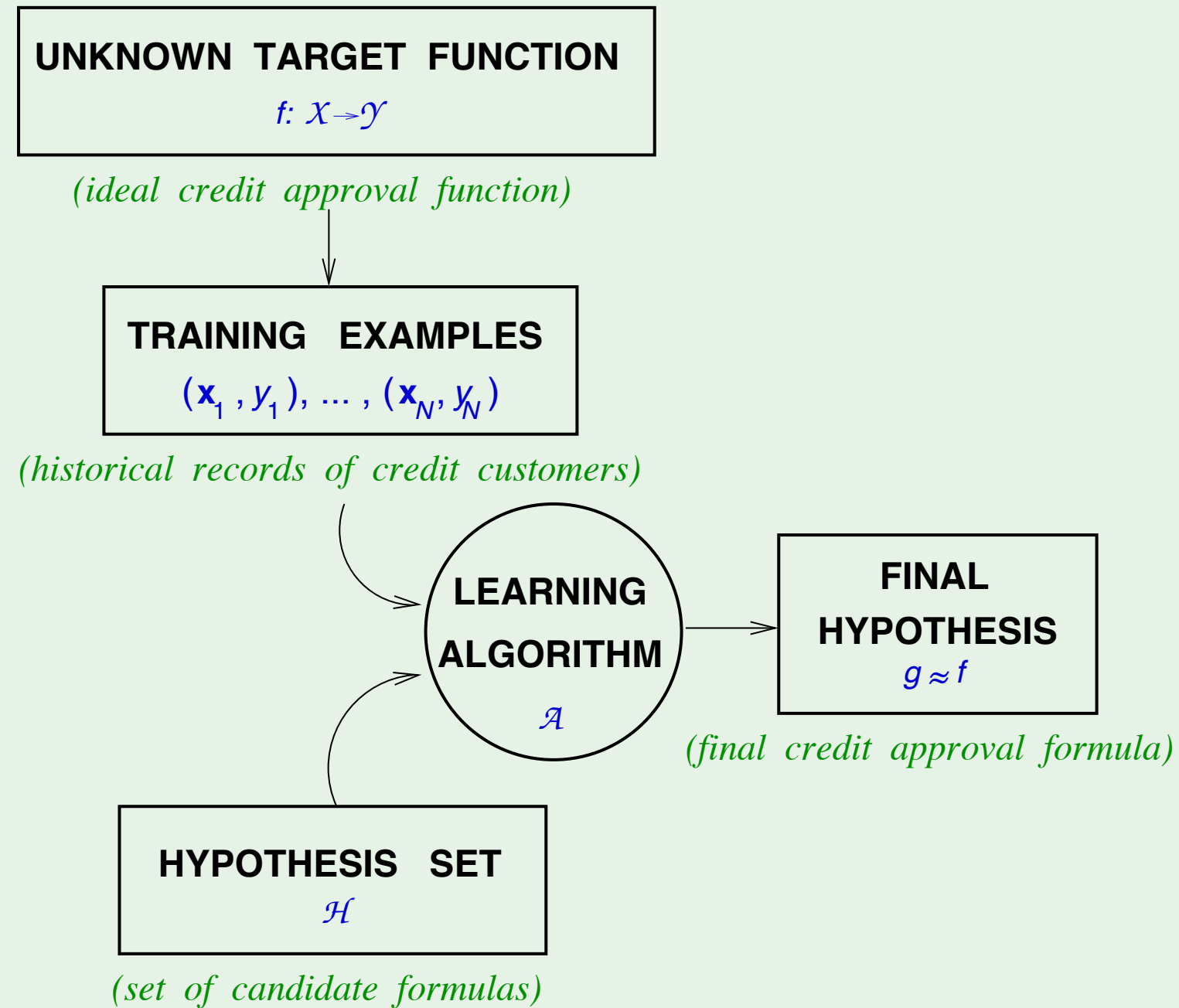# Components of learning

## Formalization:

- Input: $\mathbf{x}$     *(customer application)*

- Output: $y$     *(good/bad customer?)*

- Target function: $f : \mathcal{X} \rightarrow \mathcal{Y}$     *(ideal credit approval formula)*

- Data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \cdots , (\mathbf{x}_N, y_N)$     *(historical records)*

$$\downarrow \quad \downarrow \quad \downarrow$$

- Hypothesis: $g : \mathcal{X} \rightarrow \mathcal{Y}$     *(formula to be used)*

UNKNOWN TARGET FUNCTION
$f: \mathcal{X} \to \mathcal{Y}$

*(ideal credit approval function)*

TRAINING EXAMPLES
$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

*(historical records of credit customers)*

LEARNING ALGORITHM
$\mathcal{A}$

FINAL HYPOTHESIS
$g \approx f$

*(final credit approval formula)*

HYPOTHESIS SET
$\mathcal{H}$

*(set of candidate formulas)*

# Solution components

The 2 solution components of the learning problem:

- The Hypothesis Set

$$\mathcal{H} = \{h\} \qquad g \in \mathcal{H}$$

- The Learning Algorithm

Together, they are referred to as the *learning model*.

# Error measures

What does "$h \approx f$" mean?

Error measure: $E(h, f)$

Almost always *pointwise definition*: $\mathrm{e}\big(h(\mathbf{x}), f(\mathbf{x})\big)$

Examples:

Squared error: $\quad \mathrm{e}\big(h(\mathbf{x}), f(\mathbf{x})\big) = \big(h(\mathbf{x}) - f(\mathbf{x})\big)^2$

Binary error: $\quad \mathrm{e}\big(h(\mathbf{x}), f(\mathbf{x})\big) = [\![h(\mathbf{x}) \neq f(\mathbf{x})]\!]$

# From pointwise to overall

Overall error $E(h, f) =$ average of pointwise errors $\mathrm{e}\left(h(\mathbf{x}), f(\mathbf{x})\right)$.

In-sample error:

$$E_{\text{in}}(h) = \frac{1}{N}\sum_{n=1}^{N}\mathrm{e}\left(h(\mathbf{x}_n), f(\mathbf{x}_n)\right)$$

Out-of-sample error:

$$E_{\text{out}}(h) = \mathbb{E}_{\mathbf{x}}\left[\mathrm{e}\left(h(\mathbf{x}), f(\mathbf{x})\right)\right]$$

# A simple hypothesis set - the 'perceptron'

For input $\mathbf{x} = (x_1, \cdots, x_d)$   'attributes of a customer'

$$\text{Approve credit if} \quad \sum_{i=1}^{d} w_i x_i > \text{threshold},$$

$$\text{Deny credit if} \quad \sum_{i=1}^{d} w_i x_i < \text{threshold}.$$

This linear formula $h \in \mathcal{H}$ can be written as

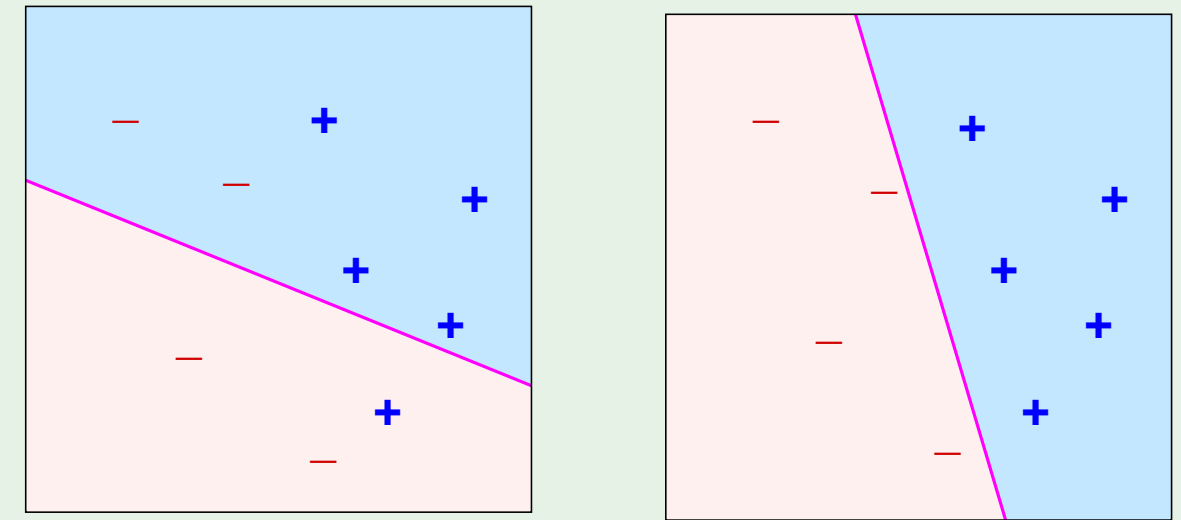$$h(\mathbf{x}) = \text{sign}\left(\left(\sum_{i=1}^{d} w_i x_i\right) - \text{threshold}\right)$$

$$h(\mathbf{x}) = \text{sign}\left(\left(\sum_{i=1}^{d} w_i\, x_i\right) + w_0 \right)$$

Introduce an artificial coordinate $x_0 = 1$:

$$h(\mathbf{x}) = \text{sign}\left(\sum_{i=0}^{d} w_i\, x_i\right)$$
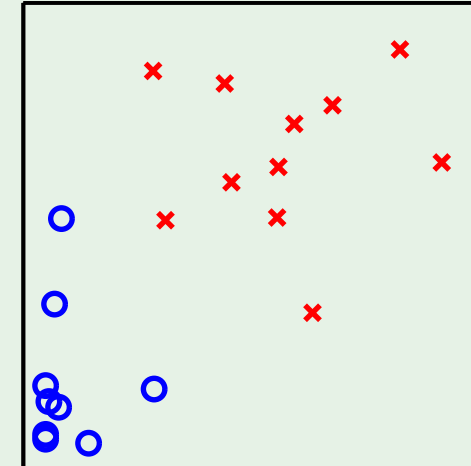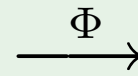
In vector form, the perceptron implements
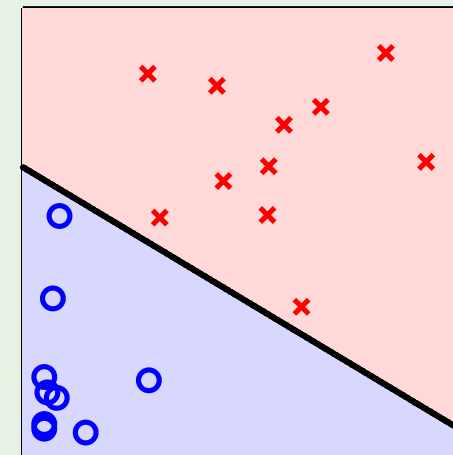
$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^\mathsf{T}\mathbf{x})$$
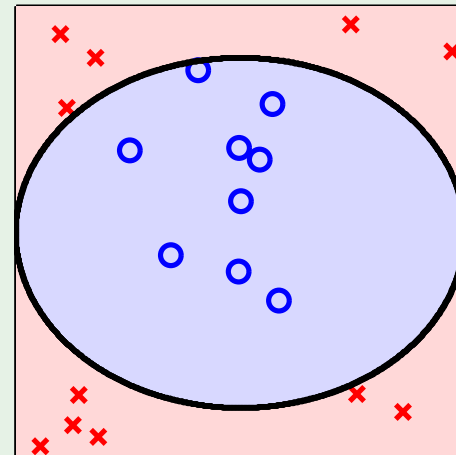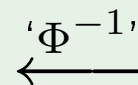


'linearly separable' data

**1.** Original data
$\mathbf{x}_n \in \mathcal{X}$

$\Phi$

**2.** Transform the data
$\mathbf{z}_n = \Phi(\mathbf{x}_n) \in \mathcal{Z}$

**4.** Classify in $\mathcal{X}$-space
$g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})) = \mathrm{sign}(\tilde{\mathbf{w}}^\mathsf{T}\Phi(\mathbf{x}))$

'$\Phi^{-1}$'

**3.** Separate data in $\mathcal{Z}$-space
$\tilde{g}(\mathbf{z}) = \mathrm{sign}(\tilde{\mathbf{w}}^\mathsf{T}\mathbf{z})$

# What transforms to what

$$\mathbf{x} = (x_0, x_1, \cdots, x_d) \quad \xrightarrow{\Phi} \quad \mathbf{z} = (z_0, z_1, \cdots \cdots \cdots, z_{\tilde{d}})$$

$$\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N \quad \xrightarrow{\Phi} \quad \mathbf{z}_1, \mathbf{z}_2, \cdots, \mathbf{z}_N$$

$$y_1, y_2, \cdots, y_N \quad \xrightarrow{\Phi} \quad y_1, y_2, \cdots, y_N$$

$$\text{No weights in } \mathcal{X} \qquad\qquad \tilde{\mathbf{w}} = (w_0, w_1, \cdots \cdots \cdots, w_{\tilde{d}})$$

$$g(\mathbf{x}) \qquad = \qquad \text{sign}(\tilde{\mathbf{w}}^{\mathsf{T}} \Phi(\mathbf{x}))$$
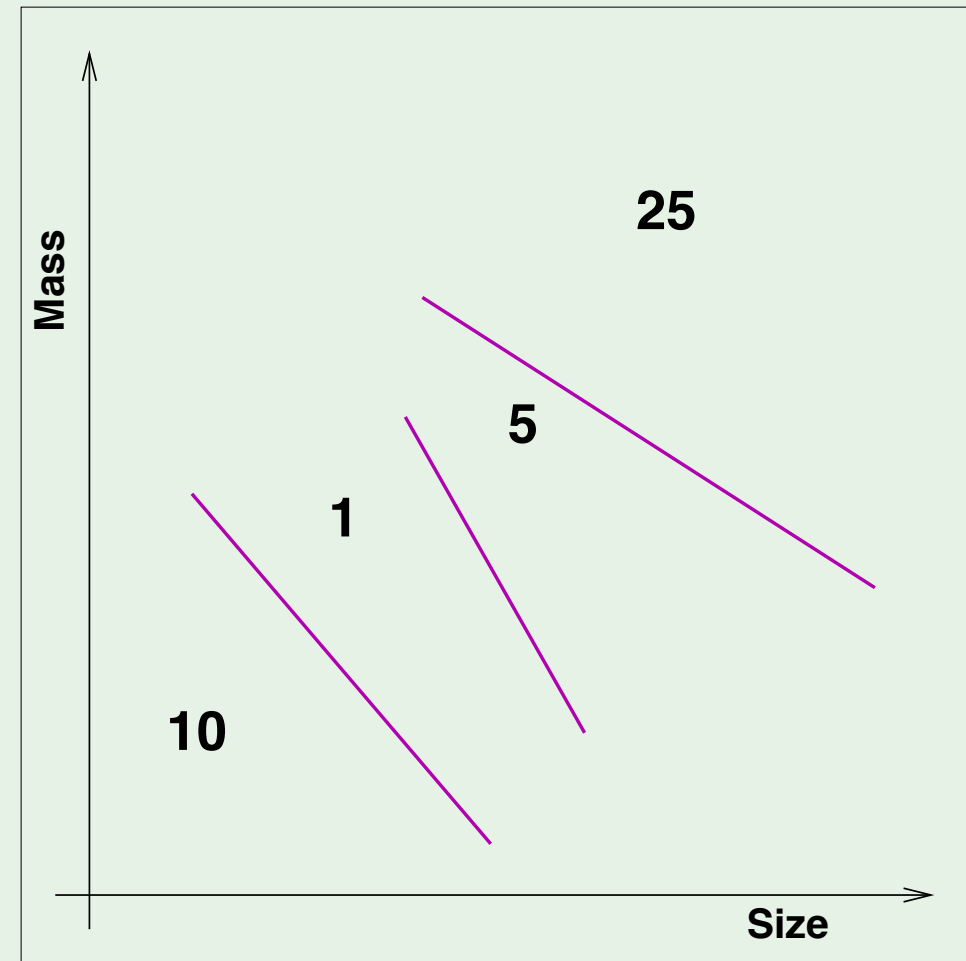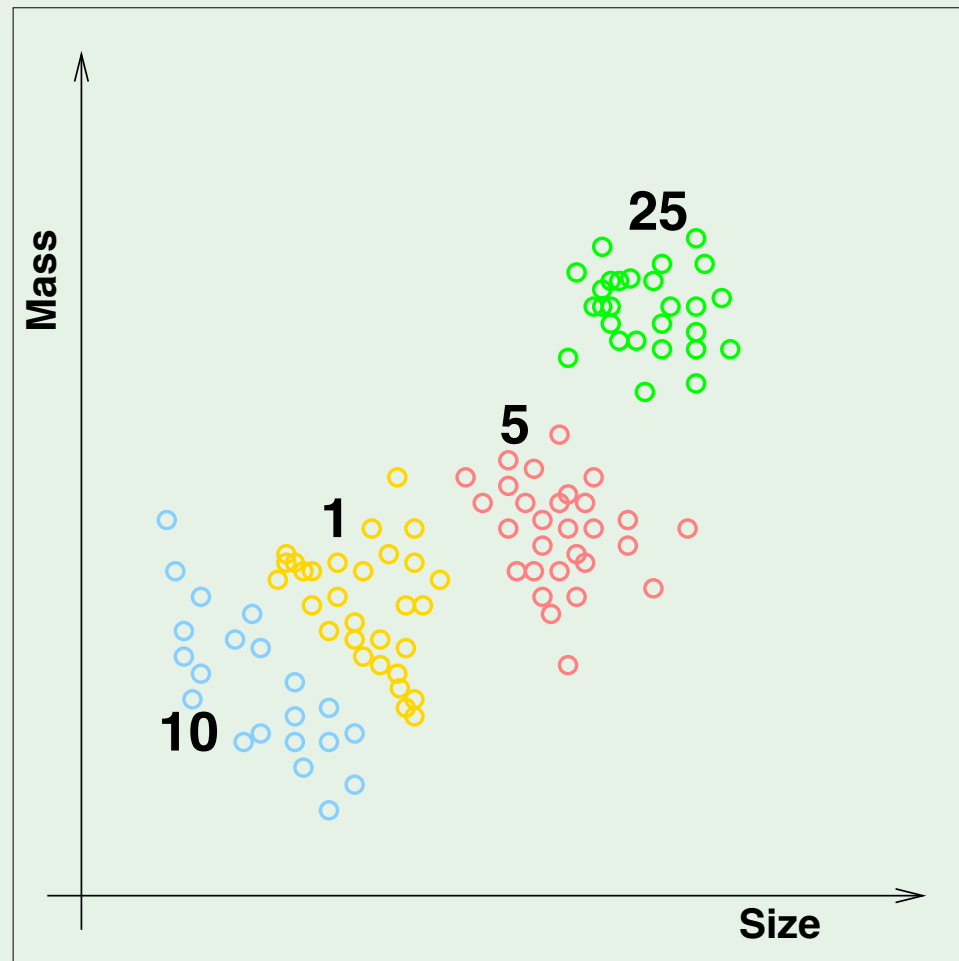
# Basic premise of learning

*"using a set of observations to uncover an underlying process"*

broad premise $\implies$ many variations

- Supervised Learning

- Unsupervised Learning

- Reinforcement Learning

# Supervised learning

Example from vending machines – **coin recognition**

# Unsupervised learning

Instead of $($input,correct output$)$, we get $($input, ? $)$