

Probabilistic Machine Learning

Assignment 1

Group 12

01/12/2024

1 Bishop 2.8

2 We will start by showing that

$$\mathbb{E}[x] = \mathbb{E}_y[\mathbb{E}_x[x \mid y]] \quad (1)$$

3 We assume that $p(x, y)$ is well defined and the expectations are integrateable.

4 We then have

$$\begin{aligned} \mathbb{E}[x] &= \int_x xp(x)dx && \text{(By definition)} \\ &= \int_x x \left(\int_y p(x, y)dy \right) dx \\ &= \int_y \left(\int_x xp(x, y)dx \right) dy \\ &= \int_y p(y) \left(\int_x \frac{xp(x, y)}{p(y)} dx \right) dy \\ &= \int_y p(y) \mathbb{E}_x[x \mid y] dy && \text{(Def. of conditional expectation)} \\ &= \mathbb{E}_y[\mathbb{E}_x[x \mid y]] && \text{(Def. of expectation)} \end{aligned}$$

5 We will now show that

$$\text{var}[x] = \mathbb{E}_y[\text{var}_x[x \mid y]] + \text{var}_y[\mathbb{E}_x[x \mid y]]$$

6 We have that

$$\begin{aligned} \text{var}[x] &= \mathbb{E}[x^2] - (\mathbb{E}[x])^2 && \text{(By definition)} \\ &= \mathbb{E}_y[\mathbb{E}_x[x^2 \mid y]] - (\mathbb{E}_y[\mathbb{E}_x[x \mid y]])^2 && \text{(Eq. 1)} \\ &= \mathbb{E}_y[\text{var}_x[x \mid y] + \mathbb{E}_x[x \mid y]] - (\mathbb{E}_y[\mathbb{E}_x[x \mid y]])^2 && \text{(Def. of variance)} \\ &= \mathbb{E}_y[\text{var}_x[x \mid y]] + \mathbb{E}_y[x \mid y] - (\mathbb{E}_y[\mathbb{E}_x[x \mid y]])^2 && \text{(Linearity of expectation)} \\ &= \mathbb{E}_y[\text{var}_x[x \mid y]] + \text{var}_y[\mathbb{E}_x[x \mid y]] && \text{(Def. of conditional variance)} \end{aligned}$$

7 ■

8 Bishop 8.9

9 We denote the Markov Blanket for some node x as $MB(x)$. A path from some
10 node $v \in \setminus MB(x) \cup x$ to x must pass through either a parent p , or a child, c of
11 x .

12 If it passes through p the path will either meet head-to-tail or tail-to-tail at p
13 and will therefore be blocked since p is conditioned on.

14 If the path passes through c it will either pass through some child of c or through
15 some co-parent of c . In the case where it passes through a child, it will meet
16 head-to-tail where the path then will be blocked since we condition on c . In the
17 case it passes through a co-parent it won't be blocked by c since it then will
18 meet head-to-head. Let denote this co-parent of c as cp . The path will then
19 meet cp either head-to-tail or tail-to-tail. Since we condition on cp the path will
20 therefore be blocked.

21 ■

22 Bishop 8.11

23 First we want to calculate $P(F = 0 \mid D = 0)$. We have from the book that
24 $P(G = 0 \mid F = 0) = 0.81$. We start by calculating the numerator in Bayes'
25 Theorem by

$$P(F = 0)P(D = 0 \mid F = 0) = P(F = 0) \sum_{G \in \{0,1\}} P(D = 0 \mid G)P(G \mid F = 0) = 0.0748$$

26 We have also from the book that $P(G = 0) = 0.315$. We can therefore calculate
27 the denominator by

$$P(D = 0) = \sum_{G \in \{0,1\}} P(D = 0 \mid G)P(G) = 0.352$$

28 We therefore have

$$P(F = 0 \mid D = 0) = \frac{P(D = 0 \mid F = 0)P(F = 0)}{P(D = 0)} = 0.2125$$

29 We now observe that $B = 0$. We then have

$$\begin{aligned} P(F = 0 \mid D = 0, B = 0) &= \frac{P(D = 0 \mid F = 0, B = 0)P(F = 0)}{P(D = 0 \mid B = 0)} \\ &= \frac{\sum_{G \in \{0,1\}} P(D = 0 \mid G)P(G \mid B = 0, F = 0)P(F = 0)}{\sum_{G \in \{0,1\}} \sum_{F \in \{0,1\}} P(D = 0 \mid G)P(G \mid B = 0, F)P(F)} \\ &\simeq 0.0587 \end{aligned}$$

30 We therefore see that the probability that tank is empty has decreased from
 31 0.2125 to 0.0587. This agrees with our intuition since we now know that the
 32 battery is flat which explains away the observation. The graphical model of the
 33 problem equates to the one we see in Fig 8.54. We also see how observing a
 34 descendant of G unblocks the path between B and F making them conditional
 35 dependent on D .

36 Week 1 Programming exercise

37 We start by sampling 20 points from $x \sim N(0, I), x \in \mathbb{R}^2$ and respective labels
 38 according to the distribution $p(y | x) = \mathcal{N}(y; x^T \theta, 0.1), \theta = [-1, 1]^T$. We can
 39 then compute the mean and covariance of the posterior by

$$\mu_{\text{posterior}} = \Sigma_{\text{posterior}} (X^T y / \sigma_y)$$

40 and

$$\Sigma_{\text{posterior}} = \left(\Sigma_{\text{prior}}^{-1} + \frac{1}{\sigma_y} X^T X \right)^{-1}$$

Where a plot of this distribution can be seen in Fig. 1. We can then calculate

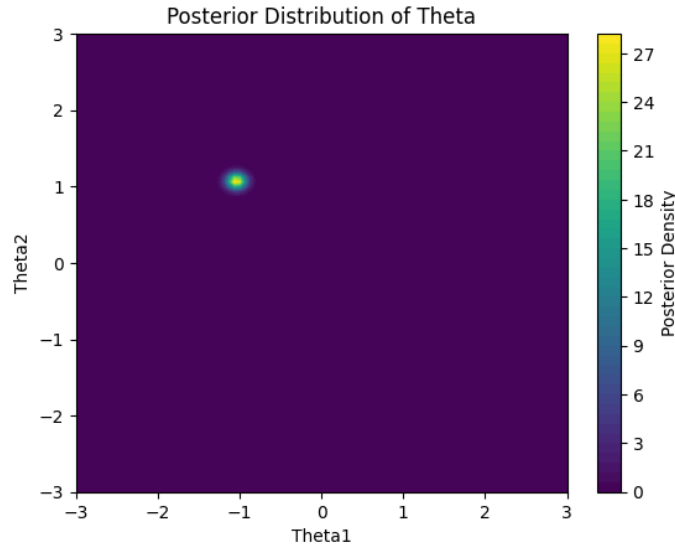


Figure 1: Posterior distribution for θ

41 the variance of the posterior predictive by $x^T \Sigma_{\theta} x + \sigma^2$. A plot of the variance
 42 can be seen on Fig. 2. . We can now repeat this but where we sample from
 43 $x \sim \mathcal{N}(0, \Sigma_x), \Sigma_X = \begin{bmatrix} 0.1 & 0 \\ 0 & 1 \end{bmatrix}$. A plot of the posterior distribution and the

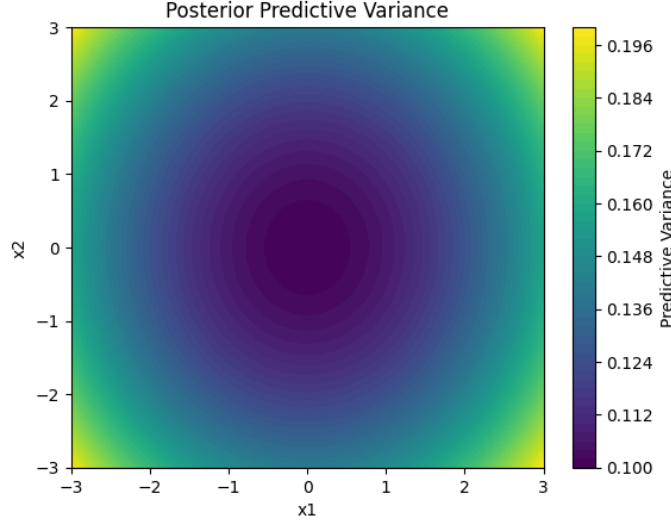


Figure 2: The variance of the posterior predictive

45 variance of the posterior predictive can be seen in Fig. 3. We can see that the
 46 posterior distribution now is wider than before. This is due to sampling with
 47 less variance along the x_1 -axis leading to less information about θ_1 and therefore
 48 a more uncertain (wider) estimate. If we where to introduce less noise for the
 49 labels the posterior distribution would become more narrow reflecting a higher
 50 certainty in parameter estimation. The predictive variance would also become
 51 reduced. .

52 Bishop 9.10

$$\begin{aligned}
 p(x_b|x_a) &= \frac{p(x_b, x_a)}{p(x_a)} && \text{Bayes theoreme} \\
 &= \frac{\sum_{k=1}^K \pi_k p(x_b, x_a|k)}{p(x_a)} && \text{using } p(x) = \sum_{k=1}^K p(x|k) \\
 &= \frac{\sum_{k=1}^K \pi_k p(x_b|x_a, k) p(x_a|k)}{\sum_{j=1}^K \pi_j p(x_a|k)} && \text{chain rule} \\
 &= \sum_{k=1}^K \frac{\pi_k p(x_a|k)}{\sum_{j=1}^K \pi_j p(x_a|k)} p(x_b|x_a, k) && \text{rearranging} \\
 &= \sum_{k=1}^K \pi'_k p(x_b|x_a, k) && \pi'_k = \frac{\pi_k p(x_a|k)}{\sum_{j=1}^K \pi_j p(x_a|k)}
 \end{aligned}$$

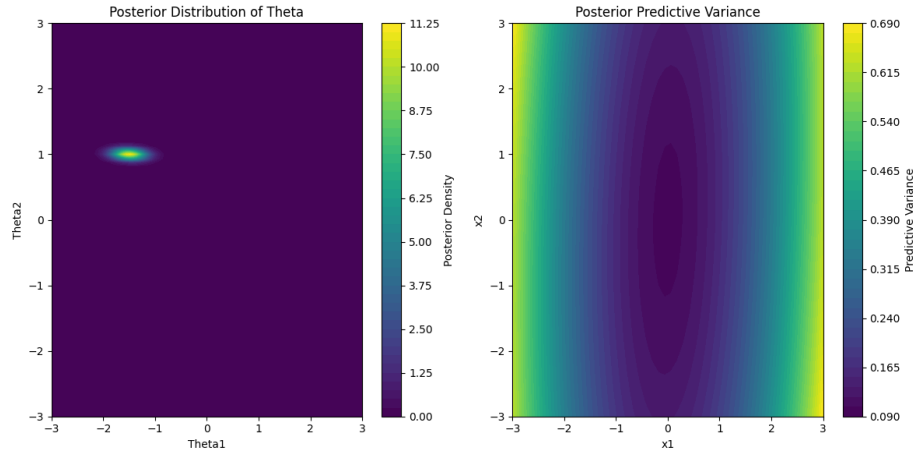


Figure 3: The variance of the posterior predictive

53 Hence we have shown expressions for mixing coefficients π'_k and component
54 densities $p(x_b|x_a, k) = \frac{p(x_b, x_a|k)}{p(x_a|k)}$

55 Bishop 10.4

56 We can find the minimization of the KL divergence between p and q by do-
 57 ing the derivative wrt μ and Σ and setting equal to zero. This works since
 58 $p(x)$ is fixed which implies $\text{KL}(p||q)$ is convex meaning local minimum is global
 59 minimum. Formally we want to find these two quantities i.e. $\frac{d}{d\mu} \text{KL}(p||q) =$
 60 $0 \wedge \frac{d}{d\Sigma} \text{KL}(p||q) = 0$. We start by deriving an expression usable for both expres-
 61 sions and then find $\frac{\partial}{\partial\mu} \log q(x) \wedge \frac{\partial}{\partial\Sigma} \log q(x)$

$$\begin{aligned} \text{KL}(p||q) &= \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx && \text{definition of } KL(p||q) \\ &= \int p(x) \log p(x) - p(x) \log q(x) dx && \frac{\log p}{\log q} = \log p - \log q \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial\mu} \log q(x) &= \frac{\partial}{\partial\mu} \log[(2\pi)^{-k/2} \det(\Sigma)^{-1/2} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))] \\ &= \frac{\partial}{\partial\mu} \left[-\frac{k}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma)) - \frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right] && \text{log rules} \\ &= \frac{\partial}{\partial\mu} \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right] && \text{non } \mu \text{ terms are zero} \\ &= -\frac{1}{2} \left[\frac{\partial}{\partial\mu} ((x-\mu)^T \Sigma^{-1} (x-\mu)) \right] \\ &= -\frac{1}{2} \left[\frac{\partial(x-\mu)^T}{\partial\mu} \Sigma^{-1} (x-\mu) + (x-\mu)^T \Sigma^{-1} \frac{\partial(x-\mu)}{\partial\mu} \right] && \text{chain rule} \\ &= -\frac{1}{2} [-\text{I} \Sigma^{-1} (x-\mu) - (x-\mu)^T \Sigma^{-1} \text{I}] \\ &= -\frac{1}{2} [-\Sigma^{-1} (x-\mu) - (x-\mu)^T \Sigma^{-1}] && \text{I}\Sigma^{-1} = \Sigma^{-1}\text{I} = \Sigma^{-1} \\ &= -\frac{1}{2} [-2\Sigma^{-1} (x-\mu)] && (AB)^T = B^T A^T \\ &= \Sigma^{-1} (x-\mu) \end{aligned}$$

62

From this point we find the derivative w.r.t μ

$$\begin{aligned}
\frac{d}{d\mu} \text{KL}(p||q) &= 0 = \frac{d}{d\mu} \int p(x) \log p(x) - p(x) \log q(x) dx \\
&= \int \frac{d}{d\mu} p(x) \log p(x) - \frac{d}{d\mu} p(x) \log q(x) dx \\
&= - \int p(x) \frac{d}{d\mu} \log q(x) dx \\
&= - \int p(x) \Sigma^{-1} (x - \mu) dx \\
&= - \Sigma^{-1} \int p(x) (x - \mu) dx \\
&= - \Sigma^{-1} \left[\int p(x) x dx - \mu \int p(x) dx \right] \quad \int p(x) dx = 1 \\
&= - (\mathbb{E}_{x \sim p(x)}[x] - \mu) \\
&= - \mathbb{E}_{x \sim p(x)}[x] + \mu
\end{aligned}$$

63

This is equivalent to $\mu = \mathbb{E}_{x \sim p(x)}[x]$

64

65

Now we do the derivative wrt to the covariance:

$$\begin{aligned}
\frac{\partial}{\partial \Sigma} \log q(x) &= \frac{\partial}{\partial \Sigma} \log[(2\pi)^{-k/2} \det(\Sigma)^{-1/2} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu))] \\
&= \frac{\partial}{\partial \Sigma} \left[-\frac{k}{2} \log(2\pi) - \frac{1}{2} \log(\det(\Sigma)) - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] \quad \text{log rules} \\
&= -\frac{1}{2} \frac{\partial}{\partial \Sigma} \log(\det(\Sigma)) - \frac{1}{2} (x - \mu)^T \frac{\partial}{\partial \Sigma} \Sigma^{-1} (x - \mu) \quad \text{non sigma terms are zero} \\
&= -\frac{1}{2} \frac{1}{\det(\Sigma)} \det(\Sigma) (\Sigma^{-1})^T - \frac{1}{2} (x - \mu)^T (-\Sigma^{-1} \Sigma^{-1}) (x - \mu) \\
&= -\frac{1}{2} \Sigma^{-1} - \frac{1}{2} (x - \mu)^T (-\Sigma^{-1} \Sigma^{-1}) (x - \mu) \\
&= -\frac{1}{2} [\Sigma^{-1} + (x - \mu)^T (-\Sigma^{-1} \Sigma^{-1}) (x - \mu)] \\
&= -\frac{1}{2} [\Sigma^{-1} - \Sigma^{-1} (x - \mu) (x - \mu)^T \Sigma^{-1}]
\end{aligned}$$

⁶⁶ We substitute in $\frac{\partial}{\partial \Sigma} \log q(x)$

$$\begin{aligned}
\frac{d}{d\Sigma} \text{KL}(p||q) &= 0 \\
&= \frac{d}{d\mu} \int p(x) \log p(x) - p(x) \log q(x) dx \\
&= \int \frac{d}{d\Sigma} p(x) \log p(x) - \frac{d}{d\Sigma} p(x) \log q(x) dx \\
&= - \int p(x) \frac{d}{d\Sigma} \log q(x) dx \\
&= - \int p(x) \left[-\frac{1}{2} (\Sigma^{-1} - \Sigma^{-1}(x - \mu)(x - \mu)^T \Sigma^{-1}) \right] dx \\
&= \frac{1}{2} \int p(x) [(\Sigma^{-1} - \Sigma^{-1}(x - \mu)(x - \mu)^T \Sigma^{-1})] dx \\
&= \frac{1}{2} \Sigma^{-1} \int p(x) dx - \frac{1}{2} \int p(x) \Sigma^{-1} (x - \mu)(x - \mu)^T \Sigma^{-1} dx \\
&= \frac{1}{2} \Sigma^{-1} - \frac{1}{2} \Sigma^{-1} \int p(x) (x - \mu)(x - \mu)^T dx \Sigma^{-1} \\
&= \frac{1}{2} \Sigma^{-1} - \frac{1}{2} \Sigma^{-1} \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T] \Sigma^{-1} \\
&\iff \\
\frac{1}{2} \Sigma^{-1} &= \frac{1}{2} \Sigma^{-1} \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T] \Sigma^{-1} \\
\Sigma^{-1} &= \Sigma^{-1} \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T] \Sigma^{-1} \\
\Sigma \Sigma^{-1} &= \Sigma \Sigma^{-1} \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T] \Sigma^{-1} \\
I &= I \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T] \Sigma^{-1} \\
I \Sigma &= I \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T] \Sigma^{-1} \Sigma \\
\Sigma &= \mathbb{E}_{x \sim p(x)} [(x - \mu)(x - \mu)^T]
\end{aligned}$$

⁶⁷ We've now shown what was asked in the exercise.

Week 2 Programming exercise

Part 1

We trained a VAE that downsampled in the following way; $28 \times 28 \rightarrow 14 \times 14 \rightarrow 7 \times 7$ using convolutions. This is linearly projected through two separate linear layers representing μ and Σ . We then used these parameters to obtain (z_1, z_2) , which was projected up into 7×7 . The upsampling steps are directly inverse proportional to the downsampling steps i.e. $7 \times 7 \rightarrow 14 \times 14 \rightarrow 28 \times 28$. The model was trained for 50 epochs using Adam and $lr = 3e - 4$. The most important detail is we weighted KL by an arbitrarily set factor of 0.001 in order to avoid mode collapse which was the default behavior.

We notice the points approximately are within a disk around origo which is to be expected as our latent space representation tries to approximate $\mathcal{N}(0, I)$.

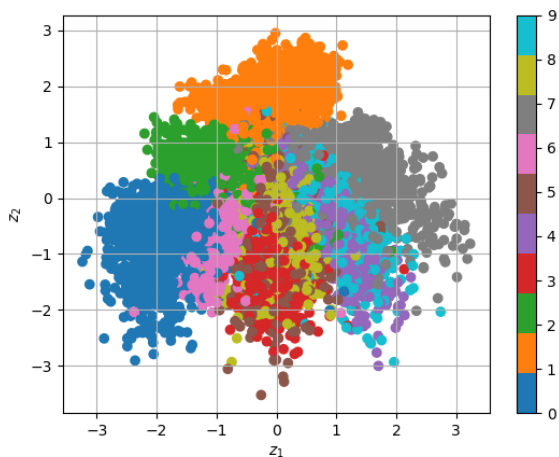


Figure 4: Encoded test samples with color-code for the different 10 digit classes.

Week 2 Programming exercise

We chose to use make a plot of $k = 15$ i.e. (15×15) Looking at the distribution $p(x|z)$, we notice it reflects the plot shown in part 1 in that members of each class are placed close to eachother. More interestingly we can better tell in this plot that classes similar to eachother are generally closer to eachother. This can be explained by the fact that similar classes share constructive components. By this meaning, we can generally deconstruct each class into a finite set of parts

89 and the classes with the biggest intersection are proportionally closer than the
 90 ones which are not.

91 Part 2

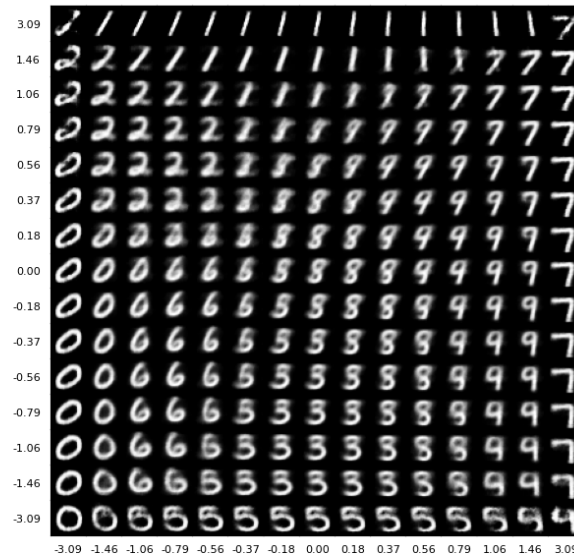


Figure 5: Visualization of latent space in terms of most likely image generations

92 A Week 1 Programming Exercise

```

1
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from scipy.stats import multivariate_normal
5
6 # Enforce reproducibility
7 np.random.seed(42)
8
9 # Parameters
10 theta_true = np.array([-1, 1])
11 sigma_y = 0.1 # Noise variance
12 prior_mean = np.zeros(2)

```

```

13 prior_cov = np.eye(2)
14 n_samples = 20
15
16 # Generate data
17 x = np.random.multivariate_normal(np.zeros(2), np.eye(2), n_samples)
18 y = x @ theta_true + np.random.normal(0, np.sqrt(sigma_y), n_samples)
19
20 # Plot data
21 plt.scatter(x[:, 0], x[:, 1], c=y, cmap='viridis', s=50, label='Data Points')
22 plt.colorbar(label='Labels (y)')
23 plt.title('Dataset for Regression')
24 plt.xlabel('x1')
25 plt.ylabel('x2')
26 plt.legend()
27 plt.show()
28
29
30 posterior_cov = np.linalg.inv(np.linalg.inv(prior_cov) + (1 / sigma_y) * (x.T @ x))
31 posterior_mean = posterior_cov @ (x.T @ y) * (1 / sigma_y)
32
33 # Grid for plotting posterior in range [-3, 3]^2
34 theta_range = np.linspace(-3, 3, 100)
35 theta_grid = np.array(np.meshgrid(theta_range, theta_range)).reshape(2, -1).T
36
37 posterior_pdf = multivariate_normal(mean=posterior_mean, cov=posterior_cov).pdf(theta_grid)
38 posterior_pdf = posterior_pdf.reshape(100, 100)
39
40 # Plot posterior
41 plt.contourf(theta_range, theta_range, posterior_pdf, levels=50, cmap='viridis')
42 plt.colorbar(label='Posterior Density')
43 plt.title('Posterior Distribution of Theta')
44 plt.xlabel('Theta1')
45 plt.ylabel('Theta2')
46 plt.show()
47
48 # Grid for x in range [-3, 3]^2
49 x_grid = np.linspace(-3, 3, 50)
50 x_test = np.array(np.meshgrid(x_grid, x_grid)).reshape(2, -1).T
51
52 predictive_mean = x_test @ posterior_mean
53 predictive_variance = np.sum((x_test @ posterior_cov) * x_test, axis=1) + sigma_y
54
55 # Plot predictive variance
56 plt.contourf(x_grid, x_grid, predictive_variance.reshape(50, 50), levels=50, cmap='viridis')
57 plt.colorbar(label='Predictive Variance')
58 plt.title('Posterior Predictive Variance')

```

```

59 plt.xlabel('x1')
60 plt.ylabel('x2')
61 plt.show()
62
63
64 sigma_x = np.array([[0.1, 0], [0, 1]])
65 x_new = np.random.multivariate_normal(np.zeros(2), sigma_x, n_samples)
66 y_new = x_new @ theta_true + np.random.normal(0, np.sqrt(sigma_y), n_samples)
67
68 # Plot data
69 plt.scatter(x_new[:, 0], x_new[:, 1], c=y_new, cmap='viridis', s=50, label='Data Points')
70 plt.colorbar(label='Labels (y)')
71 plt.title('Dataset for Regression')
72 plt.xlabel('x1')
73 plt.ylabel('x2')
74 plt.legend()
75 plt.show()
76
77
78 xT_x_new = x_new.T @ x_new
79 posterior_cov_new = np.linalg.inv(np.linalg.inv(prior_cov) + (1 / sigma_y) * xT_x_new)
80 posterior_mean_new = posterior_cov_new @ (x_new.T @ y_new) * (1 / sigma_y)
81
82 # Compute density of posterior
83 posterior_pdf_new = multivariate_normal(mean=posterior_mean_new, cov=posterior_cov_new).pdf
84 posterior_pdf_new = posterior_pdf_new.reshape(100, 100)
85
86 predictive_variance_new = np.sum((x_test @ posterior_cov_new) * x_test, axis=1) + sigma_y
87
88 fig, axes = plt.subplots(1, 2, figsize=(12, 6))
89
90 # Plot predictive variance
91 ax1 = axes[0]
92 contour1 = ax1.contourf(theta_range, theta_range, posterior_pdf_new, levels=50, cmap='viridis')
93 fig.colorbar(contour1, ax=ax1, label='Posterior Density')
94 ax1.set_title('Posterior Distribution of Theta')
95 ax1.set_xlabel('Theta1')
96 ax1.set_ylabel('Theta2')
97
98
99 # Plot new predictive variance
100 ax2 = axes[1]
101 contour2 = ax2.contourf(x_grid, x_grid, predictive_variance_new.reshape(50, 50), levels=50,
102 fig.colorbar(contour2, ax=ax2, label='Predictive Variance')
103 ax2.set_title('Posterior Predictive Variance')
104 ax2.set_xlabel('x1')

```

```

105 ax2.set_ylabel('x2')
106
107 # Adjust layout and show the plot
108 plt.tight_layout()
109 plt.show()
110

```

93 B Week 2 Programming exercise

```

1
2 import math
3 import torch
4 import torch.nn as nn
5 import numpy
6 import matplotlib.pyplot as plt
7 from torchvision import datasets, transforms
8 from tqdm import tqdm
9
10 train_loader = torch.utils.data.DataLoader(
11     datasets.MNIST('../data', train=True, download=True,
12                    transform=transforms.Compose([
13                        transforms.ToTensor(),
14                        transforms.Normalize((0.1307,), (0.3081,)),
15                        lambda x: x>0,
16                        lambda x: x.float(),
17                    ])),
18     batch_size=64, shuffle=True)
19 test_loader = torch.utils.data.DataLoader(
20     datasets.MNIST('../data', train=False, transform=transforms.Compose([
21         transforms.ToTensor(),
22         transforms.Normalize((0.1307,), (0.3081,)),
23         lambda x: x>0,
24         lambda x: x.float(),
25     ])),
26     batch_size=64, shuffle=True)
27
28
29 class VAE(nn.Module):
30     def __init__(self, latent_dim = 2):
31         super().__init__()
32
33         self.latent_dim = latent_dim
34
35
36     #The encoder serves as  $q(z/x)$ 

```

```

37         self.encoder = nn.Sequential(nn.Conv2d(in_channels=1, out_channels=32, kernel_size=3,
38                                             nn.BatchNorm2d(32),
39                                             nn.ReLU(),
40
41                                             nn.Conv2d(in_channels=32, out_channels=64, kernel_size=3,
42                                             nn.BatchNorm2d(64),
43                                             nn.ReLU(),
44
45
46                                             #flatten before applying reparametrization trick
47                                             nn.Flatten(start_dim=1) #bx64x7x7 -> bx128*7*7
48                                         )
49
50         self.mu = nn.Linear(64*7*7, latent_dim)
51         self.log_var = nn.Linear(64*7*7, latent_dim)
52
53         self.project_up = nn.Sequential(nn.Linear(2, 64*7*7),
54                                         nn.BatchNorm1d(64*7*7),
55                                         nn.ReLU()
56                                         )
57         self.decoder = nn.Sequential(nn.ConvTranspose2d(in_channels=64, out_channels=32, kernel_size=3,
58                                                         nn.BatchNorm2d(32),
59                                                         nn.ReLU(),
60
61                                                         nn.ConvTranspose2d(in_channels=32, out_channels=1, kernel_size=3,
62                                                         nn.Sigmoid()
63                                                         )
64
65
66     def encode(self, x):
67         h = self.encoder(x) #compute hidden state
68         mu = self.mu(h) #feed h into variable mu
69         log_var = self.log_var(h) #feed h into variable log_var
70         return mu, log_var
71
72     def reparameterization(self, mu, log_var):
73         epsilon = torch.randn_like(mu) #creating same shape as mu but drawn from unit gaussian
74         sigma = torch.exp(0.5*log_var)
75         z = mu+sigma*epsilon
76         return z
77
78     def decode(self, z):
79         z_projected = self.project_up(z).reshape(-1,64,7,7)
80         x_reconstructed = self.decoder(z_projected)
81         return x_reconstructed
82

```

```

83
84     def forward(self,x):
85
86         mu, log_var = self.encode(x)
87
88         z = self.reparameterization(mu, log_var)
89
90
91         x_reconstructed = self.decode(z)
92
93
94         return x_reconstructed, mu, log_var
95
96     def kl_divergence(mu, log_var):
97         pre_batch_average = 0.5*((torch.exp(log_var)+mu**2-1-log_var).sum(dim=1)) #summing over
98         return pre_batch_average.mean()
99
100
101     device = torch.device("cuda")
102     def train_epoch(model, device, optimizer, loss_function, lr, dataloader):
103         acc_loss = 0
104         total_samples = 1
105         model.train()
106
107
108         for (image_batch,_) in dataloader:
109             image_batch = image_batch.to(device)
110
111             prediction_batch, mu, log_var = model(image_batch)
112             reconstruction_loss = loss_function(prediction_batch, image_batch)
113             min_elbo_loss = reconstruction_loss + 0.001*kl_divergence(mu, log_var)
114
115             # print("reconstruction_loss",reconstruction_loss)
116             # print("kl_div",kl_divergence(mu, log_var), 0.001*kl_divergence(mu, log_var))
117             # print("elbo", min_elbo_loss,"\n")
118
119             acc_loss +=min_elbo_loss.item()*len(image_batch)
120             total_samples += len(image_batch)
121
122             optimizer.zero_grad()
123             min_elbo_loss.backward()
124             optimizer.step()
125
126
127         return acc_loss/total_samples
128

```

```

129 def val_epoch(model, device, loss_function, lr, dataloader):
130     acc_loss = 0
131     total_samples = 1
132     model.eval()
133     with torch.no_grad():
134         for (image_batch, _) in dataloader:
135             image_batch = image_batch.to(device)
136
137             prediction_batch, mu, log_var = model(image_batch)
138
139             reconstruction_loss = loss_function(prediction_batch, image_batch)
140             min_elbo_loss = reconstruction_loss + kl_divergence(mu, log_var)
141             acc_loss += min_elbo_loss.item() * len(image_batch)
142             total_samples += len(image_batch)
143
144     return acc_loss / total_samples
145 def train_model(model, device, loss_function, lr, epochs, train_dataloader, val_dataloader):
146
147     optimizer = torch.optim.AdamW(model.parameters(), lr = lr)
148     loss_list = []
149     for e in tqdm(range(epochs)):
150
151         train_epoch_loss = train_epoch(model, device, optimizer, loss_function, lr, train_dataloader)
152         val_epoch_loss = val_epoch(model, device, loss_function, lr, val_dataloader)
153         print(train_epoch_loss, val_epoch_loss)
154
155         loss_list.append((train_epoch_loss, val_epoch_loss))
156     return model, loss_list
157
158
159
160 model, loss_list = train_model(VAE(2).to(device), device, nn.MSELoss(reduction="mean").to(device), lr, epochs, train_dataloader, val_dataloader)
161
162 def visualize_prediction_distribution(model):
163     model.eval()
164
165     normal = torch.distributions.Normal(loc=0, scale=1)
166
167     z_values = torch.linspace(0.001, 0.999, 15)
168     z1 = normal.icdf(z_values)
169     z2 = normal.icdf(z_values).sort(descending=True)[0]
170     z1_grid, z2_grid = torch.meshgrid(z1, z2)
171     z_pairs = torch.dstack([z1_grid, z2_grid]).transpose(0, 1)
172
173     fig, axes = plt.subplots(15, 15, figsize=(7, 7))
174     plt.subplots_adjust(wspace=0, hspace=0)

```



```

175
176
177     z1 = z1.numpy()
178     z2 = z2.numpy()
179
180
181     for i in range(len(axes)):
182         for j in range(len(axes)):
183
184             axes[i,j].imshow(model.decode(z_pairs[i,j].to(device).unsqueeze(0)).detach().cpu())
185             #ax.axis("off")
186
187             # Customize ticks only for outer subplots
188             if i == len(axes) - 1: # Bottom row for z1 (x-axis)
189                 axes[i,j].set_xticks([14]) # Place a tick in the center (arbitrary)
190                 axes[i,j].set_xticklabels([f"{z1[j]:.2f}"], fontsize=8)
191             else:
192                 axes[i,j].set_xticks([]) # Remove ticks for inner subplots
193
194             if j == 0: # First column for z2 (y-axis)
195                 axes[i,j].set_yticks([14]) # Place a tick in the center (arbitrary)
196                 axes[i,j].set_yticklabels([f"{z2[i]:.2f}"], fontsize=8, rotation=0)
197             else:
198                 axes[i,j].set_yticks([]) # Remove ticks for inner subplots
199
200             # Remove axis frame to match the visual style
201             axes[i,j].tick_params(left=False, bottom=False)
202
203             #plt.tight_layout()
204             plt.savefig("PML_assignment1_plot2_fpk297.png")
205             plt.show()
206
207
208 visualize_prediction_distribution(model)
209
210 normal = torch.distributions.Normal(loc=0, scale=1)
211
212 z_values = torch.linspace(0.001, 0.999, 15)
213 z1 = normal.icdf(z_values)
214 z2 = normal.icdf(z_values).sort(descending=True)[0]
215 z1_grid, z2_grid = torch.meshgrid(z1,z2)
216 z_pairs = torch.dstack([z1_grid,z2_grid]).transpose(0,1)
217 z_pairs.shape
218
219
220 def plot_latent_space(model, dataloader):

```

```

221     z_points = []
222     label_points = []
223     for x, y in dataloader:
224         x, y = x.to(device), y.to(device)
225
226         mu, log_var = model.encode(x)
227         z = model.reparameterization(mu, log_var)
228
229         z_points.append(z)
230         label_points.append(y)
231
232
233     z_points = torch.cat(z_points, dim=0)
234     label_points = torch.cat(label_points, dim=0)
235
236
237     plt.scatter(z_points[:,0].detach().cpu().numpy(), z_points[:,1].detach().cpu().numpy(),
238               plt.xlabel("$z_1$")
239               plt.ylabel("$z_2$")
240               plt.grid()
241               plt.colorbar()
242               plt.savefig("PML_assignment1_plot1_fpk297.png")
243               plt.show()
244 plot_latent_space(model, test_loader)
245

```