

Imperial College
London

May 16th, 2019

Introduction to Generative Models

Olivier Dubrule and Lukas Mosser

Imperial College
London

Objectives of the Day

- Recall what unsupervised learning is, with PCA as simplest linear case
- Autoencoders and latent vectors for image parametrization
- Introduce transposed convolution
- Variational Autoencoders for coding and generating random images

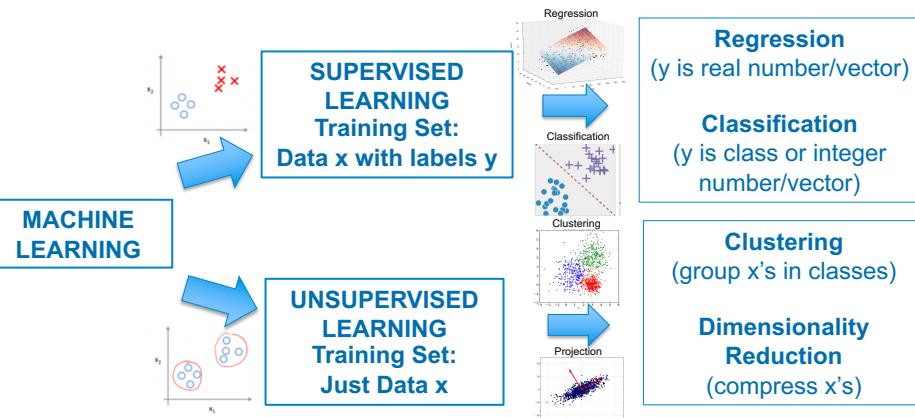
Imperial College
London

Introduction to Generative Models

1. PCA and Autoencoders
2. Transposed Convolution
3. Variational Autoencoders
4. Generative Models Generalization

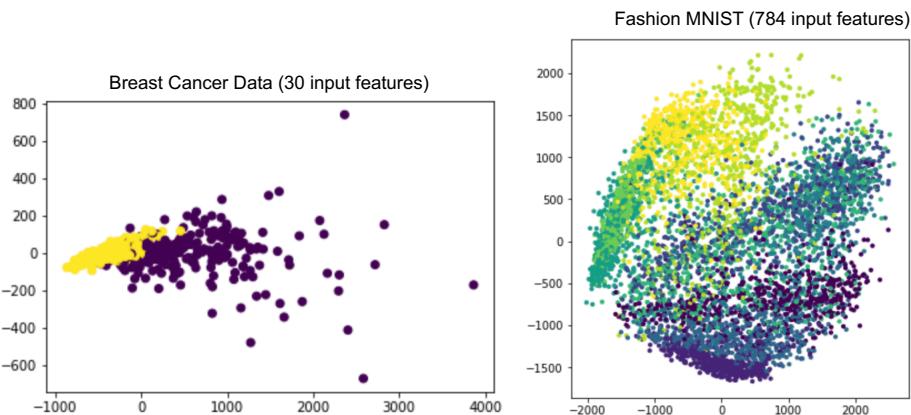
Imperial College
London

Supervised vs Unsupervised Learning



Imperial College
London

PCA Dimensionality Reduction



Imperial College
London

The PCA Method is Linear Dimensionality Reduction

If the first p eigenvectors are denoted as $(\varphi_k)_{k=1,\dots,p}$, each data vector $x^{(i)}$ of dimension n can be approximated by a linear combination of the eigenvectors:

$$x^{(i)} = \sum_{k=1}^p z_k^{(i)} \varphi_k$$

So, instead of being defined by its n original features, each $x^{(i)}$ is now defined by its p coordinates $z_k^{(i)}$ in the basis $(\varphi_k)_{k=1,\dots,p}$, where $p \ll n$: this is Dimensionality Reduction.

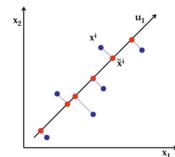
Imperial College
London

Dimensionality Reduction with Latent Vectors

Idea

Represent each image $(x^{(i)})_{i=1,\dots,m}$ where each x_i is of high dimension n (the number of features), by a set of latent vectors $(z^{(i)})_{i=1,\dots,m}$, where the vectors $z^{(i)}$ are of much smaller dimension p than that n of the images.

PCA projects only on a linear manifold, we wish to generalize to non-linear manifolds.

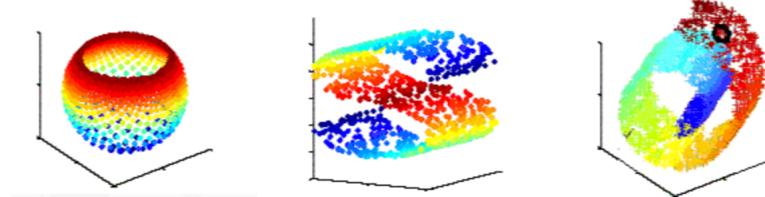


Potential applications

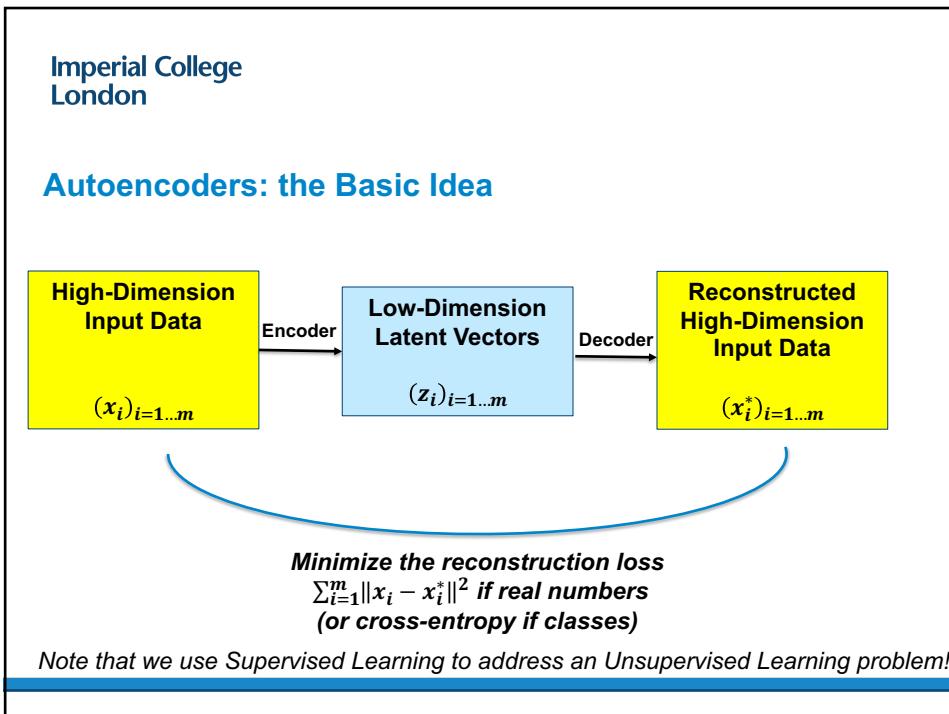
Simple parametrization of images, classification or regression based on the feature vectors $z^{(i)}$ rather than the raw data $x^{(i)}$, visualization, data compression, inversion....

Imperial College
London

Why do we Need a Non-Linear Approach ?



Examples of data that cannot be separated by projecting on a linear manifold.



Imperial College London

Autoencoders Historical Example on MNIST

Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

Dimensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer “encoder” network

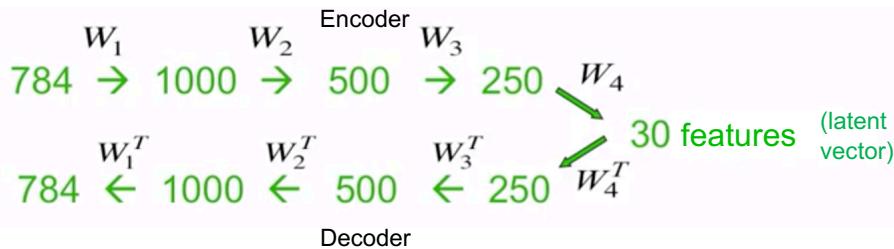
28 JULY 2006 VOL 313 SCIENCE www.sciencemag.org

From Hinton and Salakhutdinov, Science, Science, July 2006
<https://www.youtube.com/watch?v=Ex1Ur85AVs>

Imperial College
London

Autoencoder Historical Example on MNIST

Structure of Hinton and Salakhutdinov's autoencoders (loss function is cross-entropy)

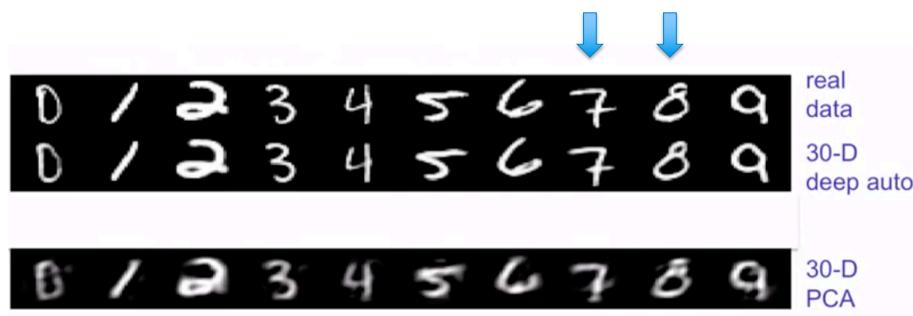


Each MNIST image represented now in 30 dimensions instead of originally in $28 \times 28 = 784$ dimensions.

From Hinton and Salakhutdinov, Science, Science, July 2006
<https://www.youtube.com/watch?v=Ex1Ur85AVs>

Imperial College
London

Autoencoder Historical Example on MNIST



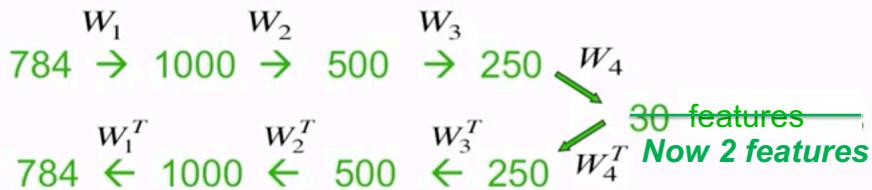
The reconstructed numbers (second line) look better than the originals, because the autoencoder model does not have enough capacity to represent the subtle variations/noise in the data.

From Hinton and Salakhutdinov, Science, Science, July 2006
<https://www.youtube.com/watch?v=Ex1Ur85AVs>

Imperial College
London

Autoencoder Historical Example on MNIST

Now using a 2-dimensional instead of 30-dimensional code



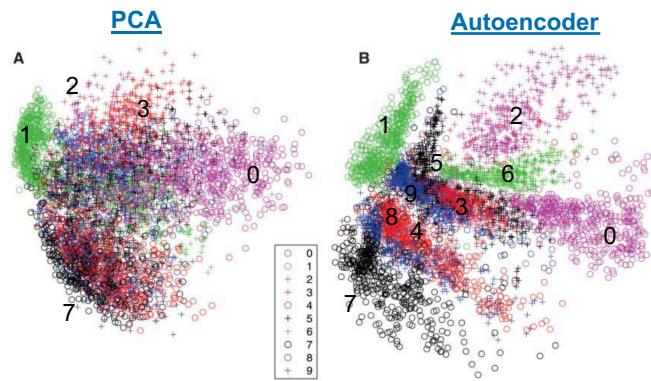
Each MNIST number represented now in 2 dimensions
instead of originally in $28 \times 28 = 784$ dimensions!

From Hinton and Salakhutdinov, Science, Science, July 2006
<https://www.youtube.com/watch?v=Ex1Ur85AVs>

Imperial College
London

Autoencoder Historical Example on MNIST

Now representing 28x28 MNIST images by just two coordinates



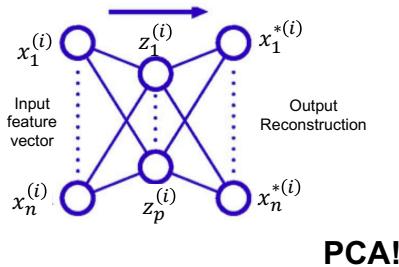
(A) The 2-D codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784- 1000-500-250-2 autoencoder.

From Hinton and Salakhutdinov, Science, Science, July 2006

Imperial College
London

Question

What result autoencoder will produce if there is just one hidden layer and no non-linear activation function in the encoder and decoder networks?



PCA!

$$z = Wx$$

$$x^* = W'z$$

Minimize:
$$\sum_{i=1}^m \|x_i - x_i^*\|^2$$

https://www.cs.toronto.edu/~urtasun/courses/CSC411/14_pca.pdf

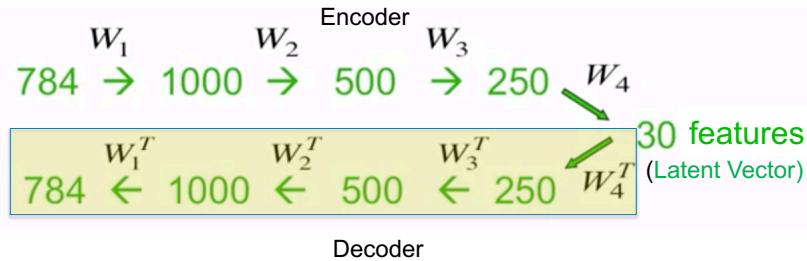
Imperial College
London

Introduction to Generative Models

1. PCA and Autoencoders
2. Transposed Convolution
3. Variational Autoencoders
4. Generative Models Generalization

Imperial College
London

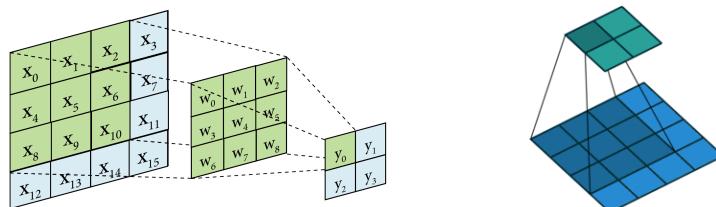
Transposed Convolution Filters



Imperial College
London

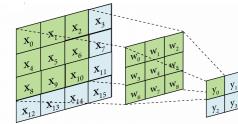
Encoder: Convolution Reduces Dimensionality

Simple 2-D Convolution Example (stride 1, no-padding)



$$y_0 = w_0x_0 + w_1x_1 + w_2x_2 + w_3x_4 + w_4x_5 + w_5x_6 + w_6x_8 + w_7x_9 + w_8x_{10} + b$$

Can we write the Convolution Operation in Matrix Form?



Encoder: Writing Convolution as a Matrix Operation

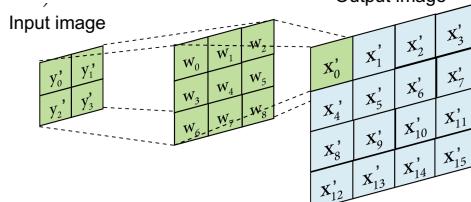
$$\begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} w_0 & w_1 & w_2 & 0 & w_3 & w_4 & w_5 & 0 & w_6 & w_7 & w_8 & 0 & 0 & 0 & 0 & 0 \\ 0 & w_0 & w_1 & w_2 & 0 & w_3 & w_4 & w_5 & 0 & w_6 & w_7 & w_8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & w_0 & w_1 & w_2 & 0 & w_3 & w_4 & w_5 & 0 & w_6 & w_7 & w_8 & 0 \\ 0 & 0 & 0 & 0 & 0 & w_0 & w_1 & w_2 & 0 & w_3 & w_4 & w_5 & 0 & w_6 & w_7 & w_8 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \\ x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \\ x_{15} \end{pmatrix}$$

Convolution can be written in matrix form $y = Wx$ with x reshaped as a 16×1 vector, y reshaped as a 4×1 vector, and a 4×16 matrix W .

https://github.com/vdumoulin/conv_arithmetic

<https://towardsdatascience.com/up-sampling-with-transposed-convolution-9ae4f2df52d0>

Decoder: Transposed Convolution to Increase Dimensionality



y' reshaped into 4×1 vector, x' reshaped into 16×1 vector and $W^T 16 \times 4$ transposed of W

$$x' = W^T y'$$

$$\begin{pmatrix} x'_0 \\ x'_1 \\ x'_2 \\ x'_3 \\ x'_4 \\ x'_5 \\ x'_6 \\ x'_7 \\ x'_8 \\ x'_9 \\ x'_{10} \\ x'_{11} \\ x'_{12} \\ x'_{13} \\ x'_{14} \\ x'_{15} \end{pmatrix} = \begin{pmatrix} w_0 & 0 & 0 & 0 \\ w_1 & w_0 & 0 & 0 \\ w_2 & w_1 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ w_3 & 0 & w_0 & 0 \\ w_4 & w_3 & w_1 & w_0 \\ w_5 & w_4 & w_2 & w_1 \\ 0 & w_5 & 0 & w_2 \\ w_6 & 0 & w_3 & 0 \\ w_7 & w_6 & w_4 & w_3 \\ w_8 & w_7 & w_5 & w_4 \\ 0 & w_8 & 0 & w_5 \\ 0 & 0 & w_6 & 0 \\ 0 & 0 & w_7 & w_6 \\ 0 & 0 & w_8 & w_7 \\ 0 & 0 & 0 & w_8 \end{pmatrix} \times \begin{pmatrix} y'_0 \\ y'_1 \\ y'_2 \\ y'_3 \end{pmatrix}$$

https://github.com/vdumoulin/conv_arithmetic

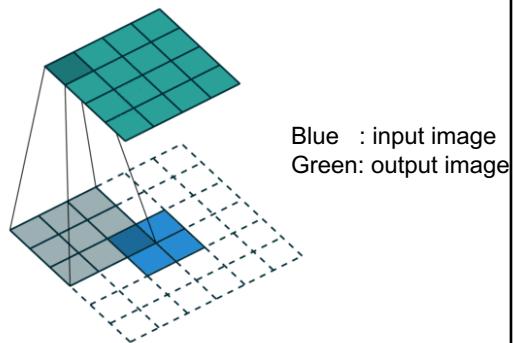
Imperial College
London

Transposed Convolution can also be Calculated by Convolution

Apply padding = 2 to array

$$\begin{matrix} y'_0 & y'_1 \\ y'_2 & y'_3 \end{matrix}$$

And convolve by filter W with
stride 1.

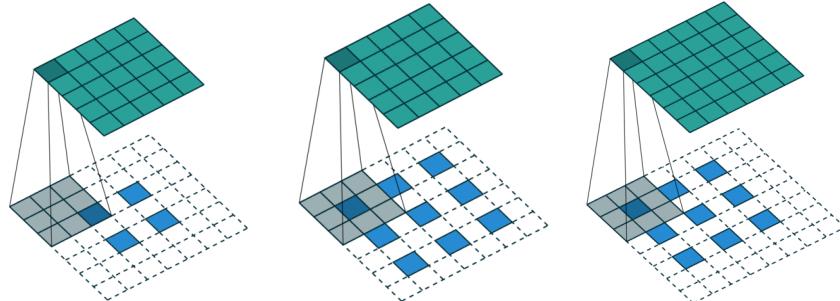


See afternoon exercise

https://github.com/vdumoulin/conv_arithmetic

Imperial College
London

Transposed Convolution can also be Calculated by Convolution



Blue : input image
Green: output image

Issue of checker board artefacts: <https://distill.pub/2016/deconv-checkerboard/>

https://github.com/vdumoulin/conv_arithmetic

Imperial College
London

Which Approach to Identify Transposed Convolution Operator?

1. Identify first the associated Direct Convolution
2. Calculate the Matrix associated with the Direct Convolution
3. Obtain Transposed Convolution from Transposed Matrix
4. Identify the Direct Convolution mimicking the Transposed one

<https://arxiv.org/pdf/1603.07285.pdf>

Imperial College
London

Introduction to Generative Models

- 1. PCA and Autoencoders**
- 2. Transposed Convolution**
- 3. Variational Autoencoders**
- 4. Generative Models Generalization**

Imperial College
London

Which Use for Autoencoders?



Throw away the calibrated Decoder and use the Encoder for Dimensionality Reduction!

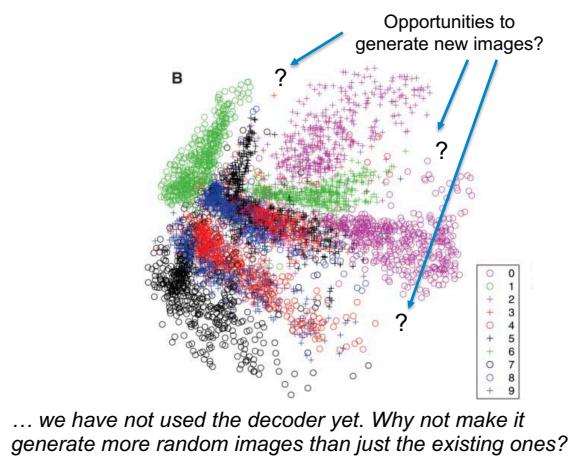
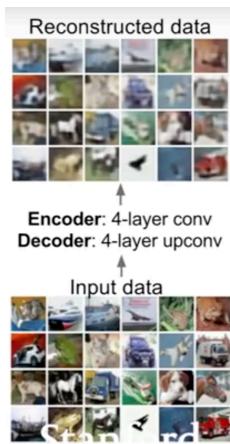
Every Image of MNIST Dataset now represented by only 30 features instead of by 784 pixel values.

Possible Application: a Supervised Learning exercise on the MNIST Dataset can use the 30 features as input features, instead of the 784 pixel values. Can avoid overfitting in case there are many input images but not many of them are labelled.

From Hinton and Salakhutdinov, Science, Science, July 2006
<https://www.youtube.com/watch?v=Ex1Ur85AVs>

Imperial College
London

Autoencoders can Only Encode and Decode Existing Images



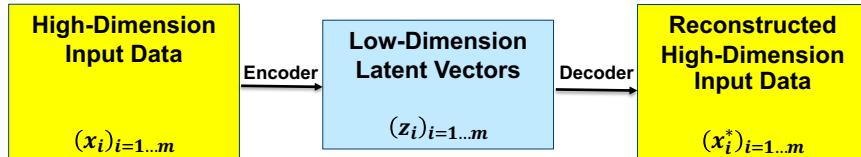
... we have not used the decoder yet. Why not make it generate more random images than just the existing ones?

<https://www.youtube.com/watch?v=5WoIgGTWV54>

From Hinton and Salakhutdinov, Science, Science, July 2006

Imperial College
London

Variational Autoencoders (VAEs): Generate New Images which were not in Training Dataset



Variational autoencoders will be used to generate new images, that is images that statistically “look like” or belong to the same general pdf as those from the training set .

Make z and x random, instead of keeping them deterministic!

Imperial College
London

The Variational Autoencoder Model (in case of images)

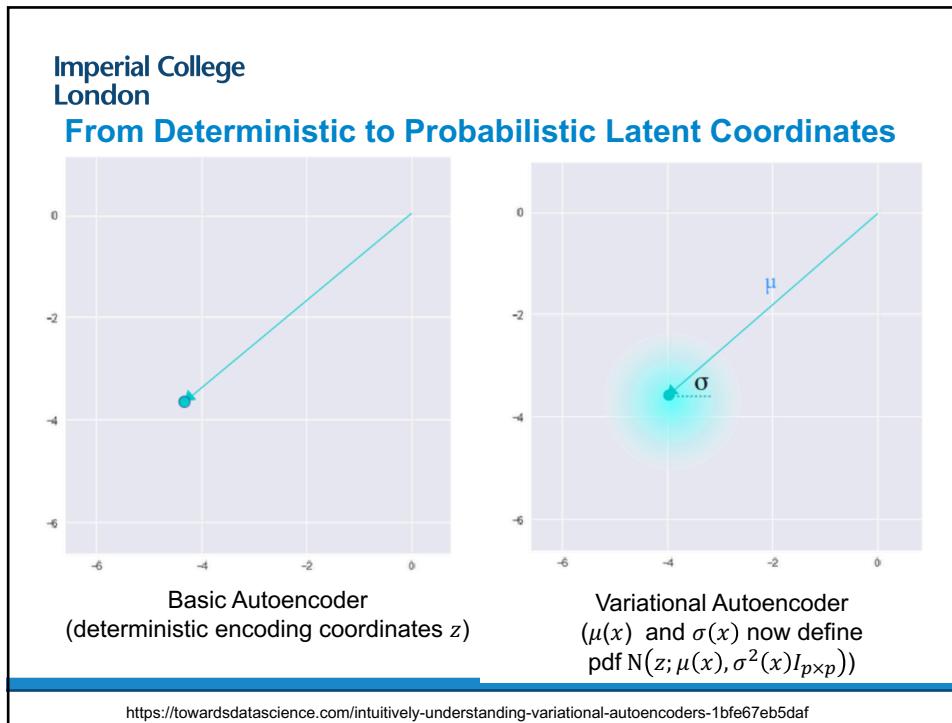
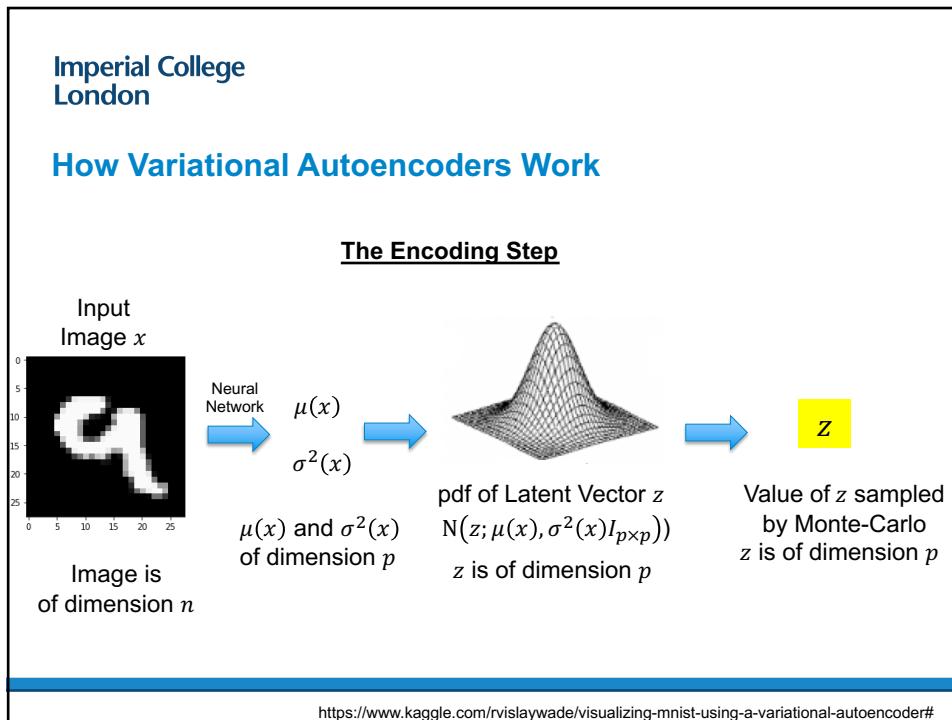
Training Data

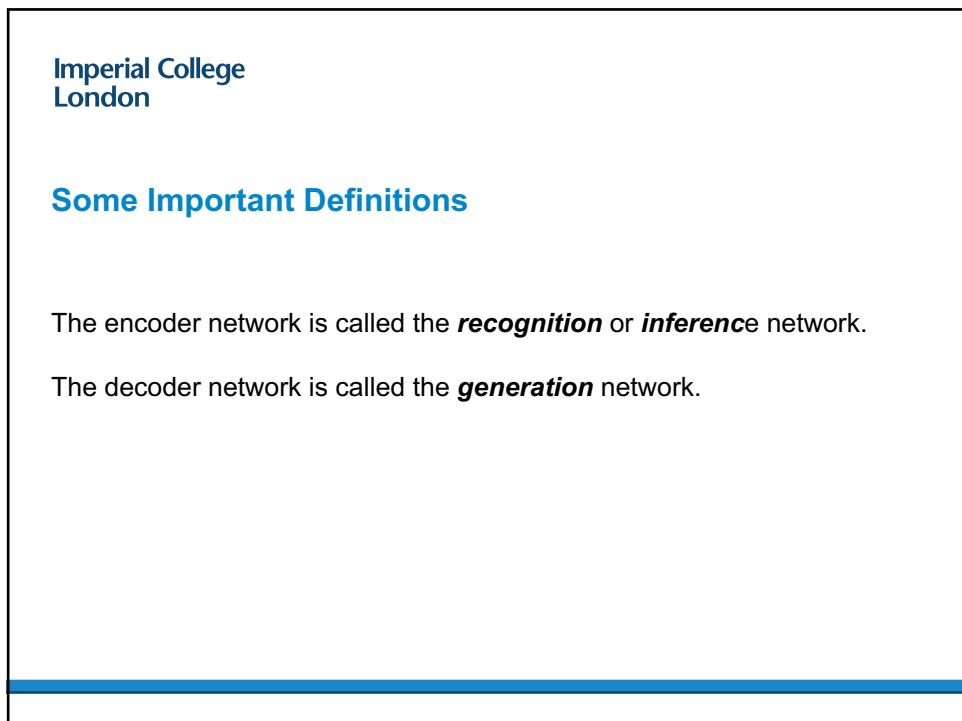
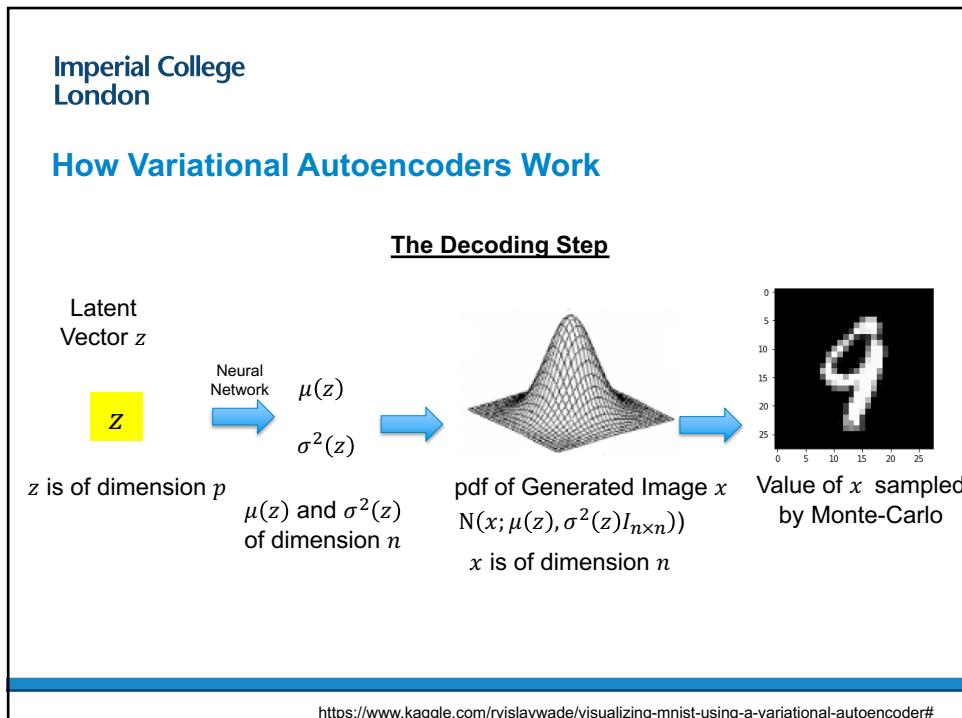


The set of images x (each image of dimension n) result from the application of a (neural network) function to a “latent vector” of much smaller dimension p than n . This latent vector is assumed to follow a standardized MultiGaussian distribution $N(z; 0, I_{p \times p})$.

How can we derive the parameters of this neural network?

By maximizing the likelihood of the Training Set samples $(x^i)_{i=1,\dots,m}$!





Imperial College
London

How Variational Autoencoders Work

Call $q_\phi(z/x)$ the multivariate normal diagonal pdf $N(z; \mu(x), \sigma^2(x)I_{p \times p})$

Call $p_\theta(x/z)$ the multivariate normal diagonal pdf $N(x; \mu(z), \sigma^2(z)I_{n \times n})$

$\mu(x)$ and $\sigma(x)$, which have the dimension p of z , are calculated from x using a neural network. The parameters of this neural network are the ϕ 's of $q_\phi(z/x)$.

$\mu(z)$ and $\sigma(z)$, which have the dimension n of x , are calculated from z using a neural network. The parameters of this neural network are the θ 's of $p_\theta(x/z)$.

These parameters must be calibrated using the m data $x^{(i)}$ by maximization of a likelihood function.

Imperial College
London

How Variational Autoencoders Work (Kingma & Welling, 2014)

For one single data point $x^{(i)}$ the exact log-likelihood $\log p(x^{(i)})$ cannot be calculated, only a lower bound can

$$\log p(x^{(i)}) \geq E_{z \sim q_\phi(z/x^{(i)})}(\log p_\theta(x^{(i)}/z)) - \text{KL}(q_\phi(z/x^{(i)}) \| q(z))$$

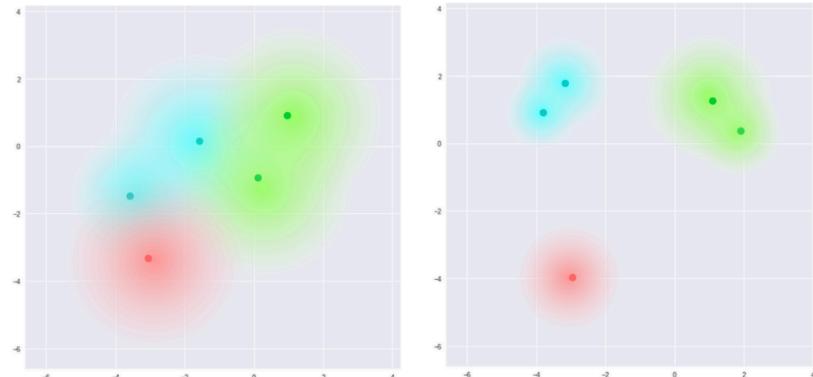
The right-hand side is called the Evidence Lower Bound (ELBO).
The ELBO is negative, as the sum of two negative terms.

The first term of the ELBO is a **reconstruction likelihood**, quantifying the quality of the « round trip » between the training image $x^{(i)}$ and its associated latent vector pdf. The second one is a **KL divergence** that aims to « spread » the range of latent vectors associated to the image $x^{(i)}$.

For a mini-batch of m images, the ELBO is summed over the mini-batch data $(x^{(i)})_{i=1..m}$.

Imperial College
London

If Only First (Likelihood) Term of Objective Function is Used

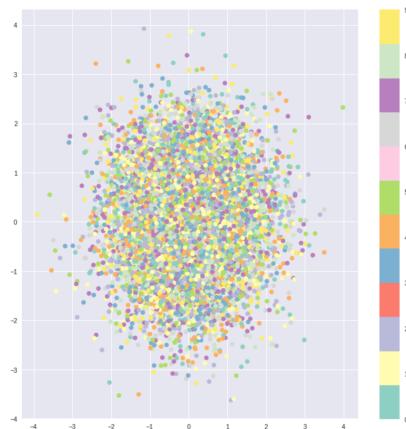


We want a continuous latent space.....But the Likelihood Term tends to create clusters

<https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>

Imperial College
London

If Only KL Divergence Term of Objective Function is Used

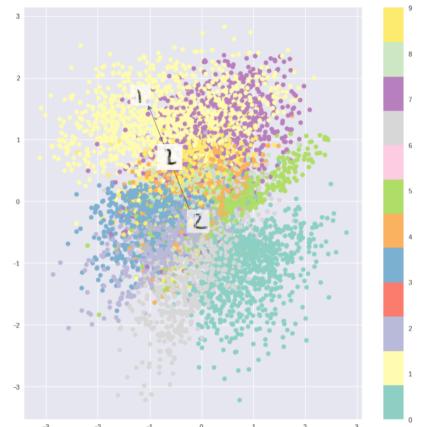


Now the space is nicely continuous but there is no more difference between coded digits

<https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>

Imperial College
London

If Both Terms of Objective Function are Used



<https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>

Imperial College
London

How Variational Autoencoders Work (Kingma & Welling, 2014)

Calculating $-KL(q_\phi(z/x^{(i)})||q(z))$

This is the KL divergence between two Gaussian distributions, each of dimension p :

$$q_\phi(z/x^{(i)}) \text{ is } N(z; \mu(x^{(i)}), \sigma^2(x^{(i)})I_{p \times p}) \text{ or } N(z; \mu^{(i)}, \sigma^{(i)2})$$

$$q(z) \text{ is } N(z; 0_p, I_{p \times p})$$

Using the formula giving the KL divergence of two multivariate Gaussians we get

$$-KL(q_\phi(z/x^{(i)})||q(z)) = \frac{1}{2} \sum_{j=1}^p \left[1 + \log \left((\sigma_j^{(i)})^2 \right) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right]$$

Where p is the dimension of the latent vector

Imperial College
London

Recall: KL divergence of 2 Multivariate Gaussians

Suppose the two distributions are: $N(\mu_1, \Sigma_1)$ and $N(\mu_2, \Sigma_2)$, each of dimension p .

The KL divergence between the two distributions is:

$$KL = \frac{1}{2} \left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - p + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right]$$

In our case the first distribution is that of the conditional Gaussian of mean $\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}$ and diagonal variance-covariance $\Sigma = \begin{pmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_p^2 \end{pmatrix}$ and the second one is the standardized multivariate normal: $\mu_2 = 0$ and $\Sigma_2 = I_{p \times p}$. Thus, in our case:

$$KL = \frac{1}{2} \left[\log \frac{1}{\prod_{i=1}^p \sigma_i^2} - p + \sum_{i=1}^p \sigma_i^2 + \sum_{i=1}^p \mu_i^2 \right] = \frac{1}{2} \left[- \sum_{i=1}^p (1 + \log(\sigma_i^2)) + \sum_{i=1}^p (\sigma_i^2 + \mu_i^2) \right]$$

Imperial College
London

How Variational Autoencoders Work (Kingma & Welling, 2014)

Calculating $E_{z \sim q_\phi(z/x^{(i)})}(\log p_\theta(x^{(i)}/z))$ for one data point $x^{(i)}$

This expression looks complex but simply quantifies the likelihood associated with the « round trip » between the image x_i and the space of associated latent vectors z . This is an expectation, which can be calculated by Monte Carlo sampling and averaging over a large set of L simulated values of z .

$$E_{z \sim q_\phi(z/x^{(i)})}(\log p_\theta(x^{(i)}/z)) = \frac{1}{L} \sum_{l=1}^L \log p_\theta(x^{(i)}/z^{(i,l)})$$

In order to express this as a function of parameters $\mu_i = \mu(x_i)$ and $\sigma_i^2 = \sigma^2(x_i)$, we write:

$$z^{(i,l)} = \mu^{(i)} + \sigma^{(i)} \odot \varepsilon^{(l)} \text{ where } \varepsilon^{(l)} \text{ is } N(\varepsilon; 0, I_{p \times p})$$

Imperial College
London

Likelihood of Uncorrelated Multivariate Gaussian

$$f(x) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

If $\Sigma = \sigma^2 I_n$, this simplifies into:

$$f(x) = (2\pi)^{-\frac{n}{2}} \frac{1}{\prod_{j=1}^n \sigma_j} \exp \left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{x_j - \mu_j}{\sigma_j} \right)^2 \right]$$

$$\log p_\theta(x^{(i)} / z^{i,l}) = -\frac{n}{2} \log(2\pi) - \sum_{j=1}^n \log(\sigma_j(z^{i,l})) - \frac{1}{2} \sum_{j=1}^n \left(\frac{x_j^{(i)} - \mu_j(z^{i,l})}{\sigma_j(z^{i,l})} \right)^2$$

<https://stats.stackexchange.com/questions/351549/maximum-likelihood-estimators-multivariate-gaussian>

Imperial College
London

Summarizing how we calculate a VAE

We have defined the ELBO as the right-hand side of :

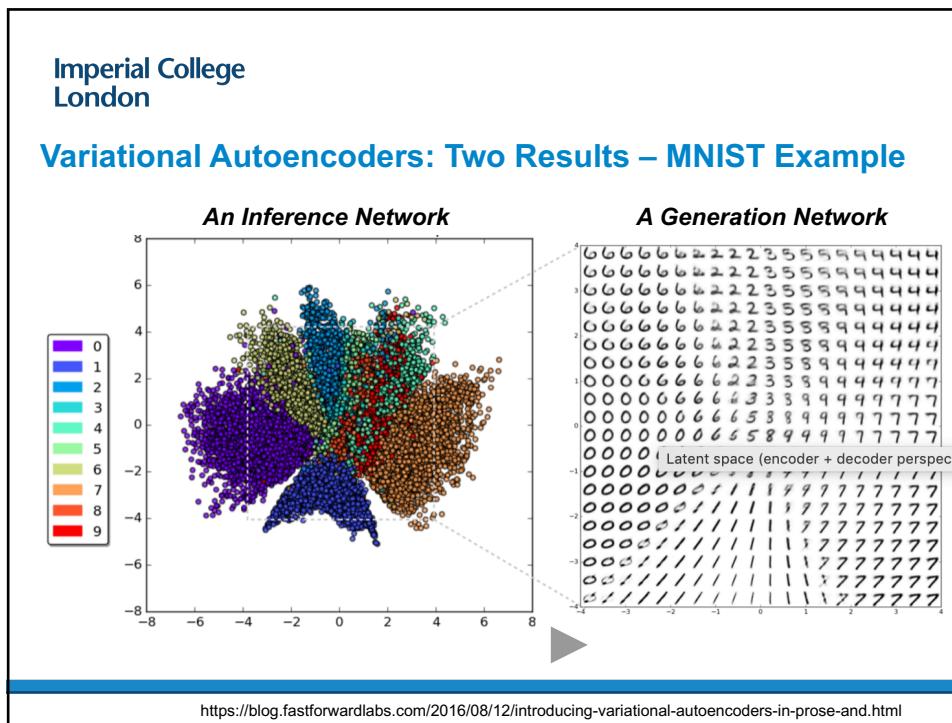
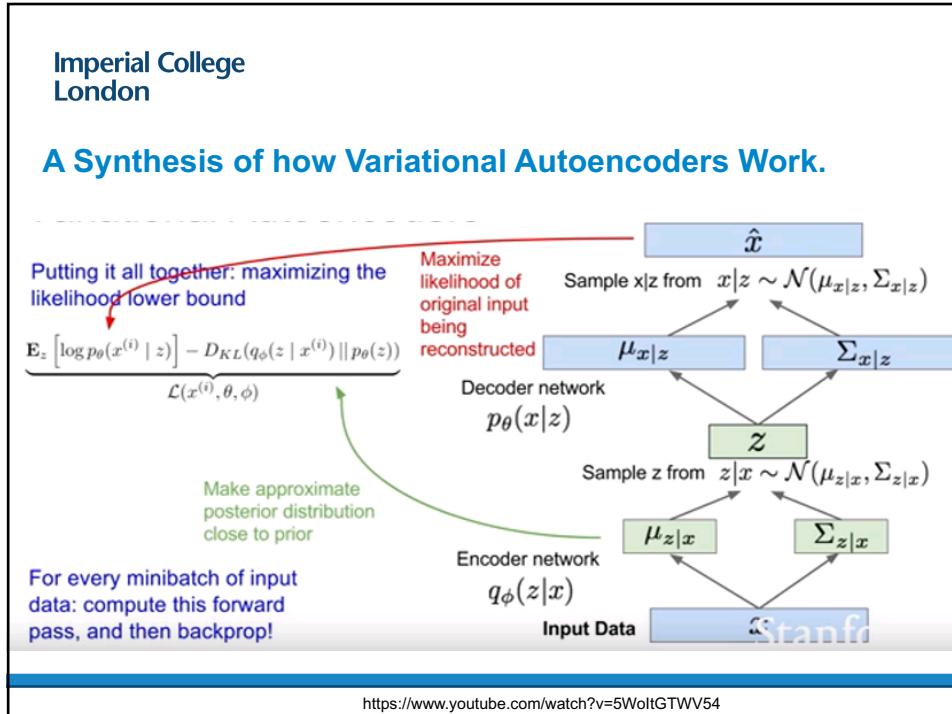
$$\log p(x^{(i)}) \geq E_{q(z)}(\log p_\theta(x^{(i)} / z)) - \text{KL}(q_\phi(z/x^{(i)}) \| q(z))$$

$E_{q(z)}(\log p_\theta(x^{(i)} / z))$ is a function of the vectors $\mu(z_i), \sigma^2(z_i), \mu(x_i), \sigma^2(x_i)$

$-\text{KL}(q_\phi(z/x^{(i)}) \| q(z))$ is a function of the vectors $\mu(x_i), \sigma^2(x_i)$

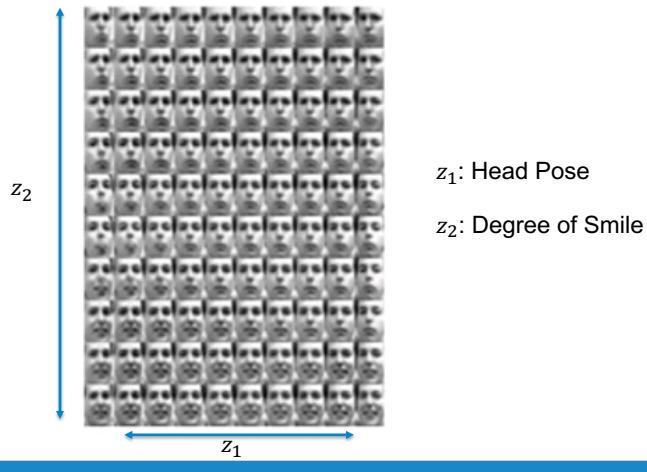
The ELBO, sum of these two terms, is thus a function of $\mu(z_i), \sigma^2(z_i), \mu(x_i), \sigma^2(x_i)$.

The parameters ϕ and θ of the neural networks respectively calculating $\mu(x), \sigma^2(x)$ as a function of x and $\mu(z), \sigma^2(z)$ as a function of z are obtained by maximizing the ELBO on the training set $(x_i)_{i=1..m}$.



Imperial College
London

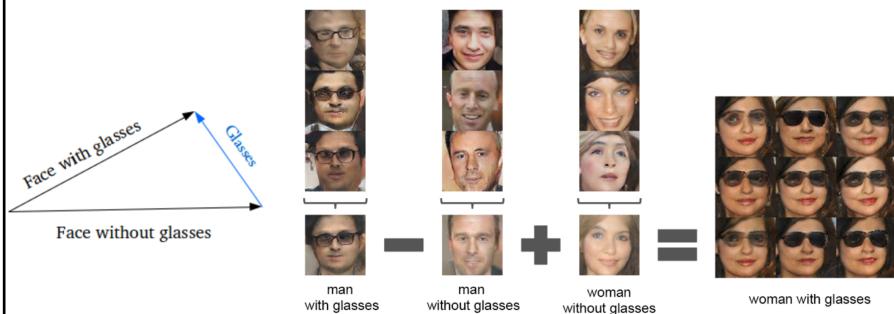
Interpreting the Latent Vector Features or Coordinates



Kingma and Welling, 2014

Imperial College
London

Arithmetic Operations on Latent Vectors



Arithmetic operations on images (from
willwhitney.github.io/gan-article/)

Imperial College
London

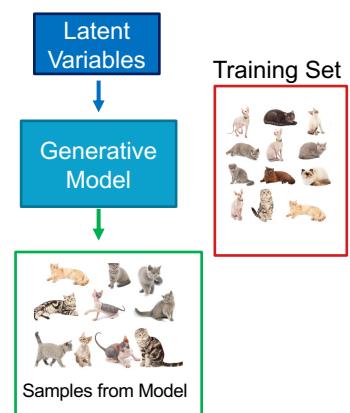
Introduction to Generative Models

1. PCA and Autoencoders
2. Transposed Convolution
3. Variational Autoencoders
4. **Generative Models Generalization**

Imperial College
London

(Deep) Generative Methods

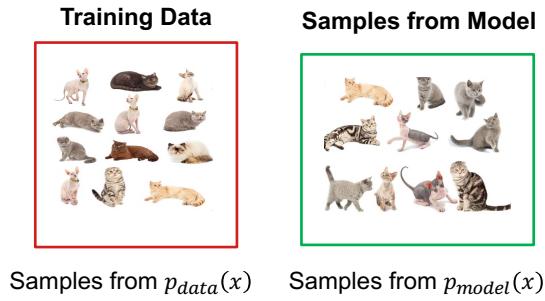
- Task: Draw samples from unknown pdf given a set of samples
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
 - Many more ...



Imperial College
London

Why Generative Models?

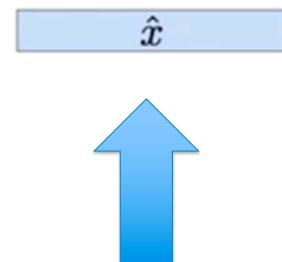
Generic Synthetic Images, Music, Text....



Goal: learn $p_{model}(x)$ from $p_{data}(x)$ in order to generate more samples “looking like” the data

Imperial College
London

VAEs: Decoder for Generating Samples



Sample z from $z \sim \mathcal{N}(0, I)$

Imperial College
London

Objective of Generative Networks Based on Latent Vectors

Apply a transformation G to random latent vector z (usually a multivariate Gaussian $N(0, I)$), such that the pdf $p_{model}(x)$ of $x = G(z)$ is close to $p_{data}(x)$



Will need tools to compare pdfs: KL Divergence, Maximum Mean Discrepancy, Discriminator...

Imperial College
London

Compare two pdfs: the Kullback-Leibler (KL) Divergence

If the distributions $p(x)$ and $q(x)$ are continuous:

$$D_{KL}(p\|q) = \int_{-\infty}^{+\infty} p(x) \log p(x) dx - \int_{-\infty}^{+\infty} p(x) \log q(x) dx$$

Two Fundamental Properties

$D_{KL}(p\|q)$ is positive, and 0 if $p(x)$ and $q(x)$ identical
Asymmetry: $D_{KL}(p\|q) \neq D_{KL}(q\|p)$

Imperial College
London

Compare two pdfs: the MMD

Consider two sample datasets:

$$X = (x^{(1)}, x^{(2)}, \dots, x^{(m)}) \text{ and } X' = (x'^{(1)}, x'^{(2)}, \dots, x'^{(m')})$$

The maximum mean discrepancy (MMD) compares the two samples by comparing the sample means in a so-called feature space:

$$MMD(X, X') = \left\| \sum_{i=1}^m \phi(x^{(i)}) - \sum_{i'=1}^{m'} \phi(x'^{(i')}) \right\|^2$$

which can be written as a function of a kernel function

$$MMD(X, X') = \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(x^{(i)}, x^{(j)}) + \frac{1}{m'^2} \sum_{i'=1}^{m'} \sum_{j'=1}^{m'} k(x'^{(i')}, x'^{(j')}) - \frac{2}{mn} \sum_{i=1}^m \sum_{j'=1}^{m'} k(x^{(i)}, x'^{(j')})$$

Imperial College
London

Compare two pdfs: Generative Adversarial Networks

If $x = G(z)$ is a sample from $p_{model}(x)$, a Discriminator function $D(x)$ is applied to x and calculates a number between 0 and 1. The closer $D(x)$ is to 1, the more likely is x to be also a sample of $p_{data}(x)$.
Both $D(x)$ and $G(z)$ are neural networks.

Imperial College
London

Conclusion

- Autoencoders a Generalization of PCAs to non-Linear Manifolds
- Variational Autoencoders give a Probabilistic Spin to Autoencoders
- Variational Autoencoders produce both an Inference and a Generative Network
- Introduction to other Generative Neural Networks and the general issue of approximating one pdf $p_{data}(x)$ by another pdf $p_{model}(x)$