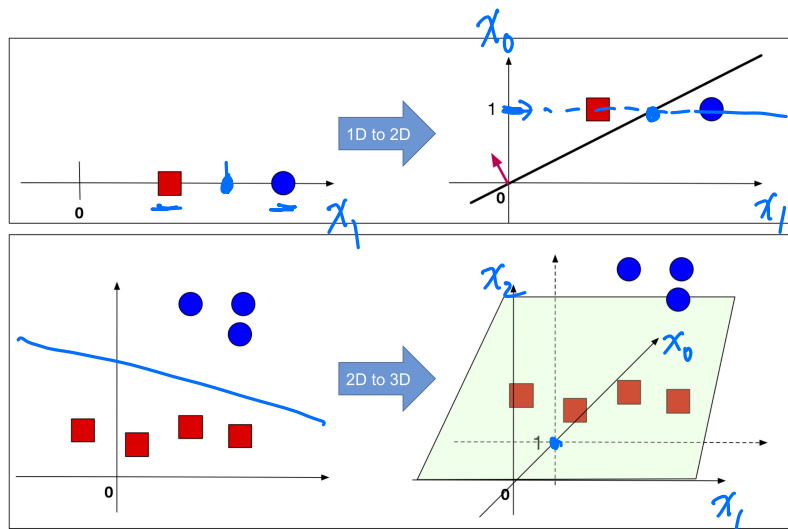
$$\text{sign}(z) = \begin{cases} 1 & \text{if } z > 0 \\ -1 & \text{if } z < 0 \end{cases}$$

$$\vec{w}^T \vec{x}' = -b$$

$$\vec{w}^T \vec{x}'' = -b$$

$$\Rightarrow \vec{w}^T (\vec{x}' - \vec{x}'') = 0$$





$$\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

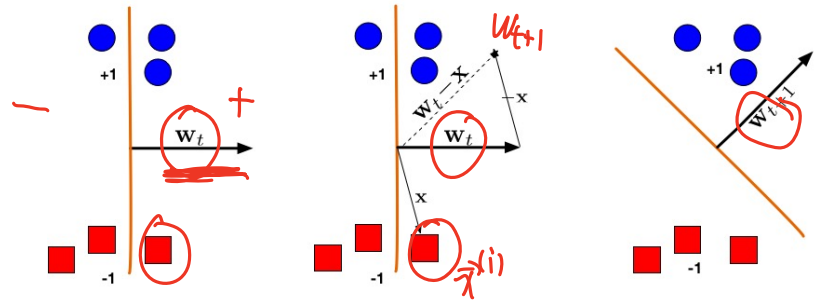
$$y^{(i)} = \{-1, 1\}$$

$$\underline{\underline{y^{(i)} (\vec{\theta}^T \vec{x}^{(i)}) > 0}} \Leftrightarrow \vec{x}^{(i)} \text{ classified correctly}$$

$$\text{Given: Find } \vec{w}^T \vec{x} + b = 0$$

Assume $b=0$

➤ Training the Perceptron



$\mathcal{D} = (\vec{x}^{(i)}, y^{(i)})$ for $i = 1 \dots n$.

Start with initial $\vec{\theta} = \vec{0}$

For $i = 1, \dots, n$

Repeat $\vec{\theta}^{next} := \vec{\theta} + \alpha (y^{(i)} - h_{\theta}(\vec{x}^{(i)})) \vec{x}^{(i)}$

- $w_{t+1} = w_t - \alpha \vec{x}^{(i)}$ if $(\vec{x}^{(i)}, \underline{y^{(i)} = -1})$
and $\underline{w_t^T \vec{x}^{(i)} > 0}$
- $w_{t+1} = w_t + \alpha \vec{x}^{(i)}$ if $y^{(i)} = 1, \text{ predict } -$

The perceptron updates its weights only on misclassified points.


```

def perceptron_sgd(X, Y):
    w = np.zeros(len(X[0])) #Initialize the weight vector for the perceptron with zeros
    eta = 1 #Set the learning rate to 1
    epochs = 20 #Set the number of epochs

    for t in range(epochs):
        for i, x in enumerate(X):
            if (np.dot(X[i], w)*Y[i]) <= 0:
                w = w + eta*X[i]*Y[i]

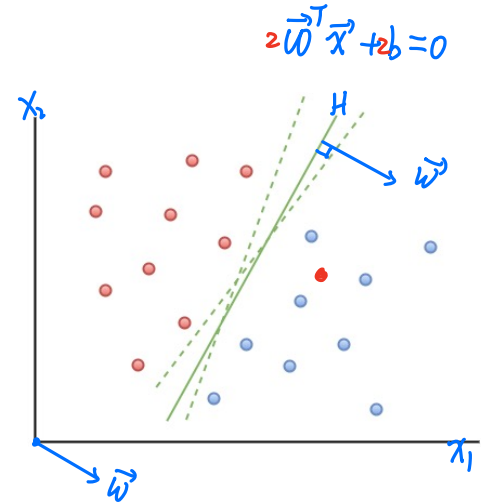
    return w

w = perceptron_sgd(X,y)
print(w)

```

The perceptron is a form of stochastic gradient decent on the loss function

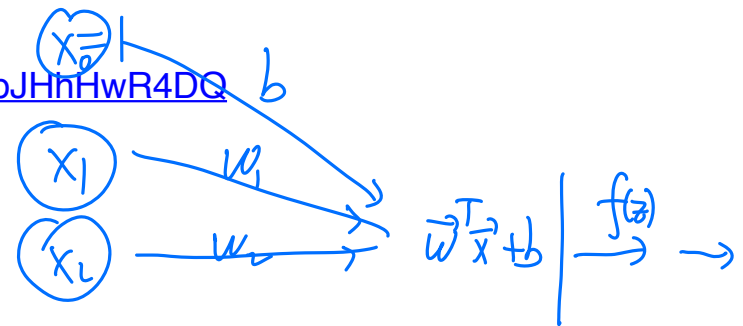
$$J(\vec{\theta}) = - \sum_{i=1}^n \left(y^{(i)} - h_{\theta}(\vec{x}^{(i)}) \right) (\vec{\theta}^T \vec{x}^{(i)})$$



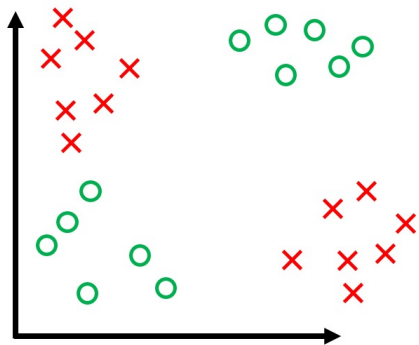
- When the data is separable, there are many solutions and which one is found depends on the starting value.
- The finite number of steps can be large, practically, if the gap is small the time to find it is large.
- When the data are not separable, the algorithm does not converge, and instead falls into a cycle.

Video illustration for perceptron:

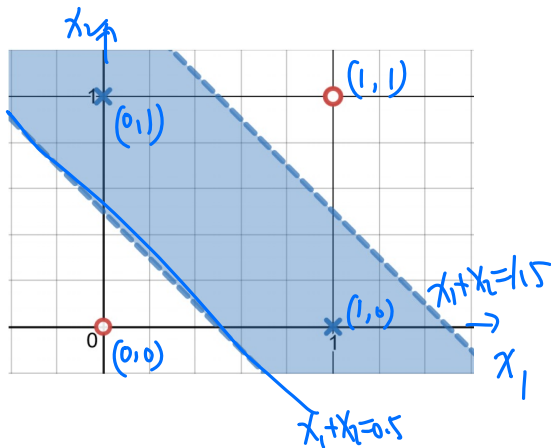
<https://www.youtube.com/watch?v=xpJHhHwR4DQ>



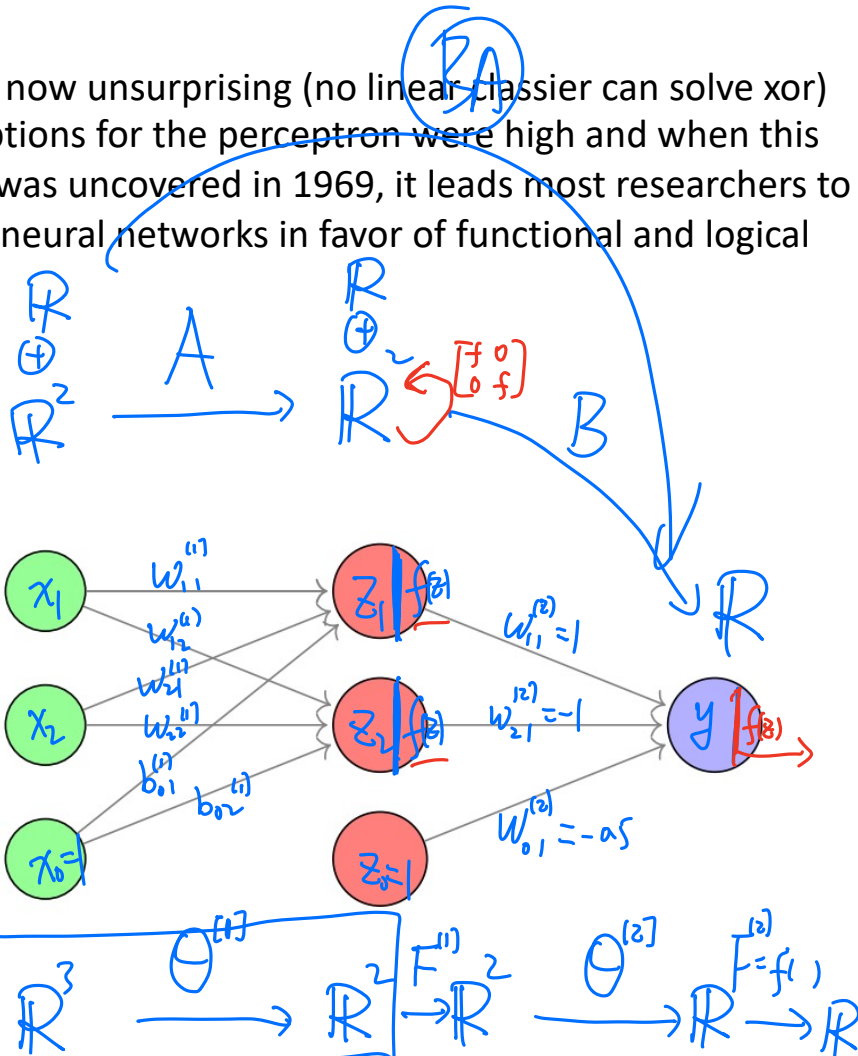
Famous example of a simple non-linearly separable data set, the XOR problem (Minsky 1969):



$f(z) = \text{ReLu}(z)$ or $\text{sign}(z)$



Although now unsurprising (no linear classifier can solve xor) the exceptions for the perceptron were high and when this problem was uncovered in 1969, it leads most researchers to abandon neural networks in favor of functional and logical methods.



$$f(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z < 0 \end{cases}$$

$$\vec{x} = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix}$$

$$\oplus \quad \oplus \\ \mathbb{R} \rightarrow \mathbb{R}$$

Ex:

$$\Theta^{(1)} = \begin{bmatrix} -0.5 & 1 & 1 \\ -1.5 & 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \Theta^{(1)} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$$F^{(1)} = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & id \end{bmatrix}$$

$$\Theta^{(2)} = [-0.5 \quad 1 \quad -1]$$

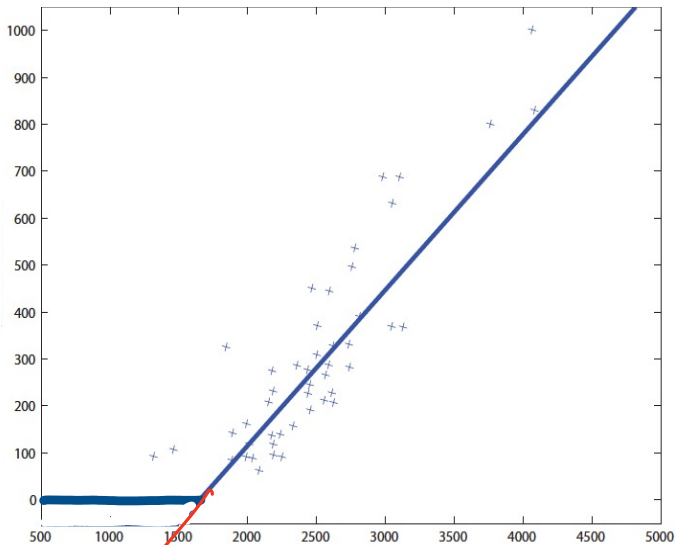
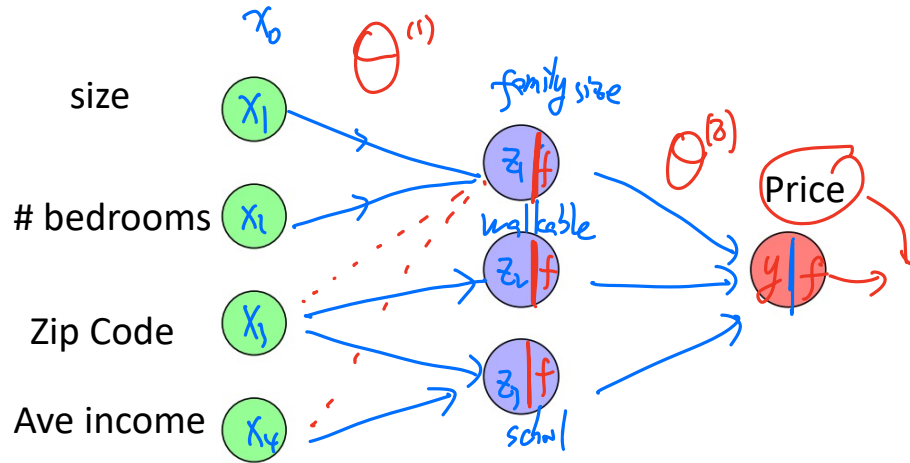
model: $h(\vec{x}) = F^{(2)} \cdot \Theta^{(2)} \cdot F^{(1)} \cdot \Theta^{(1)}(\vec{x})$

$$h(\vec{x}) = \begin{cases} 1 & \text{if } 0.5 \leq x_1 + x_2 \leq 1.5 \\ 0 & \text{if others} \end{cases}$$

➤ House price example

e.g.

$$f(z) = \text{ReLU}(z) = \begin{cases} z & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

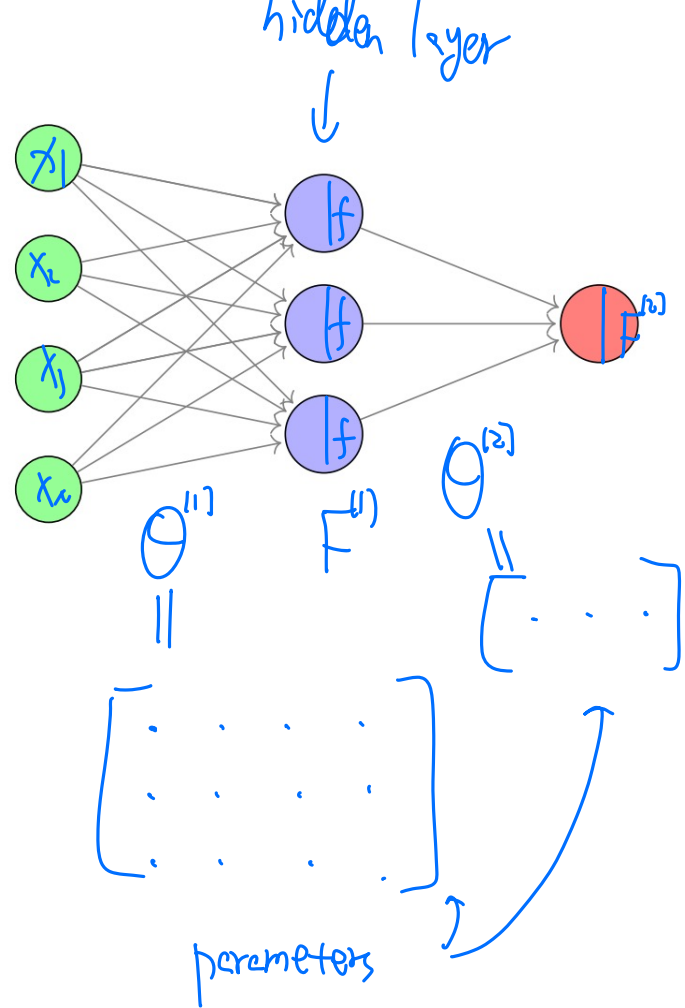


1 1 1 1

➤ Fully-connected neural networks

$$F^{(1)} = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & f \end{bmatrix} \text{ non-linear } \rightarrow$$

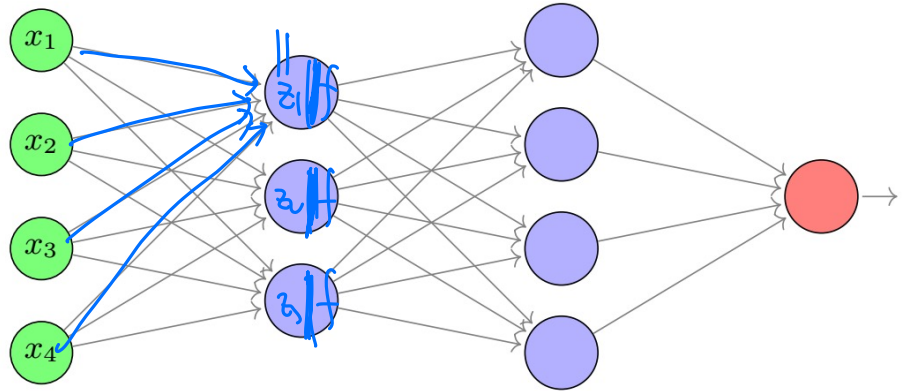
$$h(\vec{x}) = F^{(2)} \circ \theta^{(2)} \circ F^{(1)} \circ \theta^{(1)}(\vec{x})$$



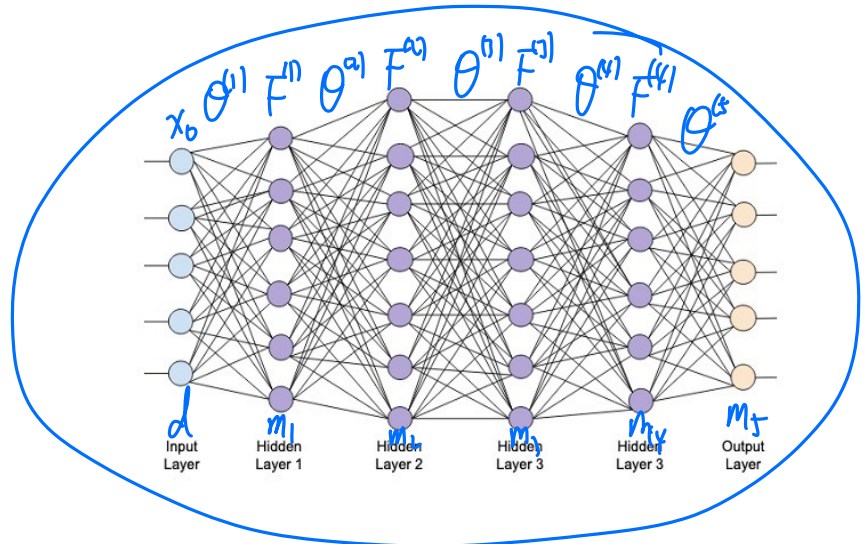
➤ Multi-layers Neural Networks.

Deep

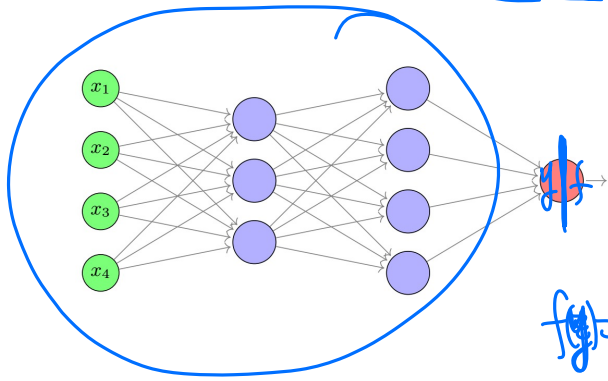
$$x_1 w_1 + x_2 w_2 + x_3 w_3 + b_4 x_4$$



$$\begin{aligned} \# \text{ parameters} &= m_1 \times (d+1) + \\ & m_2 \times (m_1+1) + \\ & m_3 \times (m_2+1) + \\ & m_4 \times (m_3+1) \end{aligned}$$

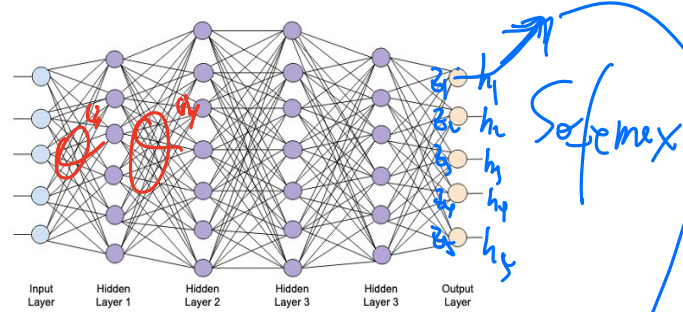


2 classes



$$f(y) = \frac{1}{1 + e^{-y}}$$

5-classes



$$h_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_5}}$$

① Data : $(\vec{x}^{(i)}, y^{(i)})$ or (X, \vec{y})

② Model : $h(\vec{x}) = F^{(4)} \cdot \Theta^{(4)} \cdot F^{(3)} \cdot \Theta^{(3)} \cdot F^{(2)} \cdot \Theta^{(2)} \cdot F^{(1)} \cdot \Theta^{(1)} \cdot \vec{x}$

③ Cost function $\left\{ \begin{array}{l} \text{① Regression: } J(\theta) = \|h(x) - \vec{y}\|_2^2 \\ \text{② Classifier (Cross-entropy)} \end{array} \right.$

④ Optimize θ

$$J(\theta) = -\frac{1}{n} \ln(P(\text{Data}))$$

$$= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(y^{(i)}=k) \ln(h_k(\vec{x}^{(i)}))$$

$$\theta^{\text{opt}} := \underset{\theta}{\operatorname{arg\,min}} J(\theta)$$

• Solution: ① $\nabla J(\theta) = 0$ \times

② Gradient descent:

③ Newton's method \times

$$\theta^{\text{next}} := \theta - \alpha \nabla J(\theta)$$

➤ **Back-propagation** (Reverse autodifferentiation)

In 1986, (*Learning representations by back-propagating errors, Nature, 323(9): 533-536*) D. E. Rumelhart popularized the idea of **back propagation** to compute gradients. It is not a learning method, but a computational trick. It is actually a simple implementation of chain rule of derivatives.

BP algorithms as stochastic gradient descent algorithms (Robbins–Monro 1950; Kiefer- Wolfowitz 1951) with Chain rules of Gradient maps

Goal: Minimize the loss function $J = H_5 \circ (H_4 \circ (H_3 \circ (H_2 \circ (H_1))))$

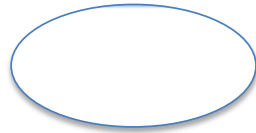
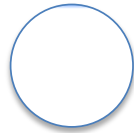
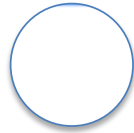
Need to calculate the gradient.

$$\nabla J(\theta) = \frac{\partial H_5}{\partial H_4} \cdot \frac{\partial H_4}{\partial H_3} \cdot \frac{\partial H_3}{\partial H_2} \cdot \frac{\partial H_2}{\partial H_1} \cdot \frac{\partial H_1}{\partial x}$$

➤ The Chain Rule:

➤ Computation Graphs

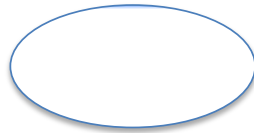
Function



Forward mode:



Backward mode:



Neural Network by Python:

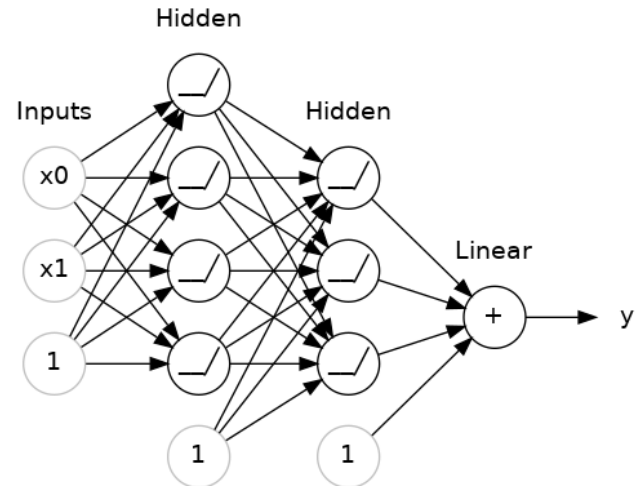
TensorFlow2: <https://www.tensorflow.org/>

Keras on TensorFlow: <https://keras.io/examples/>

PyTorch

```
from tensorflow import keras
from tensorflow.keras import layers

model = keras.Sequential([
    # the hidden ReLU layers
    layers.Dense(units=4, activation='relu', input_shape=[2]),
    layers.Dense(units=3, activation='relu'),
    # the linear output layer
    layers.Dense(units=1),
])
```



➤ Dropout and Batch Normalization

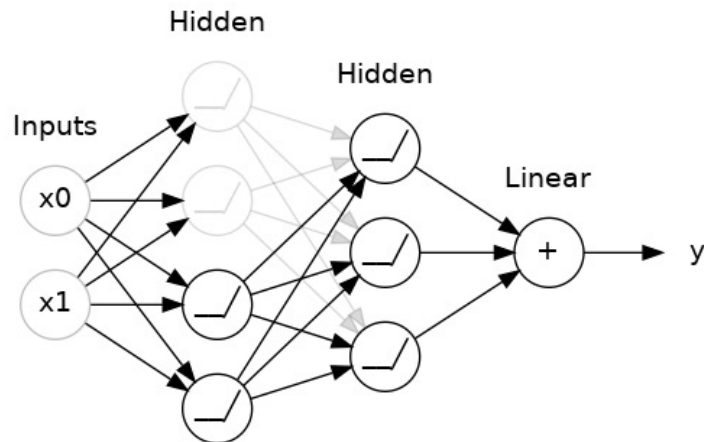
1. Dropout layer can help correct overfitting. We randomly *drop out* some fraction of a layer's input units every step of training. The weight patterns tend to be more robust.

2. Batch Normalization is something like scikit-learn's [StandardScaler](#) or [MinMaxScaler](#).

Batch normalization layer looks at each batch as it comes in, first normalizing the batch with its own mean and standard deviation, and then also putting the data on a new scale with two trainable rescaling parameters. Batch normalization, in effect, performs a kind of coordinated rescaling of its inputs.

```
from tensorflow import keras
from tensorflow.keras import layers
```

```
model = keras.Sequential([
    # the hidden ReLU layers
    layers.Dense(units=4, activation='relu', input_shape=[2]),
    layers.Dropout(0.3), # apply 30% dropout to the next layer
    layers.BatchNormalization(),
    layers.Dense(units=3, activation='relu'),
    layers.Dropout(0.3), # apply 30% dropout to the next layer
    layers.BatchNormalization(),
    # the linear output layer
    layers.Dense(units=1),
])
```



➤ Early Stopping

```
from tensorflow.keras.callbacks import EarlyStopping
from tensorflow.keras import layers, callbacks

early_stopping = EarlyStopping(
    min_delta=0.001, # minimum amount of change to count as an improvement
    patience=20, # how many epochs to wait before stopping
    restore_best_weights=True,
)
```

These parameters say: "If there hasn't been at least an improvement of 0.001 in the validation loss over the previous 20 epochs, then stop the training and keep the best model you found."

It can sometimes be hard to tell if the validation loss is rising due to overfitting or just due to random batch variation. The parameters allow us to set some allowances around when to stop.

```
model.compile(loss="sparse_categorical_crossentropy",
              optimizer="sgd",
              metrics=["accuracy"])
```

```
history = model.fit(X_train, y_train, epochs=30,
                   validation_data=(X_valid, y_valid))
```

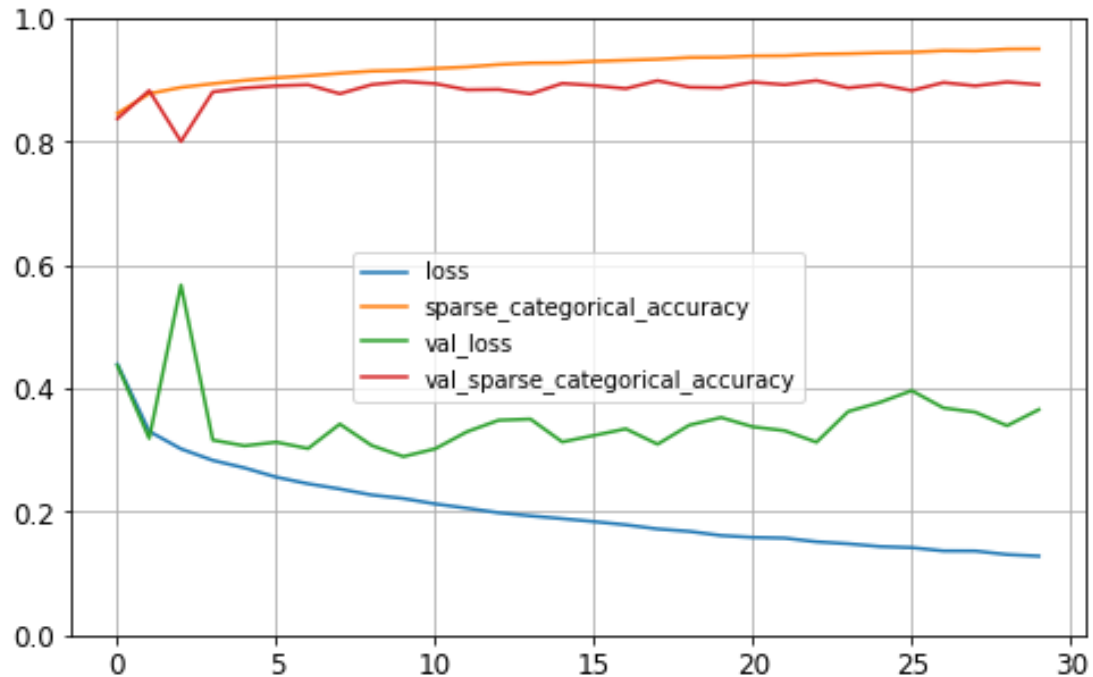
```
import pandas as pd
```

```
pd.DataFrame(history.history).plot(figsize=(8, 5))
```

```
plt.grid(True)
```

```
plt.gca().set_ylim(0, 1)
```

```
plt.show()
```



A visual proof that neural nets can compute any function

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

Play with neural network:

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

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

Online book about neural network:

<http://neuralnetworksanddeeplearning.com/chap3.html>

MIT Introduction to Deep Learning | 6.S191

https://www.youtube.com/watch?v=5tvmMX8r_OM