

EE-559 – Deep learning

7.4. Variational autoencoders

François Fleuret

<https://fleuret.org/ee559/>

Feb 9, 2020



Coming back to generating a signal, instead of training an autoencoder and modeling the distribution of Z , we can try an alternative approach:

Impose a distribution for Z and then train a decoder g so that $g(Z)$ matches the training data.

We consider the following two distributions:

- p is the distribution on $\mathcal{X} \times \mathbb{R}^d$ of a pair (X, Z) composed of an encoding state $Z \sim \mathcal{N}(0, I)$ and the output of the decoder g on it.
- q is the distribution on $\mathcal{X} \times \mathbb{R}^d$ of a pair (X, Z) composed of a sample X taken from the data distribution and the output of the encoder on it,

Our goal is that $p(X)$ mimics the data-distribution $q(X)$, that is to find g that maximizes the log-likelihood

$$\frac{1}{N} \sum_n \log p(x_n) = \hat{\mathbb{E}}_{q(X)} [\log p(X)].$$

However, with a complicated g , we can sample z and compute $g(z)$, but cannot compute $p(x)$ for a given x , and even less compute its derivatives.

The **Variational Autoencoder** proposed by Kingma and Welling (2013) relies on a tractable approximation of this log-likelihood.

Note that their framework involves **stochastic** encoder f , and decoder g , whose outputs depend on both their inputs and additional randomness.

Remember that $q(X)$ is the data distribution, and $q(Z | X = x)$ is the distribution of the latent encoding $f(x)$. We want to maximize

$$\mathbb{E}_{q(X)} [\log p(X)],$$

and we can show that

$$-\mathbb{E}_{q(X)} [\log p(X)] \leq \mathbb{E}_{q(X)} [\mathbb{D}_{\text{KL}}(q(Z | X) \| p(Z))] - \mathbb{E}_{q(X, Z)} [\log p(X | Z)].$$

So it makes sense to minimize this latter quantity.

So the final loss is

$$\mathcal{L} = \mathbb{E}_{q(X)} [\mathbb{D}_{\text{KL}}(q(Z | X) \| p(Z))] - \mathbb{E}_{q(X, Z)} [\log p(X | Z)].$$

with

- $q(X)$ is the data distribution
- $p(Z) = \mathcal{N}(0, I)$.

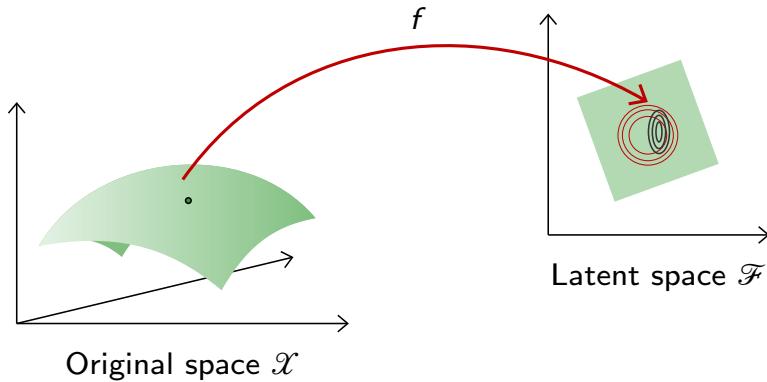
Kingma and Welling propose that both the encoder f and decoder g map to a Gaussian with diagonal covariance. Hence they map to twice the dimension (e.g. $f(x) = (\mu^f(x), \sigma^f(x))$) and

- $q(Z | X = x) \sim \mathcal{N}(\mu^f(x), \text{diag}(\sigma^f(x)))$
- $p(X | Z = z) \sim \mathcal{N}(\mu^g(z), \text{diag}(\sigma^g(z)))$.

The first term of \mathcal{L} is the average of

$$\mathbb{D}_{\text{KL}} \left(\underbrace{q(Z | X = x)}_{\mathcal{N}(\mu^f(x), \sigma^f(x))} \| \underbrace{p(Z)}_{\mathcal{N}(0, I)} \right) = -\frac{1}{2} \sum_d \left(1 + 2 \log \sigma_d^f(x) - (\mu_d^f(x))^2 - (\sigma_d^f(x))^2 \right).$$

over the x_n s.



The first term of \mathcal{L} is the average of

$$\mathbb{D}_{\text{KL}} \left(\underbrace{q(Z | X = x)}_{\mathcal{N}(\mu^f(x), \sigma^f(x))} \| \underbrace{p(Z)}_{\mathcal{N}(0, I)} \right) = -\frac{1}{2} \sum_d \left(1 + 2 \log \sigma_d^f(x) - (\mu_d^f(x))^2 - (\sigma_d^f(x))^2 \right).$$

over the x_n s.

This can be implemented as

```
param_f = model.encode(input)
mu_f, logvar_f = param_f.split(param_f.size(1)//2, 1)

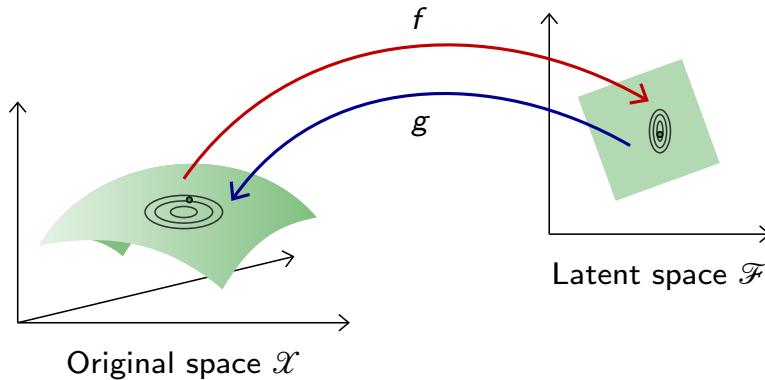
kl = - 0.5 * (1 + logvar_f - mu_f.pow(2) - logvar_f.exp())
kl_loss = kl.sum() / input.size(0)
```

As Kingma and Welling (2013), we use a constant variance of 1 for the decoder, so the second term of \mathcal{L} becomes the average of

$$-\log p(X = x \mid Z = z) = \frac{1}{2} \sum_d (x_d - \mu_d^g(z))^2 + \text{cst}$$

over the x_n , with one z_n sampled for each, i.e.

$$z_n \sim \mathcal{N}(\mu^f(x_n), \sigma^f(x_n)), \quad n = 1, \dots, N.$$



As Kingma and Welling (2013), we use a constant variance of 1 for the decoder, so the second term of \mathcal{L} becomes the average of

$$-\log p(X = x \mid Z = z) = \frac{1}{2} \sum_d (x_d - \mu_d^g(z))^2 + \text{cst}$$

over the x_n , with one z_n sampled for each, i.e.

$$z_n \sim \mathcal{N}(\mu^f(x_n), \sigma^f(x_n)), \quad n = 1, \dots, N.$$

This can be implemented as

```
std_f = torch.exp(0.5 * logvar_f)
z = torch.empty_like(mu_f).normal_() * std_f + mu_f
output = model.decode(z)

fit = 0.5 * (output - input).pow(2)
fit_loss = fit.sum() / input.size(0)
```

We had for the standard autoencoder

```
z = model.encode(input)
output = model.decode(z)
loss = 0.5 * (output - input).pow(2).sum() / input.size(0)
```

and putting everything together we get for the VAE

```
param_f = model.encode(input)
mu_f, logvar_f = param_f.split(param_f.size(1)//2, 1)

kl = - 0.5 * (1 + logvar_f - mu_f.pow(2) - logvar_f.exp())
kl_loss = kl.sum() / input.size(0)

std_f = torch.exp(0.5 * logvar_f)
z = torch.empty_like(mu_f).normal_() * std_f + mu_f
output = model.decode(z)

fit = 0.5 * (output - input).pow(2)
fit_loss = fit.sum() / input.size(0)

loss = kl_loss + fit_loss
```

During inference we do not sample, and instead use μ^f and μ^g as prediction.

Note in particular the **re-parameterization trick**:

```
z = torch.empty_like(mu_f).normal_() * std_f + mu_f
output = model.decode(z)
```

Implementing the sampling of z that way allows to compute the gradient w.r.t f 's parameters without any particular property of `normal_()`.

Original

7 2 1 0 4 1 4 9 5 9 0 6
 9 0 1 5 9 7 8 4 9 6 6 5
 4 0 7 4 0 1 3 1 3 4 7 2

Autoencoder reconstruction ($d = 32$)

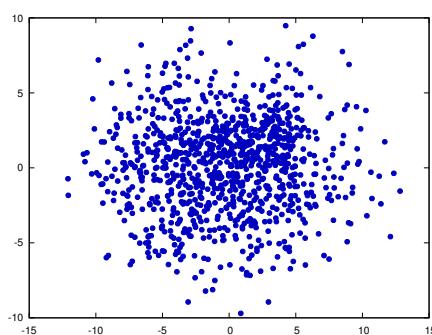
7 2 1 0 4 1 4 9 5 9 0 6
 9 0 1 5 9 7 8 4 9 6 6 5
 4 0 7 4 0 1 3 1 3 4 7 2

Variational Autoencoder reconstruction ($d = 32$)

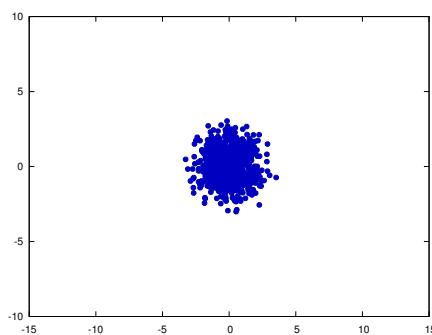
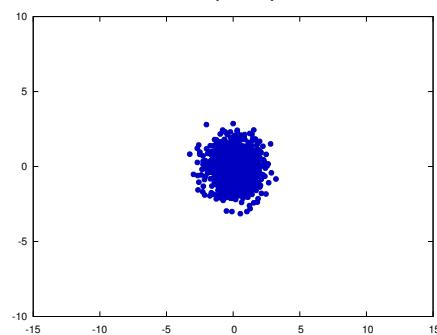
7 2 1 0 4 1 4 9 5 9 0 6
 9 0 1 5 9 7 8 4 9 6 6 5
 4 0 7 4 0 1 3 1 3 4 7 2

We can look at two latent features to check that they are Normal for the VAE.

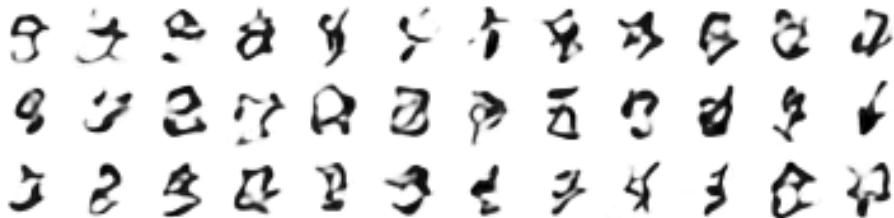
Autoencoder



Variational autoencoder

 $\mathcal{N}(0, 1)$ 

Autoencoder sampling ($d = 32$)



Variational Autoencoder sampling ($d = 32$)

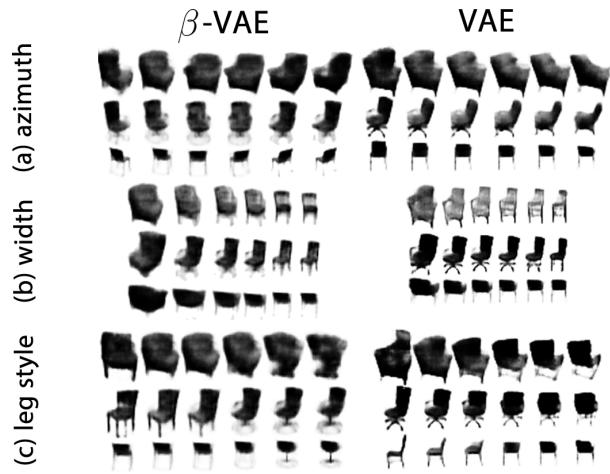


By making the embedding $\sim \mathcal{N}(0, 1)$, the VAE makes its components independent, which often results in “disentangled” representations.

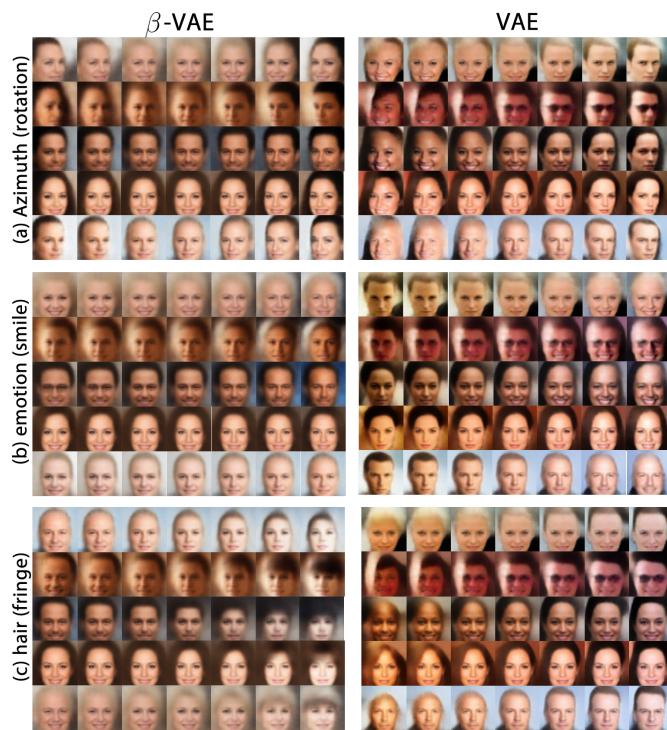
This effect can be reinforced with a greater weight of the KL term

$$\mathcal{L} = \beta \mathbb{E}_{q(X)} \left[\mathbb{D}_{\text{KL}}(q(Z | X) \| p(Z)) \right] - \mathbb{E}_{q(X, Z)} \left[\log p(X | Z) \right],$$

resulting in the β -VAE proposed by Higgins et al. (2017).



(Higgins et al., 2017)



(Higgins et al., 2017)

References

- I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. **beta-vae: Learning basic visual concepts with a constrained variational framework**. In International Conference on Learning Representations (ICLR), 2017.
- D. P. Kingma and M. Welling. **Auto-encoding variational bayes**. CoRR, abs/1312.6114, 2013.