

13. Learning

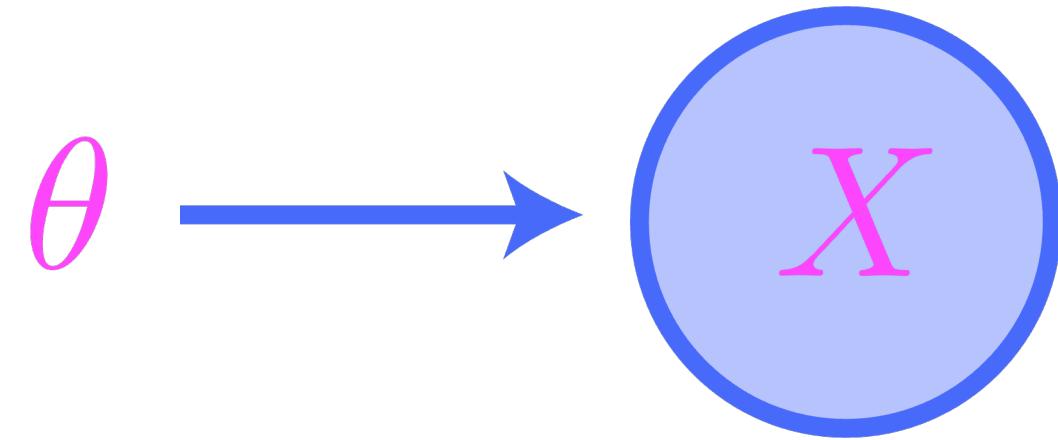
The Distance Criterion for Parameter Choice. Given two model distributions within the same family of probabilistic models, choose the model distribution whose *distance* from the empirical distribution of the data is smaller.

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 $x_1, x_2, \dots, x_m \in \{0, 1\}$ be an observed dataset corresponding to a Bernoulli random variable $X \sim \text{Ber}(\theta)$ with unknown θ . Let P_θ 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 θ .
2. Minimize the cross entropy $H_{\hat{P}}(P_\theta)$ with respect to θ .
3. Minimize the data surprisal function $\mathcal{I}(\theta; x_1, \dots, x_m)$ with respect to θ .
4. Maximize the data likelihood function $\mathcal{L}(\theta; x_1, \dots, x_m)$ with respect to θ .

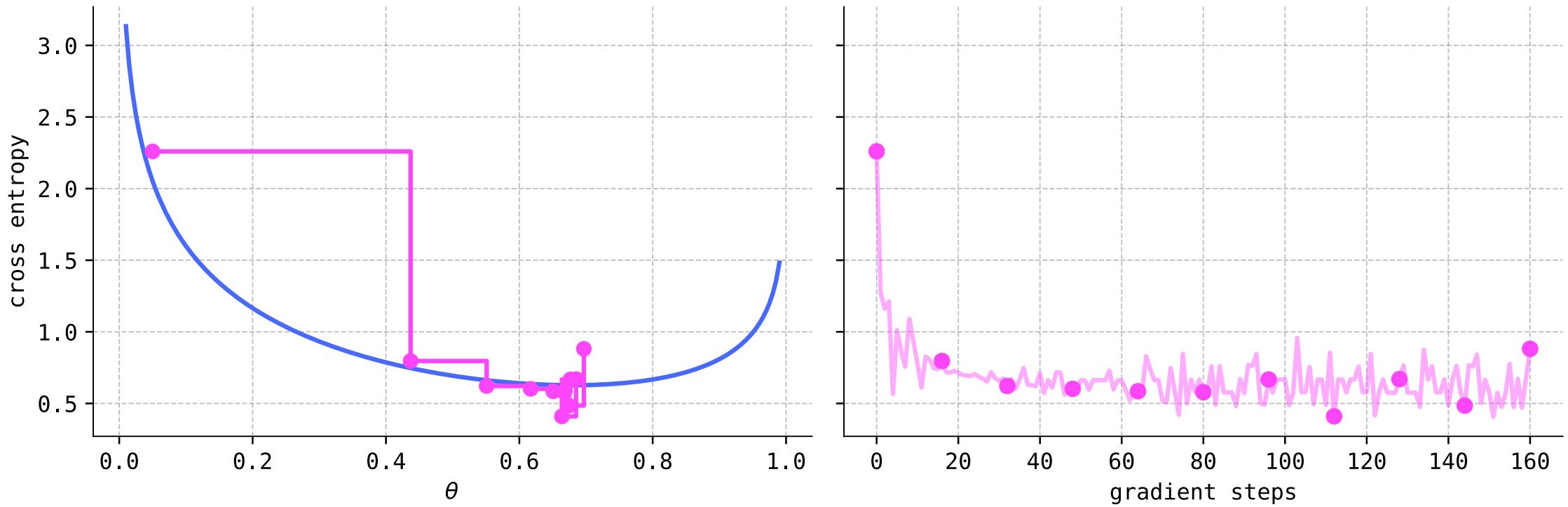
1. Minimizing the KL divergence between the empirical and model distributions has an immediate and concrete interpretation as minimizing the “distance” between these two distributions.
2. As a function of θ , the cross entropy $J(\theta) = H_{\hat{P}}(P_\theta)$ may be viewed as a stochastic objective 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 $x_1, x_2, \dots, x_m \in \{0, 1\}$ be an observed dataset corresponding to a Bernoulli random variable $X \sim \text{Ber}(\theta)$ with unknown θ . Then the (unique) maximum likelihood estimate $\hat{\theta}_{\text{MLE}}^*$ is the ratio $\sum x/m$.

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

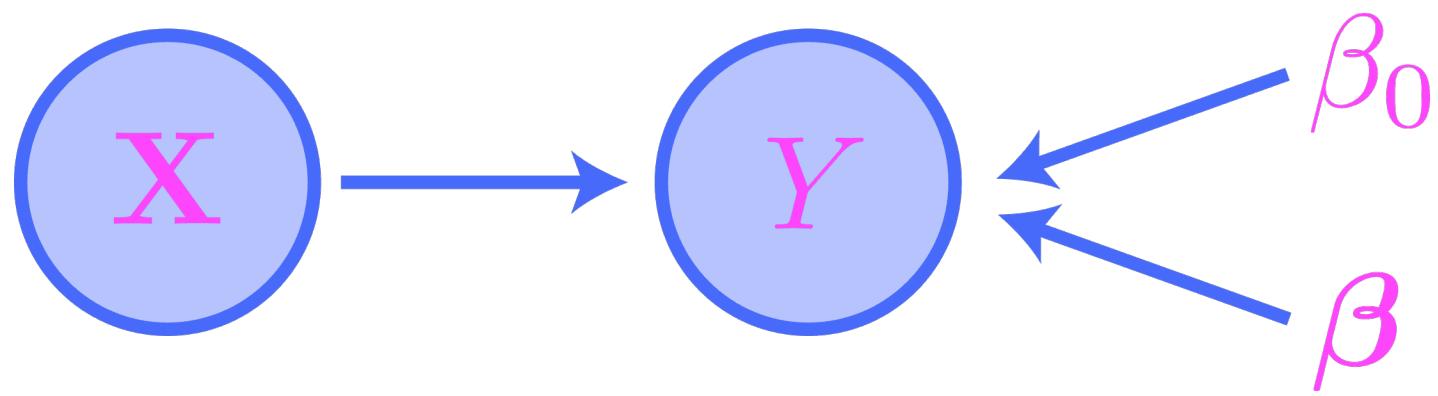




Problem Prompt

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 σ^2 , and let

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

be an observed dataset. Supposing

$$\mathbf{x}_i^\top = (x_{0i}, x_{i1}, \dots, x_{in}) = (1, x_{i1}, \dots, x_{in})$$

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

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

Provided that the $(n + 1) \times (n + 1)$ square matrix $\boldsymbol{\mathcal{X}}^\top \boldsymbol{\mathcal{X}}$ is invertible, maximum likelihood estimates for the parameters β_0 and $\boldsymbol{\beta}$ are given by

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

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 β_0 and β_1 are given by

$$(\beta_1)_{\text{MLE}}^* = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2},$$
$$(\beta_0)_{\text{MLE}}^* = \bar{y} - (\beta_1)_{\text{MLE}}^* \bar{x},$$

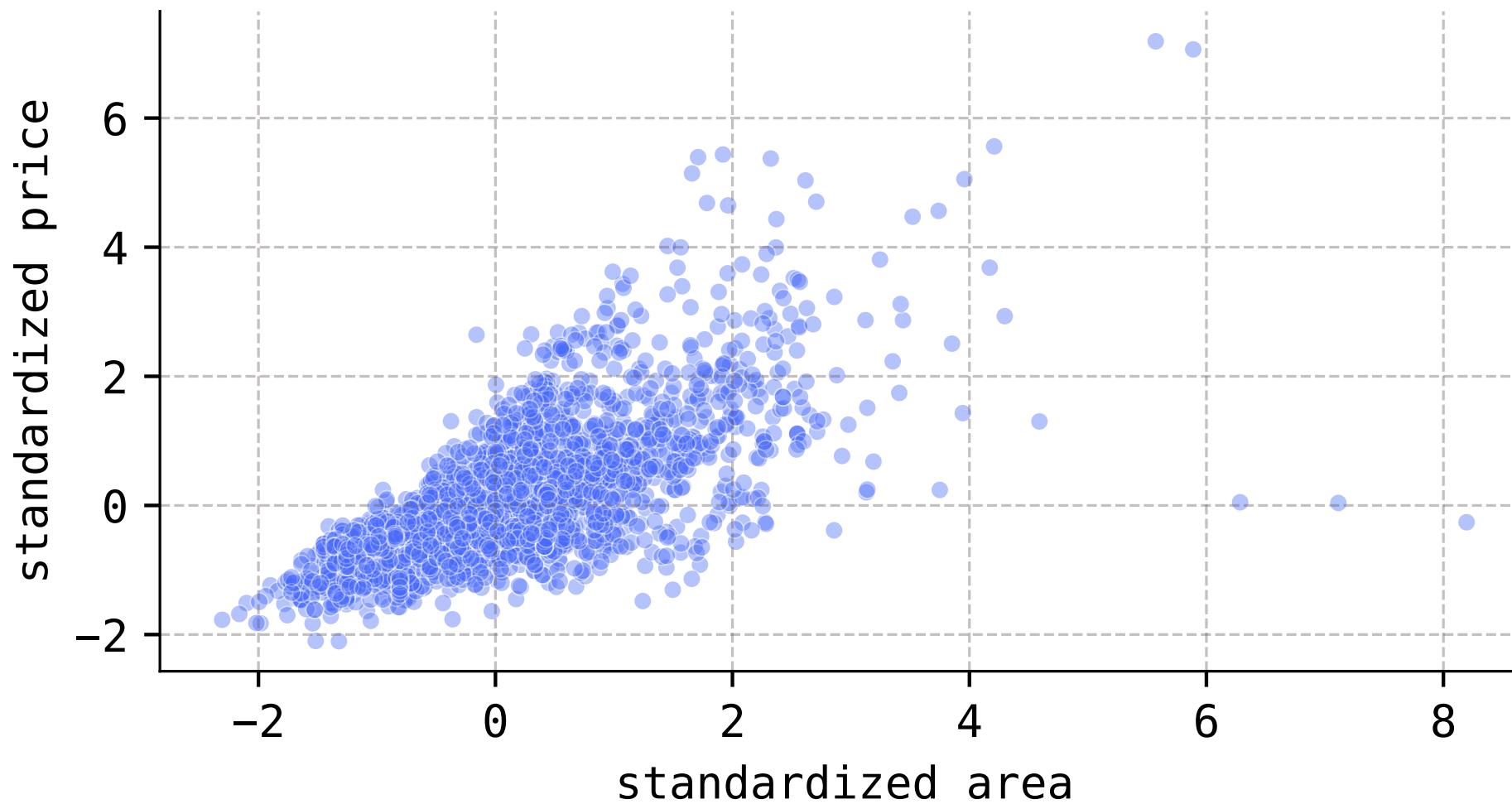
where $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$ and $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$ are the empirical means.



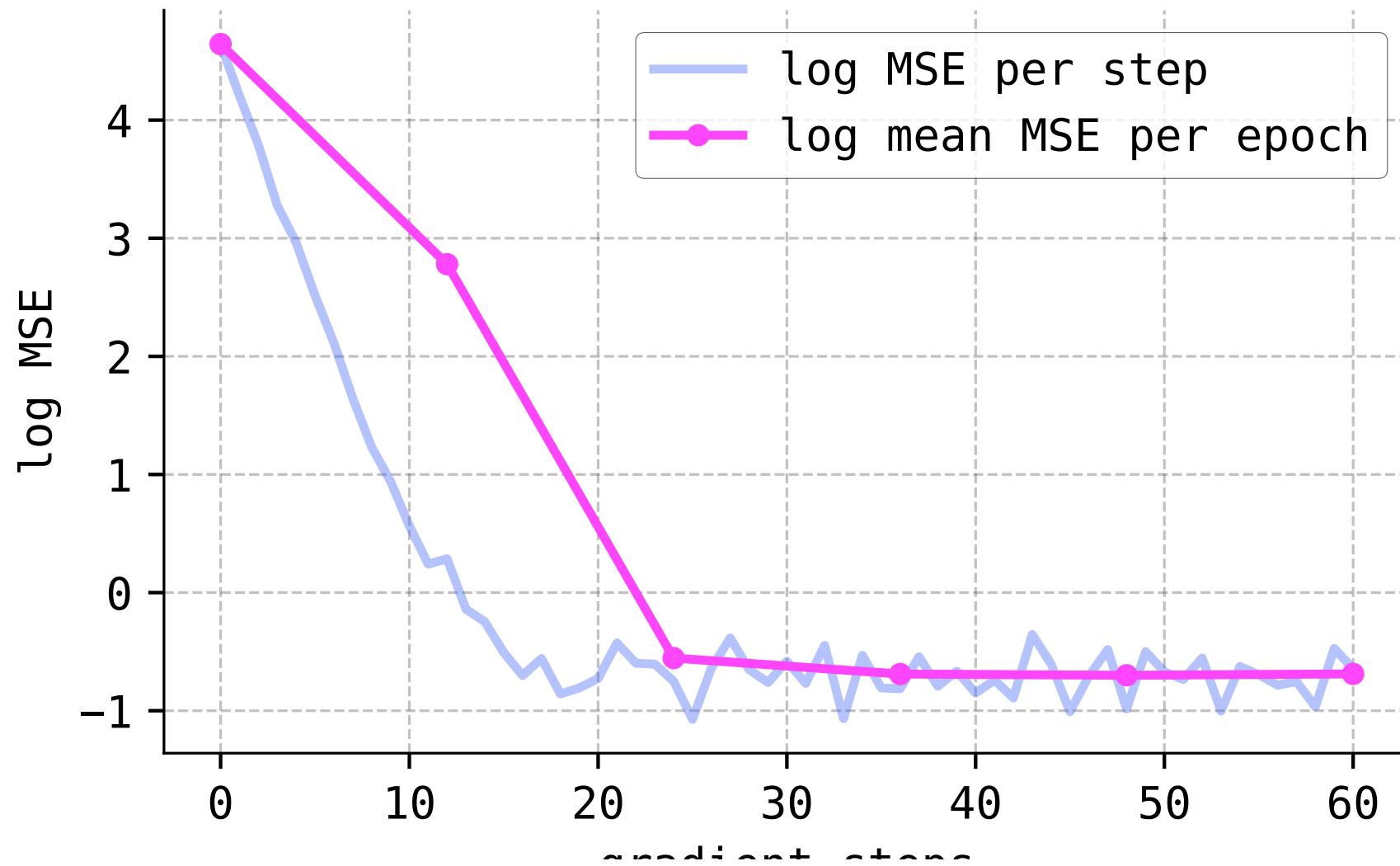
Problem Prompt

Do problem 2 on the worksheet.

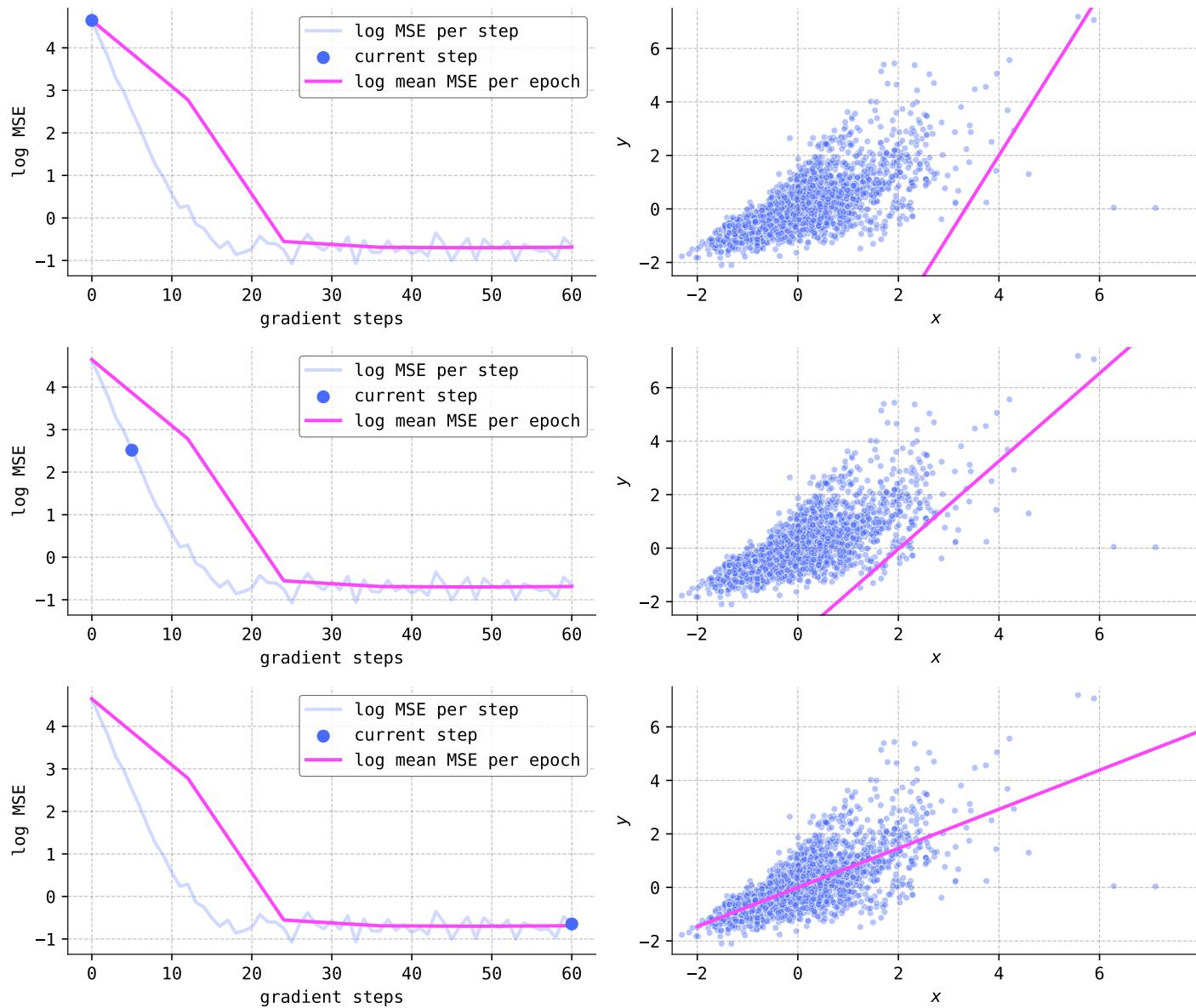
data for linear regression



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



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



13.4. MLE for logistic regression

🔔 **Theorem 13.9 (Surprisal functions of logistic regression models)**

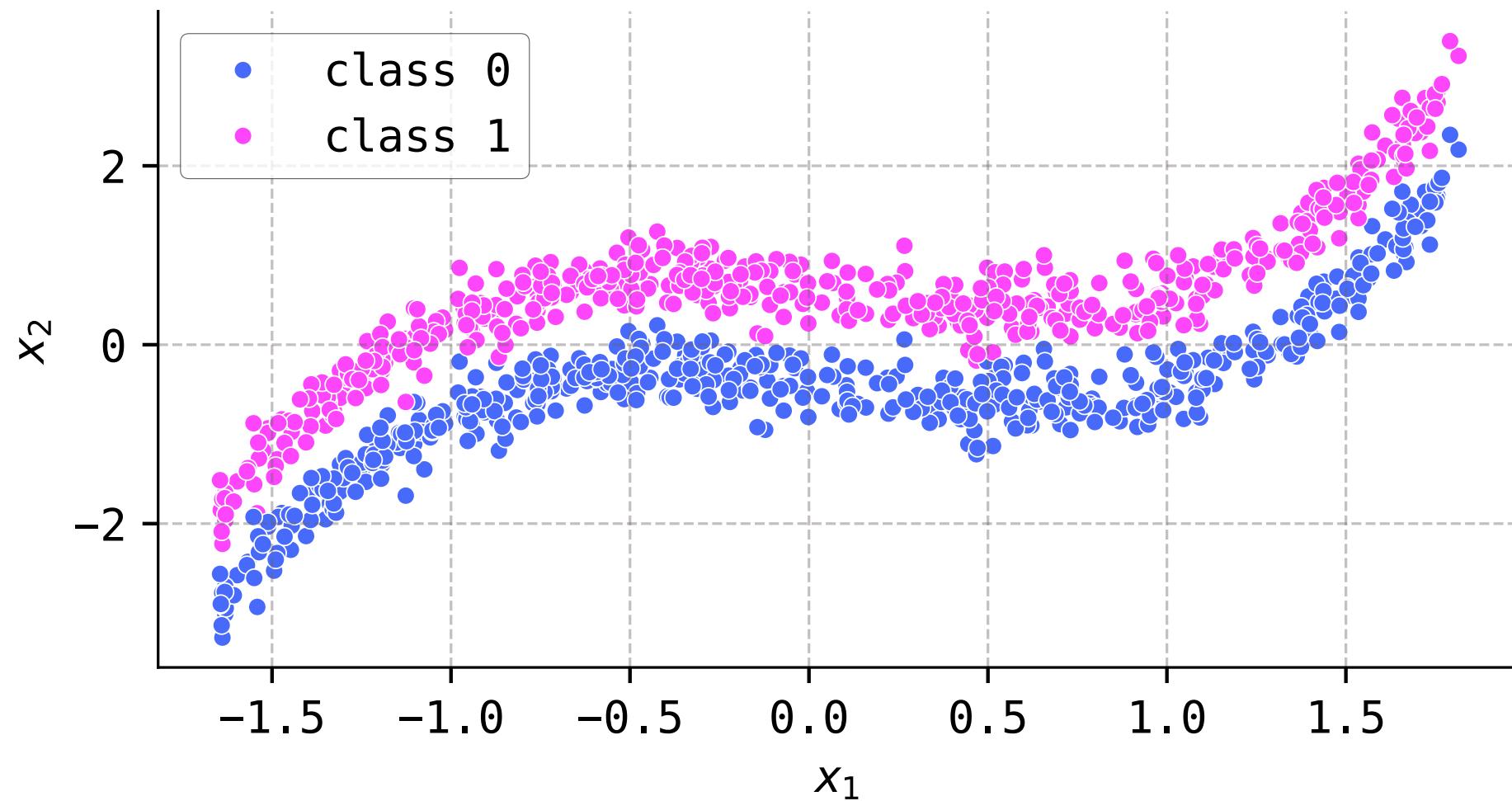
Consider a logistic regression model with predictor vector \mathbf{X} , response variable Y , and link function at Y given by

$$Y \mid \mathbf{X} \sim \text{Ber}(\phi) \quad \text{where} \quad \phi = \sigma(\beta_0 + \mathbf{x}^\top \boldsymbol{\beta}),$$

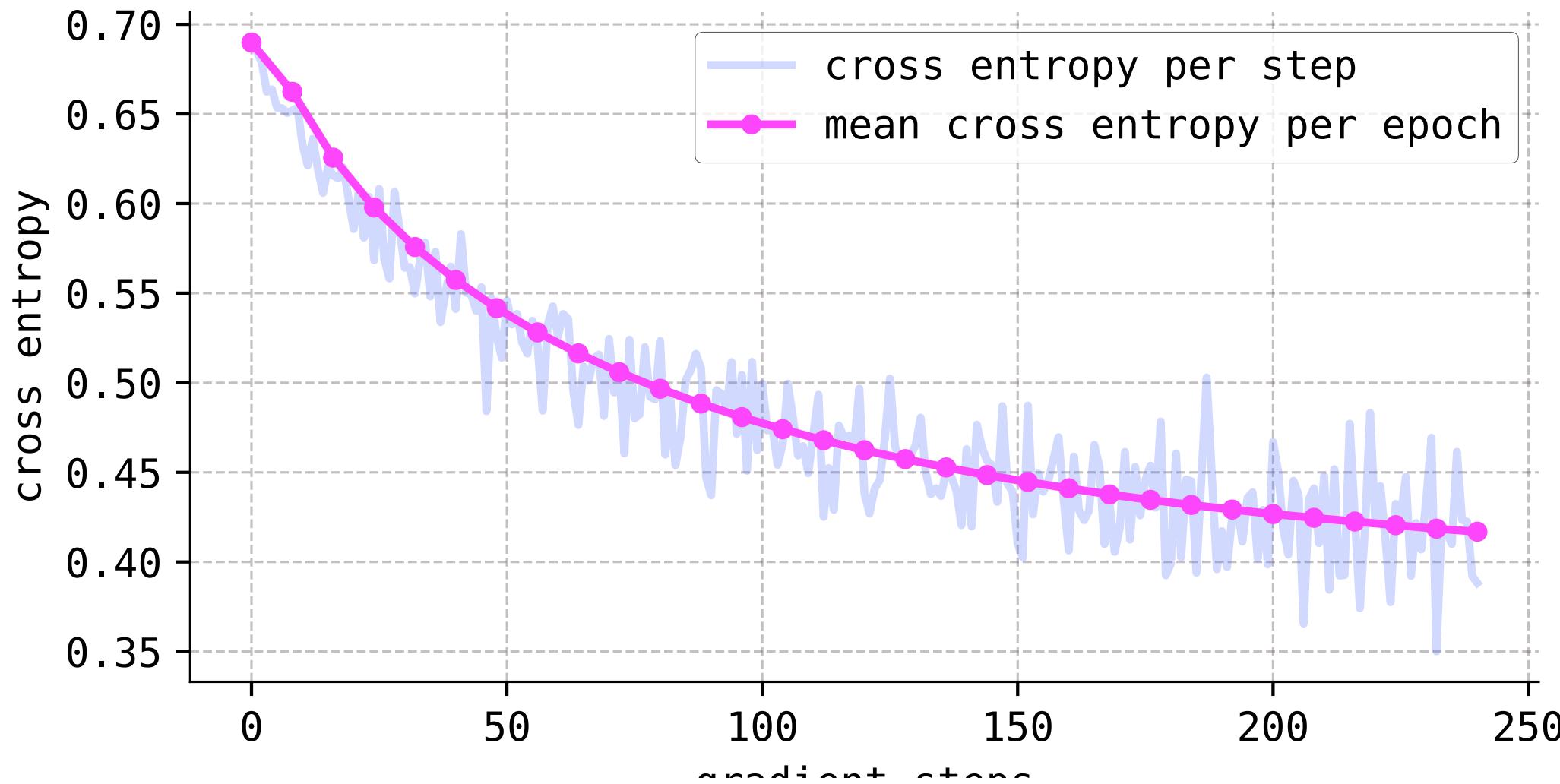
where σ is the sigmoid function. Then the model surprisal function is given by

$$\mathcal{I}_{\text{model}}(\boldsymbol{\theta}) = -y \log \phi - (1 - y) \log(1 - \phi).$$

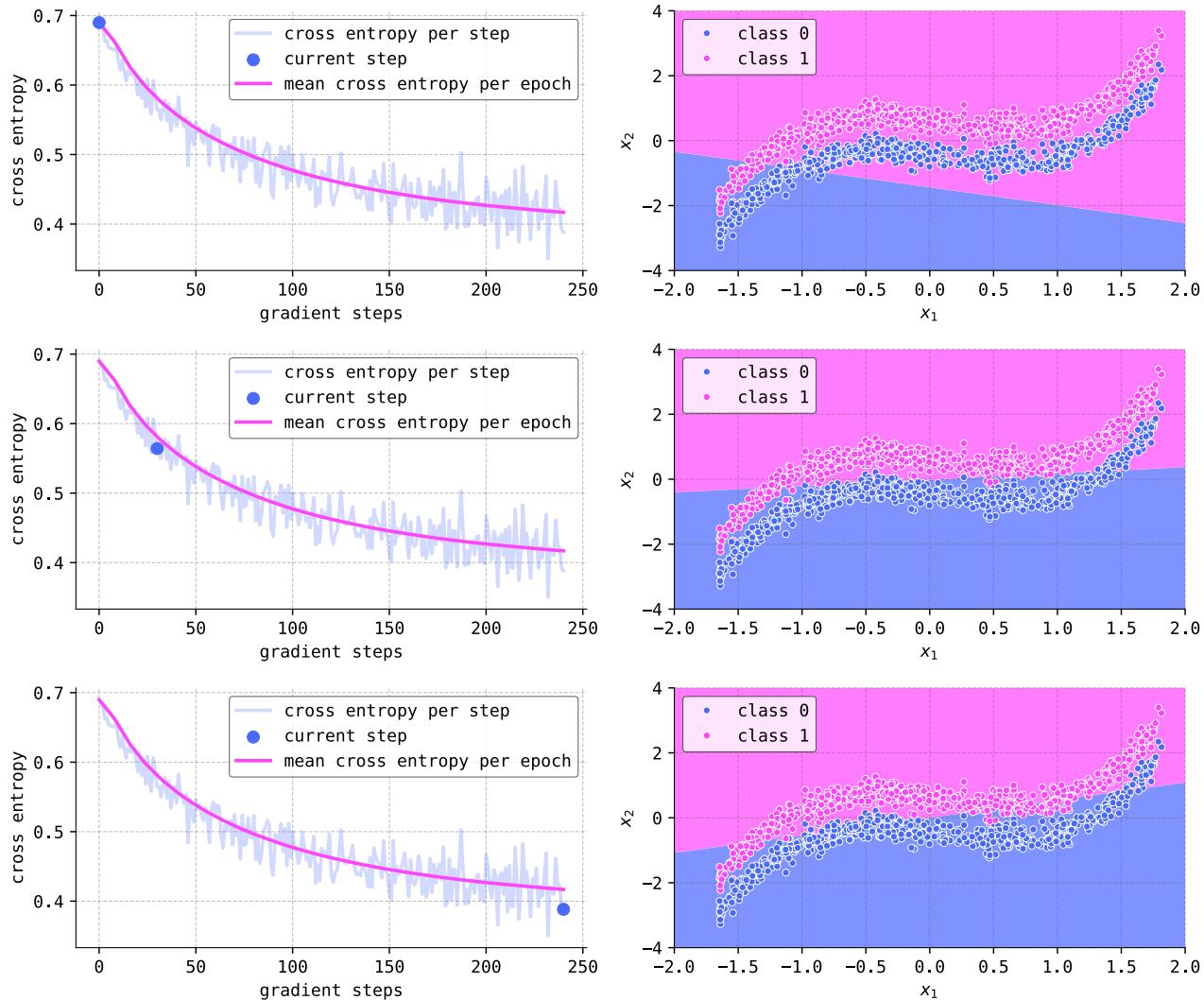
data for logistic regression



SGD for logistic regression model
 $\alpha = 0.1, \beta = 0, k = 128, N = 30$



stochastic gradient descent for logistic regression
 $\alpha = 0.1$, $\beta = 0$, $N = 30$



13.5. MLE for neural networks

Theorem 13.11 (Surprisal functions of neural network models)

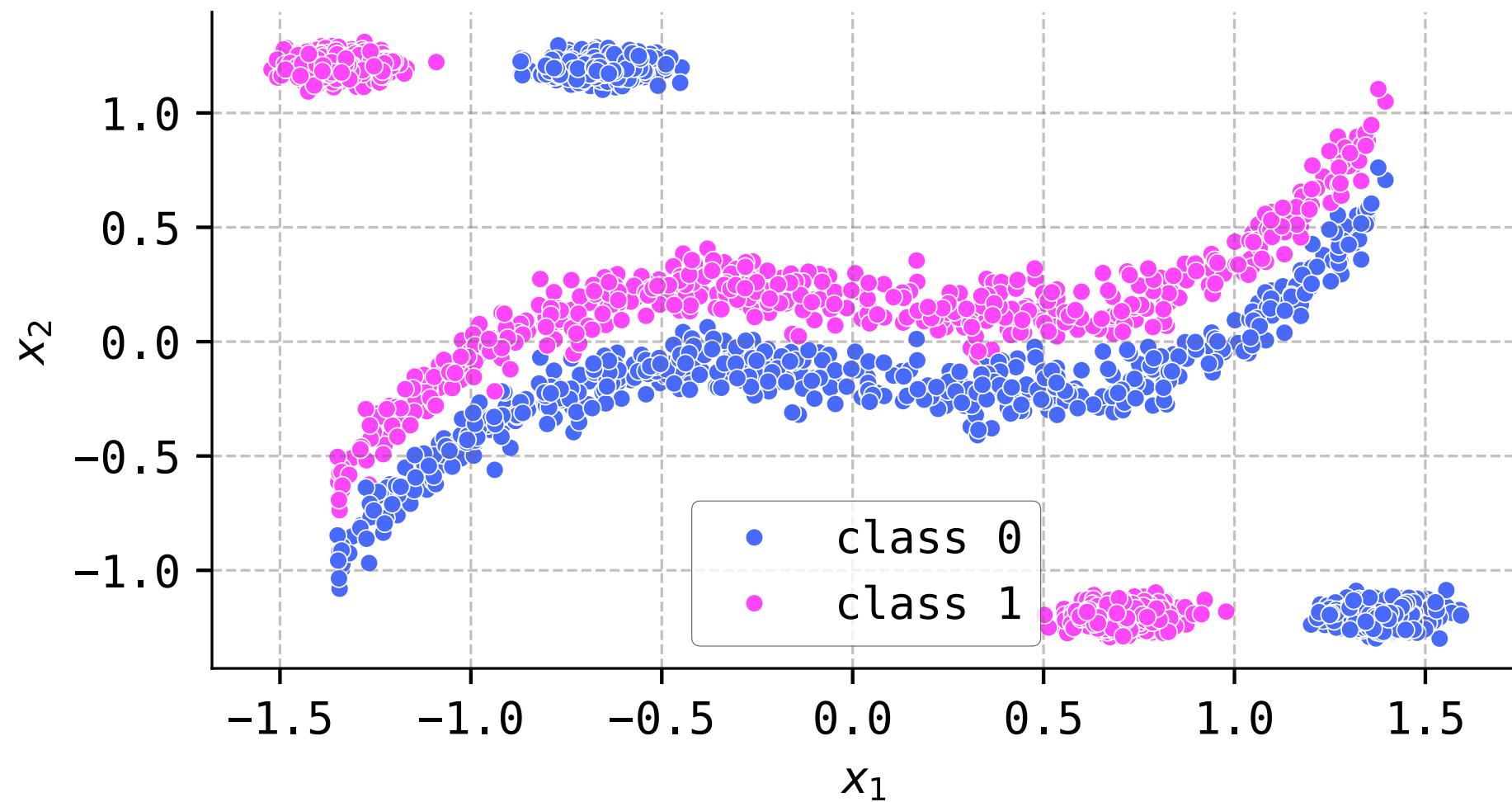
Consider a neural network model with a predictor vector \mathbf{X} , response variable Y , and link functions given by

$$\phi = \sigma(\mathbf{a}^\top \mathbf{w}_2 + b_2) \quad \text{and} \quad \mathbf{a}^\top = \rho(\mathbf{x}^\top \mathbf{W}_1 + \mathbf{b}_1^\top), \quad (13.12)$$

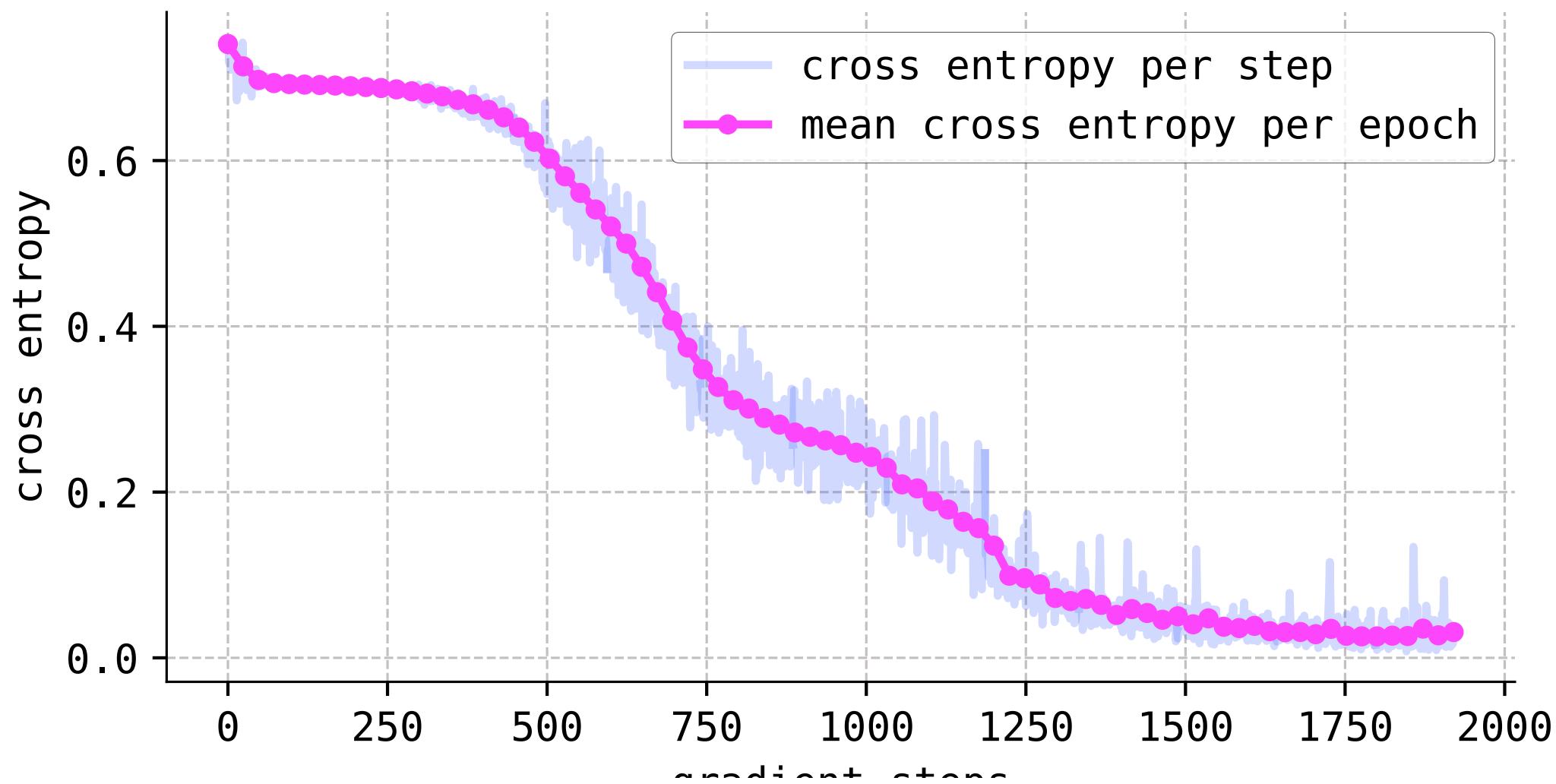
where σ is the sigmoid function and ρ the ReLU function. Then the model surprisal function is given by

$$\mathcal{I}_{\text{model}}(\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2) = -y \log \phi - (1 - y) \log(1 - \phi). \quad (13.13)$$

data for neural network model



SGD for neural network model
 $\alpha = 0.1$, $\beta = 0$, $k = 128$, $N = 80$



stochastic gradient descent for neural network model
 $\alpha = 0.1$, $\beta = 0$, $k = 128$, $N = 80$

