# Neural Network: architecture
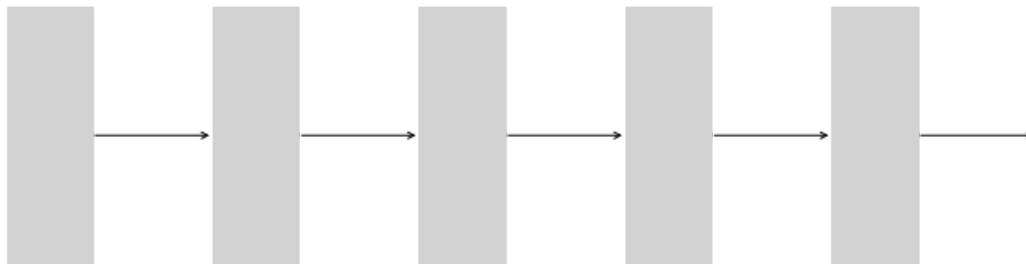
A Neural Network (Sequential architecture) is composed of

- sequence of Layers
    - layer $\ll$ transforms its input $\backslash\mathbf{y}_{(\ll-1)}$ to output $\backslash\mathbf{y}_{\backslash\mathrm{llp}}$
        - through a transformation: operation parameterized by weights $\backslash\mathbf{W}_{\backslash\mathrm{llp}}$

In [6]: `fig_tf_seq`

Out[6]:

- initial layer $0$ is Input layer:
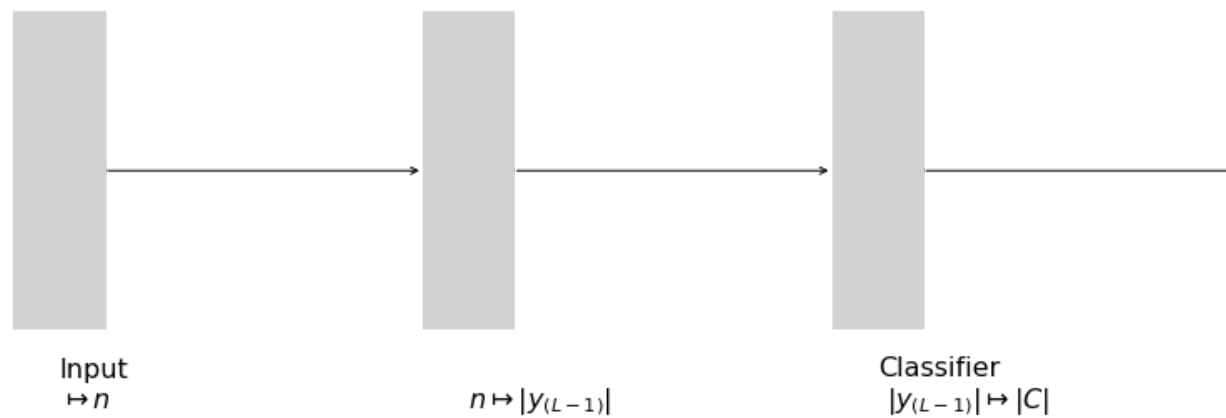  - outputs the network's inputs $\mathbf{x}$

$$\mathbf{y}_{(0)} = \mathbf{x}$$

- final layer $L$ transforms its input $\mathbf{y}_{(L-1)}$ to prediction $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \mathbf{y}_{(L)}$$

  - transformation of layer $L$ usually: Regression or Classification

Head layer: Classification or Regression

In [7]: `fig_tf_test`

Out[7]:



Input
$\mapsto n$

$n \mapsto |y_{(L-1)}|$

Classifier
$|y_{(L-1)}| \mapsto |C|$

In the above diagram

- the central box represents a sequence of 1 or more layers
- a sub-network
- just for brevity

Non-head layers (central box)

- transform raw features into synthetic features
- for consumption by Head layer

Head Layer (last box)

- Linear Regression or Classification
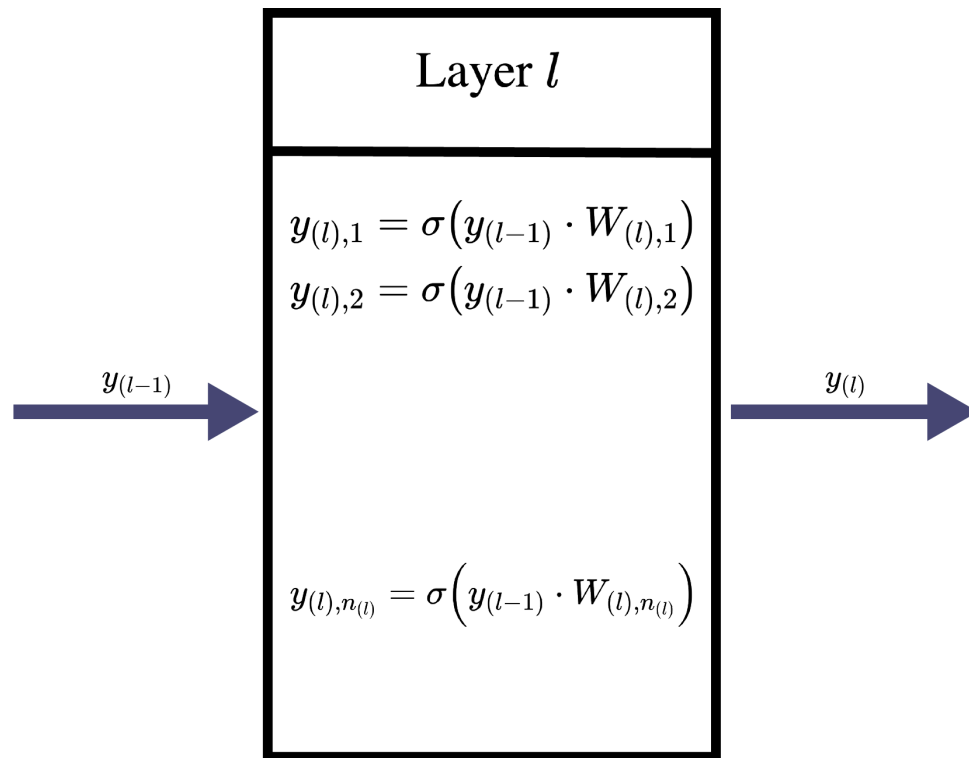
Network computes an *unknown* function

- from input of length $n$ (the raw input)
- to a vector of length $\left| \backslash \mathbf{y}_{(L-1)} \right|$

whose purpose is

- to make the final layer $L$ (e.g., the Classifier)
- create predictions (outputs)
- that have a very-low loss $\backslash \mathrm{loss}$
    - a good "approximation" of the function described by the training data

# Dense (Fully Connected) Layer

**Dense layer $\ll$ with $n_{\backslash \mathbf{llp}}$ units and sigmoid activation**

$$y_{(l-1)} \rightarrow$$

## Layer $l$

$$y_{(l),1} = \sigma\big(y_{(l-1)} \cdot W_{(l),1}\big)$$
$$y_{(l),2} = \sigma\big(y_{(l-1)} \cdot W_{(l),2}\big)$$

$$y_{(l),n_{(l)}} = \sigma\Big(y_{(l-1)} \cdot W_{(l),n_{(l)}}\Big)$$

$$y_{(l)} \rightarrow$$

Output $\mathbf{y}_{(\ll-1)}$ of layer $\ll -1$

- matched against $n_{\backslash\mathrm{llp}}$ patterns
$$\mathbf{W}_{(\ll,1)}, \cdots, \mathbf{W}_{(\ll,n_{\backslash\mathrm{llp}})}$$
- passed through a sigmoid activation function

$n_{\backslash\mathrm{llp}}$ logistic regressions

- Does $\mathbf{y}_{(\ll-1)}$ match each of the $n_{\backslash\mathrm{llp}}$ patterns

# Special cases

Usually for Head layer.

- Sigmoid activation with $n_{\backslash llp} = 1$
    - Binary logistic regression
- Softmax activation
    - Multinomial logistic regression
- Linear (or no) activation, $n_{\backslash llp} = 1$
    - Linear Regression
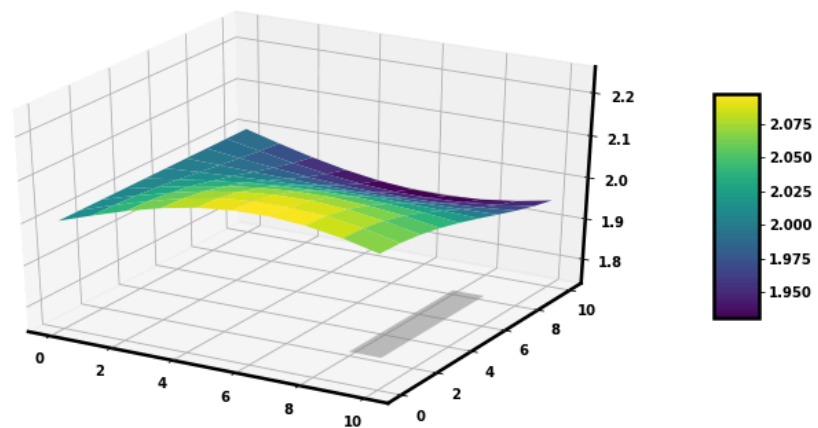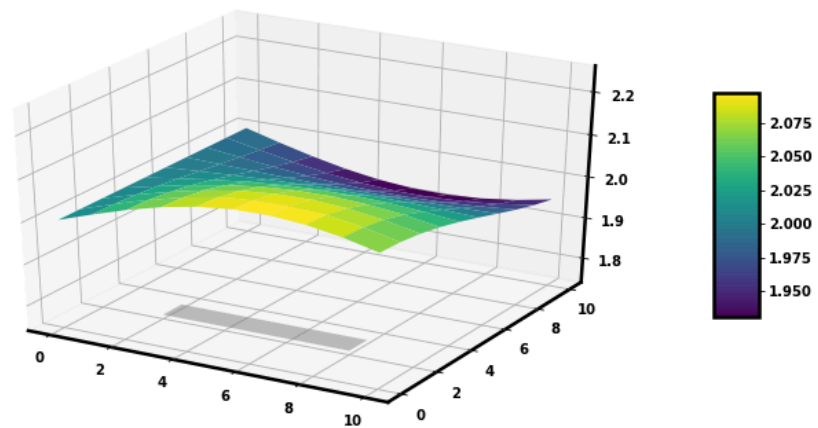
# Importance of training dataset

The Neural network thus computes a *function* from $\mathbf{x}$ to $\hat{\mathbf{y}}$.

The function *mimics* the training data

$$\langle \mathbf{X}, \mathbf{y} \rangle = [\mathbf{x}^{\backslash ip}, \mathbf{y}^{\backslash ip} | 1 \leq i \leq m]$$

where each *example* $\langle \mathbf{x}^{\backslash ip}, \mathbf{y}^{\backslash ip} \rangle$

- describes the mapping of the function on input $\mathbf{x}^{\backslash ip}$ to output $\mathbf{y}^{\backslash ip}$
- e.g., from input features $\mathbf{x}^{\backslash ip}$ to
  - continuous value $\mathbf{y}^{\backslash ip}$
  - discrete class $\mathbf{y}^{\backslash ip}$
    - really: output is a probability vector over finite set $C$ of discrete classes

The Neural Network is trained ("learns") to mimic the training data

- by solving for the weights $\mathbf{W}_{\backslash \mathrm{llp}}$ of each layer $1 \leq \ll \leq L$
- that minimize a loss function

$$\backslash \mathrm{loss} = \sum_{i-1}^{m} \backslash \mathrm{loss}^{\backslash \mathrm{ip}}$$

- where $\backslash \mathrm{loss}^{\backslash \mathrm{ip}}$ if a function of
  - how much prediction $\backslash \hat{\mathbf{y}}^{\backslash \mathrm{ip}}$ deviates from true target/label $\backslash \mathbf{y}^{\backslash \mathrm{ip}}$

The minimization procedure is usually a variant of *Gradient Descent*

The challenge is that the operation of each layer $1 \leq \ll < L$

- is usually not *interpretable*
- we can describe *how* it transforms $\backslash y_{(\ll -1)}$ to $\backslash y_{\backslash \text{llp}}$
- but not *why* it is performing the transformation
    - objective (describe) rather than subjection (why)

So the sub-network in the diagram

```python
In [9]: print("Done")
```

Done