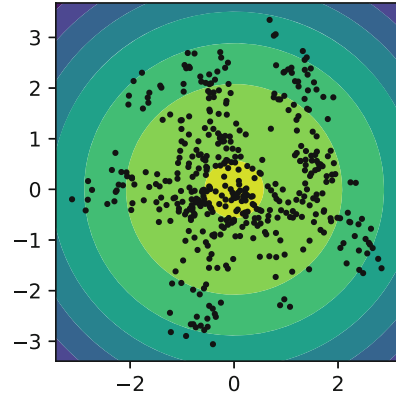


Fig. 4.8 An example of the standard Gaussian prior (contours) and the samples from the aggregated posterior (black dots)



The code for the standard Gaussian prior is presented below:

```

1 class StandardPrior(nn.Module):
2     def __init__(self, L=2):
3         super(StandardPrior, self).__init__()
4
5         self.L = L
6
7         # params weights
8         self.means = torch.zeros(1, L)
9         self.logvars = torch.zeros(1, L)
10
11     def get_params(self):
12         return self.means, self.logvars
13
14     def sample(self, batch_size):
15         return torch.randn(batch_size, self.L)
16
17     def log_prob(self, z):
18         return log_standard_normal(z)

```

Listing 4.6 A standard Gaussian prior class

4.4.1.2 Mixture of Gaussians

If we take a closer look at the aggregated posterior, we immediately notice that it is a mixture model, and a mixture of Gaussians, to be more precise. Therefore, we can use the Mixture of Gaussians (MoG) prior with K components:

$$p_{\lambda}(\mathbf{z}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{z} | \mu_k, \sigma_k^2), \quad (4.47)$$

where $\lambda = \{\{w_k\}, \{\mu_k\}, \{\sigma_k^2\}\}$ are trainable parameters.

Similarly to the standard Gaussian prior, we trained a small VAE with the mixture of Gaussians prior (with $K = 16$) and a two-dimensional latent space. In Fig. 4.9, we present samples from the encoder for the test data (black dots) and the contour plot for the MoG prior. Comparing to the standard Gaussian prior, the MoG prior fits better the aggregated posterior, allowing to *patch* holes.

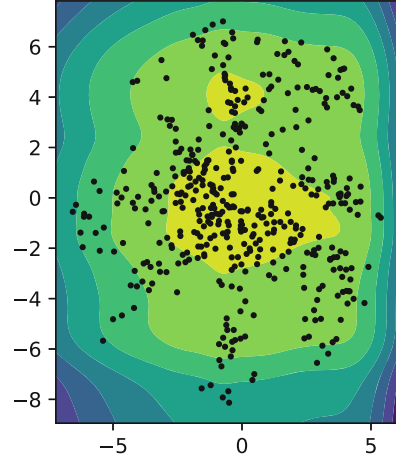
An example of the code is presented below:

```

1 class MoGPrior(nn.Module):
2     def __init__(self, L, num_components):
3         super(MoGPrior, self).__init__()
4
5         self.L = L
6         self.num_components = num_components
7
8         # params
9         self.means = nn.Parameter(torch.randn(num_components,
10 self.L)*multiplier)
11         self.logvars = nn.Parameter(torch.randn(num_components,
12 self.L))
13
14         # mixing weights
15         self.w = nn.Parameter(torch.zeros(num_components, 1, 1))
16
17     def get_params(self):
18         return self.means, self.logvars
19
20     def sample(self, batch_size):
21         # mu, lof_var
22         means, logvars = self.get_params()
23
24         # mixing probabilities
25         w = F.softmax(self.w, dim=0)
26         w = w.squeeze()
27
28         # pick components
29         indexes = torch.multinomial(w, batch_size, replacement=
30 True)
31
32         # means and logvars
33         eps = torch.randn(batch_size, self.L)
34         for i in range(batch_size):
35             indx = indexes[i]
36             if i == 0:
37                 z = means[[indx]] + eps[[i]] * torch.exp(logvars
38 [[indx]])
39             else:
40                 z = torch.cat((z, means[[indx]] + eps[[i]] *
41 torch.exp(logvars[[indx]])), 0)
42         return z
43
44     def log_prob(self, z):
45         # mu, lof_var

```

Fig. 4.9 An example of the MoG prior (contours) and the samples from the aggregated posterior (black dots)



```

41     means, logvars = self.get_params()
42
43     # mixing probabilities
44     w = F.softmax(self.w, dim=0)
45
46     # log-mixture-of-Gaussians
47     z = z.unsqueeze(0) # 1 x B x L
48     means = means.unsqueeze(1) # K x 1 x L
49     logvars = logvars.unsqueeze(1) # K x 1 x L
50
51     log_p = log_normal_diag(z, means, logvars) + torch.log(w)
52     # K x B x L
53     log_prob = torch.logsumexp(log_p, dim=0, keepdim=False) #
54     B x L
55
56     return log_prob

```

Listing 4.7 A Mixture of Gaussians prior class

4.4.1.3 VampPrior: Variational Mixture of Posterior Prior

In [23], it was noticed that we can improve on the MoG prior and approximate the aggregated posterior by introducing *pseudo-inputs*:

$$p_{\lambda}(\mathbf{z}) = \frac{1}{N} \sum_{k=1}^K q_{\phi}(\mathbf{z}|\mathbf{u}_k), \quad (4.48)$$

where $\lambda = \{\phi, \{\mathbf{u}_k^2\}\}$ are trainable parameters and $\mathbf{u}_k \in \mathcal{X}^D$ is a pseudo-input. Notice that ϕ is a part of the trainable parameters. The idea of pseudo-input is to

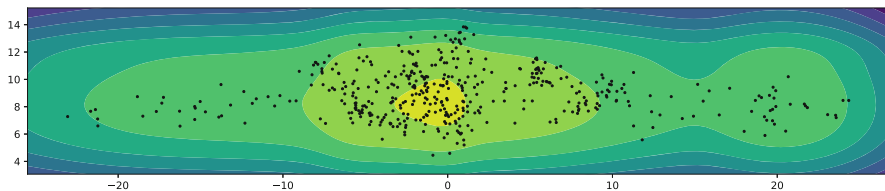


Fig. 4.10 An example of the VampPrior (contours) and the samples from the aggregated posterior (black dots)

randomly initialize objects that mimic observable variables (e.g., images) and learn them by backpropagation.

This approximation to the aggregated posterior is called the **variational mixture of posterior prior**, VampPrior for short. In [23] you can find some interesting properties and further analysis of the VampPrior. The main drawback of the VampPrior lies in initializing the pseudo-inputs; however, it serves as a good proxy to the aggregated posterior that improves the generative quality of the VAE, e.g., [10, 50, 51].

Alemi et al. [10] presented a nice connection of the VampPrior with the information-theoretic perspective on the VAE. They further proposed to introduce learnable probabilities of the components:

$$p_{\lambda}(\mathbf{z}) = \sum_{k=1}^K w_k q_{\phi}(\mathbf{z}|\mathbf{u}_k), \quad (4.49)$$

to allow the VampPrior to select more relevant components (i.e., pseudo-inputs).

As in the previous cases, we train a small VAE with the VampPrior (with $K = 16$) and a two-dimensional latent space. In Fig. 4.10, we present samples from the encoder for the test data (black dots) and the contour plot for the VampPrior. Similar to the MoG prior, the VampPrior fits better the aggregated posterior and has fewer holes. In this case, we can see that the VampPrior allows the encoders to spread across the latent space (notice the values).

An example of an implementation of the VampPrior is presented below:

```

1 class VampPrior(nn.Module):
2     def __init__(self, L, D, num_vals, encoder, num_components,
3         data=None):
4         super(VampPrior, self).__init__()
5
6         self.L = L
7         self.D = D
8         self.num_vals = num_vals
9
10        self.encoder = encoder
11
12        # pseudo-inputs
13        u = torch.rand(num_components, D) * self.num_vals

```

```

13     self.u = nn.Parameter(u)
14
15     # mixing weights
16     self.w = nn.Parameter(torch.zeros(self.u.shape[0], 1, 1))
17     # K x 1 x 1
18
19 def get_params(self):
20     # u->encoder->mu, lof_var
21     mean_vampprior, logvar_vampprior = self.encoder.encode(
22     self.u) #(K x L), (K x L)
23     return mean_vampprior, logvar_vampprior
24
25 def sample(self, batch_size):
26     # u->encoder->mu, lof_var
27     mean_vampprior, logvar_vampprior = self.get_params()
28
29     # mixing probabilities
30     w = F.softmax(self.w, dim=0) # K x 1 x 1
31     w = w.squeeze()
32
33     # pick components
34     indexes = torch.multinomial(w, batch_size, replacement=
35     True)
36
37     # means and logvars
38     eps = torch.randn(batch_size, self.L)
39     for i in range(batch_size):
40         indx = indexes[i]
41         if i == 0:
42             z = mean_vampprior[[indx]] + eps[[i]] * torch.exp
43             (logvar_vampprior[[indx]])
44         else:
45             z = torch.cat((z, mean_vampprior[[indx]] + eps[[i
46             ]] * torch.exp(logvar_vampprior[[indx]])), 0)
47     return z
48
49 def log_prob(self, z):
50     # u->encoder->mu, lof_var
51     mean_vampprior, logvar_vampprior = self.get_params() # (K
52     x L) & (K x L)
53
54     # mixing probabilities
55     w = F.softmax(self.w, dim=0) # K x 1 x 1
56
57     # log-mixture-of-Gaussians
58     z = z.unsqueeze(0) # 1 x B x L
59     mean_vampprior = mean_vampprior.unsqueeze(1) # K x 1 x L
60     logvar_vampprior = logvar_vampprior.unsqueeze(1) # K x 1
61     x L
62
63     log_p = log_normal_diag(z, mean_vampprior,
64     logvar_vampprior) + torch.log(w) # K x B x L
65     log_prob = torch.logsumexp(log_p, dim=0, keepdim=False) #
66     B x L

```