

# Sheet 1

---

Group 20

November 30, 2024

## 1 Introduction Probability Theory

**Bishop 2.8** Given random variables  $X$  and  $Y$ , where  $\mathbb{E}_{X|Y}[X]$  is what Bishop refers to as  $\mathbb{E}_X[X|Y]$  and  $\mathbb{V} = \text{var}$ ,

(a)

$$\begin{aligned}\mathbb{E}_Y[\mathbb{E}_{X|Y}[X]] &= \int_y \left( \int_x (x) p(X = x|Y = y) dx \right) p(Y = y) dy \\ &= \int_y \left( \int_x (x) p(X = x, Y = y) dx \right) dy \\ &= \int_y \int_x (x) p(X = x, Y = y) dx dy \\ &= \int_x \int_y (x) p(Y = y, X = x) dy dx \\ &= \int_x x p(X = x) dx \\ &= \mathbb{E}_X[X]\end{aligned}$$

(b) Suppose  $Z = X^2$ , from above we know

$$\mathbb{E}_Y[\mathbb{E}_{X|Y}[X^2]] = \mathbb{E}_Y[\mathbb{E}_{Z|Y}[Z]] = \mathbb{E}_Z[Z] = \mathbb{E}_X[X^2], \text{ so:}$$

$$\begin{aligned}\mathbb{E}_Y[\mathbb{V}_{X|Y}[X]] + \mathbb{V}_Y[\mathbb{E}_{X|Y}[X]] &= \mathbb{E}_Y[\mathbb{E}_{X|Y}[X^2] - \mathbb{E}_{X|Y}[X]^2] + \mathbb{E}_Y[\mathbb{E}_{X|Y}[X]^2] - \mathbb{E}_Y[\mathbb{E}_{X|Y}[X]]^2 \\ &= \mathbb{E}_Y[\mathbb{E}_{X|Y}[X^2]] - \mathbb{E}_Y[\mathbb{E}_{X|Y}[X]^2] + \mathbb{E}_Y[\mathbb{E}_{X|Y}[X]^2] - \mathbb{E}_Y[\mathbb{E}_{X|Y}[X]]^2 \\ &= \mathbb{E}_X[X^2] - \mathbb{E}_X[X]^2 \\ &= \mathbb{V}_X[X]\end{aligned}$$

**Bishop 8.9** Suppose we have a model similar to Figure 1, if we define  $A$  to be node  $x$ ,  $C$  to be the Markov blanket of  $x$  and  $B$  to be all other nodes, we can see that  $A$  is independent of  $B$  when conditioned on  $C$ , as any path from a node in  $B$  to  $A$  must either:

- (a) Be a parent of parent of  $x$ , meaning it is blocked by the head-tail link in the parent of  $x$  ( $C_1$ )
- (b) Be a different child of a parent of  $x$ , meaning it is blocked by the tail-tail link in the parent of  $x$  ( $C_1$ )
- (c) Be a parent of a co-parent of  $x$ , meaning it is blocked by the head-tail link in the co-parent of  $x$  ( $C_2$ )
- (d) Be a child of a co-parent of  $x$  that is not a child of  $x$ , meaning it is blocked by the tail-tail link in the co-parent of  $x$  ( $C_2$ )
- (e) Be a child of a child of  $x$ , meaning it is blocked by the head-tail link in the child of  $x$  ( $C_3$ )

We can see that any other path would be invalid, as a path:

- (x) Directly from a node to  $x$  would imply that this node is a parent, and therefore must be part of  $C$
- (y) Directly from  $x$  to a node would imply that this node is a child, and therefore must be part of  $C$
- (z) Directly from a node to a child of  $x$  would imply that this node is a co-parent, and therefore must be part of  $C$

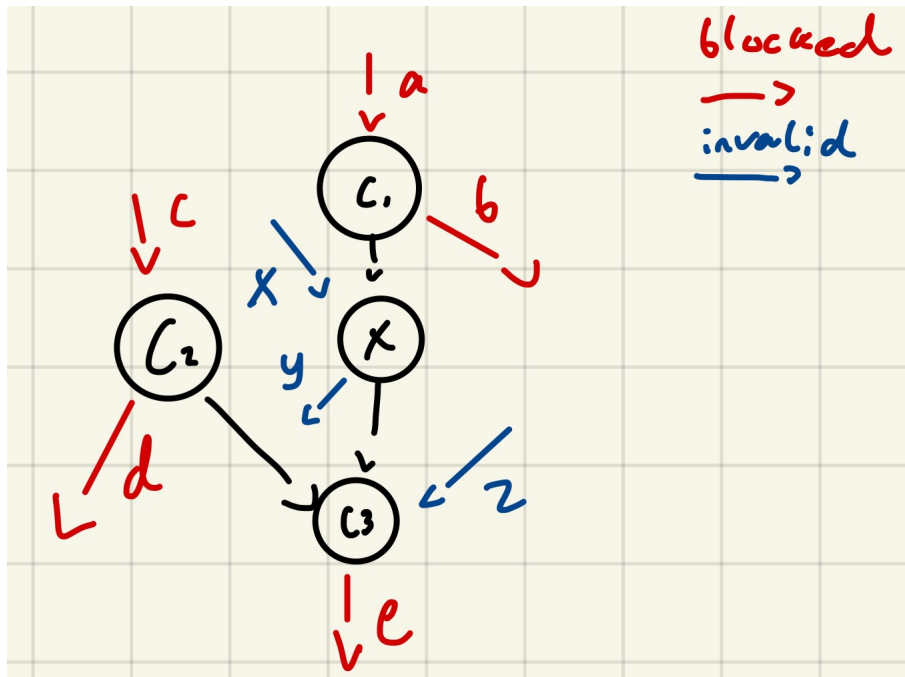


Figure 1: Small illustration of Markov blanket independence

**Bishop 8.11** See equations and Figure 2

$$\begin{aligned}
 p(F = 0|D = 0) &= \frac{p(F = 0, D = 0)}{p(D = 0)} \\
 &= \frac{\sum_b \sum_g p(B = b, F = 0, G = g, D = 0)}{\sum_b \sum_f \sum_g p(B = b, F = f, G = g, D = 0)} \\
 &= \frac{\sum_b \sum_g p(B = b)p(F = 0)p(G = g|B = b, F = 0)p(D = 0|G = g)}{\sum_b \sum_f \sum_g p(B = b)p(F = f)p(G = g|B = b, F = f)p(D = 0|G = g)} \\
 &\approx 0.2125 \\
 p(F = 0|D = 0, B = 0) &= \frac{p(F = 0, D = 0, B = 0)}{p(D = 0, B = 0)} \\
 &= \frac{\sum_g p(B = 0, F = 0, G = g, D = 0)}{\sum_f \sum_g p(B = 0, F = f, G = g, D = 0)} \\
 &= \frac{\sum_g p(B = 0)p(F = 0)p(G = g|B = 0, F = 0)p(D = 0|G = g)}{\sum_f \sum_g p(B = 0)p(F = f)p(G = g|B = 0, F = f)p(D = 0|G = g)} \\
 &\approx 0.1096
 \end{aligned}$$

```

1 prob_b = np.array([0.1, 0.9])
2 prob_f = np.array([0.1, 0.9])
3 prob_g_given_bf = np.array([[[0.9, 0.1], [0.8, 0.2]], [[0.8, 0.2], [0.2, 0.8]]])
4 prob_d_given_g = np.array([[0.9, 0.1], [0.1, 0.9]])
5
6 prob_bfgd = (
7     1 # [B, F, G, D]
8     * prob_b.reshape(2, 1, 1, 1)
9     * prob_f.reshape(1, 2, 1, 1)
10    * prob_g_given_bf.reshape(2, 2, 2, 1)
11    * prob_d_given_g.reshape(1, 1, 2, 2)
12 )
✓ [14] < 10 ms

1 prob_fd = prob_bfgd.sum((0, 2), keepdims=True)
2 prob_f_given_d = prob_fd / prob_fd.sum(1, keepdims=True)
3 print(f"P(F=0 | D=0) = {prob_f_given_d[:, 0, :, 0].flat[0]:.2%}")
✓ [15] < 10 ms

P(F=0 | D=0) = 21.25%

1 prob_bfd = prob_bfgd.sum(2, keepdims=True)
2 prob_f_given_bd = prob_bfd / prob_bfd.sum(1, keepdims=True)
3 print(f"P(F=0 | B=0, D=0) = {prob_f_given_bd[0, 0, :, 0].flat[0]:.2%}")
✓ [16] < 10 ms

P(F=0 | B=0, D=0) = 10.96%

```

Figure 2: Calculation of probability values for 8.11

**W1 Programming** See Figure 3 for the graphs, see Sec. 3 (end of document) for the code.

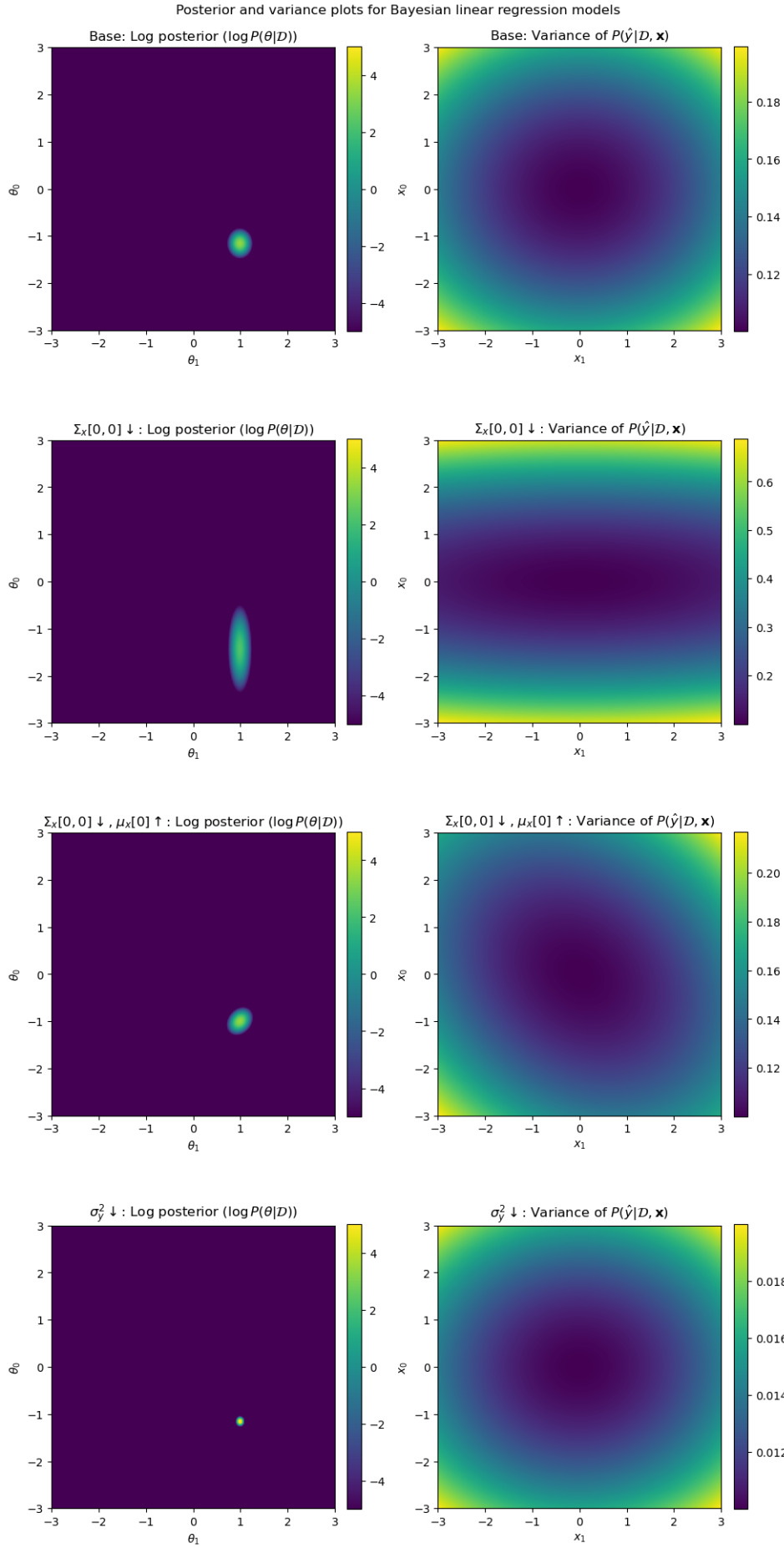


Figure 3: Parameter posterior and posterior-predictive variance plots for Bayesian linear regression using generated datasets. The first row is from the base dataset, the second row has variance 0.1 for  $x_0$ , the third row has variance 0.1 and mean 1 for  $x_0$ , while the final row has variance 0.01 for the target distribution. Note that, while the left column has identical hue scales for the heatmaps, the right column does not.

We can see that, compared to the base distribution, reducing the variance of  $x_0$  makes the variance for  $\theta_0$  increase significantly. This somewhat counter-intuitive result comes from the fact that the distributions are zero-mean, meaning that if a component has low variance, then it will have little impact on the target variable: if all the elements of  $x_0$  are  $\approx 0$ , the value for  $\theta_0$  has a lower impact on  $y$ , making it harder to estimate. If we use the same lowered variance with a nonzero mean, we see that the parameter posterior variance shrinks significantly (as  $x_0$  now always impacts  $y$ ), and that the posterior becomes diagonal (reflecting that, with the non-zero means, it is possible to trade one variable off for another).

For the posterior predictive variance plots, we see that compared to the base distribution, the version with lower  $x_0$  variance shows higher overall variance, as well as a clear increase in variance when  $|x|$  increases. This is because the direction of  $\theta_0$  is uncertain, meaning that having larger values of  $x$  in this direction will make our prediction less certain. Increasing the mean of  $x_0$  again lowers the variance, while also changing the direction of increasing predictive posterior variance to align with the direction of parameter posterior uncertainty.

Finally, decreasing  $\sigma_y^2$  does not change the direction of uncertainty in either of the two plots, but significantly reduces the overall variance: the primary source of uncertainty in the base version was not in the estimation of  $\theta$  but due to the randomness in  $y$ , so reducing this randomness greatly reduces the uncertainty of the overall system.

## 2 Mixture Models and PPCA

**Bishop 9.10** If your component distribution(s) allow for tractable inference of  $p(x_b|x_a, k)$  and  $p(x_a|k')$ , then

$$\begin{aligned}
 p(x_b|x_a) &= \sum_{k=1}^K p(x_b, k|x_a) \\
 &= \sum_{k=1}^K p(x_b|x_a, k)p(k|x_a) \\
 &= \sum_{k=1}^K p(x_b|x_a, k) \frac{p(x_a|k)p(k)}{p(x_a)} \\
 &= \sum_{k=1}^K \pi_k \frac{p(x_a|k)}{p(x_a)} p(x_b|x_a, k) \\
 &= \sum_{k=1}^K \pi_k \frac{p(x_a|k)}{\sum_{k'} \pi_{k'} p(x_a|k')} p(x_b|x_a, k)
 \end{aligned}$$

would give you a mixture distribution with coefficients  $\pi_k^{cond} = \pi_k \frac{p(x_a|k)}{\sum_{k'} \pi_{k'} p(x_a|k')}$  and component densities  $C_{k,x_a}(x_b) = p(x_b|x_a, k)$ .

**Bishop 10.4**  $\mathbb{E}$  is expectation,  $\mathbb{V}$  is variance,  $H$  is entropy,  $\Sigma$  is symmetric positive definite:

$$\begin{aligned}
\text{KL}[q||p] &= - \int_x p(x) \ln \left( \frac{q(x)}{p(x)} \right) dx \\
&= - \int_x p(x) \ln \mathcal{N}(x; \mu, \Sigma) dx - \int_x p(x) \ln p(x) \\
&= - \mathbb{E}_{x \sim p(x)} \left[ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) - \ln \left( \sqrt{\det 2\pi \Sigma} \right) \right] - H_{x \sim p(x)}[x] \\
&= \frac{1}{2} \mathbb{E}_{x \sim p(x)} [(x - \mu)^T \Sigma^{-1} (x - \mu)] + \mathbb{E}_{x \sim p(x)} \left[ \ln \left( \sqrt{\det 2\pi \Sigma} \right) \right] - H_{x \sim p(x)}[x] \\
&= \frac{1}{2} \mathbb{E}_{x \sim p(x)} [(x - \mu)^T \Sigma^{-1} (x - \mu)] + \ln \left( \sqrt{\det 2\pi \Sigma} \right) - H_{x \sim p(x)}[x] \\
\nabla_\mu \text{KL}[q||p] &= \nabla_\mu \left( \frac{1}{2} \mathbb{E}_{x \sim p(x)} [(x - \mu)^T \Sigma^{-1} (x - \mu)] + \ln \left( \sqrt{\det 2\pi \Sigma} \right) - H_{x \sim p(x)}[x] \right) \\
&= \frac{1}{2} \nabla_\mu \mathbb{E}_{x \sim p(x)} [(x - \mu)^T \Sigma^{-1} (x - \mu)] \\
&= \frac{1}{2} \mathbb{E}_{x \sim p(x)} [\nabla_\mu ((x - \mu)^T \Sigma^{-1} (x - \mu))] \\
&= \frac{1}{2} \mathbb{E}_{x \sim p(x)} [-2 \Sigma^{-1} (x - \mu)] \\
&= -\Sigma^{-1} (\mathbb{E}_{x \sim p(x)} [x] - \mu)
\end{aligned} \tag{86}$$

Suppose  $\nabla_\mu \text{KL}[q||p] = 0$ , then

$$0 = -\Sigma^{-1} (\mathbb{E}_{x \sim p(x)} [x] - \mu)$$

$$0 = \mathbb{E}_{x \sim p(x)} [x] - \mu$$

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

$$\begin{aligned}
\nabla_\Sigma \text{KL}[q||p] &= \nabla_\Sigma \left( \frac{1}{2} \mathbb{E}_{x \sim p(x)} [(x - \mu)^T \Sigma^{-1} (x - \mu)] + \ln \left( \sqrt{\det 2\pi \Sigma} \right) - H_{x \sim p(x)}[x] \right) \\
&= \frac{1}{2} \nabla_\Sigma \mathbb{E}_{x \sim p(x)} [(x - \mu)^T \Sigma^{-1} (x - \mu)] + \nabla_\Sigma \ln \left( \sqrt{\det 2\pi \Sigma} \right) \\
&= \frac{1}{2} \mathbb{E}_{x \sim p(x)} [\nabla_\Sigma ((x - \mu)^T \Sigma^{-1} (x - \mu))] + \frac{1}{2} \nabla_\Sigma \ln (\det \Sigma) + \nabla_\Sigma \text{const} \\
&= -\frac{1}{2} \mathbb{E}_{x \sim p(x)} [\Sigma^{-1} (x - \mu) (x - \mu)^T \Sigma^{-1}] + \frac{1}{2} \Sigma^{-1}
\end{aligned} \tag{61, 57}$$

Suppose  $\nabla_\Sigma \text{KL}[q||p] = 0$ , then

$$0 = -\frac{1}{2} \mathbb{E}_{x \sim p(x)} [\Sigma^{-1} (x - \mu) (x - \mu)^T \Sigma^{-1}] + \frac{1}{2} \Sigma^{-1}$$

$$\Sigma^{-1} = \Sigma^{-1} \mathbb{E}_{x \sim p(x)} [(x - \mu) (x - \mu)^T] \Sigma^{-1}$$

$$\Sigma \Sigma^{-1} \Sigma = \Sigma \Sigma^{-1} \mathbb{V}_{x \sim p(x)} [x] \Sigma^{-1} \Sigma$$

$$\Sigma = \mathbb{V}_{x \sim p(x)} [x] = \text{Cov}_{x \sim p(x)} [x, x]$$

So, at the stationary point (the minimum of the KL divergence), the mean is the sample mean, and the covariance is the sample covariance.

**W2 Programming** See Figure 4 for samples from the Gaussian mixture model, and Figure 5 for an approximation of the conditional distribution for the same mixture. See Sec. 4 for the code that generated these figures.

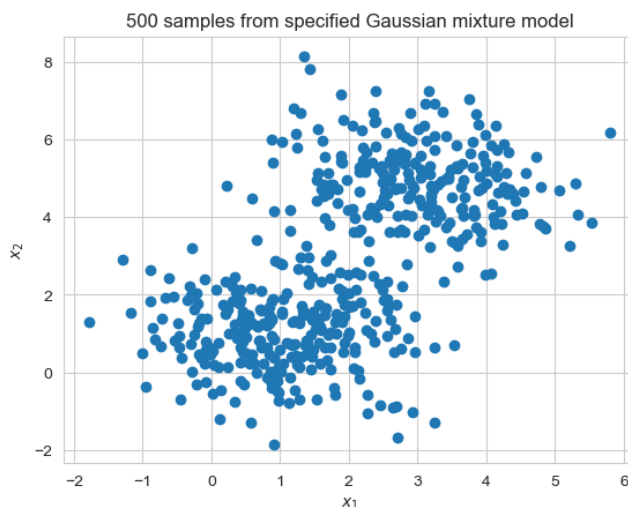


Figure 4: Samples for Gaussian mixture model of two components, with identity covariance and means  $(1, 1)$  and  $(3, 5)$ .

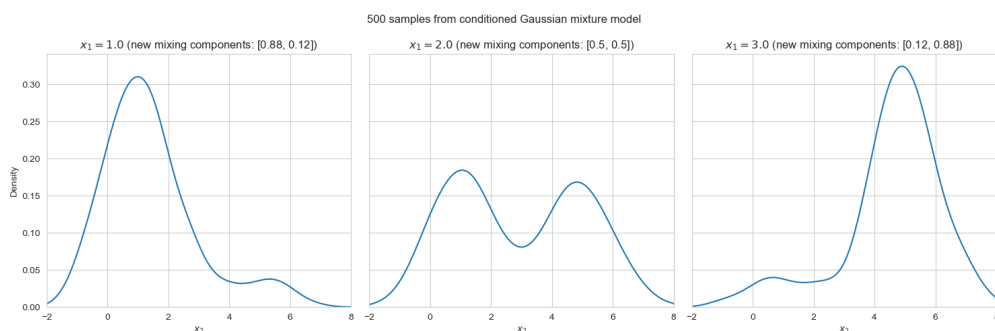


Figure 5: Conditional distributions of  $x_2$  given  $x_1$  from same Gaussian mixture model, for  $x_1 = 1$ ,  $x_1 = 2$  and  $x_1 = 3$  (500 samples, smoothed by `sns.kdeplot`).

See Figure 6 for the distribution of latents for three sizes of VAE: a small one with a single transformation layer of size 10 ( $\text{img} \rightarrow 10 \rightarrow 2 \rightarrow 10 \rightarrow \text{img}$ ), a medium one with transformation layers of sizes 100 and 10 ( $\text{img} \rightarrow 100 \rightarrow 10 \rightarrow 2 \rightarrow 10 \rightarrow 100 \rightarrow \text{img}$ ) and a large one with transformation layers 300, 100 and 10. The small and medium VAEs were trained for 10 epochs, while the large one was trained for 20 epochs on the MNIST training set. Note that increasing the size of the VAE increases the ability of the model to separate the digits in the latent space.

See Figure 7 for a decoding of the latent space for the three VAEs. Note that increasing the size of the model, increases the number of clearly legible digits from 4 in the small model (0, 1, 7, 9) to all 10 in the large model latent space. Again, see Sec. 4 for the code that generated these figures.

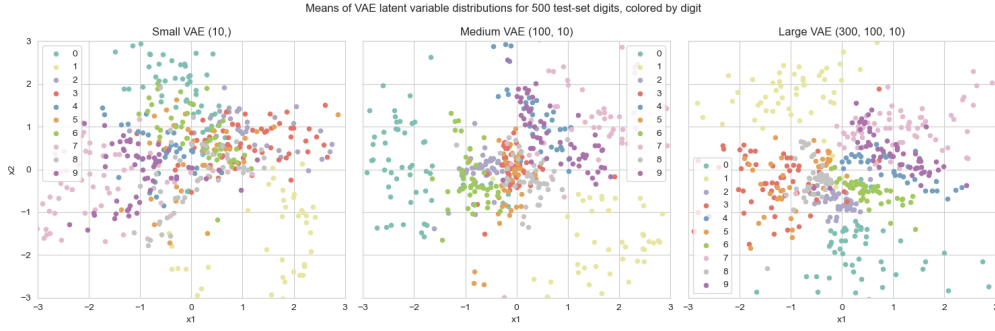


Figure 6: Distribution of latent variables for small-, medium- and large-sized VAEs, with datapoints sourced from test-set and coloured by the true label.

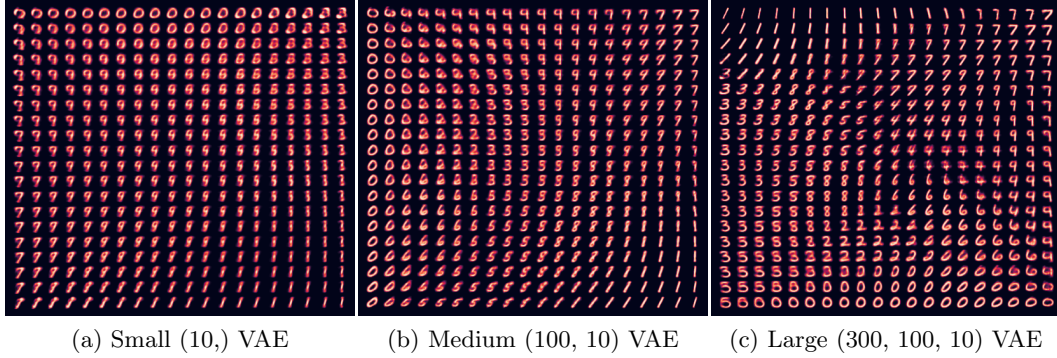


Figure 7: Decoded representations of latent variables, where latents were sampled evenly using the inverse cumulative density function, for three different sizes of VAEs.

### 3 Code W1 Programming

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4
5 # I'm not sure if we're allowed to use library functions for the distributions,
6 # so here are some multivariate distribution functions, re-implemented.
7 @np.vectorize(signature="(d),(d,d)->(r)")
8 def custom_multivariate_normal(mean: np.ndarray, a_transform: np.ndarray) -> np.
    ndarray:
9     n_dim = a_transform.shape[1]
10     e_noise = np.random.normal(0, 1, n_dim)
11     return mean + a_transform @ e_noise
12
13
14 @np.vectorize
15 def custom_single_normal(mean: float, var: float) -> float:
16     return np.random.normal(0, 1) * np.sqrt(var) + mean
17
18
19 def posterior_distribution_params(
20     data: np.ndarray, targets: np.ndarray, target_variance: float
21 ) -> tuple[np.ndarray, np.ndarray]:
22     num, features = data.shape[:2]
23     means = (
24         data.T
25         @ np.linalg.inv(target_variance * np.identity(num) + data @ data.T)
26         @ targets
27     ).flatten()
28     covariance = (
29         np.identity(features)

```



```

30     - data.T
31     @ np.linalg.inv(target_variance * np.identity(num) + data @ data.T)
32     @ data
33 )
34 return means, covariance
35
36
37 @np.vectorize(signature="(),(2),(2,2)->()")
38 def custom_2d_pdf(
39     x0: float, x1: float, means: np.ndarray, covariance: np.ndarray
40 ) -> float:
41     difference = (np.array([x0, x1]) - means).reshape((-1, 1))
42     d_dim = 2
43     power = -0.5 * difference.T @ np.linalg.inv(covariance) @ difference
44     exp_part = float(np.exp(power)[0, 0])
45     normalize_part = (2 * np.pi) ** (d_dim / 2) * np.sqrt(np.linalg.det(covariance))
46     return exp_part / normalize_part
47
48
49 def plot_param_posterior(
50     means: np.ndarray,
51     covariance: np.ndarray,
52     res: int = 200,
53     extra_title: str = "",
54     ax: plt.Axes = None,
55 ):
56     axis = np.linspace(-3, 3, axis=0, num=res).reshape((-1, 1))
57     pdf = custom_2d_pdf(axis, axis.T, means.flatten(), covariance)
58     pdf[pdf == 0] = 1e-10
59     log_pdf = np.log(pdf)
60     if ax is None:
61         _, ax = plt.subplots()
62     img = ax.imshow(log_pdf[::-1, :], vmin=-5, vmax=5)
63     plt.colorbar(img, fraction=0.046, pad=0.04)
64     img.set_extent((-3, 3, -3, 3))
65     ax.autoscale(False)
66     ax.set_ylabel(r"$\theta_0$")
67     ax.set_xlabel(r"$\theta_1$")
68     if extra_title:
69         extra_title += ": "
70     ax.set_title(extra_title + r"Log posterior ($\log P(\theta|\mathcal{D})$)")
71
72
73 @np.vectorize(signature="(),(2),(2,2)->()")
74 def custom_2d_posterior_predictive_variance(
75     x0: float, x1: float, target_var: float, post_covar: np.ndarray
76 ) -> float:
77     x_vec = np.array([x0, x1]).reshape((2, 1))
78     return target_var + x_vec.T @ post_covar @ x_vec
79
80
81 def plot_pp_var(
82     target_var: float,
83     post_covar: np.ndarray,
84     res: int = 200,
85     extra_title: str = "",
86     ax=None,
87 ):
88     axis = np.linspace(-3, 3, axis=0, num=res).reshape((-1, 1))
89     test_pdf = custom_2d_posterior_predictive_variance(
90         axis, axis.T, target_var, post_covar
91     )
92     if ax is None:
93         _, ax = plt.subplots()
94     img = ax.imshow(test_pdf[::-1, :])
95     plt.colorbar(img, fraction=0.046, pad=0.04)
96     img.set_extent((-3, 3, -3, 3))

```

```

97     ax.autoscale(False)
98     ax.set_ylabel(r"$x_0$")
99     ax.set_xlabel(r"$x_1$")
100     if extra_title:
101         extra_title += ": "
102     ax.set_title(extra_title + r"Variance of  $\hat{y}$ ,  $\mathbf{x}$ ")
103
104
105 def plot_experiment(
106     x_mean: np.ndarray | list,
107     x_var_diag: np.ndarray | list,
108     target_var: float,
109     axs=None,
110     seed: int = 42,
111     t: str = "",
112 ):
113     np.random.seed(seed)
114     if axs is None:
115         _, axs = plt.subplots(ncols=2, figsize=(10, 5))
116     x_vals = custom_multivariate_normal(
117         np.repeat(np.atleast_2d(x_mean), 20, axis=0), np.diag(x_var_diag) ** 0.5
118     )
119     theta = np.array([-1, 1]).reshape((2, 1))
120     y_vals = np.vectorize(custom_single_normal)(x_vals @ theta, 0.1)
121     post_mean, post_covar = posterior_distribution_params(x_vals, y_vals, target_var)
122     plot_param_posterior(post_mean, post_covar, ax=axs[0], extra_title=t)
123     plot_pp_var(target_var, post_covar, ax=axs[1], extra_title=t)
124
125
126 _, axss = plt.subplots(nrows=4, ncols=2, figsize=(10, 20))
127 plot_experiment([0, 0], [1, 1], 0.1, axs=axss[0], t="Base")
128 plot_experiment([0, 0], [0.1, 1], 0.1, axs=axss[1], t=r"$\Sigma_x[0,0]\downarrow$")
129 plot_experiment(
130     [1, 0],
131     [0.1, 1],
132     0.1,
133     axs=axss[2],
134     t=r"$\Sigma_x[0,0]\downarrow$, $\mu_x[0]\uparrow$",
135 )
136 plot_experiment([0, 0], [1, 1], 0.01, axs=axss[3], t=r"$\sigma^2_y\downarrow$")
137 plt.suptitle("Posterior and variance plots for Bayesian linear regression models")
138 plt.tight_layout()
139 plt.show()

```

## 4 Code W2 Programming

```

1 from itertools import islice
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import polars as pl
6 import seaborn as sns
7 import torch
8 import torch.distributions as dist
9 import torch.nn as nn
10 import torch.nn.functional as func
11 from torch.utils.data import DataLoader
12 from torchvision import datasets, transforms
13 from tqdm import tqdm, trange
14
15 torch.random.manual_seed(42)
16 mixture = dist.MixtureSameFamily(
17     dist.Categorical(torch.tensor([0.5, 0.5], dtype=torch.float32)),
18     dist.MultivariateNormal(
19         torch.tensor([[1, 1], [3, 5]], dtype=torch.float32),

```

```

20     torch.stack([torch.eye(2), torch.eye(2)]),
21 ),
22 )
23
24 sample = np.asarray(mixture.sample((500,)))
25 plt.scatter(sample[:, 0], sample[:, 1])
26 plt.xlabel("$x_1$")
27 plt.ylabel("$x_2$")
28 plt.title("500 samples from specified Gaussian mixture model")
29 plt.show()
30
31
32 def condition_first(
33     gmm: dist.MixtureSameFamily, cond: torch.Tensor
34 ) -> dist.MixtureSameFamily:
35     (c,) = cond.shape
36     cat: dist.Categorical = gmm.mixture_distribution
37     norm: dist.MultivariateNormal = gmm.component_distribution
38     s_11 = norm.covariance_matrix[:, :c, :c]
39     s_21 = norm.covariance_matrix[:, c:, :c]
40     s_22 = norm.covariance_matrix[:, c:, c:]
41     m_1 = norm.mean[:, :c]
42     m_2 = norm.mean[:, c:]
43     s_21_inv_s11 = s_21 @ torch.linalg.inv(s_11)
44     cond_cat = dist.Categorical(
45         logits=(
46             cat.logits
47             + dist.MultivariateNormal(loc=m_1, covariance_matrix=s_11).log_prob(cond
48         )
49     )
50     cond_norm = dist.MultivariateNormal(
51         loc=m_2 + torch.einsum("knc,kc->kn", s_21_inv_s11, cond - m_1),
52         covariance_matrix=s_22 - torch.einsum("knc,kNc->knN", s_21_inv_s11, s_21),
53     )
54     return dist.MixtureSameFamily(cond_cat, cond_norm)
55
56
57 def plot_conditioned(gmm: dist.MixtureSameFamily, cond: torch.Tensor, ax=None):
58     gmm_cond = condition_first(gmm, cond)
59     samples = np.asarray(gmm_cond.sample((500,)))
60     if ax is None:
61         _, ax = plt.subplots()
62     sns.kdeplot(samples, ax=ax, legend=False)
63     ax.set_xlabel("$x_2$")
64     ax.set_xlim((-2, 8))
65     ax.set_title(
66         f"$x_1$={cond[0]}$ "
67         f"(new mixing components: {[round(x, 2) for x in gmm_cond.
68             mixture_distribution.probs.tolist()]})"
69     )
70
71 torch.random.manual_seed(42)
72 _, axs = plt.subplots(ncols=3, figsize=(15, 5), sharey=True)
73 plot_conditioned(mixture, torch.tensor([1], dtype=torch.float32), ax=axs[0])
74 plot_conditioned(mixture, torch.tensor([2], dtype=torch.float32), ax=axs[1])
75 plot_conditioned(mixture, torch.tensor([3], dtype=torch.float32), ax=axs[2])
76 plt.suptitle("500 samples from conditioned Gaussian mixture model")
77 plt.tight_layout()
78 plt.show()
79
80 train_loader = DataLoader(
81     datasets.MNIST(
82         "../data",
83         train=True,
84         download=True,

```

```

85         transform=transforms.Compose(
86             [
87                 transforms.ToTensor(),
88                 transforms.Normalize((0.1307,), (0.3081,)),
89                 lambda x: x > 0,
90                 lambda x: x.float(),
91             ]
92         ),
93     ),
94     batch_size=50,
95     shuffle=True,
96 )
97 test_loader = DataLoader(
98     datasets.MNIST(
99         "../data",
100         train=False,
101         transform=transforms.Compose(
102             [
103                 transforms.ToTensor(),
104                 transforms.Normalize((0.1307,), (0.3081,)),
105                 lambda x: x > 0,
106                 lambda x: x.float(),
107             ]
108         ),
109     ),
110     batch_size=50,
111     shuffle=True,
112 )
113
114
115 class VAE(nn.Module):
116     def __init__(
117         self,
118         output_dim: int,
119         transform_dims: list[int],
120         latent_dim: int,
121         # An experiment to see if multivariate encodings helped (not really)
122         multivariate: bool = False,
123     ):
124         super().__init__()
125         self.latent_dim = latent_dim
126         trans_enc = []
127         for dim in transform_dims:
128             trans_enc.append(nn.LazyLinear(dim))
129             trans_enc.append(nn.ReLU())
130         self.trans_enc = nn.Sequential(*trans_enc)
131         self.enc_mean = nn.LazyLinear(latent_dim)
132         self.enc_log_var = nn.LazyLinear(
133             latent_dim * (latent_dim + 1) // 2 if multivariate else latent_dim
134         )
135         dec = []
136         for dim in reversed(transform_dims):
137             dec.append(nn.LazyLinear(dim))
138             dec.append(nn.ReLU())
139         dec.append(nn.LazyLinear(output_dim))
140         dec.append(nn.Sigmoid())
141         self.dec = nn.Sequential(*dec)
142         self.multivariate = multivariate
143
144     def encode(self, x: torch.Tensor) -> dist.Distribution:
145         trans = self.trans_enc(x.view(x.shape[0], -1))
146         mean = self.enc_mean(trans)
147         var_vals = torch.exp(self.enc_log_var(trans))
148         if self.multivariate:
149             var = torch.empty((x.shape[0], self.latent_dim, self.latent_dim))
150             idx_u = torch.tril_indices(self.latent_dim, self.latent_dim)
151             var[:, idx_u[0], idx_u[1]] = var_vals

```

```

152         var.mT[:, idx_u[0], idx_u[1]] = var_vals
153         # hacky psd transform, the 0.1*eye providing a margin of error
154         var = torch.einsum("bij,bik->bjk", var, var) + torch.eye(2) * 0.1
155         return dist.MultivariateNormal(mean, var)
156     else:
157         return dist.Normal(mean, var_vals)
158
159     def decode(self, x: torch.Tensor) -> torch.Tensor:
160         return self.dec(x)
161
162     def forward(self, x) -> tuple[torch.Tensor, dist.Distribution]:
163         norm = self.encode(x)
164         return self.decode(norm.rsample()), norm
165
166
167     def loss_fn(
168         images: torch.Tensor,
169         reconstructions: torch.Tensor,
170         distributions: dist.Distribution,
171     ) -> torch.Tensor:
172         recon_loss = func.binary_cross_entropy(
173             reconstructions, images.view_as(reconstructions), reduction="sum"
174         )
175         if isinstance(distributions, dist.Normal):
176             diverg_loss = -0.5 * torch.sum(
177                 1
178                 + torch.log(distributions.variance)
179                 - distributions.mean.pow(2)
180                 - distributions.variance
181             )
182         elif isinstance(distributions, dist.MultivariateNormal):
183             m_diff = 1 - distributions.mean
184             batch_trace, batch_det = torch.vmap(torch.trace), torch.vmap(torch.det)
185             diverg_loss = 0.5 * (
186                 torch.einsum("bd,bd->b", m_diff, m_diff)
187                 + batch_trace(distributions.covariance_matrix)
188                 + torch.log(batch_det(distributions.covariance_matrix))
189                 - distributions.mean.shape[1]
190             ).sum(0)
191         else:
192             raise NotImplementedError
193
194         return recon_loss + diverg_loss
195
196
197 test_model = VAE(784, [100, 10], 2, multivariate=False)
198 test_optim = torch.optim.Adam(test_model.parameters())
199 test_img = next(iter(train_loader))[0]
200 test_recon, test_dist = test_model(test_img)
201 print(test_dist)
202 print(test_recon.shape)
203 loss_fn(test_img, test_recon, test_dist)
204
205 device = "cpu"
206
207
208     def train_epoch(
209         model: VAE, optimizer: torch.optim.Optimizer, epoch: int, show_bar: bool = True
210     ) -> float:
211         model.train()
212         train_loss = 0
213         bar = tqdm(
214             train_loader,
215             total=len(train_loader.dataset) // train_loader.batch_size,
216             desc=f"Epoch {epoch}",
217             leave=False,
218             position=1,

```

```

219         disable=not show_bar,
220     )
221     for batch_idx, (data, _) in enumerate(bar):
222         data = data.to(device)
223         optimizer.zero_grad()
224         reconstructions, distribution = model(data)
225         loss = loss_fn(data, reconstructions, distribution)
226         loss.backward()
227         train_loss += loss.item()
228         optimizer.step()
229         if batch_idx % 100 == 0:
230             bar.set_postfix(loss=loss.item())
231     return train_loss
232
233
234 def train_loop(
235     epochs: int,
236     model: VAE = None,
237     optimizer: torch.optim.Optimizer = None,
238     seed: int = 42,
239     sub_bar: bool = False,
240     transform_dims: list[int] = None,
241     multivariate=False,
242 ) -> VAE:
243     torch.random.manual_seed(seed)
244     if transform_dims is None:
245         transform_dims = [10]
246     if model is None:
247         model = VAE(784, transform_dims, 2, multivariate=multivariate)
248     model = model.to(device)
249     model(next(iter(train_loader))[0].to(device))
250     if device in ("cpu", "cuda"):
251         model.compile()
252     if optimizer is None:
253         optimizer = torch.optim.Adam(model.parameters())
254     bar = trange(epochs, unit="epoch", desc=f"Training {transform_dims}")
255     for epoch in bar:
256         train_loss = train_epoch(model, optimizer, epoch, show_bar=sub_bar)
257         bar.set_postfix(epoch_loss=train_loss)
258     return model
259
260
261 trained_model_s = train_loop(10, transform_dims=[10])
262 trained_model_m = train_loop(10, transform_dims=[100, 10])
263 trained_model_l = train_loop(20, transform_dims=[300, 60, 10])
264
265
266 def plot_model(model: VAE, seed: int = 42, grid_size: int = 10):
267     with torch.no_grad():
268         model.eval()
269         torch.random.manual_seed(seed)
270         space_1d = torch.linspace(0.01, 0.99, steps=grid_size)
271         space_2d = torch.cartesian_prod(space_1d.flip(0), space_1d).flip(1)
272         samples = dist.Normal(0, 1).icdf(space_2d).to(device)
273         images = model.decode(samples).reshape(-1, 28, 28).cpu()
274         _, axss = plt.subplots(
275             grid_size,
276             grid_size,
277             gridspec_kw={
278                 "wspace": -0.8 if grid_size > 7 else -0.1,
279                 "hspace": 0,
280                 "bottom": 0,
281                 "top": 1,
282                 "left": 0,
283                 "right": 1,
284             },
285             facecolor="#020419",

```

```

286         )
287         for ax, img in zip(axss.flat, images):
288             ax.imshow(img)
289             ax.axis("off")
290         plt.show()
291
292
293 plot_model(trained_model_s, grid_size=20)
294 plot_model(trained_model_m, grid_size=20)
295 plot_model(trained_model_l, grid_size=20)
296
297
298 def plot_classes(
299     model: VAE, batches: int = 1, title="", seed: int = 42, ax=None, limits: int = 3
300 ):
301     torch.random.manual_seed(seed)
302     res = []
303     if ax is None:
304         _, ax = plt.subplots()
305     with torch.no_grad():
306         model.eval()
307         for data, labels in islice(train_loader, batches):
308             data = data.to(device)
309             latents = model.encode(data).mean
310             res.append(
311                 pl.DataFrame(
312                     {
313                         "x1": latents[:, 0].numpy(),
314                         "x2": latents[:, 1].numpy(),
315                         "Digit": labels.numpy(),
316                     }
317                 )
318             )
319     for (dig, grp), col in zip(
320         pl.concat(res).sort("Digit").group_by("Digit", maintain_order=True),
321         plt.colormaps["Set3"].colors,
322     ):
323         grp.to_pandas().plot.scatter(
324             x="x1", y="x2", label=dig[0], color=tuple(np.array(col) * 0.9), ax=ax
325         )
326     ax.set_xlim((-limits, limits))
327     ax.set_ylim((-limits, limits))
328     ax.set_title(title)
329
330
331 n_batches = 10
332 _, axs = plt.subplots(ncols=3, tight_layout=True, figsize=(15, 5), sharey=True)
333 plot_classes(trained_model_s, batches=n_batches, title="Small VAE (10,)", ax=axs[0])
334 plot_classes(
335     trained_model_m, batches=n_batches, title="Medium VAE (100, 10)", ax=axs[1]
336 )
337 plot_classes(
338     trained_model_l, batches=n_batches, title="Large VAE (300, 100, 10)", ax=axs[2]
339 )
340 plt.suptitle(
341     "Means of VAE latent variable distributions for 500 test-set digits, colored by digit"
342 )
343 plt.show()
344
345 trained_model_multivar = train_loop(20, transform_dims=[300, 60, 10], multivariate=
    True)
346
347 plot_classes(
348     trained_model_multivar,
349     batches=n_batches,
350     title="Multivariate VAE (300, 100, 10,)",

```

```
351     limits=8,  
352 )  
353 plot_model(trained_model_multivar, grid_size=20)
```