

EE-559 – Deep learning

9. Autoencoders and generative models

François Fleuret

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

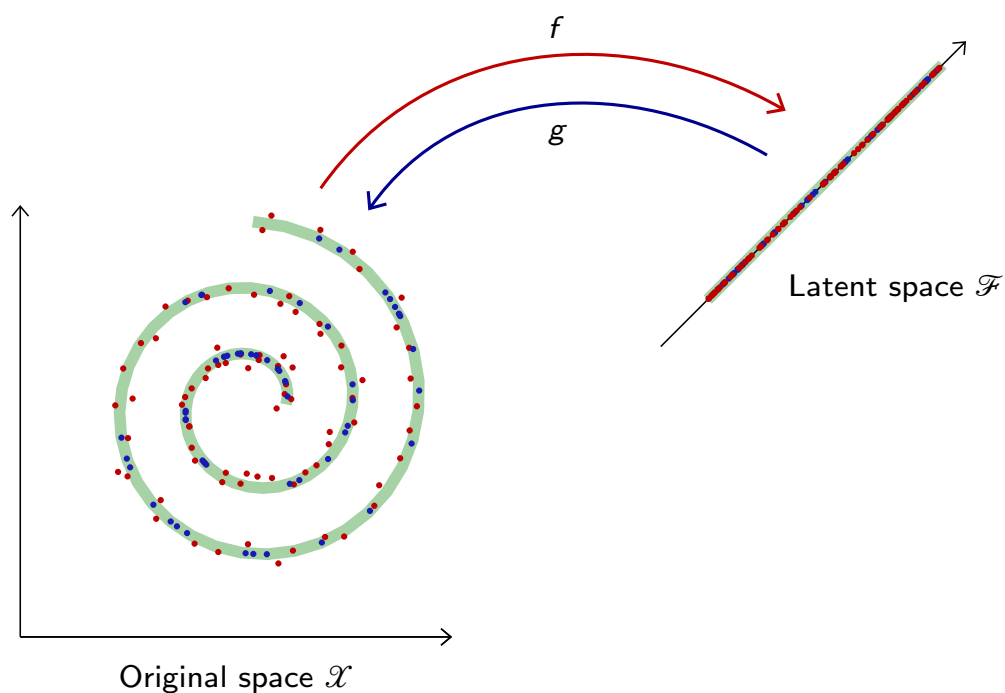
February 16, 2018



Embeddings and generative models

Many applications such as image synthesis, denoising, super-resolution, speech synthesis, compression, etc. require to go beyond classification and regression, and model explicitly a high dimension signal.

This modeling consists of finding “meaningful degrees of freedom” that describe the signal, and are of lesser dimension.

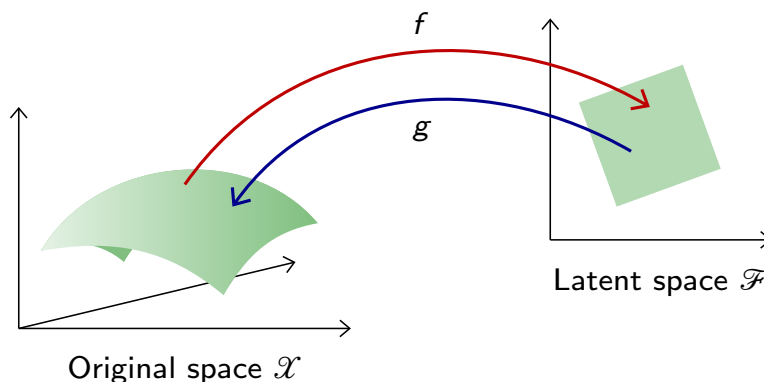


When dealing with real-world signals, this objective involves the same theoretical and practical issues as for classification or regression: defining the right class of high-dimension models, and optimizing them.

Regarding synthesis, we saw that deep feed-forward architectures exhibit good generative properties, which motivates their use explicitly for that purpose.

Autoencoders

An autoencoder combines an **encoder** f that embeds the original space \mathcal{X} into a **latent** space of lower dimension \mathcal{F} , and a **decoder** g to map back to \mathcal{X} , such that their composition $g \circ f$ is [close to] the identity on the data.



A proper autoencoder has to capture a “good” parametrization of the signal, and in particular the statistical dependencies between the signal components.

Let q^X be the data distribution over \mathcal{X} . A good autoencoder could be characterized with the MSE loss

$$\mathbb{E}_{X \sim q^X} [\|X - g \circ f(X)\|^2] \simeq 0.$$

Given two parametrized mappings $f(\cdot; w)$ and $g(\cdot; w)$, training consists of minimizing an empirical estimate of that loss

$$\hat{w}_f, \hat{w}_g = \operatorname{argmin}_{w_f, w_g} \frac{1}{N} \sum_{n=1}^N \|x_n - g(f(x_n; w_f); w_g)\|^2.$$

A simple example of such an autoencoder would be with both f and g linear, in which case the optimal solution is given by PCA. Better results can be achieved with more sophisticated classes of mappings, in particular deep architectures.

Transposed convolutions

Constructing deep generative architectures, such as the decoder of an autoencoder, requires layers to increase the signal dimension, the contrary of what we have done so far with feed-forward networks.

The generative process we saw previously was based on optimizing the input, relying on back-propagation to expend the signal from a low-dimension representation to the high-dimension signal space.

The same can be done in the forward pass with **transposed convolution layers** whose forward operation corresponds to a convolution layer backward pass.

Remember that for a fully connected layer, if x is an (column) input vector, w the weight matrix, and y the (column) output vector, we have

$$y = w x,$$

and with our tensorial notation

$$\left[\frac{\partial \ell}{\partial x} \right] = \begin{pmatrix} \frac{\partial \ell}{\partial x_1} \\ \vdots \\ \frac{\partial \ell}{\partial x_D} \end{pmatrix} = w^T \left[\frac{\partial \ell}{\partial y} \right].$$

The gradient wrt the input is the gradient wrt the output multiplied by the transposed matrix of weights.

Consider now the 1d convolution case, with a kernel κ

$$\begin{aligned} y_i &= (x \circledast \kappa)_i \\ &= \sum_a x_{i+a-1} \kappa_a \\ &= \sum_u x_u \kappa_{u-i+1}, \end{aligned}$$

from which

$$\begin{aligned} \left[\frac{\partial \ell}{\partial x} \right]_u &= \frac{\partial \ell}{\partial x_u} \\ &= \sum_i \frac{\partial \ell}{\partial y_i} \frac{\partial y_i}{\partial x_u} \\ &= \sum_i \frac{\partial \ell}{\partial y_i} \kappa_{u-i+1}. \end{aligned}$$

This final expression looks a lot like a standard convolution, except that the kernel coefficients are visited in reverse order.

If $*$ denotes this operation, we have

$$(x * \kappa)_i = \sum_a x_a \kappa_{i-a+1},$$

which is actually the usual convolution from signal processing.

Coming back to the backward pass of the convolution layer, we have:

$$\begin{aligned} \left[\frac{\partial \ell}{\partial x} \right]_u &= \sum_i \frac{\partial \ell}{\partial y_i} \kappa_{u-i+1} \\ &= \left(\left[\frac{\partial \ell}{\partial y} \right] * \kappa \right)_u. \end{aligned}$$

Or, more concisely, if

$$y = x \circledast \kappa$$

then

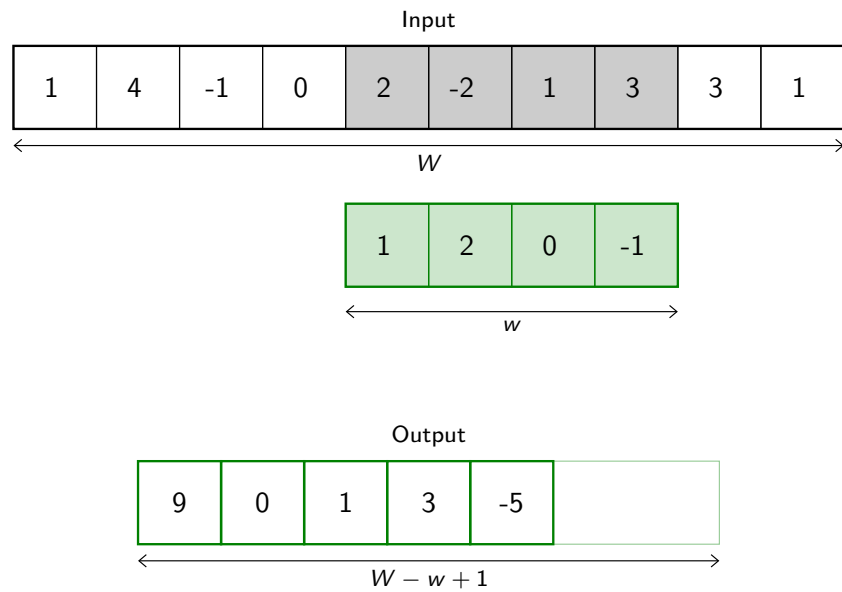
$$\left[\frac{\partial \ell}{\partial x} \right] = \left[\frac{\partial \ell}{\partial y} \right] * \kappa.$$

In the deep-learning field, since it corresponds to transposing the weight matrix of the equivalent fully-connected layer, it is called a **transposed convolution**.

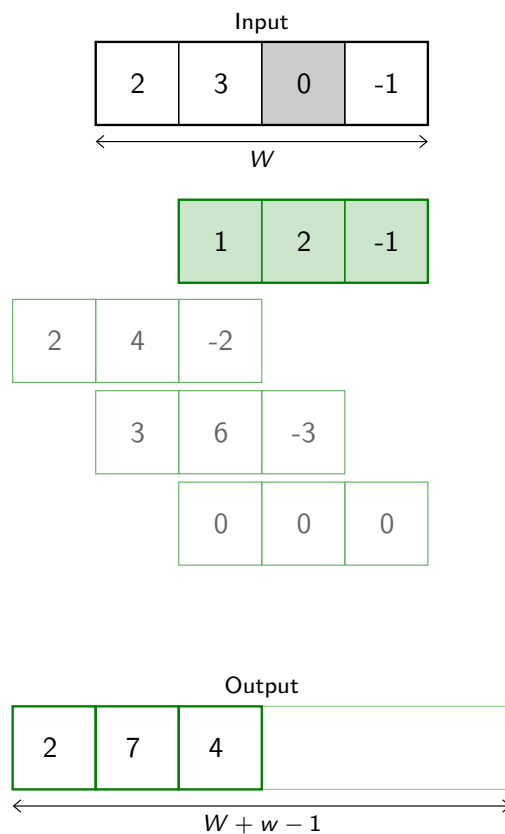
$$\begin{pmatrix} \kappa_1 & \kappa_2 & \kappa_3 & 0 & 0 & 0 & 0 \\ 0 & \kappa_1 & \kappa_2 & \kappa_3 & 0 & 0 & 0 \\ 0 & 0 & \kappa_1 & \kappa_2 & \kappa_3 & 0 & 0 \\ 0 & 0 & 0 & \kappa_1 & \kappa_2 & \kappa_3 & 0 \\ 0 & 0 & 0 & 0 & \kappa_1 & \kappa_2 & \kappa_3 \end{pmatrix}^T = \begin{pmatrix} \kappa_1 & 0 & 0 & 0 & 0 \\ \kappa_2 & \kappa_1 & 0 & 0 & 0 \\ \kappa_3 & \kappa_2 & \kappa_1 & 0 & 0 \\ 0 & \kappa_3 & \kappa_2 & \kappa_1 & 0 \\ 0 & 0 & \kappa_3 & \kappa_2 & \kappa_1 \\ 0 & 0 & 0 & \kappa_3 & \kappa_2 \\ 0 & 0 & 0 & 0 & \kappa_3 \end{pmatrix}$$

While a convolution can be seen as a series of inner products, a transposed convolution can be seen as a weighted sum of translated kernels.

Convolution layer

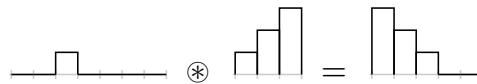


Transposed convolution layer

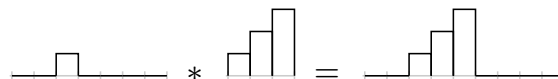


`torch.nn.functional.conv_transpose1d` implements the operation we just described. It takes as input a batch of multi-channel samples, and produces a batch of multi-channel samples.

```
>>> x = Variable(Tensor([[[0, 0, 1, 0, 0, 0, 0]]]))
>>> k = Variable(Tensor([[[1, 2, 3]]]))
>>> F.conv1d(x, k)
Variable containing:
(0 ,.,.) =
  3  2  1  0  0
[torch.FloatTensor of size 1x1x5]
```



```
>>> F.conv_transpose1d(x, k)
Variable containing:
(0 ,.,.) =
  0  0  1  2  3  0  0  0  0
[torch.FloatTensor of size 1x1x9]
```



The class `torch.nn.ConvTranspose1d` implements that operation into a `torch.nn.Module`.

```
>>> x = Variable(Tensor([[[ 2, 3, 0, -1]]]))
>>> m = nn.ConvTranspose1d(1, 1, kernel_size=3)
>>> m.bias.data.zero_()

0
[torch.FloatTensor of size 1]

>>> m.weight.data.copy_(Tensor([ 1, 2, -1 ]))

(0 ,.,.) =
  1  2 -1
[torch.FloatTensor of size 1x1x3]

>>> y = m(x)
>>> y
Variable containing:
(0 ,.,.) =
  2  7  4 -4 -2  1
[torch.FloatTensor of size 1x1x6]
```

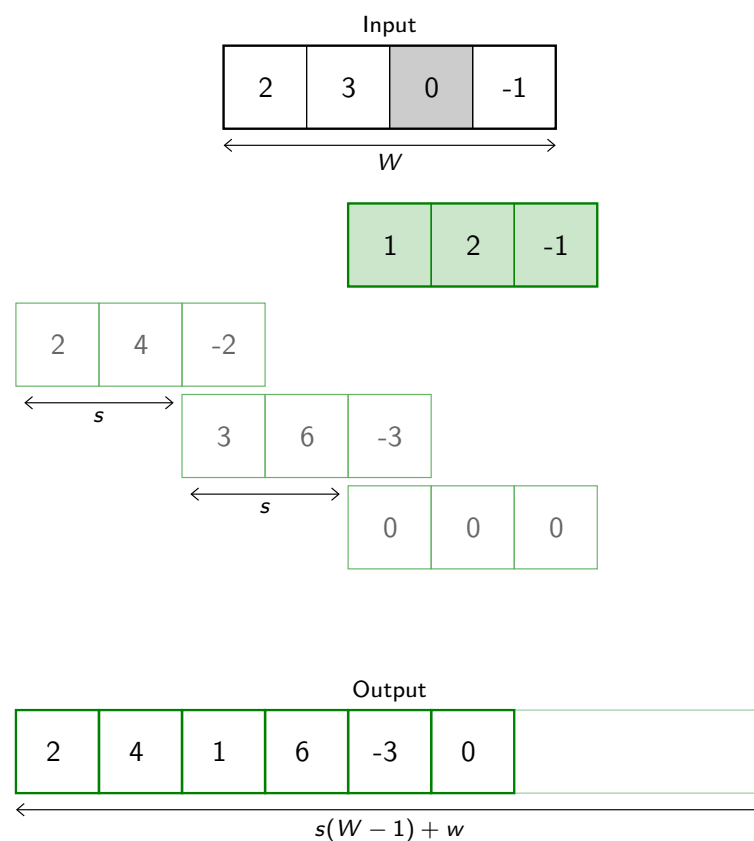
Transposed convolutions also have a **dilation** parameter that behaves as for convolution and expands the kernel size without increasing the number of parameters by making it sparse.

They also have a **stride** and **padding** parameters, however, due to the relation between convolutions and transposed convolutions:



While for convolutions **stride** and **padding** are defined in the input map, for transposed convolutions these parameters are defined in the output map, and the latter modulates a cropping operation.

Transposed convolution layer (stride = 2)



The composition of a convolution and a transposed convolution of same parameters keep the signal size [roughly] unchanged.



A convolution with a stride greater than one may ignore parts of the signal, so its composition with the corresponding transposed convolution may generate a smaller map than the original one.

With an input signal of size W , a kernel of size w , a stride of s , and \div the Euclidean division

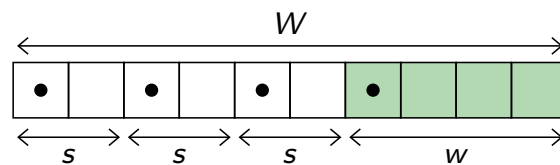
- a convolution layer generates a signal of size $(W - w) \div s + 1$,
- a transposed convolution layer generates a signal of size $(W - 1)s + w$.

If we compose a convolution layer and a transposed convolution layer of same kernel size and stride, to get back the original size, we need

$$((W - w) \div s)s + w = W.$$

So, composing a 1d convolution of kernel size w and stride s and the transposed convolution of same parameters maintains the signal size W , if and only if

$$\exists q \in \mathbb{N}, W = w + s q.$$



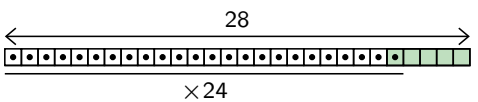
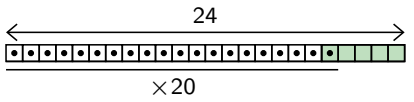
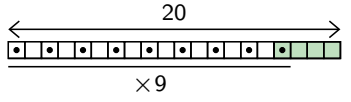
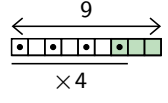
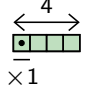
Deep Autoencoders

A deep autoencoder combines an encoder composed of convolutional layers, and a decoder composed of the reciprocal transposed convolution layers.

To run a simple example on MNIST, we consider the following model, where dimension reduction is obtained through filter sizes and strides > 1 , avoiding max-pooling layers.

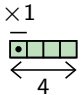
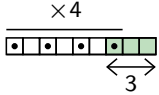
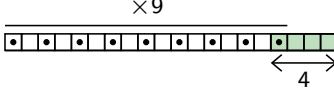
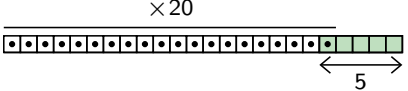
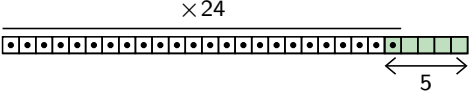
```
AutoEncoder (
  (encoder): Sequential (
    (0): Conv2d(1, 32, kernel_size=(5, 5), stride=(1, 1))
    (1): ReLU (inplace)
    (2): Conv2d(32, 32, kernel_size=(5, 5), stride=(1, 1))
    (3): ReLU (inplace)
    (4): Conv2d(32, 32, kernel_size=(4, 4), stride=(2, 2))
    (5): ReLU (inplace)
    (6): Conv2d(32, 32, kernel_size=(3, 3), stride=(2, 2))
    (7): ReLU (inplace)
    (8): Conv2d(32, 8, kernel_size=(4, 4), stride=(1, 1))
  )
  (decoder): Sequential (
    (0): ConvTranspose2d(8, 32, kernel_size=(4, 4), stride=(1, 1))
    (1): ReLU (inplace)
    (2): ConvTranspose2d(32, 32, kernel_size=(3, 3), stride=(2, 2))
    (3): ReLU (inplace)
    (4): ConvTranspose2d(32, 32, kernel_size=(4, 4), stride=(2, 2))
    (5): ReLU (inplace)
    (6): ConvTranspose2d(32, 32, kernel_size=(5, 5), stride=(1, 1))
    (7): ReLU (inplace)
    (8): ConvTranspose2d(32, 1, kernel_size=(5, 5), stride=(1, 1))
  )
)
```

Encoder

Tensor sizes / operations	
$1 \times 28 \times 28$	
<code>nn.Conv2d(1, 32, kernel_size=5, stride=1)</code>	
$32 \times 24 \times 24$	
<code>nn.Conv2d(32, 32, kernel_size=5, stride=1)</code>	
$32 \times 20 \times 20$	
<code>nn.Conv2d(32, 32, kernel_size=4, stride=2)</code>	
$32 \times 9 \times 9$	
<code>nn.Conv2d(32, 32, kernel_size=3, stride=2)</code>	
$32 \times 4 \times 4$	
<code>nn.Conv2d(32, 8, kernel_size=4, stride=1)</code>	
$8 \times 1 \times 1$	

Decoder

Tensor sizes / operations

	$8 \times 1 \times 1$	
<code>nn.ConvTranspose2d(8, 32, kernel_size=4, stride=1)</code>		
	$32 \times 4 \times 4$	
<code>nn.ConvTranspose2d(32, 32, kernel_size=3, stride=2)</code>		
	$32 \times 9 \times 9$	
<code>nn.ConvTranspose2d(32, 32, kernel_size=4, stride=2)</code>		
	$32 \times 20 \times 20$	
<code>nn.ConvTranspose2d(32, 32, kernel_size=5, stride=1)</code>		
	$32 \times 24 \times 24$	
<code>nn.ConvTranspose2d(32, 1, kernel_size=5, stride=1)</code>		
	$1 \times 28 \times 28$	

Training is achieved with MSE and Adam

```

model = AutoEncoder(embedding_dim, nb_channels)
mse_loss = nn.MSELoss()

if torch.cuda.is_available():
    model.cuda()
    mse_loss.cuda()

optimizer = optim.Adam(model.parameters(), lr = 1e-3)

for epoch in range(args.nb_epochs):
    for input, _ in iter(train_loader):
        if torch.cuda.is_available(): input = input.cuda()
        input = Variable(input)
        output = model(input)
        loss = mse_loss(output, input)
        model.zero_grad()
        loss.backward()
        optimizer.step()

```

X (original samples)

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

$g \circ f(X)$ (CNN, $d = 8$)

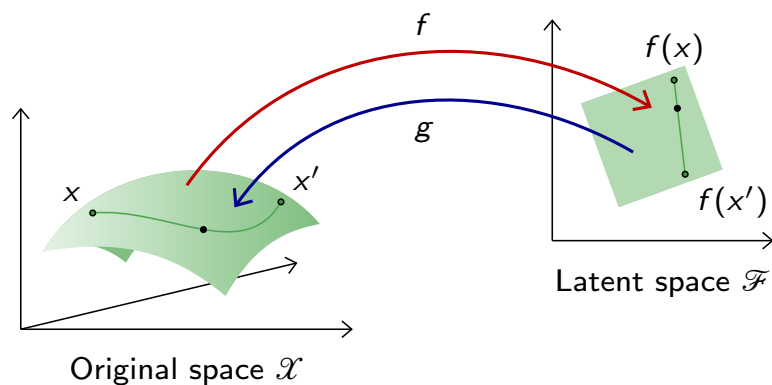
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 5 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

$g \circ f(X)$ (PCA, $d = 8$)

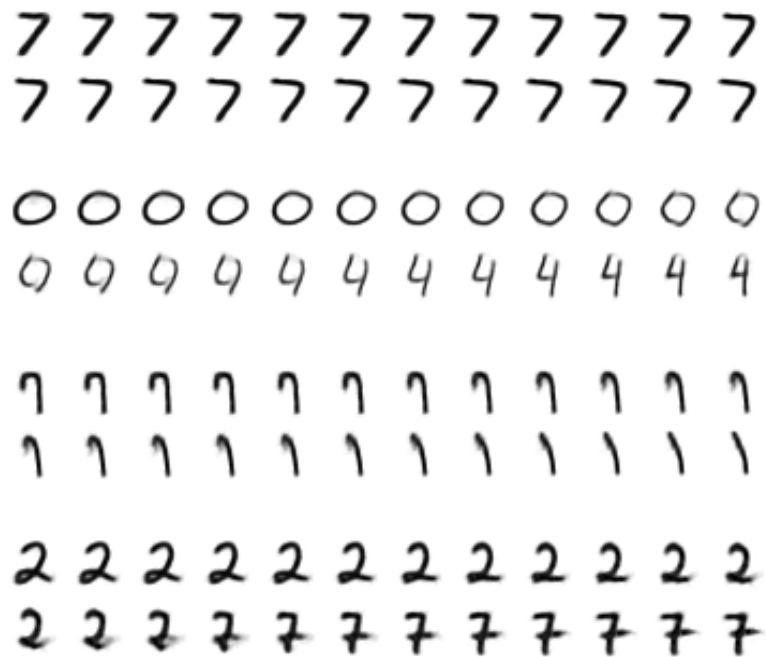
7 2 1 0 4 1 4 9 5 9 0 6
9 0 1 0 9 7 3 4 9 6 6 5
4 0 7 4 0 1 3 1 3 4 7 2

To get an intuition of the latent representation, we can pick two samples x and x' at random and interpolate samples along the line in the latent space

$$\forall x, x' \in \mathcal{X}^2, \alpha \in [0, 1], \xi(x, x', \alpha) = g((1 - \alpha)f(x) + \alpha f(x')).$$



Autoencoder interpolation ($d = 8$)

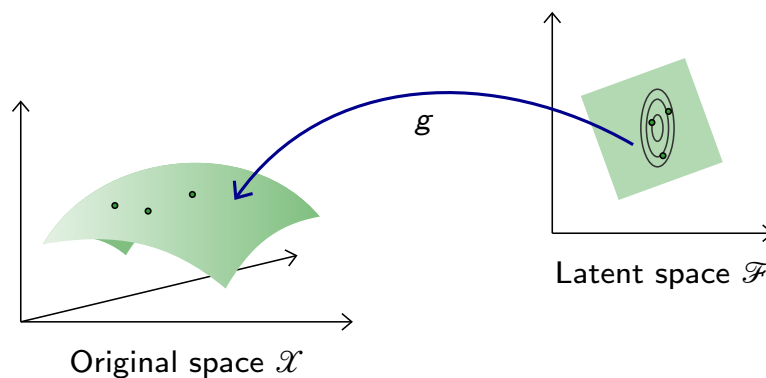


And we can assess the generative capabilities of the decoder g by introducing a [simple] density model q^Z over the latent space \mathcal{Z} , sample there, and map the samples into the image space \mathcal{X} with g .

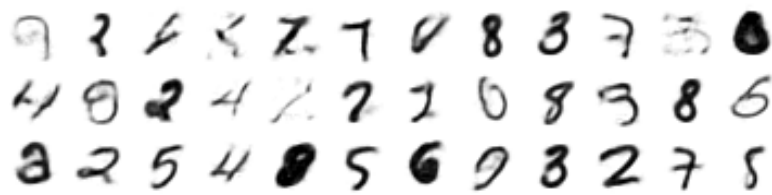
We can for instance use a Gaussian model with diagonal covariance matrix.

$$f(X) \sim \mathcal{N}(\hat{m}, \hat{\Delta})$$

where \hat{m} is a vector and $\hat{\Delta}$ a diagonal matrix, both estimated on training data.



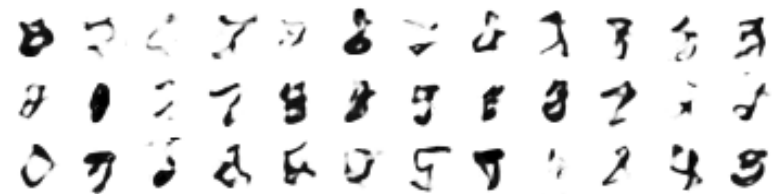
Autoencoder sampling ($d = 8$)



Autoencoder sampling ($d = 16$)



Autoencoder sampling ($d = 32$)



These results are unsatisfying, because the density model used on the latent space \mathcal{Z} is too simple and inadequate.

Building a “good” model amounts to our original problem of modeling an empirical distribution, although it may now be in a lower dimension space.

Denoising Autoencoders

Vincent et al. (2010) interpret the autoencoder in a probabilistic framework as a way of building an encoder that maximizes the mutual information between the input and the latent state.

Let X be a sample, θ the encoder parameter, Z the latent representation of X , and $q_\theta(x, z)$ the distribution of (X, Z) .

We have

$$\begin{aligned}\arg\max_{\theta} \mathbb{I}_{(X,Z) \sim q_\theta}(X, Z) &= \arg\max_{\theta} \underbrace{\mathbb{H}_{(X,Z) \sim q_\theta}(X)}_{\text{cst}} - \mathbb{H}_{(X,Z) \sim q_\theta}(X | Z) \\ &= \arg\max_{\theta} -\mathbb{H}_{(X,Z) \sim q_\theta}(X | Z) \\ &= \arg\max_{\theta} \mathbb{E}_{(X,Z) \sim q_\theta} \left[\log q_\theta(X | Z) \right].\end{aligned}$$

Unfortunately, while computing $q_\theta(Z | X)$ is easy, $q_\theta(X | Z)$ is not.

For any distribution p we have

$$\mathbb{E}_{(X,Z) \sim q_\theta} \left[\log q_\theta(X | Z) \right] \geq \mathbb{E}_{(X,Z) \sim q_\theta} \left[\log p(X | Z) \right],$$

which intuitively means that the amount of information captured by f is greater than the amount captured by f composed with any [noisy] mapping.

So we can in particular design and optimize a parametrized model p_η that we have reasons to believe will [almost] reach the bound.

If we consider the following model for p

$$p_\eta(x | Z = z) = \exp \left(-\frac{\|x - g(z; \eta)\|^2}{2\sigma^2} \right)$$

where g is deterministic, and if we also consider a deterministic f , we get

$$\begin{aligned} \mathbb{E}_{(X,Z) \sim q_\theta} \left[\log p_\eta(X | Z) \right] &= -\mathbb{E}_{(X,Z) \sim q_\theta} \left[\frac{\|X - g(Z; \eta)\|^2}{2\sigma^2} \right] \\ &= -\mathbb{E}_{(X,Z) \sim q_\theta} \left[\frac{\|X - g(f(X; \theta); \eta)\|^2}{2\sigma^2} \right]. \end{aligned}$$

If optimizing η makes the bound tight, the final loss is the reconstruction error

$$\operatorname{argmax}_{\theta} \mathbb{I}(X, Z) \simeq \operatorname{argmin}_{\theta} \left(\min_{\eta} \hat{\mathbb{E}} \left[\|X - g(f(X; \theta); \eta)\|^2 \right] \right).$$

This abstract view of the encoder as “maximizing information” justifies its use to build generic encoding layers.

In the perspective of building a good feature representation, just retaining information is not enough, otherwise the identity would be a good choice.

Reducing dimension, or forcing sparsity is a way to push the model to maximize retained information in a constraint coding space.

In their work, Vincent et al. propose to degrade the signal with noise before feeding it to the encoder, but to keep the MSE to the original signal. This forces the encoder to retain meaningful structures.

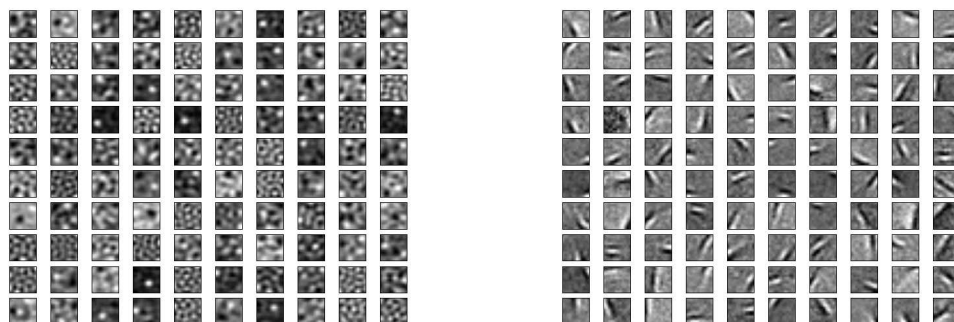


Figure 6: Weight decay vs. Gaussian noise. We show typical filters learnt from natural image patches in the over-complete case (200 hidden units). *Left*: regular autoencoder with weight decay. We tried a wide range of weight-decay values and learning rates: filters never appeared to capture a more interesting structure than what is shown here. Note that some local blob detectors are recovered compared to using no weight decay at all (Figure 5 right). *Right*: a denoising autoencoder with additive Gaussian noise ($\sigma = 0.5$) learns Gabor-like local oriented edge detectors. Clearly the filters learnt are qualitatively very different in the two cases.

(Vincent et al., 2010)

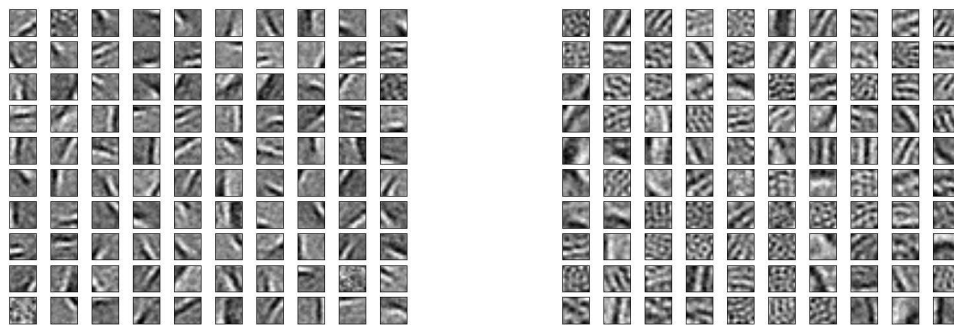


Figure 7: Filters obtained on natural image patches by denoising autoencoders using other noise types. *Left*: with 10% salt-and-pepper noise, we obtain oriented Gabor-like filters. They appear slightly less localized than when using Gaussian noise (contrast with Figure 6 right). *Right*: with 55% zero-masking noise we obtain filters that look like oriented gratings. For the three considered noise types, denoising training appears to learn filters that capture meaningful natural image statistics structure.

(Vincent et al., 2010)

Vincent et al. build deep MLPs whose layers are initialized successively as encoders trained within a noisy autoencoder.

A final classifying layer is added and the full structure can be fine-tuned.

Data Set	SVM_{rbf}	DBN-1	SAE-3	DBN-3	SDAE-3 (v)
<i>MNIST</i>	1.40 ± 0.23	1.21 ± 0.21	1.40 ± 0.23	1.24 ± 0.22	1.28 ± 0.22 (25%)
<i>basic</i>	3.03 ± 0.15	3.94 ± 0.17	3.46 ± 0.16	3.11 ± 0.15	2.84 ± 0.15 (10%)
<i>rot</i>	11.11 ± 0.28	14.69 ± 0.31	10.30 ± 0.27	10.30 ± 0.27	9.53 ± 0.26 (25%)
<i>bg-rand</i>	14.58 ± 0.31	9.80 ± 0.26	11.28 ± 0.28	6.73 ± 0.22	10.30 ± 0.27 (40%)
<i>bg-img</i>	22.61 ± 0.37	16.15 ± 0.32	23.00 ± 0.37	16.31 ± 0.32	16.68 ± 0.33 (25%)
<i>bg-img-rot</i>	55.18 ± 0.44	52.21 ± 0.44	51.93 ± 0.44	47.39 ± 0.44	43.76 ± 0.43 (25%)
<i>rect</i>	2.15 ± 0.13	4.71 ± 0.19	2.41 ± 0.13	2.60 ± 0.14	1.99 ± 0.12 (10%)
<i>rect-img</i>	24.04 ± 0.37	23.69 ± 0.37	24.05 ± 0.37	22.50 ± 0.37	21.59 ± 0.36 (25%)
<i>convex</i>	19.13 ± 0.34	19.92 ± 0.35	18.41 ± 0.34	18.63 ± 0.34	19.06 ± 0.34 (10%)
<i>tzanetakis</i>	14.41 ± 2.18	18.07 ± 1.31	16.15 ± 1.95	18.38 ± 1.64	16.02 ± 1.04 (0.05)

(Vincent et al., 2010)

Variational Autoencoders

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.

If p^X is the distribution of $X = g(Z)$, we should ideally look for the g that maximizes the [empirical] log-likelihood

$$\frac{1}{N} \sum_n \log p^X(x_n).$$

Unfortunately, while we can sample z and compute $g(z)$, we cannot compute $p^X(x)$ 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.

Their framework considers a **stochastic** encoder f , and decoder g , whose outputs depend on their inputs as usual but with some remaining randomness.

We consider the following two distributions:

- 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,
- 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 on it.

We will use notations such as q^X for marginals, and $q^{Z|X=x}$ for conditionals.

As we said, ideally we would maximize

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

We can equivalently maximize

$$\begin{aligned} \mathcal{S} &= \hat{\mathbb{E}}_{X \sim q^X} [\log p^X(X) - \log q^X(X)] \\ &= \hat{\mathbb{E}}_{(X, Z) \sim q} \left[\log p(X, Z) + \log \frac{p^X(X)}{p(X, Z)} - \log q(X, Z) - \log \frac{q^X(X)}{q(X, Z)} \right] \\ &= -\hat{\mathbb{D}}_{KL}(q \| p) + \hat{\mathbb{E}}_{X' \sim q^X} [\hat{\mathbb{D}}_{KL}(q^{Z|X=X'} \| p^{Z|X=X'})]. \end{aligned}$$

While the second term cannot be properly estimated, we can use the first as a lower “variational” bound, and rewrite it

$$\begin{aligned} \mathcal{S} &\geq -\hat{\mathbb{D}}_{KL}(q \| p) \\ &= -\hat{\mathbb{E}}_{X' \sim q^X} [\hat{\mathbb{D}}_{KL}(q^{Z|X=X'} \| p^Z)] + \hat{\mathbb{E}}_{(X, Z') \sim q} [\log p^{X|Z=Z'}(X)]. \end{aligned}$$

This can be seen as optimizing at the same time the g we want, and an f such that $q^{Z|X=x}$ gives better zs than the prior to get x back.

Kingma and Welling use Gaussians with diagonal covariance for $q^{Z|X}$ and $p^{X|Z}$.

So in practice the encoder maps a data point from the signal space \mathbb{R}^c to [the parameters of] a Gaussian in the latent space \mathbb{R}^d

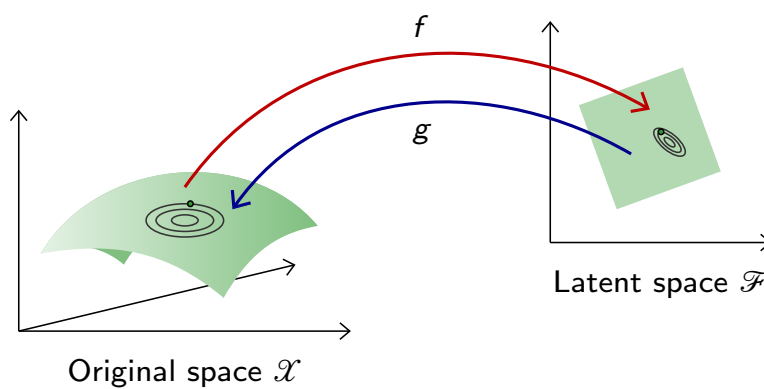
$$f : \mathbb{R}^c \rightarrow \mathbb{R}^{2d}$$

$$x \mapsto (\mu_1^f, \dots, \mu_d^f, \sigma_1^f, \dots, \sigma_d^f),$$

and the decoder maps a latent value from \mathbb{R}^d to [the parameters of] a Gaussian in the signal space \mathbb{R}^c

$$g : \mathbb{R}^d \rightarrow \mathbb{R}^{2c}$$

$$z \mapsto (\mu_1^g, \dots, \mu_c^g, \sigma_1^g, \dots, \sigma_c^g).$$



We have to minimize

$$\mathcal{L} = \hat{\mathbb{E}}_{X' \sim q^X} \left[\hat{\mathbb{D}}_{KL} \left(q^{Z|X=X'} \parallel p^Z \right) \right] - \hat{\mathbb{E}}_{(X, Z') \sim q} \left[\log p^{X|Z=Z'}(X) \right].$$

Since $q^{Z|X}$ and p^Z are Gaussian, we have

$$\hat{\mathbb{D}}_{KL} \left(q^{Z|X=x} \parallel p^Z \right) = \frac{1}{2} \sum_d \left(1 + 2 \log \sigma_d^f(x) - \left(\mu_d^f(x) \right)^2 - \left(\sigma_d^f(x) \right)^2 \right).$$

And with

$$z_l^n \sim \mathcal{N} \left(\mu^f(x_n), \sigma^f(x_n) \right), \quad n = 1, \dots, N, \quad l = 1, \dots, L,$$

we have

$$-\hat{\mathbb{E}}_{(X, Z') \sim q} \left[\log p^{X|Z=Z'}(X) \right] \simeq \frac{1}{2} \sum_{n=1}^N \sum_{l=1}^L \sum_c \frac{(x_{n,d} - \mu_c^g(z_l^n))^2}{2 (\sigma_c^g(z_l^n))^2}$$

Kingma and Welling point out that using $L = 1$ is enough.

```
class VariationalAutoEncoder:
    def __init__(self, dim_x, dim_z, dim_hidden):
        self.dim_x = dim_x
        self.dim_z = dim_z
        self.dim_hidden = dim_hidden

        self.encoder_f = nn.Sequential(
            nn.Linear(self.dim_x, self.dim_hidden),
            nn.ReLU(),
            nn.Linear(self.dim_hidden, self.dim_z * 2)
        )

        self.decoder_g = nn.Sequential(
            nn.Linear(self.dim_z, self.dim_hidden),
            nn.ReLU(),
            nn.Linear(self.dim_hidden, self.dim_x * 2)
        )
```

```

def log_g_x_given_z(self, z, x):
    mu, logvar = self.decoder_g(z).split(self.dim_x, 1)
    u = - (x - mu).pow(2) / (2 * logvar.exp()) - logvar.mul(0.5)
    log_p = u.sum(1) - 0.5 * math.log(2 * math.pi)
    return log_p

def kl_f_given_x_to_g(self, x):
    mu, logvar = self.encoder_f(x).split(self.dim_z, 1)
    u = - 0.5 * (1 + logvar - mu.pow(2) - logvar.exp())
    kl = u.sum(1)
    return kl

def sample_from_f_given_x(self, x):
    mu, logvar = self.encoder_f(x).split(self.dim_z, 1)
    std = logvar.mul(0.5).exp()
    u = Variable(mu.data.new(mu.size()).normal_())
    z = u * std + mu
    return z

```

```

def train(self, x, nb_epochs, bs, lr):
    optimizer = optim.Adam(chain(self.encoder_f.parameters(),
                                  self.decoder_g.parameters()),
                             lr = lr)

    for k in range(0, nb_epochs):
        for xb in x.split(bs):
            loss_kl = self.kl_f_given_x_to_g(xb).mean()

            zb = self.sample_from_f_given_x(xb)
            loss_E = - self.log_g_x_given_z(zb, xb).mean()

            loss = loss_kl + loss_E

            optimizer.zero_grad()
            loss.backward()
            optimizer.step()

```

```

dim, nb = 2, 250

x = toy_moon_set(nb)
x = Variable((x - x.mean(0)) / x.std())

file = open('toy-vae-train.dat', 'w')
for k in range(x.size(0)):
    file.write('{:f} {:f}\n'.format(x.data[k, 0], x.data[k, 1]))
file.close()

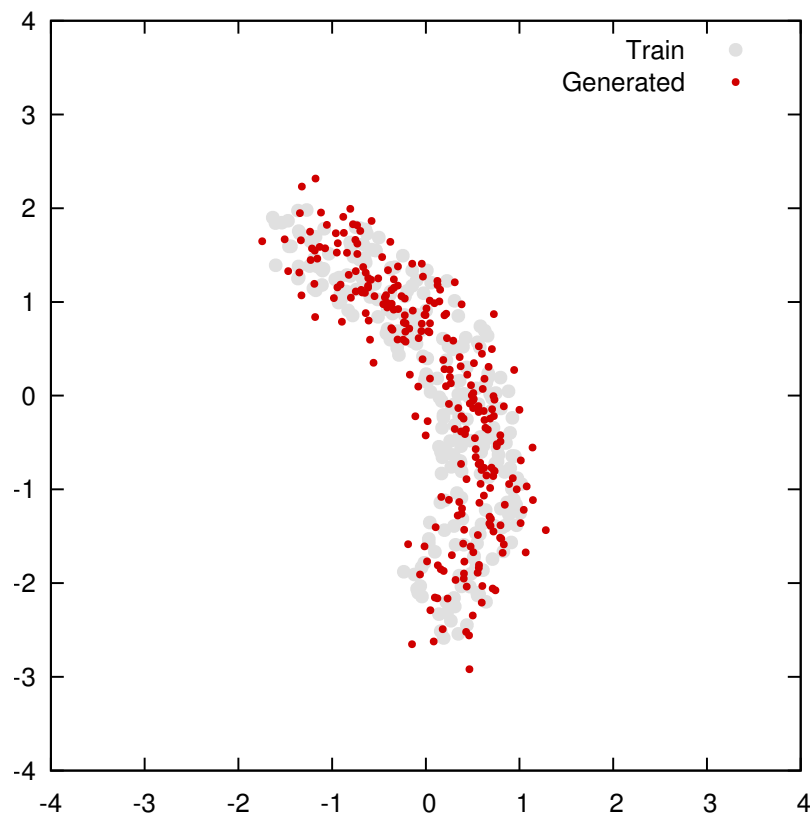
vae = VariationalAutoEncoder(dim_x = dim, dim_z = dim, dim_hidden = 100)

vae.train(x, nb_epochs = 500, bs = 10, lr = 1e-2)

z = Variable(torch.Tensor(nb, vae.dim_z).normal_())
mu, logvar = vae.decoder_g(z).split(vae.dim_x, 1)
std = logvar.mul(0.5).exp()
u = Variable(mu.data.new(mu.size()).normal_())
h = mu + u * std

file = open('toy-vae-synth.dat', 'w')
for k in range(h.size(0)):
    file.write('{:f} {:f}\n'.format(h.data[k, 0], h.data[k, 1]))
file.close()

```



For MNIST, we keep our convolutional structure, but the encoder now maps to twice the number of dimensions, which corresponds to the μ^f s and σ^f s, and we use a fixed variance for the decoder.

We use Adam for training and the loss estimate for the standard autoencoder

```
output = model(input)
loss = mse_loss(output, input)
```

becomes

```
param = model.encode(input)
mu, logvar = param.split(param.size(1)//2, 1)
logvar = logvar + math.log(0.01)
std = logvar.mul(0.5).exp()

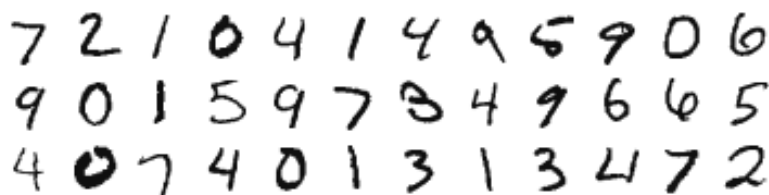
kl = - 0.5 * (1 + logvar - mu.pow(2) - logvar.exp())
kl = kl.mean()

u = Variable(mu.data.new(mu.size()).normal_())
z = u * std + mu
output = model.decode(z)

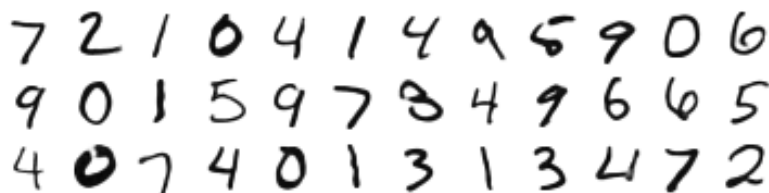
loss = mse_loss(output, input) + 0.5 * kl
```

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

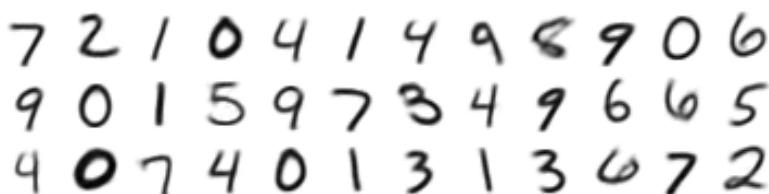
Original



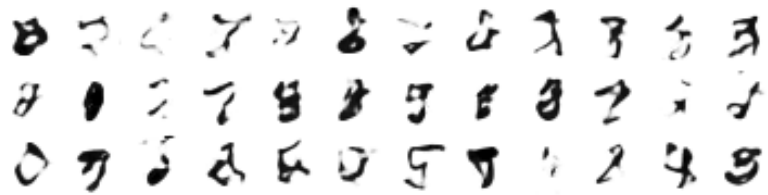
Autoencoder reconstruction ($d = 32$)



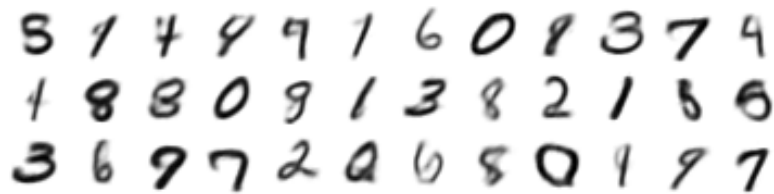
Variational Autoencoder reconstruction ($d = 32$)



Autoencoder sampling ($d = 32$)



Variational Autoencoder sampling ($d = 32$)



Non-Volume Preserving network

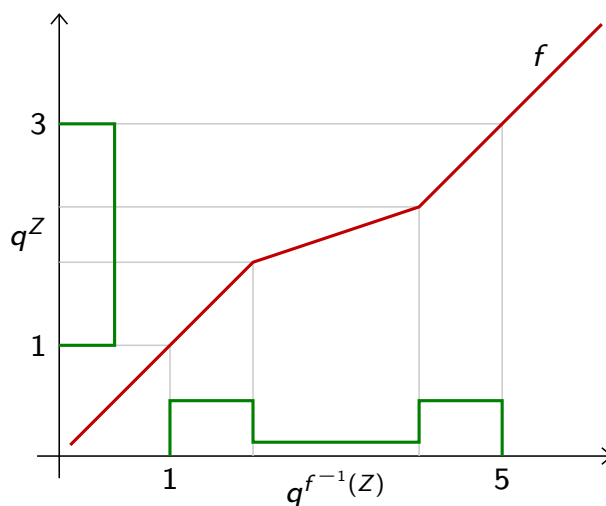
The motivation for the Variational Autoencoder was that although computing $g(z)$ is easy, maximizing $p^X(x)$ is not.

Some methods rely on network architectures for f , such that with $Z \sim \mathcal{N}(0, I)$ the distribution of $f^{-1}(Z)$ fits the data, which boils down to maximizing

$$\sum_n \log q^{f^{-1}(Z)}(x_n).$$

A standard result of probability theory is that if f is continuous, invertible and [almost everywhere] differentiable, then

$$\forall x, \quad q^{f^{-1}(Z)}(x) = q^Z(f(x)) |J_f(x)|.$$



Remember that if f is a composition of functions

$$f = f^{(K)} \circ \dots \circ f^{(1)}$$

we have

$$J_f(x) = \prod_{k=1}^K J_{f^{(k)}} \left(f^{(k-1)} \circ \dots \circ f^{(1)}(x) \right),$$

so

$$\log |J_f(x)| = \sum_{k=1}^K \log \left| J_{f^{(k)}} \left(f^{(k-1)} \circ \dots \circ f^{(1)}(x) \right) \right|.$$

So we are finally aiming at maximizing

$$\sum_n \left(\log q^Z(f(x_n)) + \sum_{k=1}^K \log \left| J_{f^{(k)}} \left(f^{(k-1)} \circ \dots \circ f^{(1)}(x_n) \right) \right| \right),$$

and since $Z \sim \mathcal{N}(0, I)$,

$$\log q^Z(f(x_n)) = -\frac{1}{2} (\|f(x_n)\|^2 + d \log 2\pi).$$

If $f^{(k)}$ are standard layers, computing their Jacobian determinant is intractable.

Besides, to be able to use this architecture for sampling we have to compute $f^{-1}(z)$ for a $z \sim \mathcal{N}(0, I)$, so we need invertible $f^{(k)}$ s.

Dinh et al. (2014) introduced the **coupling layers** to address both issues.

The resulting Non-Volume Preserving network (NVP) is an example of a **Normalizing flow**, which is invertible and allows sampling (Rezende and Mohamed, 2015).

We use here the formalism from Dinh et al. (2016).

Given a dimension d , a Boolean vector $b \in \{0, 1\}^d$ and two mappings

$$\begin{aligned} s &: \mathbb{R}^d \rightarrow \mathbb{R}^d \\ t &: \mathbb{R}^d \rightarrow \mathbb{R}^d, \end{aligned}$$

we define a [fully connected] coupling layer as the transformation

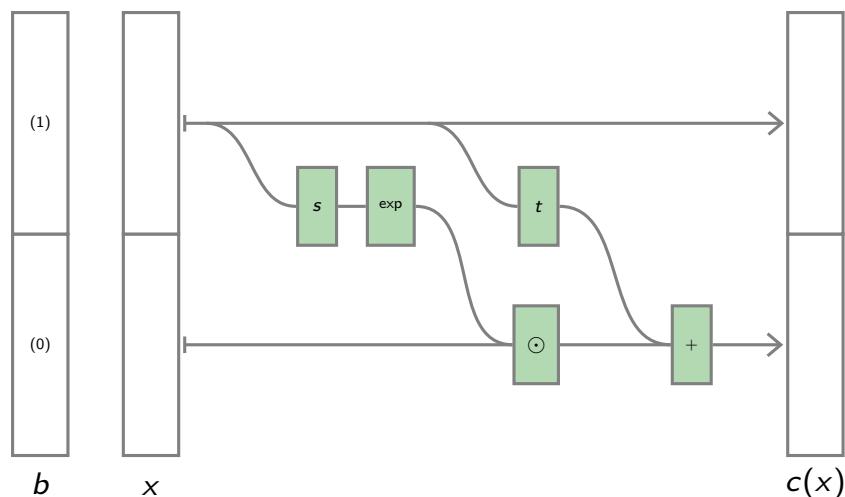
$$\begin{aligned} c &: \mathbb{R}^d \rightarrow \mathbb{R}^d \\ x &\mapsto b \odot x + (1 - b) \odot \left(x \odot \exp(s(b \odot x)) + t(b \odot x) \right) \end{aligned}$$

where \exp is component-wise, and \odot is the Hadamard component-wise product.

The expression

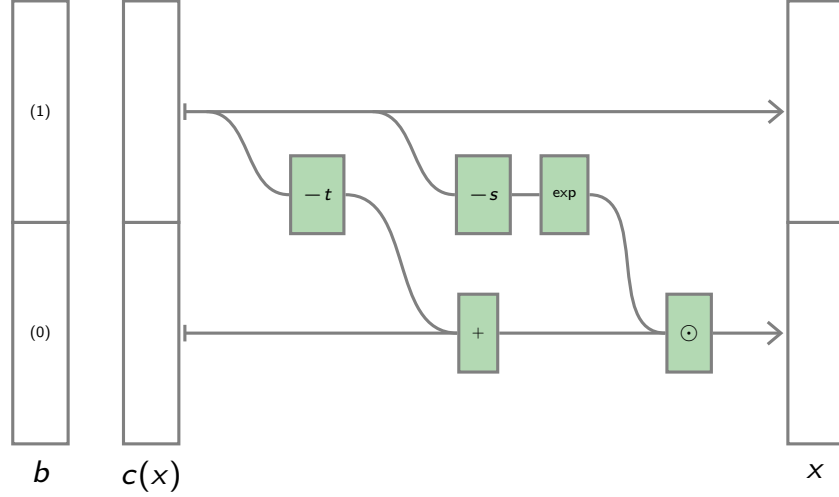
$$c(x) = b \odot x + (1 - b) \odot \left(x \odot \exp(s(b \odot x)) + t(b \odot x) \right)$$

can be understood as: forward $b \odot x$ unchanged, and apply to $(1 - b) \odot x$ an invertible transformation parametrized by $b \odot x$.



The consequence is that c is invertible, and if $y = c(x)$

$$x = b \odot y + (1 - b) \odot \left(y - t(b \odot y) \right) \odot \exp(-s(b \odot y)).$$



The second property of this mapping is the simplicity of its Jacobian determinant. Since

$$c_i(x) = b_i \odot x_i + (1 - b_i) \odot \left(x_i \odot \exp(s_i(b \odot x)) + t_i(b \odot x) \right)$$

we have, $\forall i, j, x$,

$$\begin{aligned} b_i = 1 &\Rightarrow c_i(x) = x_i \\ &\Rightarrow \frac{\partial c_i}{\partial x_j} = \delta_{i,j} \end{aligned}$$

and

$$\begin{aligned} b_i = 0 &\Rightarrow c_i(x) = x_i \exp(s_i(b \odot x)) + t_i(b \odot x) \\ &\Rightarrow \frac{\partial c_i}{\partial x_j} = \left(\delta_{i,j} + \underbrace{x_i \frac{\partial s_i(b \odot x)}{\partial x_j}}_{0 \text{ if } b_j=0} \right) \exp(s_i(b \odot x)) + \underbrace{\frac{\partial t_i(b \odot x)}{\partial x_j}}_{0 \text{ if } b_j=0} \\ &\Rightarrow \frac{\partial c_i}{\partial x_j} = \delta_{i,j} \exp(s_i(b \odot x)) + b_j \left(\dots \right). \end{aligned}$$

Hence $\frac{\partial c_i}{\partial x_j}$ can be non-zero only if $i = j$, or $(1 - b_i)b_j = 1$.

If we re-order both the rows and columns of the Jacobian to put first the non-zeros entries of b , and then the zeros, it becomes lower triangular

$$J_c(x) = \left(\begin{array}{ccc|ccc} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & (0) \\ \hline & & & \exp(s_k(x \odot b)) & & \\ & (\neq 0) & & & \ddots & \\ & & & & & \exp(s_{k'}(x \odot b)) \end{array} \right)$$

its determinant remains unchanged, and we have

$$\begin{aligned} \log |J_{f^{(k)}}(x)| &= \sum_{i: b_i=0} s_i(x \odot b) \\ &= \sum_i ((1 - b) \odot s(x \odot b))_i. \end{aligned}$$

```
dim = 6

x = Variable(Tensor(1, dim).normal_(), requires_grad = True)
b = Variable(Tensor(1, dim).zero_())
b.data.narrow(1, 0, dim//2).fill_(1.0)

s = nn.Sequential(nn.Linear(dim, dim), nn.Tanh())
t = nn.Sequential(nn.Linear(dim, dim), nn.Tanh())

c = b * x + (1 - b) * (x * s(b * x).exp() + t(b * x))

j = torch.cat([torch.autograd.grad(c_k, x, retain_graph=True)[0] for c_k in c[0]])

print(j)
```

prints

```
1.0000  0.0000  0.0000  0.0000  0.0000  0.0000
0.0000  1.0000  0.0000  0.0000  0.0000  0.0000
0.0000  0.0000  1.0000  0.0000  0.0000  0.0000
-0.8182  0.5622 -0.4035  1.4700  0.0000  0.0000
0.2610  0.1475  0.1689  0.0000  1.1358  0.0000
-0.2910  0.0287  0.0194  0.0000  0.0000  0.8508
[torch.FloatTensor of size 6x6]
```

To recap, with $f^{(k)}$, $k = 1, \dots, K$ coupling layers,

$$f = f^{(K)} \circ \dots \circ f^{(1)},$$

and $x_n^{(0)} = x_n$ and $x_n^{(k)} = f^{(k)}(x_n^{(k-1)})$, we train by maximizing

$$\mathcal{L}(f) = \sum_n \left(-\frac{1}{2} \left(\|x_n^{(K)}\|^2 + d \log 2\pi \right) + \sum_{k=1}^K \log |J_{f^{(k)}}(x_n^{(k-1)})| \right),$$

with

$$\log |J_{f^{(k)}}(x)| = \sum_i \left((1 - b^{(k)}) \odot s^{(k)} (x \odot b^{(k)}) \right)_i.$$

And to sample we just need to generate $Z \sim \mathcal{N}(0, I)$ and compute $f^{-1}(Z)$.

A coupling layer can be implemented with

```
class NVPCouplingLayer(Module):
    def __init__(self, map_s, map_t, b):
        super(NVPCouplingLayer, self).__init__()
        self.map_s = map_s
        self.map_t = map_t
        self.b = Variable(b.clone().unsqueeze(0), requires_grad = False)

    def forward(self, x_and_logdetjac):
        x, logdetjac = x_and_logdetjac
        s, t = self.map_s(self.b * x), self.map_t(self.b * x)
        logdetjac += ((1 - self.b) * s).sum(1)
        y = self.b * x + (1 - self.b) * (torch.exp(s) * x + t)
        return (y, logdetjac)

    def invert(self, y):
        s, t = self.map_s(self.b * y), self.map_t(self.b * y)
        return self.b * y + (1 - self.b) * (torch.exp(-s) * (y - t))
```

The `unsqueeze` method here adds a dimension of size 1 in front of the tensor size, so that we can use it with broadcasting on the multi-sample batch.

We can then define a complete network with one-hidden layer tanh MLPs for the s and t mappings

```
class NVPNet(Module):
    def __init__(self, dim, hdim, depth):
        super(NVPNet, self).__init__()
        b = Tensor(dim)
        self.layers = nn.ModuleList()
        for d in range(depth):
            if d%2 == 0:
                # Tag half the dimensions
                i = torch.randperm(b.numel()).narrow(0, 0, b.numel() // 2)
                b.zero_()[i] = 1
            else:
                b = 1 - b
            map_s = nn.Sequential(nn.Linear(dim, hdim), nn.Tanh(), nn.Linear(hdim, dim))
            map_t = nn.Sequential(nn.Linear(dim, hdim), nn.Tanh(), nn.Linear(hdim, dim))
            self.layers.append(NVPCouplingLayer(map_s, map_t, b))

    def forward(self, x):
        for m in self.layers: x = m(x)
        return x

    def invert(self, y):
        for m in reversed(self.layers): y = m.invert(y)
        return y
```

And the log-proba of individual samples of a batch

```
def LogProba(x_and_logdetjac):
    (x, logdetjac) = x_and_logdetjac
    log_p = logdetjac - 0.5 * x.pow(2).add(math.log(2 * math.pi)).sum(1)
    return log_p
```

Training is achieved by maximizing the mean log-proba

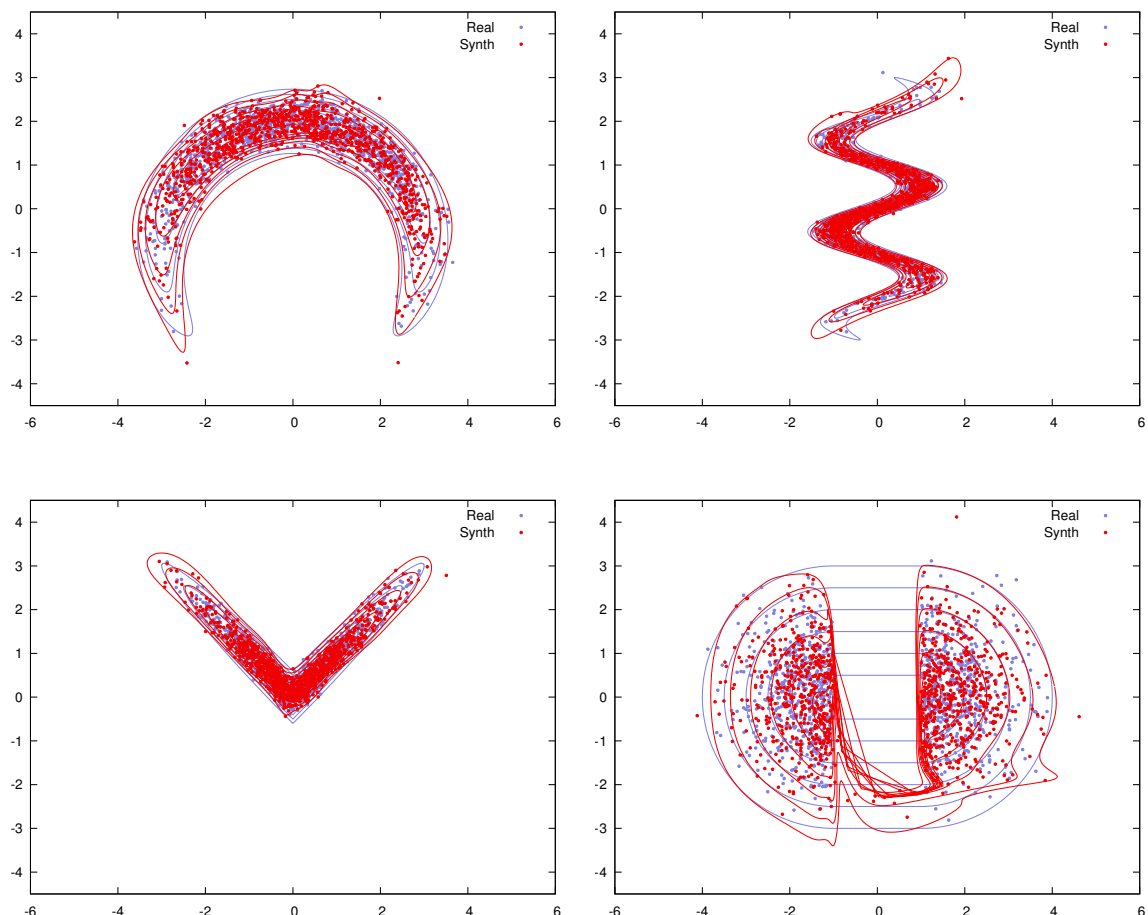
```
batch_size = 100

model = NVPNet(dim = 2, hidden_dim = 2, depth = 4)
optimizer = optim.Adam(model.parameters(), lr = 1e-2)

for e in range(args.nb_epochs):
    for b in range(0, nb_train_samples, batch_size):
        output = model((input.narrow(0, b, batch_size), 0))
        loss = - LogProba(output).mean()
        model.zero_grad()
        loss.backward()
        optimizer.step()
```

Finally, we can sample according to $q_{f^{-1}}(Z)$ with

```
z = Variable(Tensor(nb_train_samples, dim).normal_())
x = model.invert(z).data
```



Dinh et al. (2016) apply this approach to convolutional layers by using bs consistent with the activation map structure, and reducing the map size while increasing the number of channels.

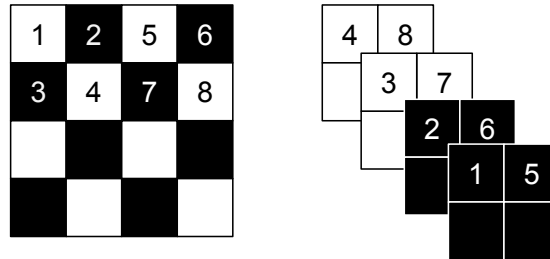


Figure 3: Masking schemes for affine coupling layers. On the left, a spatial checkerboard pattern mask. On the right, a channel-wise masking. The squeezing operation reduces the $4 \times 4 \times 1$ tensor (on the left) into a $2 \times 2 \times 4$ tensor (on the right). Before the squeezing operation, a checkerboard pattern is used for coupling layers while a channel-wise masking pattern is used afterward.

(Dinh et al., 2016)

They combine these layers by alternating masks, and branching out half of the channels at certain points to forward them unchanged.

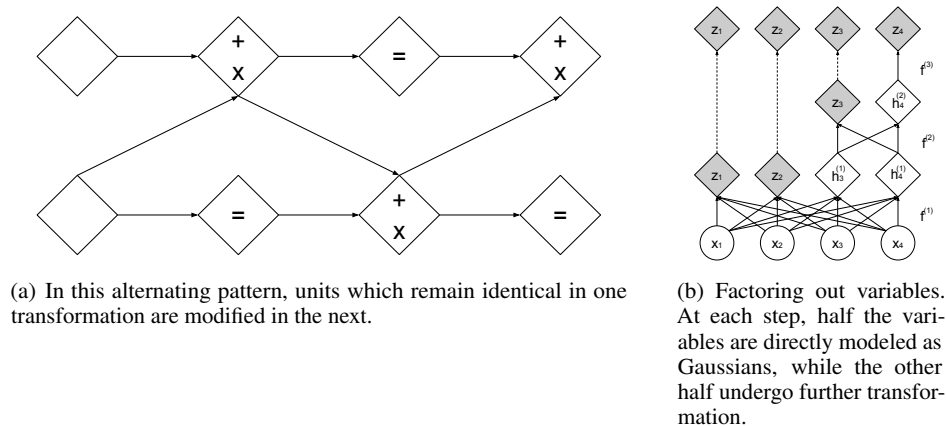


Figure 4: Composition schemes for affine coupling layers.

(Dinh et al., 2016)

The structure for generating images consists of

- $\times 2$ stages
 - $\times 3$ checkerboard coupling layers,
 - a squeezing layer,
 - $\times 3$ channel coupling layers,
 - a factor-out layer.
- $\times 1$ stage
 - $\times 4$ checkerboard coupling layers
 - a factor-out layer.

The s and t mappings get more complex in the later layers.

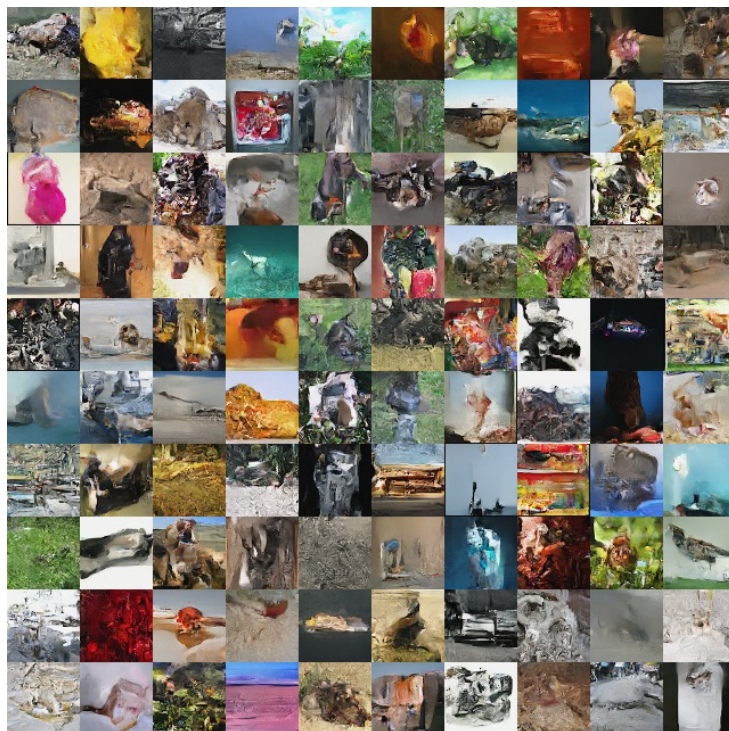


Figure 7: Samples from a model trained on *Imagenet* (64×64).

(Dinh et al., 2016)

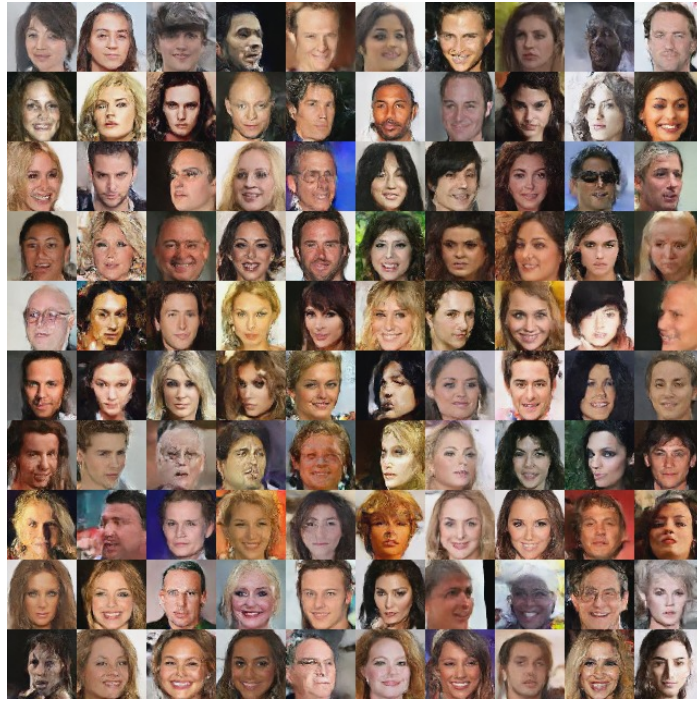


Figure 8: Samples from a model trained on *CelebA*.

(Dinh et al., 2016)



Figure 9: Samples from a model trained on *LSUN* (bedroom category).

(Dinh et al., 2016)



Figure 10: Samples from a model trained on *LSUN* (*church outdoor* category).

(Dinh et al., 2016)

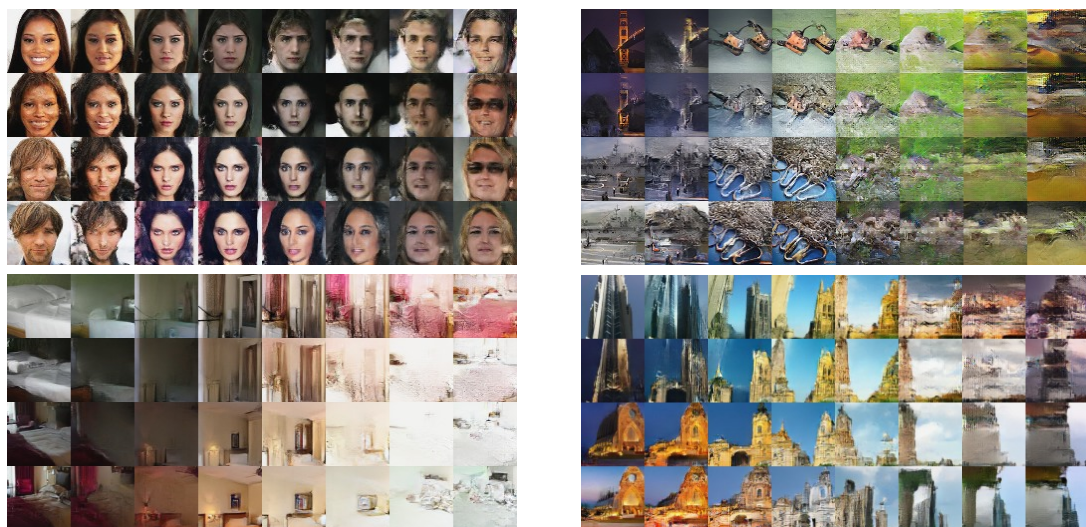


Figure 6: Manifold generated from four examples in the dataset. Clockwise from top left: CelebA, Imagenet (64×64), LSUN (tower), LSUN (bedroom).

(Dinh et al., 2016)

References

- L. Dinh, D. Krueger, and Y. Bengio. NICE: non-linear independent components estimation. *CoRR*, abs/1410.8516, 2014.
- L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using real NVP. *CoRR*, abs/1605.08803, 2016.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- D. Rezende and S. Mohamed. Variational inference with normalizing flows. *CoRR*, abs/1505.05770, 2015.
- P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research (JMLR)*, 11:3371–3408, 2010.