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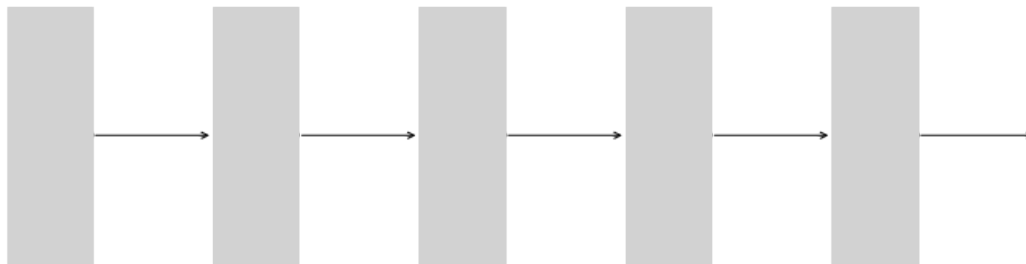
A Neural Network (Sequential architecture) is composed of

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

In [6]: `fig_tf_seq`

Out[6]:

- initial layer $0$ is Input layer:
  - outputs the network's inputs $\mathbf{x}$
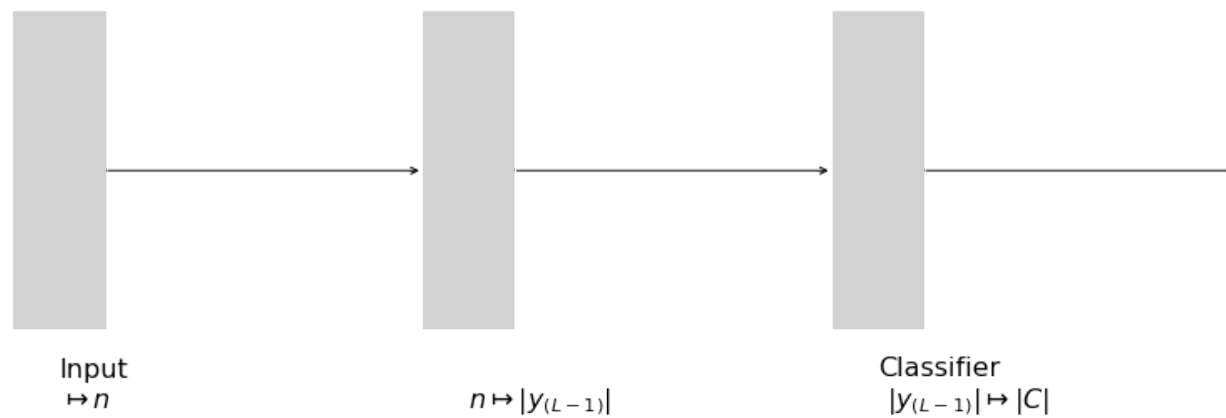$$\mathbf{y}_{(0)} = \mathbf{x}$$
- final layer $L$ transforms its input $\mathbf{y}_{(L-1)}$ to prediction $\hat{\mathbf{y}}$
$$\hat{\mathbf{y}} = \mathbf{y}_{\text{llp}}$$
  - transformation of layer $L$ usually: Regression or Classification

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
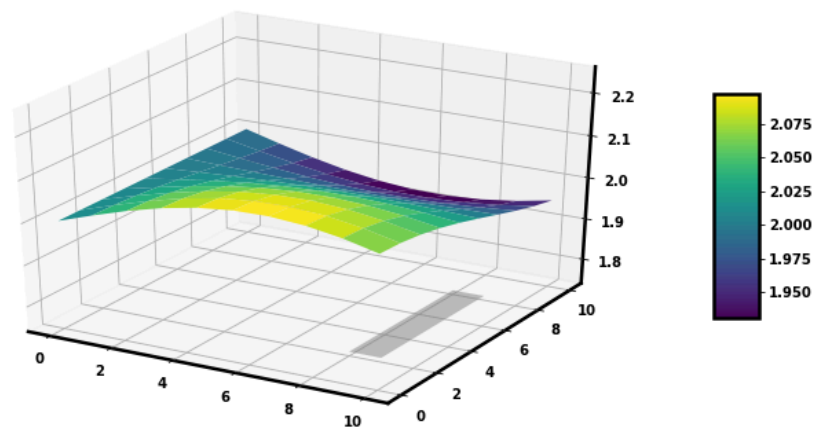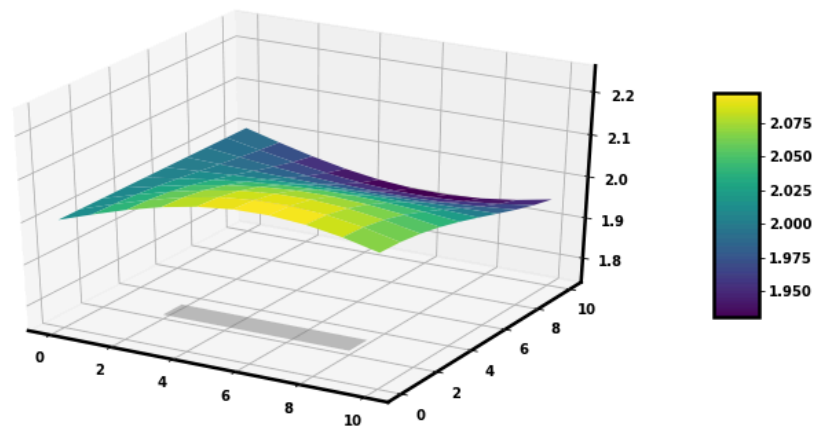- a sub-network
- just for brevity

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

The function *mimics* the training data
$$\langle \mathbf{X}, \mathbf{y} \rangle = [\mathbf{x}^{(i)}, \mathbf{y}^{(i)} | 1 \leq i \leq m]$$

where each *example* $\langle \mathbf{x}^{(i)}, \mathbf{y}^{(i)} \rangle$

- describes the mapping of the function on input $\mathbf{x}^{(i)}$ to output $\mathbf{y}^{(i)}$
- e.g., from input features $\mathbf{x}^{(i)}$ to
  - continuous value $\mathbf{y}^{(i)}$
  - discrete class $\mathbf{y}^{(i)}$
    - 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 $\backslash \mathbf{y}^{\backslash \mathrm{ip}}$

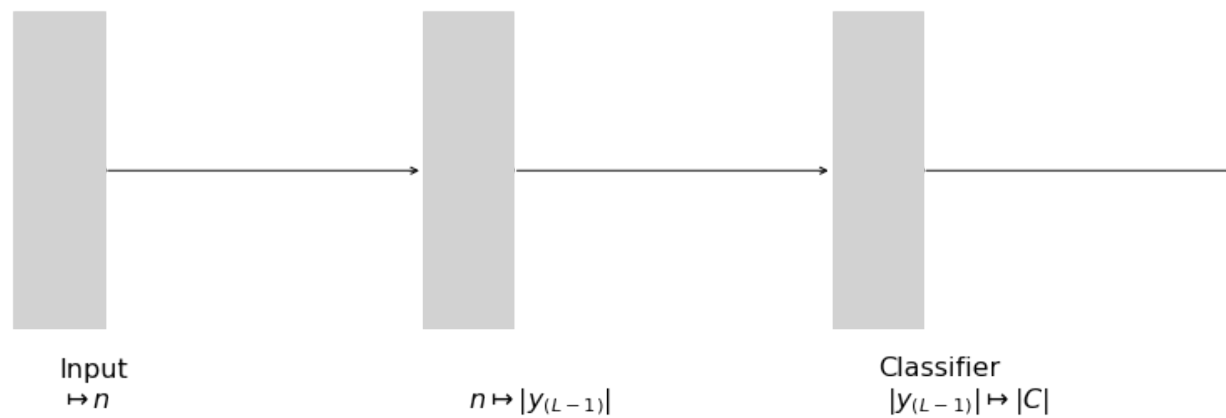The minimization procedure is usually a variant of *Gradient Descent*

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

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

So the sub-network in the diagram

In [8]: `fig_tf_test`

Out[8]:



Input
$\mapsto n$

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

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

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

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

Done