

# CS 446 MJT — Homework 6

*junzhew3*  
*haoxuans8*

Version 2

## Instructions.

- Homework is due **Tuesday, April 30, at 11:59pm**; no late homework accepted.
- Everyone must submit individually at gradescope under **hw6** and **hw6code**.
- The “written” submission at **hw6 must be typed**, and submitted in any format gradescope accepts (to be safe, submit a PDF). You may use L<sup>A</sup>T<sub>E</sub>X, markdown, google docs, MS word, whatever you like; but it must be typed!
- When submitting at **hw6**, gradescope will ask you to mark out boxes around each of your answers; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.
- We reserve the right to reduce the auto-graded score for **hw6code** if we detect funny business (e.g., rather than implementing an algorithm, you keep re-submitting the assignment to the auto-grader, eventually completing a binary search for the answers).
- There are **no regrade requests** on **hw6code**, which is the code auto-grader; however, you can re-submit and re-grade as many times as you like before the deadline! Start early and report any issues on piazza!
- Methods and functions in the template and utility code include docstrings to describe the inputs and outputs. The autograder relies on correct implementations of these methods. Follow the docstrings to avoid failing tests.

1. ***k*-means.**

Recall the *k*-means problem with  $n$  data points  $S = (\mathbf{x}_i)_{i=1}^n$  in  $\mathbb{R}^d$ ,  $k$  centers  $(\boldsymbol{\mu}_i)_{i=1}^k$  in  $\mathbb{R}^d$ , and corresponding clusters  $(S_j)_{j=1}^k$ . The objective of *k*-means is then to minimize the cost:

$$\phi_S(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k) = \sum_{i=1}^n \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2.$$

In this problem you will develop an alternate formulation of this cost function in terms of pairwise distances.

- (a) Let  $S_{\mathbf{z}} \subset S$  be the cluster induced by using a particular point  $\mathbf{z} \in \mathbb{R}^d$  as a center. Then the cost of using any point  $\mathbf{z}$  as a center is then given by

$$\phi_{S_{\mathbf{z}}}(\mathbf{z}) = \sum_{\mathbf{x} \in S_{\mathbf{z}}} \|\mathbf{x} - \mathbf{z}\|^2.$$

Let  $\boldsymbol{\mu}(S_{\mathbf{z}})$  be the sample mean of the cluster  $S_{\mathbf{z}}$ . Prove that the cost  $\phi_{S_{\mathbf{z}}}(\mathbf{z})$  is equivalent to

$$\phi_{S_{\mathbf{z}}}(\mathbf{z}) = \phi_{S_{\mathbf{z}}}(\boldsymbol{\mu}(S_{\mathbf{z}})) + |S_{\mathbf{z}}| \|\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z}\|^2.$$

- (b) Show that

$$\phi_{S_j}(\boldsymbol{\mu}_j) = \frac{1}{2|S_j|} \sum_{\mathbf{a}, \mathbf{b} \in S_j} \|\mathbf{a} - \mathbf{b}\|^2.$$

Conclude that solving the *k*-means problem is equivalent to solving

$$\min_{S_1, \dots, S_k} \sum_{j=1}^k \frac{1}{2|S_j|} \sum_{\mathbf{a}, \mathbf{b} \in S_j} \|\mathbf{a} - \mathbf{b}\|^2.$$

**Solution.**

(Your solution here.)

- (a) Assume  $S_{\mathbf{z}} = \{\mathbf{x}_1, \dots, \mathbf{x}_a\}$ ,  $|S_{\mathbf{z}}| = a$

$$\phi_{S_{\mathbf{z}}}(\mathbf{z}) = \sum_{\mathbf{x} \in S_{\mathbf{z}}} \|\mathbf{x} - \mathbf{z}\|^2 = \sum_{i=1}^a \|(\mathbf{x}_i - \boldsymbol{\mu}(S_{\mathbf{z}})) + (\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z})\|^2 \quad (1)$$

$$= \sum_{i=1}^a \|(\mathbf{x}_i - \boldsymbol{\mu}(S_{\mathbf{z}}))\|^2 + \|(\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z})\|^2 + 2(\mathbf{x}_i - \boldsymbol{\mu}(S_{\mathbf{z}})) \bullet (\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z}) \quad (2)$$

$$= \phi_{S_{\mathbf{z}}}(\boldsymbol{\mu}(S_{\mathbf{z}})) + a\|\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z}\|^2 + \sum_{i=1}^a 2(\mathbf{x}_i - \boldsymbol{\mu}(S_{\mathbf{z}})) \bullet (\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z}) \quad (3)$$

$$\sum_{i=1}^a (\mathbf{x}_i - \boldsymbol{\mu}(S_{\mathbf{z}})) \bullet (\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z}) = \sum_{i=1}^a (\mathbf{x}_i \bullet \boldsymbol{\mu}(S_{\mathbf{z}}) - \boldsymbol{\mu}(S_{\mathbf{z}}) \bullet \boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{x}_i \bullet \mathbf{z} + \boldsymbol{\mu}(S_{\mathbf{z}}) \bullet \mathbf{z}) \quad (4)$$

$$= \sum_{i=1}^a \mathbf{x}_i \bullet \boldsymbol{\mu}(S_{\mathbf{z}}) - a\boldsymbol{\mu}(S_{\mathbf{z}}) \bullet \boldsymbol{\mu}(S_{\mathbf{z}}) - \sum_{i=1}^a \mathbf{x}_i \bullet \mathbf{z} + a\boldsymbol{\mu}(S_{\mathbf{z}}) \bullet \mathbf{z} \quad (5)$$

Since  $\boldsymbol{\mu}(S_{\mathbf{z}}) = \frac{\sum_{i=1}^a \mathbf{x}_i}{a}$ , so  $\sum_{i=1}^a (\mathbf{x}_i - \boldsymbol{\mu}(S_{\mathbf{z}})) \bullet (\boldsymbol{\mu}(S_{\mathbf{z}}) - \mathbf{z}) = 0$ .

So:

$$\phi_{S_z}(z) = \phi_{S_z}(\mu(S_z)) + |S_z| \|\mu(S_z) - z\|^2.$$

(b)

$$\sum_{\mathbf{a}, \mathbf{b} \in S_j} \|\mathbf{a} - \mathbf{b}\|^2 = \sum_{\mathbf{a}, \mathbf{b} \in S_j} \left\| \left( \mathbf{a} - \mu(S_j) \right) + \left( \mu(S_j) - \mathbf{b} \right) \right\|^2 \quad (6)$$

$$= \sum_{\mathbf{a}, \mathbf{b} \in S_j} \left( \left\| \mathbf{a} - \mu(S_j) \right\|^2 + \left\| \mu(S_j) - \mathbf{b} \right\|^2 + 2 \left( \mathbf{a} - \mu(S_j) \right) \bullet \left( \mu(S_j) - \mathbf{b} \right) \right) \quad (7)$$

Since  $\sum_{i=1}^a (\mathbf{x}_i - \mu(S_z)) \bullet (\mu(S_z) - z) = 0$  is proved in (a), so:

$$\sum_{\mathbf{a}, \mathbf{b} \in S_j} 2 \left( \mathbf{a} - \mu(S_j) \right) \bullet \left( \mu(S_j) - \mathbf{b} \right) = 0 \quad (8)$$

So:

$$\sum_{\mathbf{a}, \mathbf{b} \in S_j} \|\mathbf{a} - \mathbf{b}\|^2 = 2 |S_j| \sum_{\mathbf{x} \in S_j} \left\| \mathbf{x} - \mu(S_j) \right\|^2 \quad (9)$$

So:

$$\min_{S_1, \dots, S_k} \sum_{j=1}^k \frac{1}{2 |S_j|} \sum_{\mathbf{a}, \mathbf{b} \in S_j} \|\mathbf{a} - \mathbf{b}\|^2 = \min_{S_1, \dots, S_k} \sum_{j=1}^k \sum_{\mathbf{x} \in S_j} \left\| \mathbf{x} - \mu(S_j) \right\|^2 = \min_{\mathbf{x} \in S} \sum_j \min_j \left\| \mathbf{x} - \mu_j \right\|^2 \quad (10)$$

So it is proved.

## 2. Wasserstein Distance.

Consider two discrete distributions with weights  $(\alpha_i)_{i=1}^n$  and  $(\beta_j)_{j=1}^m$  on points  $(\mathbf{x}_i)_{i=1}^n$  and  $(\mathbf{z}_j)_{j=1}^m$ . The Wasserstein distance between these two distributions (let's call them  $\mu$  and  $\nu$ ) is

$$W(\mu, \nu) = \max_{\|f\|_{\text{Lip}} \leq 1} \sum_{i=1}^n \alpha_i f(\mathbf{x}_i) - \sum_{j=1}^m \beta_j f(\mathbf{z}_j).$$

- (a) Suppose  $n = m$  and  $\alpha_i = \beta_i = 1/n$ , meaning both distributions are uniform. Show that for any permutation  $\pi$  of  $(1, \dots, n)$ .

$$W(\mu, \nu) \leq \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\|.$$

Note that this implies  $W(\mu, \nu) \leq \min_{\pi} \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\|$ .

- (b) Choose  $((\alpha_i, \mathbf{x}_i))_{i=1}^n$  and  $((\beta_j, \mathbf{z}_j))_{j=1}^m$  with  $m = n$  so that

$$0 < W(\mu, \nu) = \min_{\pi} \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\|.$$

- (c) Choose  $((\alpha_i, \mathbf{x}_i))_{i=1}^n$  and  $((\beta_j, \mathbf{z}_j))_{j=1}^m$  with  $m = n$  so that

$$0 < W(\mu, \nu) \leq \frac{1}{100} \min_{\pi} \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\|.$$

**Solution.** (Your solution here.)

- (a)

$$W(\mu, \nu) = \max_{\|f\|_{\text{Lip}} \leq 1} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f(\mathbf{z}_i)) = \max_{\|f\|_{\text{Lip}} \leq 1} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f(\mathbf{z}_{\pi(i)})) \quad (11)$$

Since  $\|f\|_{\text{Lip}} \leq 1$ ,  $\|f(\mathbf{x}_i) - f(\mathbf{z}_i)\| \leq \|\mathbf{x}_i - \mathbf{z}_i\|$ . So:

$$W(\mu, \nu) = \max_{\|f\|_{\text{Lip}} \leq 1} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f(\mathbf{z}_{\pi(i)})) \leq \max_{\|f\|_{\text{Lip}} \leq 1} \frac{1}{n} \sum_{i=1}^n \|f(\mathbf{x}_i) - f(\mathbf{z}_{\pi(i)})\| \quad (12)$$

$$= \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\| \leq \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\| \quad (13)$$

So:  $W(\mu, \nu) \leq \min_{\pi} \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\|$ .

- (b) Suppose  $\alpha_i = \beta_i = 1/n$ ,  $\mathbf{x}_i = a$  and  $\mathbf{z}_i = b$ ,  $a > b$ ,  $f(x) = x$ . So:

$$W(\mu, \nu) > 0, \max_{\|f\|_{\text{Lip}} \leq 1} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f(\mathbf{z}_{\pi(i)})) = \max_{\|f\|_{\text{Lip}} \leq 1} \frac{1}{n} \sum_{i=1}^n \|f(\mathbf{x}_i) - f(\mathbf{z}_{\pi(i)})\| \quad (14)$$

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\| = \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\| = a - b \quad (15)$$

So  $W(\mu, \nu) = \min_{\pi} \max_i \|\mathbf{x}_i - \mathbf{z}_{\pi(i)}\|$

- (c) Suppose  $\alpha_i = \beta_i = \frac{1}{100n}$  and  $\mathbf{x}_i = a$  and  $\mathbf{z}_i = b$ ,  $a > b$ ,  $f(x) = x$ .

### 3. Boosting.

In this problem we will consider boosting applied to interval classifiers on the real line. An interval classifier has the form  $h(x) := \mathbb{1}[a \leq x \leq b]$ ; let  $\mathcal{H}$  denote all such classifiers (meaning for all  $a \leq b$ ). Boosting therefore outputs a function of the form

$$g(x) = \sum_{j=1}^m \alpha_j h_j(x) = \sum_{j=1}^m \alpha_j \cdot \mathbb{1}[a_j \leq x \leq b_j].$$

For all parts of this problem let  $(x_i, y_i)_{i=1}^n$  be a data set of  $n$  points  $x_i \in \mathbb{R}$  along with associated labels  $y_i \in \{-1, 1\}$ . Assume that the  $x_i$  are in sorted order and distinct, meaning  $x_i < x_{i+1}$ .

(a) Let  $(q_1, \dots, q_n)$  be any weights on the training set, meaning  $q_i \geq 0$  and  $\sum_i q_i = 1$ . Show that

$$\min_{h \in \mathcal{H}} \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i] \leq \min \left\{ \sum_{\substack{i \in \{1, \dots, n\} \\ y_i > 0}} q_i, \sum_{\substack{i \in \{1, \dots, n\} \\ y_i < 0}} q_i \right\}.$$

**Remark.** This calculation is related to the “weak learning assumption” discussed in lecture. The only difference is these predictors map to  $\{0, 1\}$ , rather than  $\{-1, +1\}$ .

(b) Show that

$$\min_{h \in \mathcal{H}} \sum_{i=1}^n \frac{1}{n} \mathbb{1}[2h(x_i) - 1 \neq y_i] \leq \frac{n - L}{n},$$

where  $L$  is the length of the longest contiguous subsequence of examples having the same labels, meaning  $y_j = y_{j+1} = \dots = y_{j+L-1}$  for some  $j$ .

(c) Show that there exists an integer  $m$ , reals  $(\alpha_1, \dots, \alpha_m)$ , and interval classifiers  $(h_1, \dots, h_m)$  with  $h_j \in \mathcal{H}$  so that, for every  $i$ ,

$$y_i = \sum_{j=1}^m \alpha_j h_j(x_i).$$

In other words, that there exists a perfect boosted interval classifier.

**Solution.** (*Your solution here.*)

(a) for any  $(x_i, y_i)_{i=1}^n$ , There is a classifier  $h_1$  that predict  $h(x_i) = 1$  ( $a < (x_i)_{i=1}^n < b$ ), and another classifier  $h_2$  that predict  $h(x_i) = 0$  ( $(x_i)_{i=1}^n < a$  or  $(x_i)_{i=1}^n > b$ ). So:

$$\min \left\{ \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i], \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i] \right\} = \min \left\{ \sum_{\substack{i \in \{1, \dots, n\} \\ y_i > 0}} q_i, \sum_{\substack{i \in \{1, \dots, n\} \\ y_i < 0}} q_i \right\} \quad (16)$$

$h_1, h_2 \in \mathcal{H}$ , So :

$$\min_{h \in \mathcal{H}} \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i] \leq \min \left\{ \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i], \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i] \right\} \quad (17)$$

$$\min_{h \in \mathcal{H}} \sum_{i=1}^n q_i \mathbb{1}[2h(x_i) - 1 \neq y_i] \leq \min \left\{ \sum_{\substack{i \in \{1, \dots, n\} \\ y_i > 0}} q_i, \sum_{\substack{i \in \{1, \dots, n\} \\ y_i < 0}} q_i \right\} \quad (18)$$

(b) if  $y_j = y_{j+1} = \dots = y_{j+L-1} = 1$  for some  $j$ . There is a classifier  $h_1$  where  $a < x_j$  and

$b > x_{j+L-1}$  to make  $y_j, y_{j+1}, \dots, y_{j+L-1}$  all correct. So:

$$\sum_{i=1}^n \frac{1}{n} \mathbb{1}[2h_1(x_i) - 1 \neq y_i] \leq \frac{n-L}{n} \quad (19)$$

if  $y_j = y_{j+1} = \dots = y_{j+L-1} = -1$  for some  $j$ . There is a classifier  $h_2$  where  $a > x_{j+L-1}$  or

$b < x_j$  to make  $y_j, y_{j+1}, \dots, y_{j+L-1}$  all correct. So:

$$\sum_{i=1}^n \frac{1}{n} \mathbb{1}[2h_2(x_i) - 1 \neq y_i] \leq \frac{n-L}{n} \quad (20)$$

$h_1, h_2 \in \mathcal{H}$ , So :

$$\min_{h \in \mathcal{H}} \sum_{i=1}^n \frac{1}{n} \mathbb{1}[2h(x_i) - 1 \neq y_i] \leq \frac{n-L}{n},$$

(c) There are  $m$  classifiers  $(h_1, \dots, h_m)$  where  $a_1 \leq x_1 \leq b_1 \leq a_2 \leq x_2 \leq b_2 \leq \dots \leq a_m \leq x_m \leq b_m$ .

So for each classifier  $h_i$ ,  $h_i(x_i) = 1, h_i(x_j) = 0$ . We can assume  $(\alpha_1, \dots, \alpha_m) = (y_1, \dots, y_m)$ , so:

$$y_i = \sum_{j=1}^m \alpha_j h_j(x_i) = \alpha_i h_i(x_i) = y_i$$

So it is proved.

#### 4. Variational Autoencoders.

In this problem you will implement a Variational Autoencoder (VAE) to model points sampled from an unknown distribution. This will be done by constructing an encoder network and a decoder network. The encoder network  $f_{\text{enc}} : X \subset \mathbb{R}^2 \rightarrow \mathbb{R}^h \times \mathbb{R}^h$  takes as input a point  $\mathbf{x}$  from the input space and outputs parameters  $(\boldsymbol{\mu}, \boldsymbol{\xi})$  where  $\boldsymbol{\xi} = \log \boldsymbol{\sigma}^2$ . The decoder network  $f_{\text{dec}} : \mathbb{R}^h \rightarrow \mathbb{R}^2$  takes as input a latent vector  $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$  and outputs an element  $\hat{\mathbf{x}} \in \mathbb{R}^2$  that we would hope is similar to members of the input space  $X$ . You will train this model by minimizing the (regularized) empirical risk

$$\widehat{\mathcal{R}}_{\text{VAE}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\text{dec}}(f_{\text{enc}}(\mathbf{x})), \mathbf{x}) + \lambda \text{KL}(\mathcal{N}(\boldsymbol{\mu}(\mathbf{x}_i), \exp(\boldsymbol{\xi}(\mathbf{x}_i)/2)), \mathcal{N}(0, I)).$$

- (a) Let  $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\sigma}^2)$ . In your written submission show that

$$\text{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathcal{N}(0, I)) = -\frac{1}{2} \left[ h + \sum_{j=1}^h \left( \log \sigma_j^2 - \mu_j^2 - \sigma_j^2 \right) \right],$$

where  $\text{KL}(p, q) = \int p(\mathbf{x}) \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$  is the *KL divergence* between two densities  $p, q$ . You may use the fact that the KL-divergence between two  $h$ -dimensional normal distributions  $\mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  is given by

$$\text{KL}(\mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)) = \frac{1}{2} \left( \text{tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - h + \ln \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} \right).$$

- (b) Use the empirical risk discussed above to implement a VAE in the class `VAE`. Use ReLU activations between each layer, except on the last layer of the decoder use sigmoid. Use the ADAM optimizer to optimize in the `step()` function. Make use of the PyTorch library for this. Use `torch.optim.Adam()`, there is no need to implement it yourself. Please refer to the docstrings in `hw6.py` for more implementation details.
- (c) Implement the `fit` function using the `net.step()` function from the `VAE` class. See the docstrings in `hw6.py` for more details.
- (d) Fit a VAE on the data generated by `generate_data` in `hw6_utils.py`. Use a learning rate  $\eta = 0.01$ , latent space dimension  $h = 6$ , KL-divergence scaling factor  $\lambda = 5 \times 10^{-5}$ , and train for 8000 iterations. Use least squares as the loss, that is, let  $\ell(f(\mathbf{x}), \hat{\mathbf{x}}) = \|f(\mathbf{x}) - \hat{\mathbf{x}}\|_2^2$ . Include separate plots of each of the following in your written submission:
- Your empirical risk  $\widehat{\mathcal{R}}_{\text{VAE}}$  on the training data vs iteration count;
  - The data points  $(\mathbf{x})_{i=1}^n$  along with their encoded and decoded approximations  $\hat{\mathbf{x}} = f_{\text{dec}}(f_{\text{enc}}(\mathbf{x}))$ ;
  - The data points  $(\mathbf{x})_{i=1}^n$  along with their encoded and decoded approximations  $\hat{\mathbf{x}}$ , and  $n$  generated points  $f_{\text{dec}}(\mathbf{z})$  where  $\mathbf{z} \sim \mathcal{N}(0, I)$ .

After you are done training, save your neural network to a file using `torch.save(model.cpu().state_dict(), "vae.pb")`. You will submit this file to the autograder with your code submission.

- (e) What is the difference between the  $\hat{\mathbf{x}}$  and  $f_{\text{dec}}(\mathbf{z})$  in general? Why are they different in the plots?
- (f) Repeat part (d) except this time use L1 as your loss, that is let  $\ell(f(\mathbf{x}), \hat{\mathbf{x}}) = \|f(\mathbf{x}) - \hat{\mathbf{x}}\|_1 = \sum_{j=1}^2 |x_j - \hat{x}_j|$ . Again, be sure to include the plots in your written submission.
- (g) Fit a VAE with  $\lambda \in \{1, 0.01, 0.001\}$  and L1 loss on the same data again, but this time only plot (iii) from part (d). Discuss your results. Do you expect the VAE to generalize more closely to the true distribution better or worse as you increase  $\lambda$ ? Out of all of the parameters you tried including  $5 \times 10^{-5}$ , which  $\lambda$  parameter seems to give the right balance? Be sure to provide a brief justification for your choice.

**Solution.**

(a)

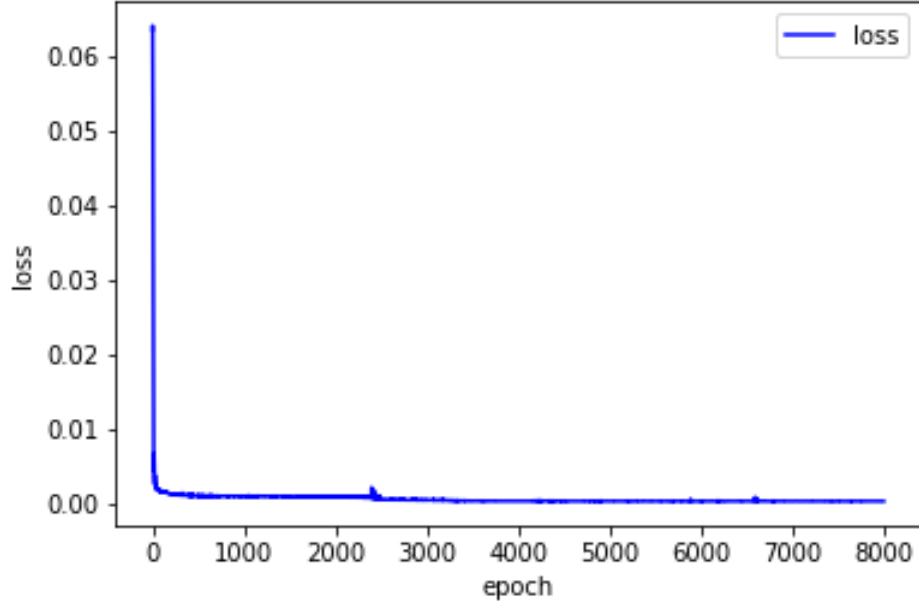
$$\text{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathcal{N}(0, I)) = -\frac{1}{2} \left( h + \ln(|\boldsymbol{\Sigma}|) - \text{tr}(\boldsymbol{\Sigma}) - \boldsymbol{\mu}^T \boldsymbol{\mu} \right) \quad (21)$$

Since  $\boldsymbol{\Sigma} = \text{diag}(\boldsymbol{\sigma}^2)$ ,  $\text{tr}(\boldsymbol{\Sigma}) = \sum_{j=1}^h \boldsymbol{\sigma}_j^2$ .

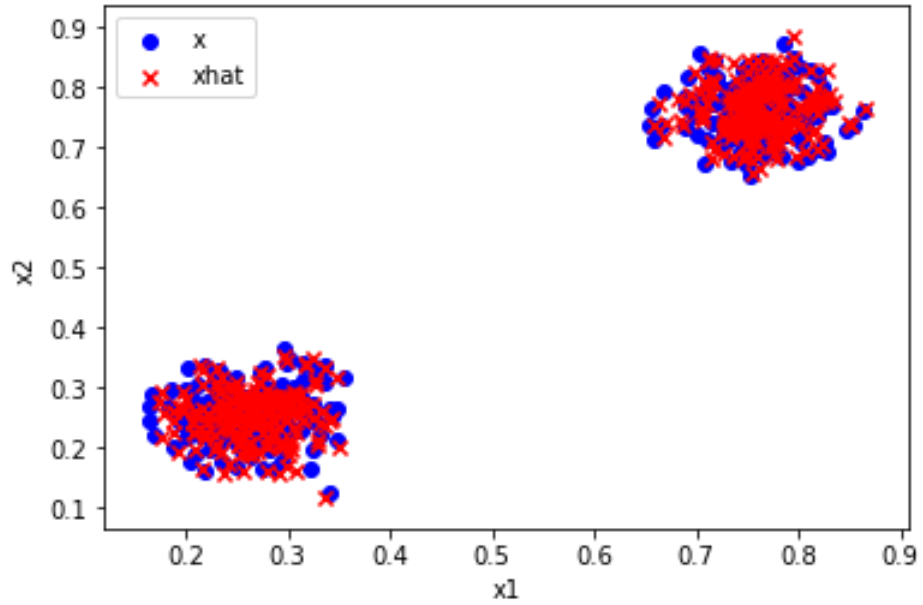
$\boldsymbol{\mu}^T \boldsymbol{\mu} = \sum_{j=1}^h \boldsymbol{\mu}_j^2$  and  $\ln(|\boldsymbol{\Sigma}|) = \ln(\boldsymbol{\sigma}_1^2 \boldsymbol{\sigma}_2^2 \cdots \boldsymbol{\sigma}_h^2) = \sum_{j=1}^h \ln(\boldsymbol{\sigma}_j^2)$ . So:

$$\text{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathcal{N}(0, I)) = -\frac{1}{2} \left[ h + \sum_{j=1}^h \left( \log \sigma_j^2 - \mu_j^2 - \sigma_j^2 \right) \right] \quad (22)$$

(d) 1.

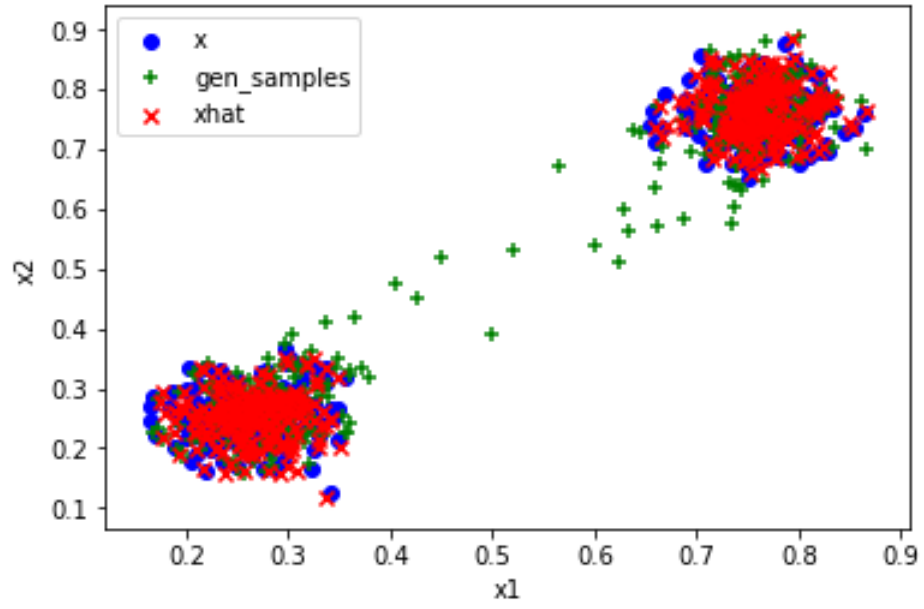


2.



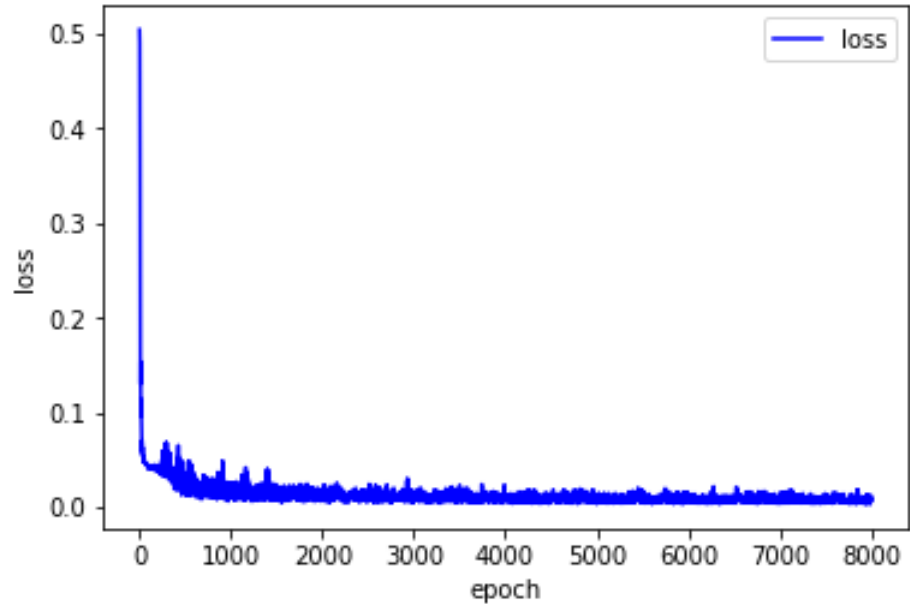


3.

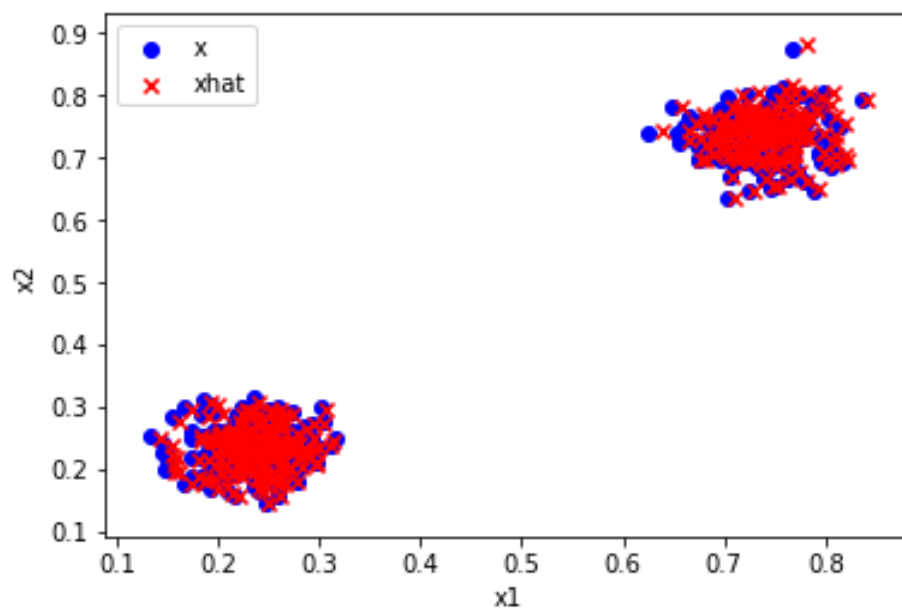


(e)  $\hat{\mathbf{x}}$  is similar to  $\mathbf{x}$ . Most points of  $f_{\text{dec}}(\mathbf{z})$  are similar to  $\mathbf{x}$ . But  $f_{\text{dec}}(\mathbf{z})$  have few points between the two distributions of  $\mathbf{x}$ .

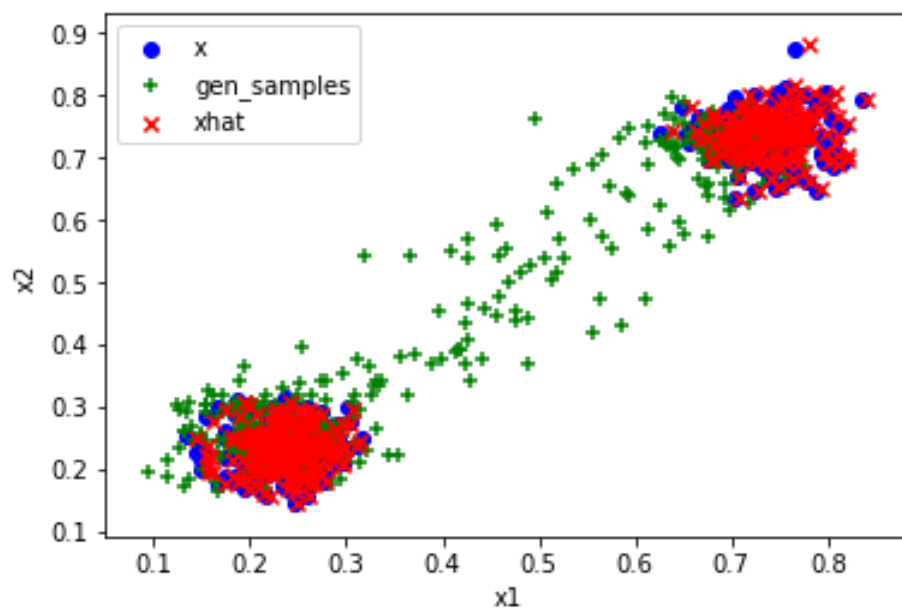
(f) 1.



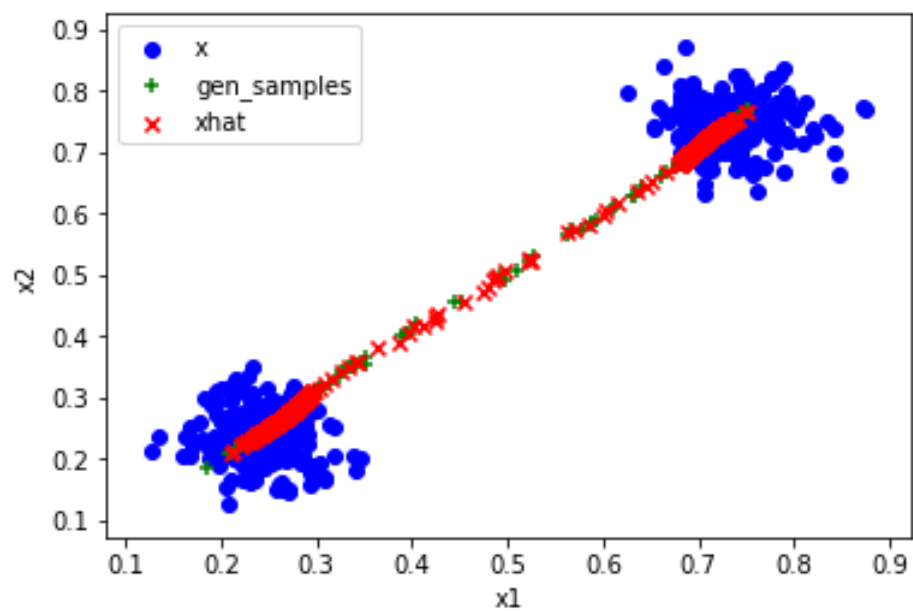
2.



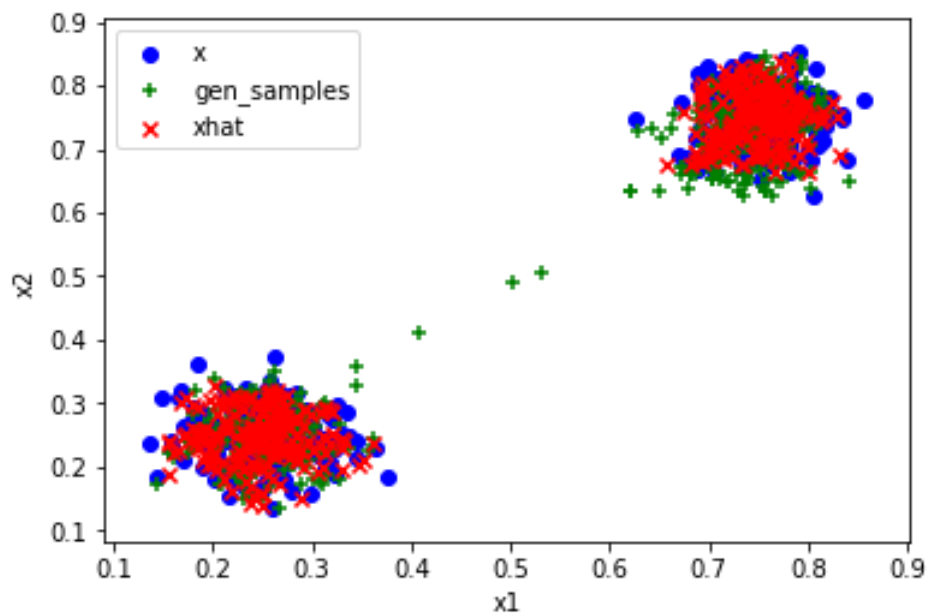
3.



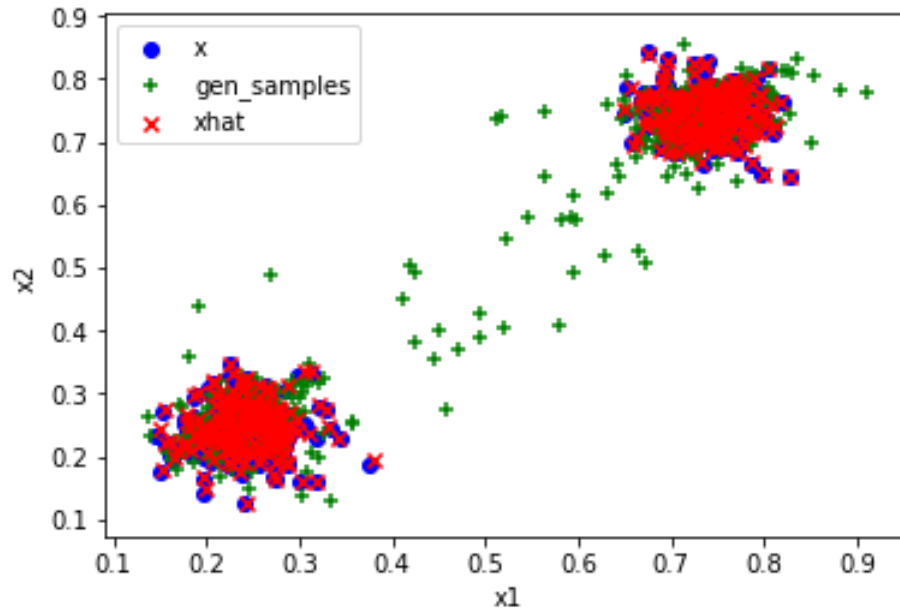
(g) 1.  $\lambda = 1$



2.  $\lambda = 0.01$



3.  $\lambda = 0.001$



The VAE generalizes  $\hat{\mathbf{x}}$  more closely to the true distribution as I decrease  $\lambda$ , but the gen-samples will not have a good distribution if  $\lambda$  is too small. The outputs will be a line if the  $\lambda$  is too large.

The distribution is best when  $\lambda$  is 0.01 since both distributions of  $\hat{\mathbf{x}}$  and gen-samples are good.

Too small  $\lambda$  will ignore the similarity between the output of encoder and Gaussian distribution. And too large  $\lambda$  will also ignore the similarity between the output and  $\mathbf{x}$ .

5. Naive Bayes (Extra credit!).

Let  $\mathbf{X} = (X_1, \dots, X_d)$  be a vector of  $d$  binary random variables whose distribution, labeled by a boolean function  $f : \{0, 1\}^d \rightarrow \{0, 1\}$ . Naive bayes proceeds by forming estimates of various probabilities, and predicting with

$$\hat{f}(\mathbf{x}) = \arg \max_y \widehat{\Pr}(Y = y) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = y).$$

- (a) Suppose  $f(\mathbf{x}) = \mathbb{1}\left(\sum_{j=1}^d x_j \geq \frac{d}{2}\right)$ , and that the various  $\widehat{\Pr}$  estimates in  $\hat{f}$  are exact. Show that the naive Bayes predictor  $\hat{f}(\mathbf{x})$  classifies perfectly in this case. For this problem you can assume  $d$  is odd.

**Hint.** Use symmetry arguments to make computing the probabilities easier.

- (b) Under the same setup from part(a), construct a boolean function  $f : \{0, 1\}^3 \rightarrow \{0, 1\}$  for which naive Bayes will be unable to correctly classify every binary vector  $\mathbf{x} \in \{0, 1\}^3$ . Be sure to verify that your construction works.

**Solution.** (Your solution here.)

- (a) Suppose  $\widehat{\Pr}(Y = 0) = \widehat{\Pr}(Y = 1) = \frac{1}{2}$ , and  $\widehat{\Pr}(X_i = 0) = \widehat{\Pr}(X_i = 1) = \frac{1}{2}$ . So:

$$\widehat{\Pr}(X_i = 1 | Y = 1) = \frac{\widehat{\Pr}(Y = 1 | X_i = 1) * \widehat{\Pr}(X_i = 1)}{\widehat{\Pr}(Y = 1)} \quad (23)$$

$$= \widehat{\Pr}(Y = 1 | X_i = 1) = 1 - \widehat{\Pr}(Y = 0 | X_i = 1) = 1 - \widehat{\Pr}(X_i = 1 | Y = 0) = \widehat{\Pr}(X_i = 0 | Y = 0) \quad (24)$$

So in the same way, we can prove  $\widehat{\Pr}(X_i = 0 | Y = 1) = \widehat{\Pr}(X_i = 1 | Y = 0)$ .

Since  $f(\mathbf{x}) = \mathbb{1}\left(\sum_{j=1}^d x_j \geq \frac{d}{2}\right)$ ,  $\widehat{\Pr}(X_i = 1 | Y = 1) \geq \widehat{\Pr}(X_i = 0 | Y = 1)$ .

So for any  $\mathbf{X}$ , if true  $Y$  is 1, we can have  $n_1 = \sum_i x_i$ ,  $n_2 = d - \sum_i x_i$ , where  $n_1 \geq n_2$ . Then,

$$\widehat{\Pr}(Y = 1) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = 1) = \widehat{\Pr}(Y = 1) \widehat{\Pr}(X_i = 1 | Y = 1)^{n_1} \widehat{\Pr}(X_i = 0 | Y = 1)^{n_2}$$

$$\widehat{\Pr}(Y = 0) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = 0) = \widehat{\Pr}(Y = 0) \widehat{\Pr}(X_i = 1 | Y = 0)^{n_1} \widehat{\Pr}(X_i = 0 | Y = 0)^{n_2}$$

$$\begin{aligned} \frac{\widehat{\Pr}(Y = 1) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = 1)}{\widehat{\Pr}(Y = 0) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = 0)} &= \frac{\widehat{\Pr}(X_i = 1 | Y = 1)^{n_1} \widehat{\Pr}(X_i = 0 | Y = 1)^{n_2}}{\widehat{\Pr}(X_i = 1 | Y = 0)^{n_1} \widehat{\Pr}(X_i = 0 | Y = 0)^{n_2}} \\ &= \frac{\widehat{\Pr}(X_i = 1 | Y = 1)^{n_1 - n_2}}{\widehat{\Pr}(X_i = 1 | Y = 0)^{n_1 - n_2}} > 1 \end{aligned}$$

So in the same way, if true  $Y$  is 0, we can prove:

$$\frac{\widehat{\Pr}(Y = 0) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = 0)}{\widehat{\Pr}(Y = 1) \prod_{i=1}^d \widehat{\Pr}(X_i = x_i | Y = 1)} > 1 \quad (25)$$

So the naive Bayes predictor  $\hat{f}(\mathbf{x})$  classifies perfectly in this case.

(b) Suppose  $f(\mathbf{x}) = \mathbb{1}\left(\sum_{j=1}^3 x_j = 2\right)$ . Assume  $\mathbf{X} = (1, 1, 1)$  whose true value is  $f(\mathbf{X}) = 0$ . So:

$$2\widehat{\Pr}(Y = 1) \prod_{i=1}^3 \widehat{\Pr}(X_i = 1|Y = 1) = \left(\frac{2}{3}\right)^3 \quad (26)$$

$$2\widehat{\Pr}(Y = 0) \prod_{i=1}^3 \widehat{\Pr}(X_i = 1|Y = 0) = \left(\frac{2}{5}\right)^3 \quad (27)$$

So the prediction is wrong. So it is proved.