

# Outline

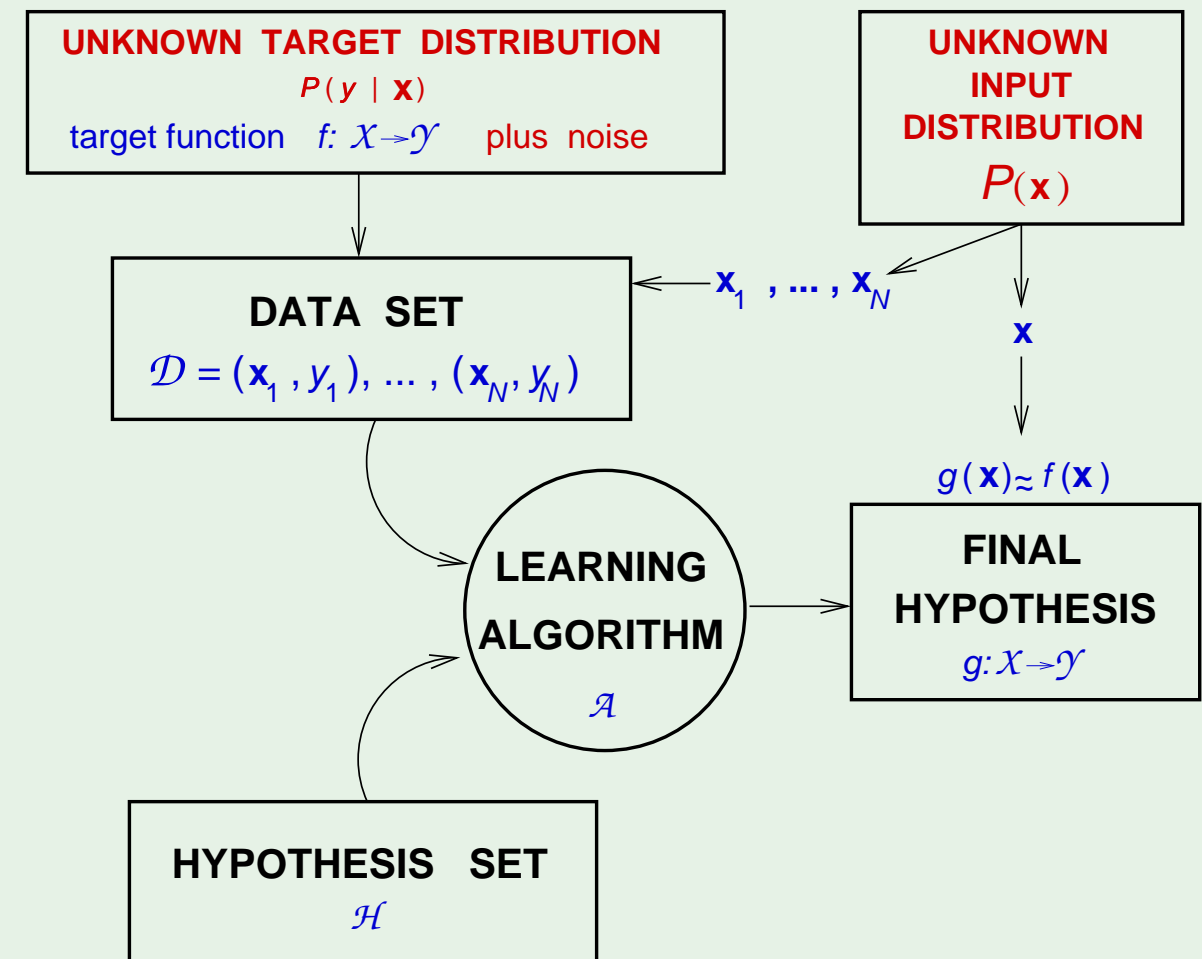
- The map of machine learning
- Bayesian learning
- Aggregation methods
- Acknowledgments

# Probabilistic approach

Extend probabilistic role to all components

$P(\mathcal{D} \mid h = f)$  decides which  $h$  (likelihood)

How about  $P(h = f \mid \mathcal{D})$  ?



# The prior

$P(h = f \mid \mathcal{D})$  requires an additional probability distribution:

$$P(h = f \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid h = f) P(h = f)}{P(\mathcal{D})} \propto P(\mathcal{D} \mid h = f) P(h = f)$$

$P(h = f)$  is the **prior**

$P(h = f \mid \mathcal{D})$  is the **posterior**

Given the prior, we have the full distribution

## Example of a prior

Consider a perceptron:  $h$  is determined by  $\mathbf{w} = w_0, w_1, \dots, w_d$

A possible prior on  $\mathbf{w}$ : Each  $w_i$  is independent, uniform over  $[-1, 1]$

This determines the prior over  $h$  -  $P(h = f)$

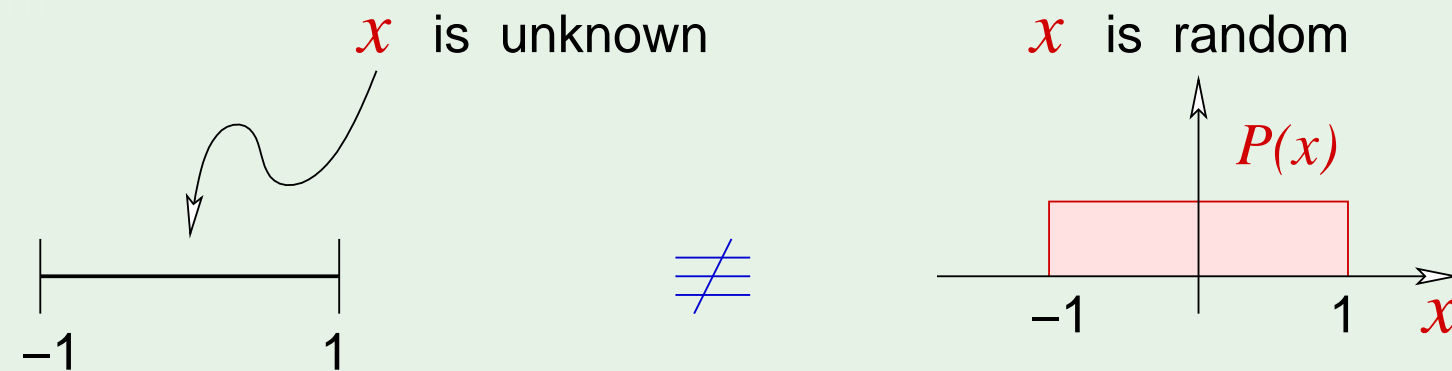
Given  $\mathcal{D}$ , we can compute  $P(\mathcal{D} \mid h = f)$

Putting them together, we get  $P(h = f \mid \mathcal{D})$

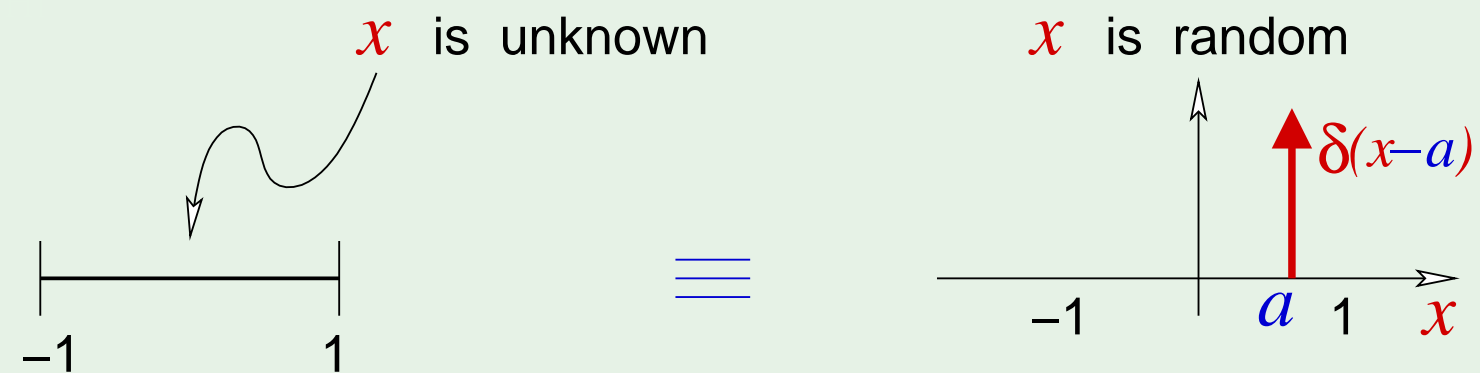
$$\propto P(h = f)P(\mathcal{D} \mid h = f)$$

# A prior is an assumption

Even the most “neutral” prior:



The true equivalent would be:



## If we knew the prior

... we could compute  $P(h = f \mid \mathcal{D})$  for every  $h \in \mathcal{H}$

$\implies$  we can find the most probable  $h$  given the data

we can derive  $\mathbb{E}(h(\mathbf{x}))$  for every  $\mathbf{x}$

we can derive the **error bar** for every  $\mathbf{x}$

we can derive everything in a principled way

# When is Bayesian learning justified?

1. The prior is **valid**

trumps all other methods

2. The prior is **irrelevant**

just a computational catalyst