

# Machine learning problem formulation

- The machine learning approach expresses NLP as an optimization problem:

$$\hat{\mathbf{Y}} = \operatorname{argmax}_{\mathbf{y} \in f(\mathbf{x}; \boldsymbol{\theta})} \Psi(\mathbf{Y}, \mathbf{X}; \boldsymbol{\theta})$$

*where*  $\mathbf{x} \in X$  is the input

$\mathbf{y} \in Y$  is the output

$\Psi(\cdot) \rightarrow \mathbb{R}$  is a function expressing the learning objective

$f(\cdot)$  is the function, or model, that maps  $\mathbf{x}$  to  $\mathbf{y}$

$\boldsymbol{\theta}$  parameterizes  $f(\cdot)$

# Maximum likelihood estimation

- Estimates the parameters,  $\theta$ , of a distribution using a likelihood function,  $\mathcal{L}(\mathbf{D}; \theta)$ , given some data,  $\mathbf{D}$ .

- We maximize the likelihood function by minimizing its -logarithm:

$$\mathcal{L}(\theta | \mathbf{D}) = P(\mathbf{D}; \theta)$$

$$= \left( \prod_{i=1}^M P(\mathbf{y}_i | \mathbf{x}_i; \theta) \right)^{\frac{1}{M}}$$

$$= \frac{1}{M} \sum_{i=1}^M \log P(\mathbf{y}_i | \mathbf{x}_i; \theta) \quad \text{Note: technically natural log, but true to within a constant}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^M -\log P(\mathbf{y}_i | \mathbf{x}_i; \theta)$$

- This expresses an optimization problem; the form of  $p(\mathbf{D}; \theta)$  dictates how we solve it
  - In deep learning, this function is a neural network; we compute its gradient w.r.t.  $\theta$ , and then estimate  $\hat{\theta}$  using stochastic gradient descent (SGD).