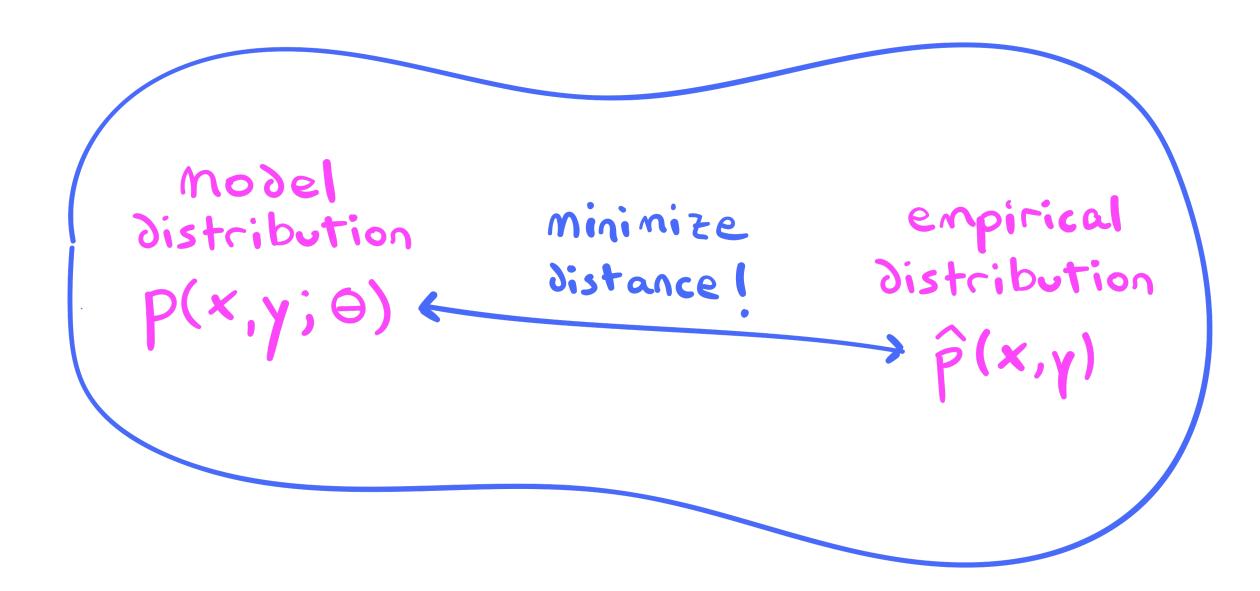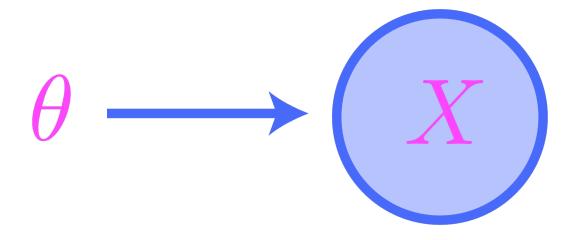# 13. Learning

**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

$$\theta \longrightarrow \boxed{X}$$

**🔔 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 $X \sim \mathcal{B}er(\theta)$ 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:
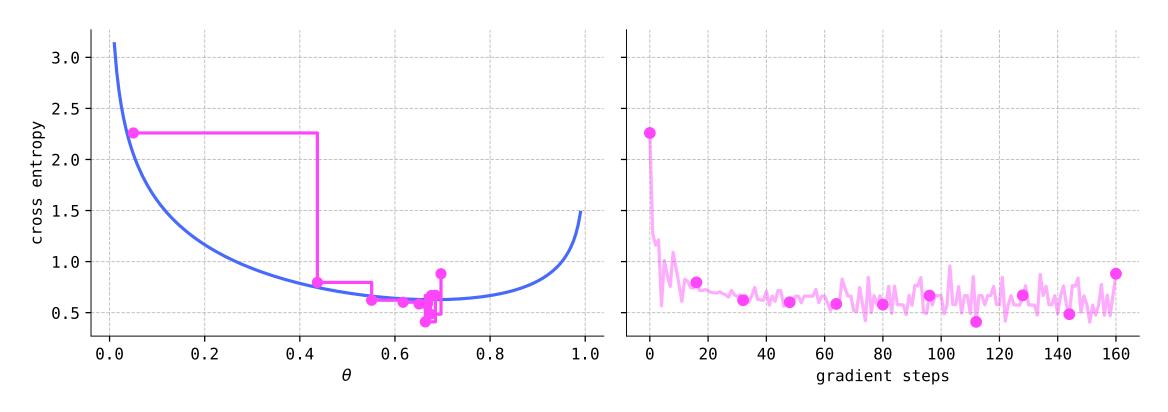
1. Minimize the KL divergence $D(\hat{P} \parallel P_\theta)$ with respect to $\theta$.
2. Minimize the cross entropy $H_{\hat{P}}(P_\theta)$ with respect to $\theta$.
3. Minimize the data surprisal function $\mathcal{I}(\theta; x_1, \ldots, x_m)$ with respect to $\theta$.
4. Maximize the data likelihood 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
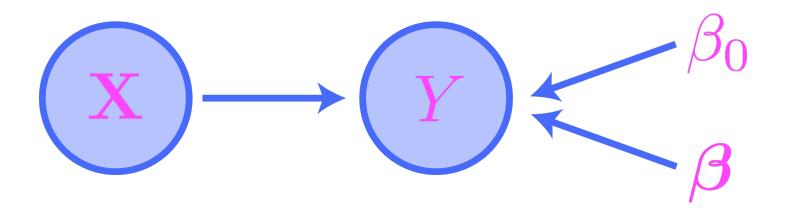$k = 8$, $\alpha = 0.01$, $\beta = 0$, $N = 10$

**🔔 Problem Prompt**

Do problem 1 on the worksheet.

# 13.3. MLE for linear regression

**🔔 Theorem 13.8 (MLEs for linear regression models with known variance)**

Consider a linear regression model with *fixed* variance $\sigma^2$, and let

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \mathbb{R}$$

be an observed dataset. Supposing

$$\mathbf{x}_i^\mathsf{T} = (x_{0i}, x_{i1}, \ldots, x_{in}) = (1, x_{i1}, \ldots, x_{in})$$

for each $i = 1, \ldots, m$, let

$$\mathbf{\mathcal{X}} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\mathsf{T} & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{x}_m^\mathsf{T} & \rightarrow \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

Provided that the $(n+1) \times (n+1)$ square matrix $\mathbf{\mathcal{X}}^\mathsf{T}\mathbf{\mathcal{X}}$ is invertible, maximum likelihood estimates for the parameters $\beta_0$ and $\boldsymbol{\beta}$ are given by

$$\boldsymbol{\theta}_{\text{MLE}}^\star = (\mathbf{\mathcal{X}}^\mathsf{T}\mathbf{\mathcal{X}})^{-1}\mathbf{\mathcal{X}}^\mathsf{T}\mathbf{y}.$$

**🔔 Corollary 13.1 (MLEs for simple linear regression models with known variance)**

Let the notation be as in Theorem 13.8, but assume that $\mathbf{X}$ is $1$-dimensional, equal to a random variable $X$. Then MLEs for the parameters $\beta_0$ and $\beta_1$ are given by
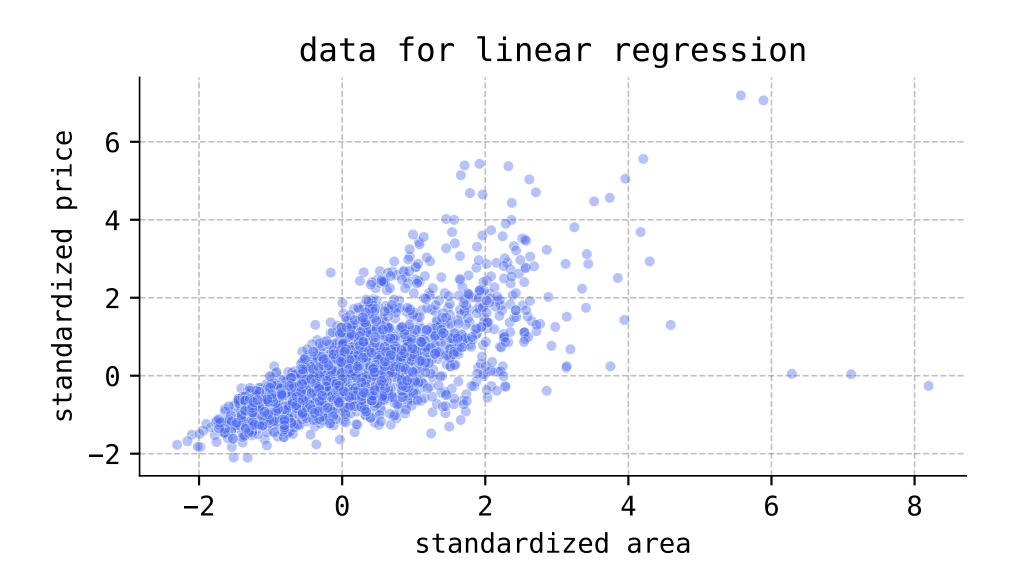
$$(\beta_1)^\star_{\mathrm{MLE}} = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2},$$

$$(\beta_0)^\star_{\mathrm{MLE}} = \bar{y} - (\beta_1)^\star_{\mathrm{MLE}}\bar{x},$$

where $\bar{x} = \frac{1}{m}\sum_{i=1}^m x_i$ and $\bar{y} = \frac{1}{m}\sum_{i=1}^m y_i$ are the empirical means.

**🔔 Problem Prompt**
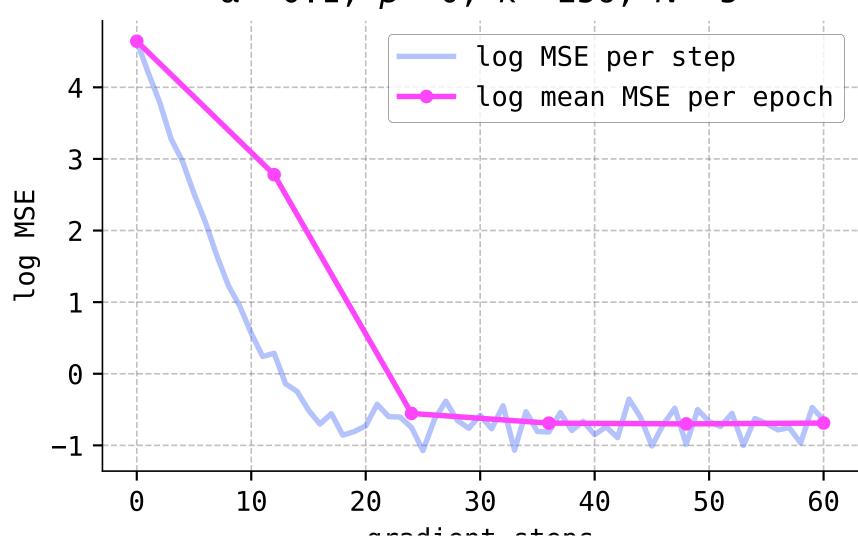
Do problem 2 on the worksheet.

data for linear regression

SGD for linear regression
$\alpha = 0.1, \ \beta = 0, \ k = 256, \ N = 5$

log MSE per step
log mean MSE per epoch

log MSE

gradient steps

stochastic gradient descent for linear regression
$\alpha = 0.1,\ \beta = 0,\ k = 256,\ N = 5$