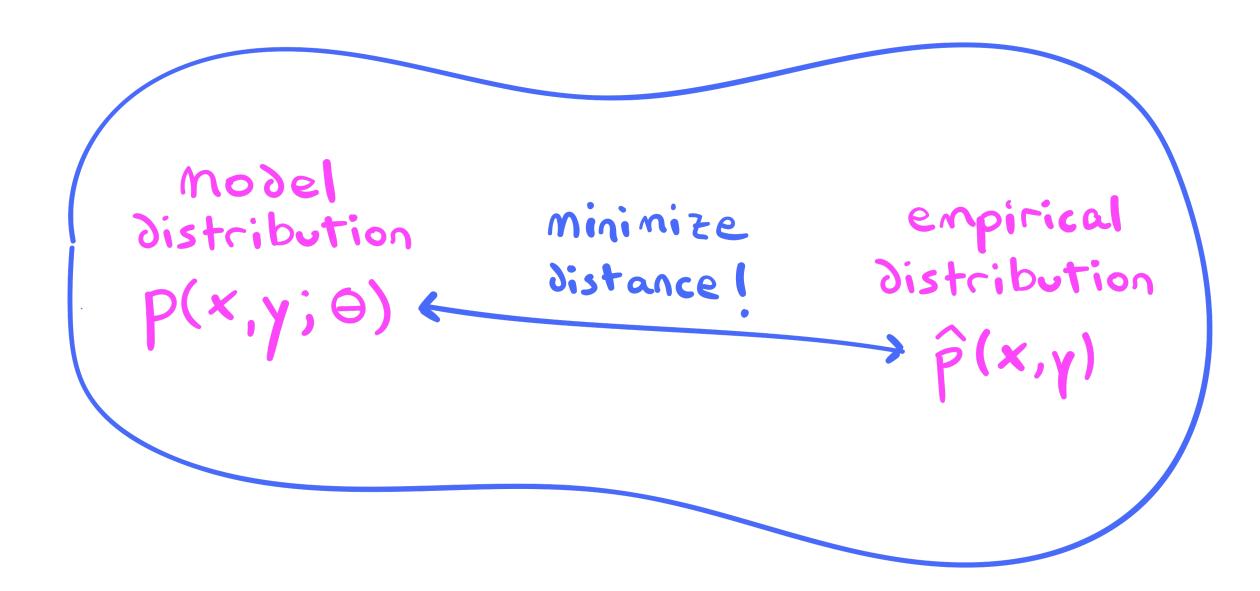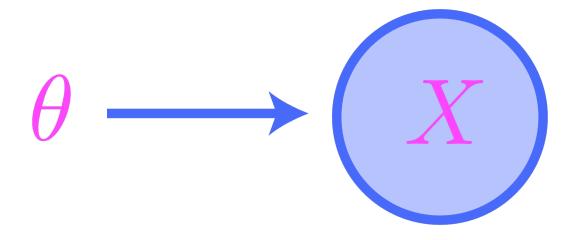# 13. Learning

**The Distance Criterion for Parameter Choice.** Given two model distributions within the same family of probabilistic models, choose the model distribution whose *distance* from the empirical distribution of the data is smaller.

model
distribution
$P(x, y; \theta)$

minimize
distance!

empirical
distribution
$\hat{P}(x, y)$

# 13.1. A first look at likelihood-based learning objectives

## 🔔 Theorem 13.1 (Equivalent learning objectives for the univariate Bernoulli model)

Let $x_1, x_2, \ldots, x_m \in \{0, 1\}$ be an observed dataset corresponding to a Bernoulli random variable $X \sim \mathcal{Ber}(\theta)$ with unknown $\theta$. Let $P_\theta$ be the model distribution of $X$ and let $\hat{P}$ be the empirical distribution of the dataset. The following optimization objectives are equivalent:

1. Minimize the KL divergence $D(\hat{P} \parallel P_\theta)$ with respect to $\theta$.
2. Minimize the cross entropy $H_{\hat{P}}(P_\theta)$ with respect to $\theta$.
3. Minimize the data surprisal function $\mathcal{I}(\theta; x_1, \ldots, x_m)$ with respect to $\theta$.
4. Maximize the data likelihood function $\mathcal{L}(\theta; x_1, \ldots, x_m)$ with respect to $\theta$.

1. Minimizing the KL divergence between the empirical and model distributions has an immediate and concrete interpretation as minimizing the "distance" between these two distributions.

2. As a function of $\theta$, the cross entropy $J(\theta) = H_{\hat{P}}(P_\theta)$ may be viewed as a stochastic objective function, since it is exactly the mean of the model surprisal function. This opens the door for applications of the stochastic gradient descent algorithm studied in [Section 11.4](#).

3. The third optimization objective seeks the model probability distribution according to which the data is *least surprising*.

4. The fourth optimization objective seeks the model probability distribution according to which the data is *most likely*.

**🔔 Theorem 13.2 (MLE for the univariate Bernoulli model)**

Let $x_1, x_2, \ldots, x_m \in \{0, 1\}$ be an observed dataset corresponding to a Bernoulli random variable $X \sim \mathcal{B}er(\theta)$ with unknown $\theta$. Then the (unique) maximum likelihood estimate $\theta^\star_{\mathrm{MLE}}$ is the ratio $\Sigma x/m$.

stochastic gradient descent for univariate Bernoulli model
$k = 8$, $\alpha = 0.01$, $\beta = 0$, $N = 10$