

# VAE with Normalizing Flow

Zhenya Liu

# Overview

---

1. VAE

2. VAE with Normalizing Flow

3. Experiment

# Variational Autoencoder

---

Recall

**Variational Autoencoder(VAE)**(Kingma & Welling, 2014) is a generative model using neural networks to train on ELBO.

# Variational Autoencoder

---

$$\begin{aligned}\log P(\mathbf{X}) &= D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X})) - \mathbb{E}_{\mathbf{Q}}[\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] \\ &= D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X})) + \mathcal{L}(Q)\end{aligned}$$

Here  $\mathcal{L}(Q) = -\mathbb{E}_{\mathbf{Q}}[\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})]$  is known as the **Evidence Lower Bound** (ELBO) of  $\log P(\mathbf{X})$ .

# Variational Autoencoder

---

Construct neural networks  $p_{\theta}(\mathbf{x}|\mathbf{z})$  (*probabilistic decoder*) for likelihood distribution and  $q_{\phi}(\mathbf{z}|\mathbf{x})$  (*probabilistic encoder*) for marginal distribution.

**Objective:** Maximize  $\log p_{\theta}(\mathbf{x})$ , and approximate  $p_{\theta}(\mathbf{z}|\mathbf{x})$  by  $q_{\phi}(\mathbf{z}|\mathbf{x})$ .

Recall:

Maximize  $\mathcal{L}(Q) \implies$  maximize  $\log P(\mathbf{X})$  and minimize  $D_{\text{KL}}(Q(\mathbf{Z})|P(\mathbf{Z}|\mathbf{X}))$ .

so the **ELBO** is a sufficient objective for jointly training  $\theta, \phi$ .

# Variational Autoencoder

---

Marginal log likelihood for individual datapoint

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}),$$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [-\log q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}) + \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = -D_{KL} \left( q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}^{(i)}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}) \right) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z}) \right]$$

The *KL* divergence term on **RHS** can always be analytically calculated, but the expectation term (reconstruction error) can be approximated by Monte Carlo sampling with reparameterization techniques.

# VAE with Normalizing Flow

---

In VAE,  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$  is defined as factorial gaussian distribution. This may not be a good approximation to the true posterier distribution.

By constructing normalizing flows, we can finally output a flexible distribution which could be a better approximation to the posterior. We embed VAE with learnable normalizing flows.

# VAE with Normalizing Flow

---

Recall the change of variables in normalizing flows:

Given initial distribution  $\mathbf{z}_0$ , and a sequence of bijective transformations, then we can get the final distribution  $\mathbf{z}_K$ ,

$$\mathbf{z}_K = f_K \circ \dots \circ f_2 \circ f_1 (\mathbf{z}_0)$$

$$\log q_K (\mathbf{z}_K) = \log q_0 (\mathbf{z}_0) - \sum_{k=1}^K \log \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|$$



# VAE with Normalizing Flow

---

We consider the **planar** flows for all transformation,

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h\left(\mathbf{w}^\top \mathbf{z} + b\right)$$

where  $\lambda = \{\mathbf{w} \in \mathbb{R}^D, \mathbf{u} \in \mathbb{R}^D, b \in \mathbb{R}\}$  are free parameters,  $h(\cdot)$  is a smooth element-wise non-linear function with derivative  $h'$ .

# VAE with Normalizing Flow

---

From planar flows, we can calculate the jacobian determinant in the linear time *w.r.t* dimension of  $\mathbf{z}$ ,

$$\psi(\mathbf{z}) = h' \left( \mathbf{w}^\top \mathbf{z} + b \right) \mathbf{w}$$
$$\left| \det \frac{\partial f}{\partial \mathbf{z}} \right| = \left| \det \left( \mathbf{I} + \mathbf{u} \psi(\mathbf{z})^\top \right) \right| = \left| 1 + \mathbf{u}^\top \psi(\mathbf{z}) \right|$$

The last equality is about calculating the determinant of  $\mathbf{I}$  plus rank 1 matrix, and the result achieved by **Matrix Determinant Lemma**.

# Matrix Determinant Lemma

---

Since

$$\begin{pmatrix} \mathbf{I} & 0 \\ \mathbf{v}^\top & 1 \end{pmatrix} \begin{pmatrix} \mathbf{I} + \mathbf{u}\mathbf{v}^\top & \mathbf{u} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{I} & 0 \\ -\mathbf{v}^\top & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{u} \\ 0 & 1 + \mathbf{v}^\top \mathbf{u} \end{pmatrix}$$

Calculating the determinant of both sides,

$$\det(\mathbf{I} + \mathbf{u}\mathbf{v}^\top) = (1 + \mathbf{v}^\top \mathbf{u})$$

# VAE with Normalizing Flow

---

Then we can conclude our normalizing flow

$$\mathbf{z}_K = f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z})$$
$$\log q_K(\mathbf{z}_K) = \log q_0(\mathbf{z}) - \sum_{k=1}^K \log \left| 1 + \mathbf{u}_k^\top \psi_k(\mathbf{z}_{k-1}) \right|$$

## Remark

---

The planar flow  $f$  may not be invertible, but we have techniques to enforce the invertibility during the implementation.

For instance, when using  $h(x) = \tanh(x)$ , one sufficient condition for  $f(\mathbf{z})$  to be invertible is  $\mathbf{w}^T \mathbf{u} \geq -1$ .

In practice, for arbitrary  $\mathbf{u}$  and  $\mathbf{w}$ , we will calculate a new  $\hat{\mathbf{u}}$ ,  
 $\hat{\mathbf{u}}(\mathbf{w}, \mathbf{u}) = \mathbf{u} + [m(\mathbf{w}^T \mathbf{u}) - (\mathbf{w}^T \mathbf{u})] \frac{\mathbf{w}}{\|\mathbf{w}\|^2}$ , where  $m(x) = -1 + \log(1 + e^x)$ .  
so that  $\mathbf{w}^T \hat{\mathbf{u}} \geq -1$ .

# VAE with Normalizing Flow

---

Back to VAE, we define the encoder  $q_\phi(\mathbf{z} \mid \mathbf{x}) := q_K(\mathbf{z}_K)$ , and get the negative ELBO as

$$\begin{aligned}\mathcal{F}(\mathbf{x}) = -\mathcal{L}(Q) &= \mathbb{E}_{q_\phi(z|x)} [\log q_\phi(\mathbf{z} \mid \mathbf{x}) - \log p(\mathbf{x}, \mathbf{z})] \\ &= \mathbb{E}_{q_0(z_0)} [\log q_K(\mathbf{z}_K) - \log p(\mathbf{x}, \mathbf{z}_K)] \\ &= \mathbb{E}_{q_0(z_0)} [\log q_0(\mathbf{z}_0)] - \mathbb{E}_{q_0(z_0)} [\log p(\mathbf{x}, \mathbf{z}_K)] \\ &\quad - \mathbb{E}_{q_0(z_0)} \left[ \sum_{k=1}^K \log \left| 1 + \mathbf{u}_k^\top \psi_k(\mathbf{z}_{k-1}) \right| \right]\end{aligned}$$

where the initial density  $q_0(\mathbf{z})$  is defined as  $\mathcal{N}(\mu, \sigma)$ , and  $\mu, \sigma$  are trained with all parameters in the encoder and normalizing flows.

# VAE with Normalizing Flow

---

In our implementation, we calculate  $\mathcal{F}(x)$  in the same way as VAE, and we use Monte Carlo sampling to approximate the KL divergence between  $q_K(\mathbf{z}_K)$  and  $p(\mathbf{z})$ .

$$\begin{aligned}\mathcal{F}(\mathbf{x}) &= -\mathcal{L}(Q) = \mathbb{E}_{q_\phi(z|x)} [\log q_\phi(\mathbf{z} | \mathbf{x}) - \log p(\mathbf{x}, \mathbf{z})] \\ &= \mathbb{E}_{q_0(z_0)} [\log q_K(\mathbf{z}_K) - \log p(\mathbf{x}, \mathbf{z}_K)] \\ &= \mathbb{E}_{q_0(z_0)} \left[ \log q_0(\mathbf{z}_0) - \sum_{k=1}^K \log \left| 1 + \mathbf{u}_k^\top \psi_k(\mathbf{z}_{k-1}) \right| - \log p(\mathbf{z}_K) \right] \\ &\quad - \mathbb{E}_{q_0(z_0)} [\log p(\mathbf{x}|\mathbf{z}_K)] \\ &= \text{KL}[q_K(\mathbf{z}_K)||p(\mathbf{z})] - \text{Reconstruction Loss}\end{aligned}$$

# Implementation

---

```
1 class PlanarFlow(nn.Module):
2     def __init__(self, z_dim = None):
3         super(PlanarFlow, self).__init__()
4         self.init_sigma = 0.01
5         self.n_features = z_dim
6         self.weights = nn.Parameter(torch.randn(1, z_dim).normal_(0,
self.init_sigma))
7         self.bias = nn.Parameter(torch.zeros(1).normal_(0, self.
init_sigma))
8         self.u = nn.Parameter(torch.randn(1, z_dim).normal_(0, self.
init_sigma))
```



# Implementation

---

```
1 class PlanarFlow(nn.Module):
2     def forward(self, zk):
3         u = self.u
4         weights = self.weights
5         bias = self.bias
6         ## transforation for invertibility
7         u_temp = (weights @ u.t.squeeze())
8         m_u_temp = -1 + F.softplus(u_temp)
9         uhat = u + (m_u_temp - u_temp) * (weights / (weights @ weights.t
10
11         z_temp = zk @ weights.t() + bias
12         new_z = zk + uhat * torch.tanh(z_temp)
13         ## now compute psi
14         psi = (1 - torch.tanh(z_temp)**2) @ weights
15         det_jac = 1 + psi @ uhat.t
16         logdet_jacobian = torch.log(torch.abs(det_jac) + 1e-8).squeeze()
17
18         return new_z, logdet_jacobian
```

# Implementation

---

```
1 class NormalizingFlows(nn.Module):
2     def __init__(self, n_flows = 1, z_dim = None, flow_type = PlanarFlow
3     ):
4         '''
5         :param z_dims: dimension of the latent variables
6         :param n_flows: how many flows we should use in term of sequence
7         of function f_k (f_k-1(f_k-2(..
8         :param flow_type: we have implemented only the Planar Flow
9         '''
10        super(NormalizingFlows, self).__init__()
11        self.n_flows = n_flows
12        self.flow_type = flow_type
13        self.z_dim = z_dim
14
15        flows_sequence = [self.flow_type(self.z_dim) for _ in range(self
16        .n_flows)]
17
18        self.flows = nn.ModuleList(flows_sequence)
```

# Implementation

---

```
1 class NormalizingFlows(nn.Module):
2
3
4     def forward(self, z):
5
6         logdet_jacobians = []
7         for k, flow in enumerate(self.flows):
8             z, logdet_j = flow(z)
9             logdet_jacobians.append(logdet_j)
10
11         z_k = z
12         logdet_jacobians = torch.stack(logdet_jacobians, dim=1)
13         return z_k, torch.sum(logdet_jacobians, 1)
```

# Implementation

---

```
1 BATCH_SIZE = 64
2 HIDDEN_LAYERS = [200,200]
3 Z_DIM = 40
4 N_FLOWS = 10
5
6 N_EPOCHS = 50
7 LEARNING_RATE = 1e-5
8
9 model = VariationalAutoencoderWithFlows(28*28, HIDDEN_LAYERS, Z_DIM,
    N_FLOWS)
10 optimizer = torch.optim.Adam(model.parameters(), lr=LEARNING_RATE)
```

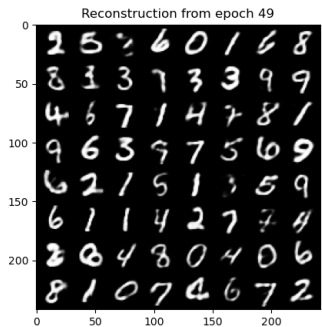
# Implementation

---

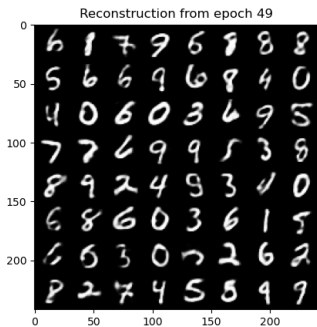
```
1 for epoch in range(N_EPOCHS):
2     for i, data in enumerate(train_loader):
3         images = data[0]
4
5         # loss is KL + BCE(reconstruction)
6         reconstruction = model(images)
7         likelihood = F.binary_cross_entropy(reconstruction, images,
reduction='sum')
8         kl = model.qz - beta * model.pz
9         bound = torch.sum(likelihood) + torch.sum(kl)
10        L = bound / len(images)
11        L.backward()
12        optimizer.step()
13        optimizer.zero_grad()
```

# Experiment

---



(a) MNIST Reconstruction of VAE



(b) MNIST Reconstruction of VAE\_NF

# References

---



Rezende, D. & Mohamed, S. Variational inference with normalizing flows. *International Conference On Machine Learning*. pp. 1530-1538 (2015)