



Deep Generative Modeling - part 1 -

Hana Sebia

hana.sebia@inria.fr

December, 2023

AIstroSight, Inria, Lyon Center

Which face is real ?



A



B



C

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn function to map
 $x \rightarrow y$

Examples: Classification,
regression, object detection,
semantic segmentation, etc.

Unsupervised Learning

Data: x

x is data, no labels!

Goal: Learn some *hidden* or
underlying structure of the data

Examples: Clustering, feature or
dimensionality reduction, etc.

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn function to map

$$x \rightarrow y$$

Examples: Classification,
regression, object detection,
semantic segmentation, etc.

Unsupervised Learning

Data: x

x is data, no labels!

Goal: Learn the *hidden* or
underlying structure of the data

Examples: Clustering, feature or
dimensionality reduction, etc.

Goal: Take as input training samples from some distribution and learn a model that represent that distribution



How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

Why Generative Models ? Debiasing

Capable of uncovering **underlying features** in a dataset



Homogeneous skin color, pose

vs



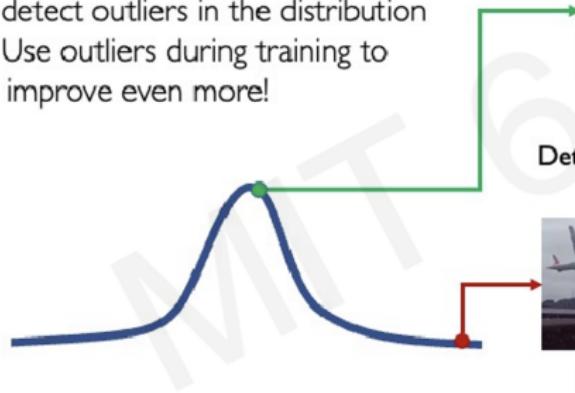
Diverse skin color, pose, illumination

How can we use this information to create fair and representative datasets?

Why Generative Models ? Outlier detection

Inria

- **Problem:** How can we detect when we encounter something new or rare?
- **Strategy:** Leverage generative models, detect outliers in the distribution
- Use outliers during training to improve even more!



95% of Driving Data:
(1) sunny, (2) highway, (3) straight road



Detect outliers to avoid unpredictable behavior when training



Edge Cases



Harsh Weather



Pedestrians

Generative AI

- Unsupervised process
- Take training samples from some unknown distribution and learn a model that can capture this distribution

Generative Model Trilemma

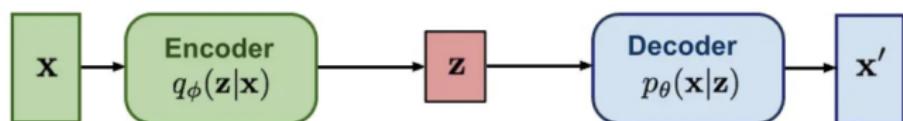
- **Quality of samples** : the generated data should be realistic and accurate compared to the actual data distribution
- **Speed of sampling** : the number of network passes required to generate a new sample
- **Mode coverage** : the model have to capture the entire distribution to ensure sample diversity

Deep Generative Models

GAN: Adversarial training

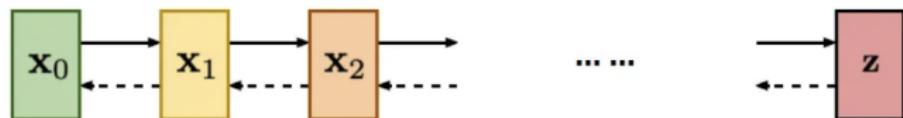


VAE: maximize variational lower bound



Diffusion models:

Gradually add Gaussian noise and then reverse



Autoencoders: background

Unsupervised approach for learning a **lower-dimensional** feature representation from unlabeled training data



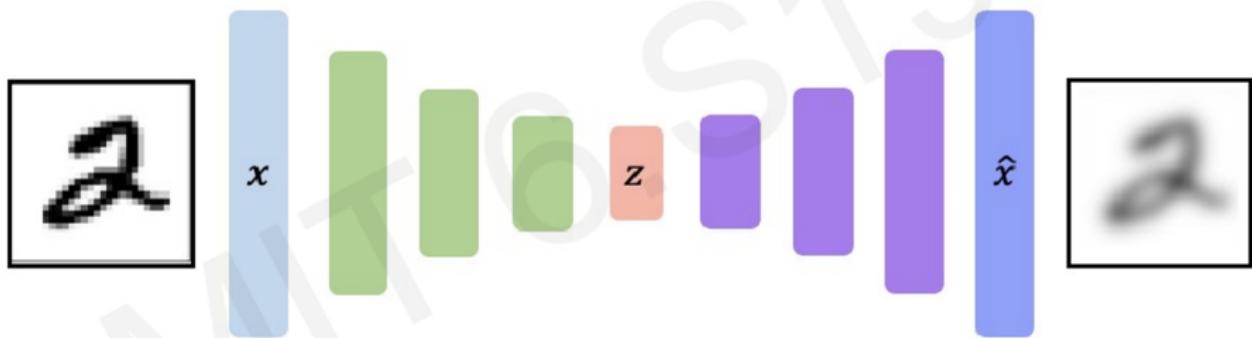
Why do we care about a low-dimensional z ?



"Encoder" learns mapping from the data, x , to a low-dimensional latent space, z

How can we learn this latent space?

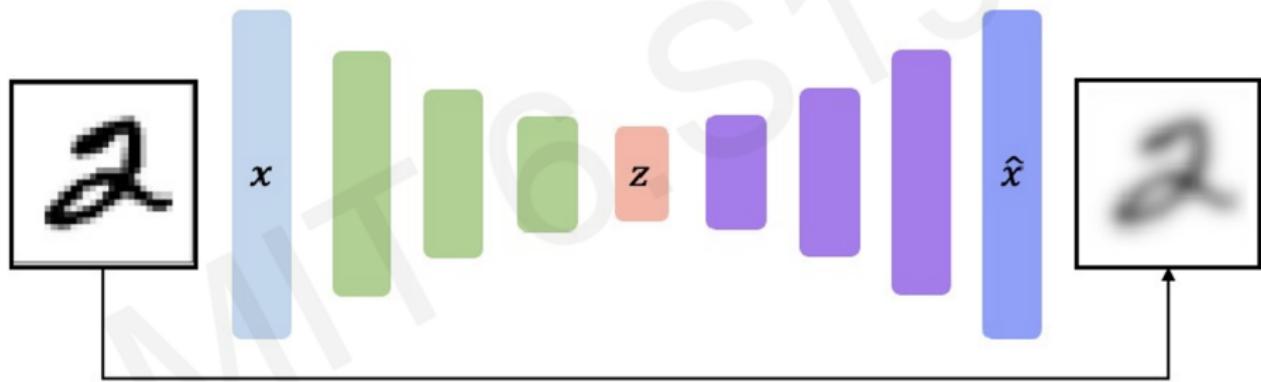
Train the model to use these features to **reconstruct the original data**



"Decoder" learns mapping back from latent space, z ,
to a reconstructed observation, \hat{x}

How can we learn this latent space?

Train the model to use these features to **reconstruct the original data**



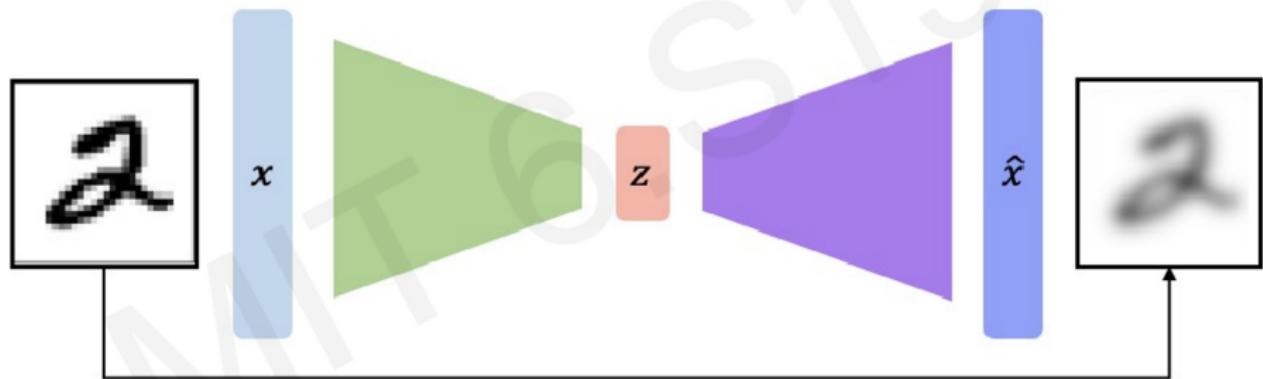
$$\mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2$$

Loss function doesn't
use any labels!!

Autoencoders: background

How can we learn this latent space?

Train the model to use these features to **reconstruct the original data**



$$\mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2$$

Loss function doesn't
use any labels!!

Dimension of latent space

Autoencoding is a form of compression!
Smaller latent space will force a larger training bottleneck

2D latent space

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

5D latent space

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

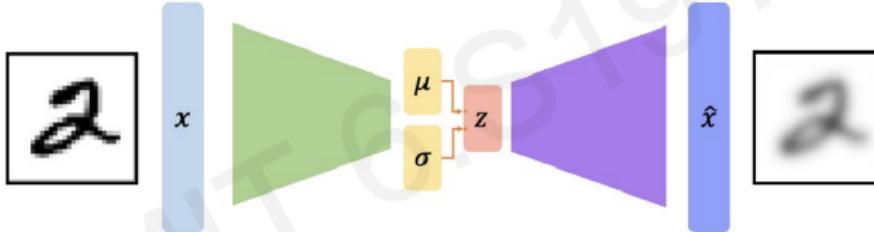
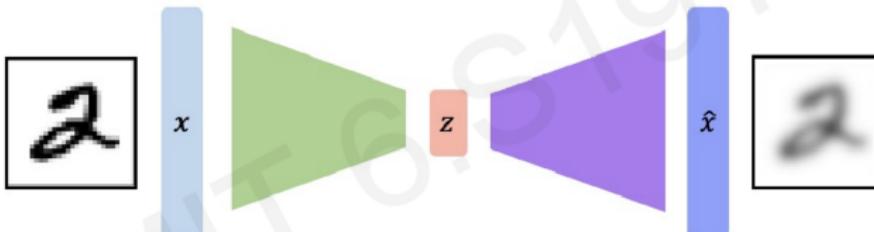
Ground Truth

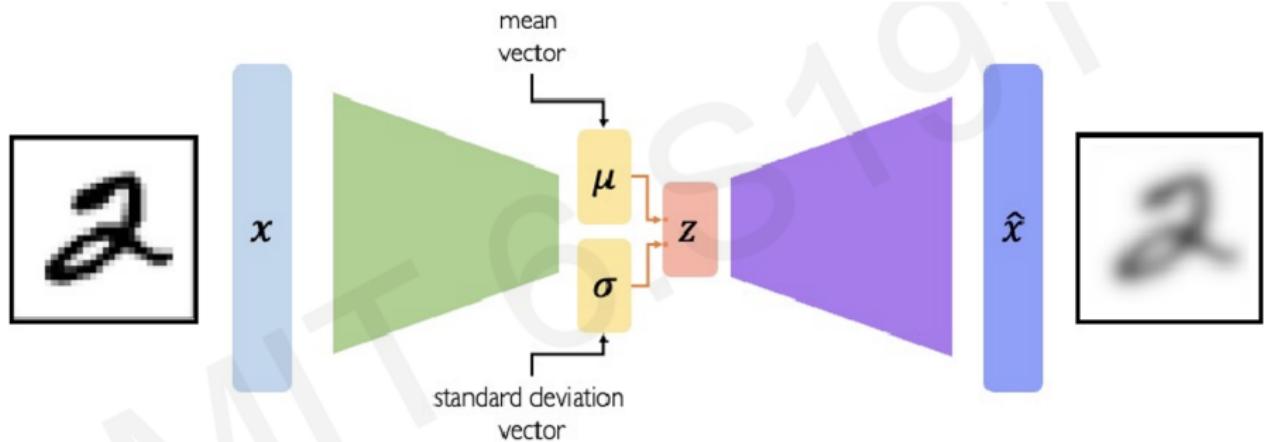
7	2	1	0	4	1	4	9	5	9
0	6	9	0	1	5	9	7	8	4
9	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	6	4	3	0
7	0	2	9	1	7	3	2	9	7
9	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9

- **Bottleneck hidden layer** forces network to learn a compressed latent space representation
- **Reconstruction loss** forces the latent representation to capture (or encode) as much *information* about the data as possible
- **Autoencoding** = Automatically **encoding** data; *Auto* = **self**-encoding

Variational Autoencoders (VAEs)

key difference between traditional autoencoders and variational autoencoders

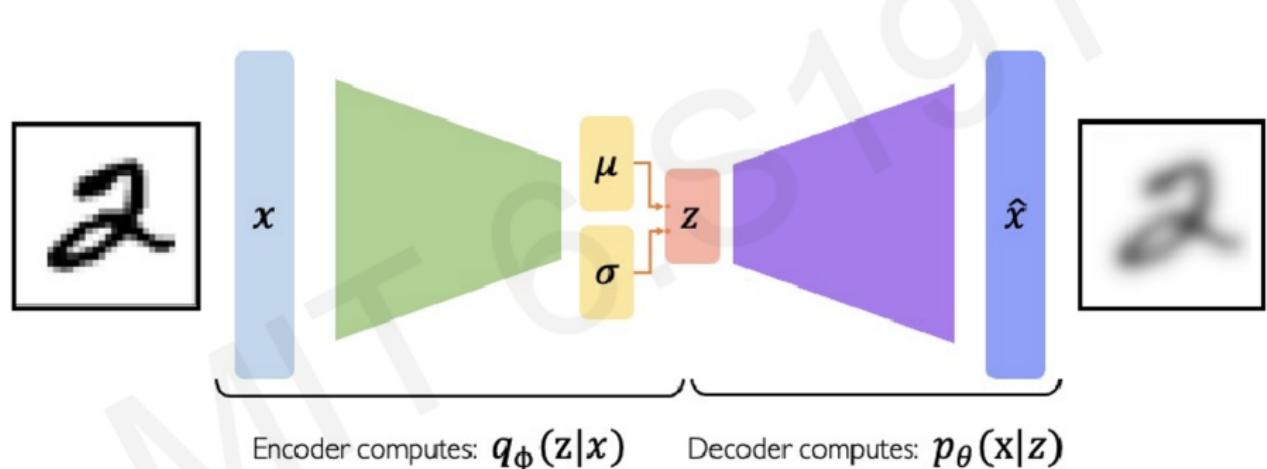




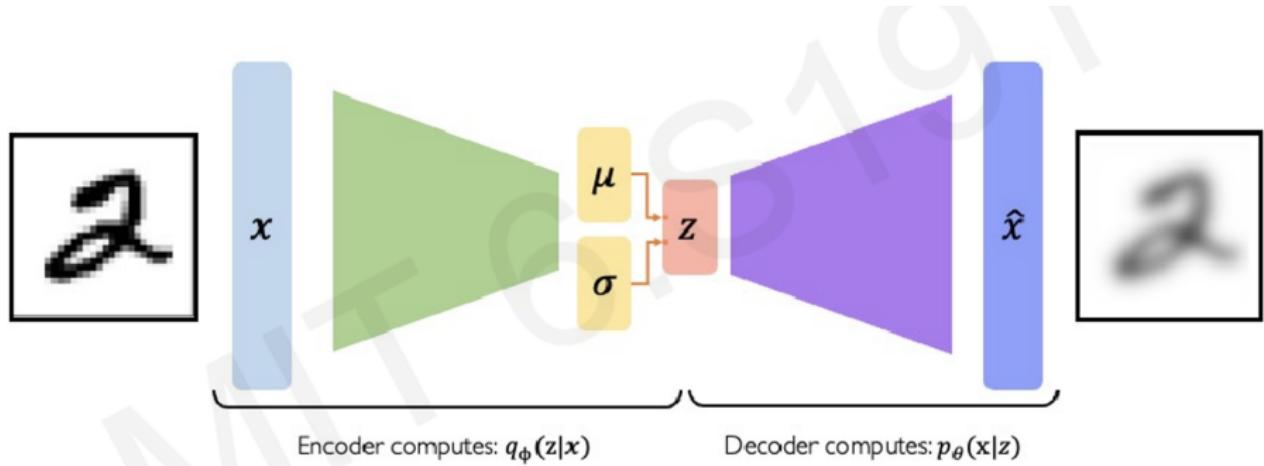
Variational autoencoders are a probabilistic twist on autoencoders!

Sample from the mean and standard deviation to compute latent sample

VAE Optimization

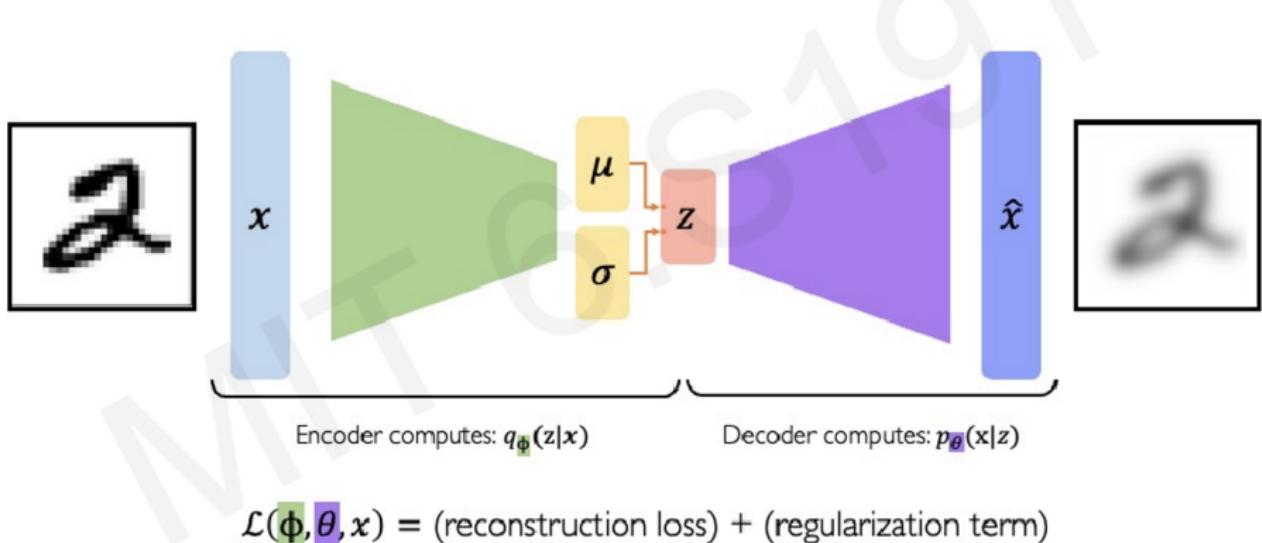


VAE Optimization

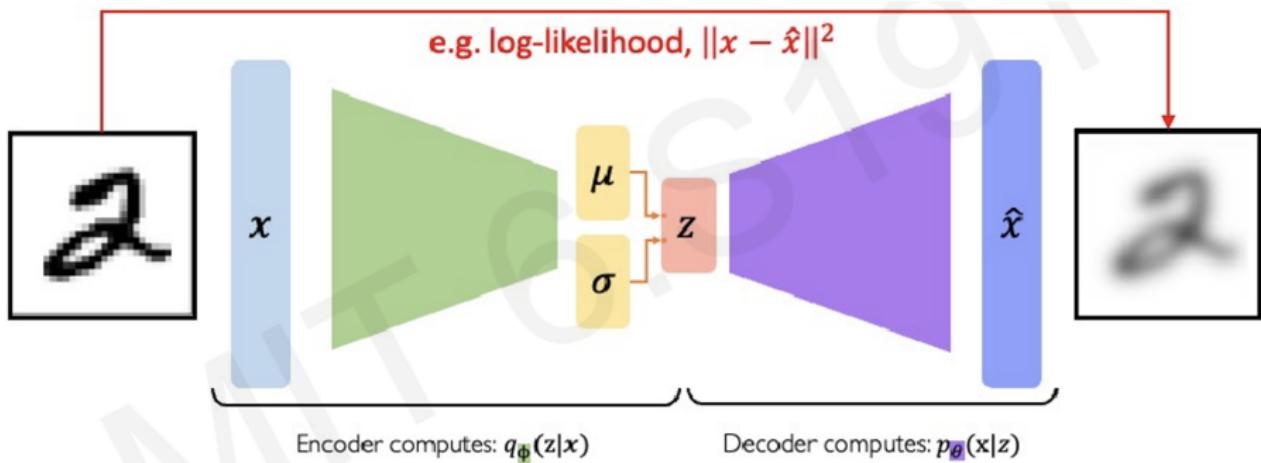


$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

VAE Optimization

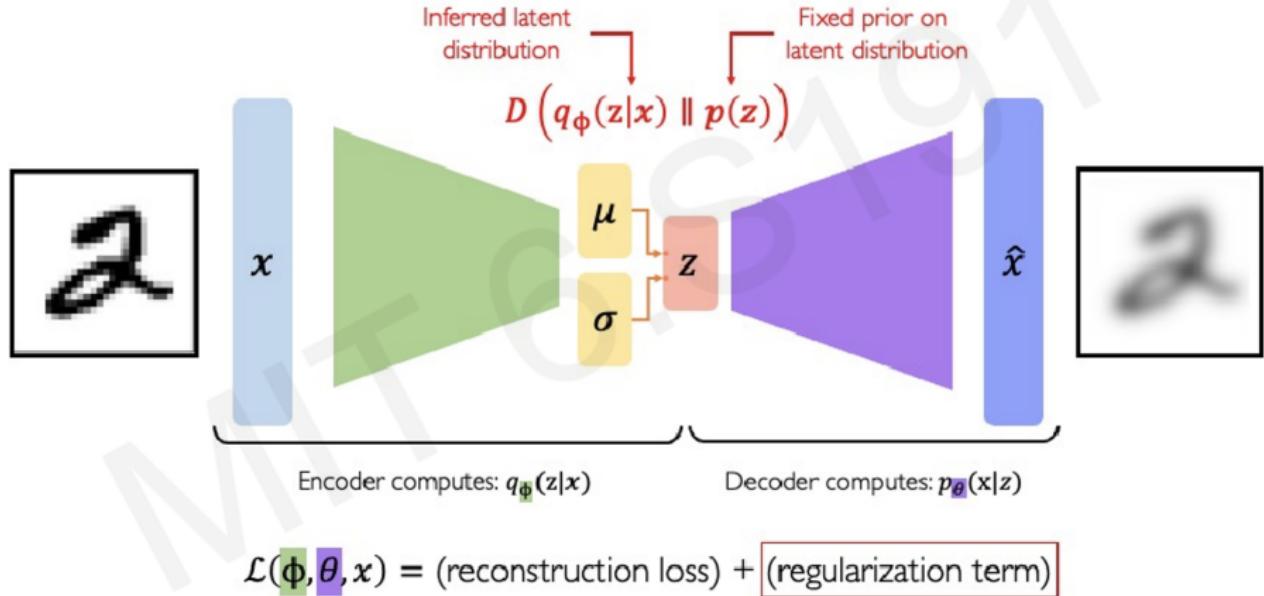


VAE Optimization



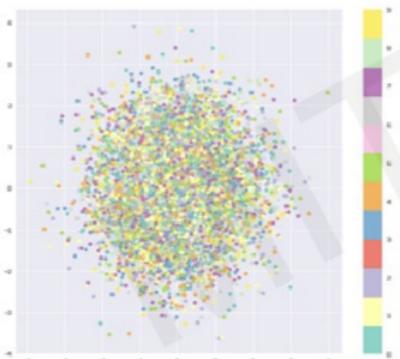
$$\mathcal{L}(\Phi, \theta, x) = \text{(reconstruction loss)} + \text{(regularization term)}$$

VAE Optimization



$$D \left(q_{\phi}(z|x) \parallel p(z) \right)$$

Inferred latent distribution Fixed prior on latent distribution



Common choice of prior – Normal Gaussian:

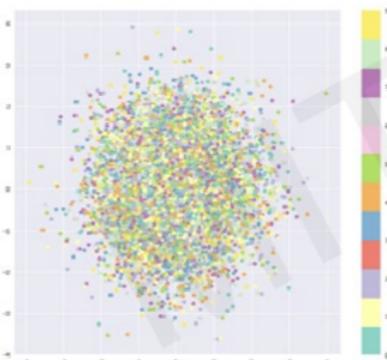
$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

- Encourages encodings to distribute encodings evenly around the center of the latent space
- Penalize the network when it tries to "cheat" by clustering points in specific regions (i.e., by memorizing the data)

$$D(q_{\phi}(z|x) \parallel p(z))$$

$$= -\frac{1}{2} \sum_{j=0}^{k-1} (\sigma_j + \mu_j^2 - 1 - \log \sigma_j)$$

KL-divergence
between the two
distributions



Common choice of prior – Normal Gaussian:

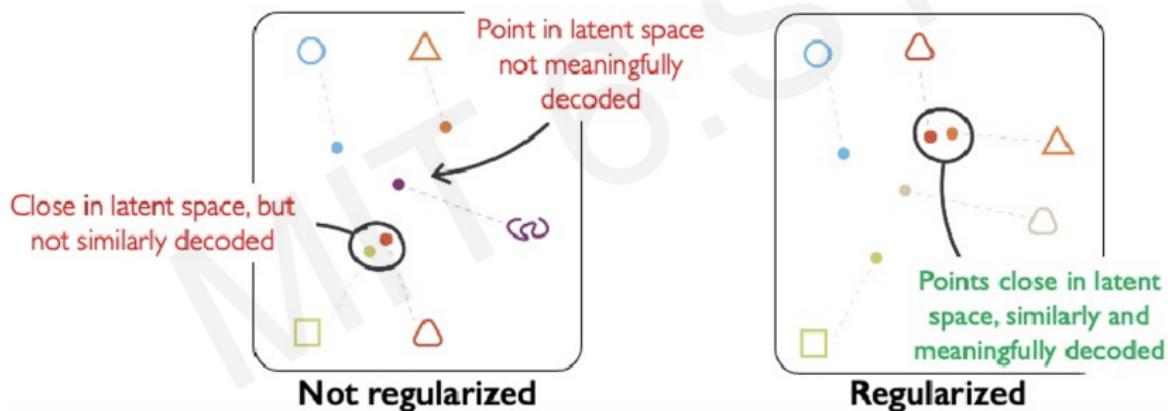
$$p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$$

- Encourages encodings to distribute encodings evenly around the center of the latent space
- Penalize the network when it tries to "cheat" by clustering points in specific regions (i.e., by memorizing the data)

What properties do we want to achieve from regularization?



- Continuity:** points that are close in latent space → similar content after decoding
- Completeness:** sampling from latent space → “meaningful” content after decoding



Intuition on regularization and Normal prior

Inria

1. **Continuity:** points that are close in latent space → similar content after decoding
2. **Completeness:** sampling from latent space → “meaningful” content after decoding

Encoding as a distribution does not
guarantee these properties!

Small variances →
Pointed distributions

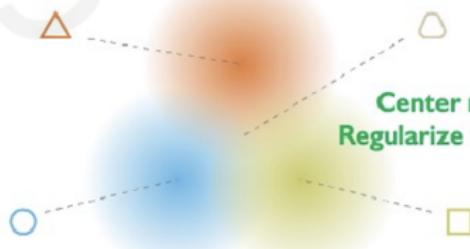
Different means →
Discontinuities



Not regularized

Normal prior →
continuity + completeness

Center means
Regularize variances



Regularized

Intuition on regularization and Normal prior

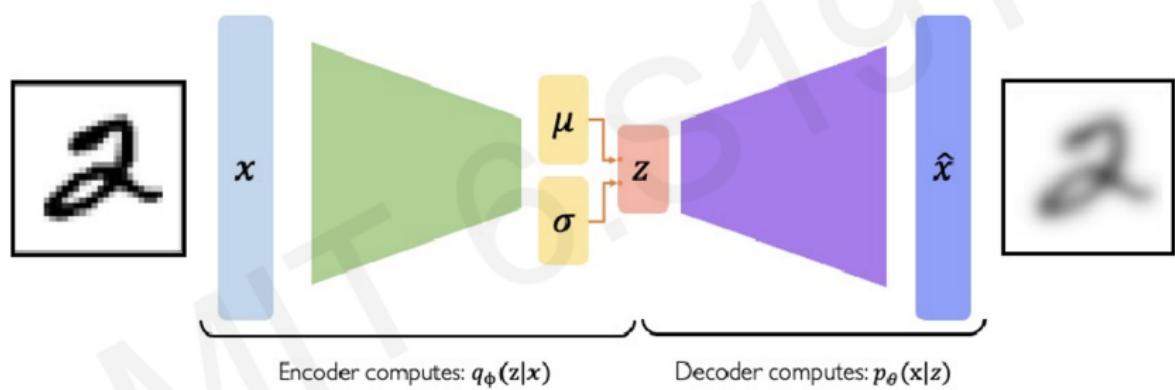
Inria

1. **Continuity:** points that are close in latent space \rightarrow similar content after decoding
2. **Completeness:** sampling from latent space \rightarrow "meaningful" content after decoding



Regularization with Normal prior helps enforce **information gradient** in the latent space.

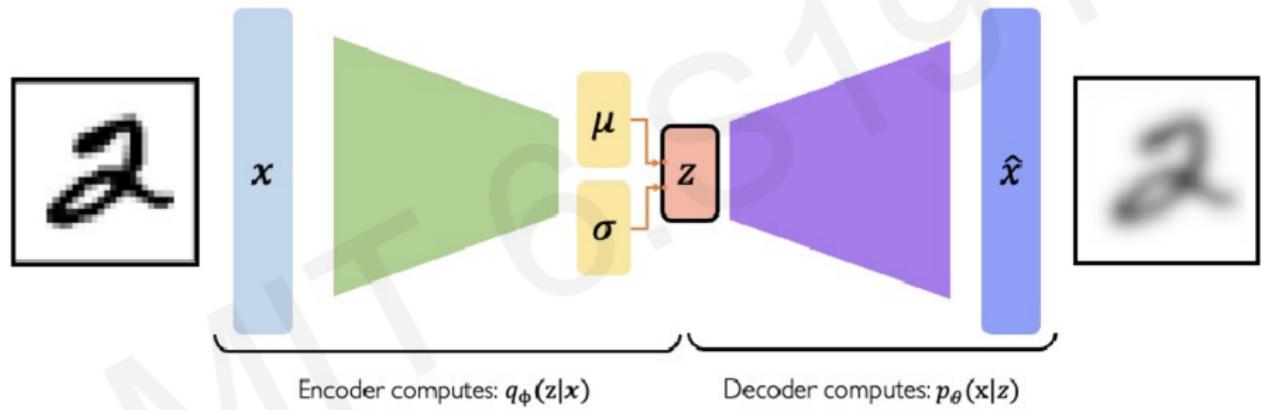
VAE computation graph



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

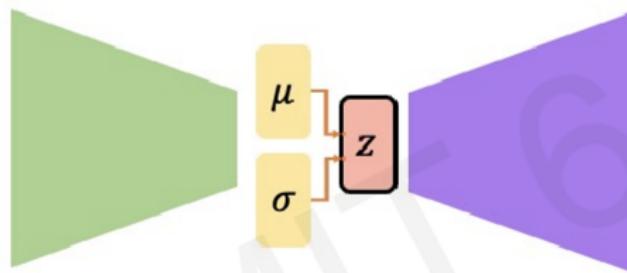
VAE computation graph

Problem: We cannot backpropagate gradients through sampling layers!



$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term})$$

Reparametrizing the sampling layer



Key Idea:

$$z \sim \mathcal{N}(\mu, \sigma^2)$$

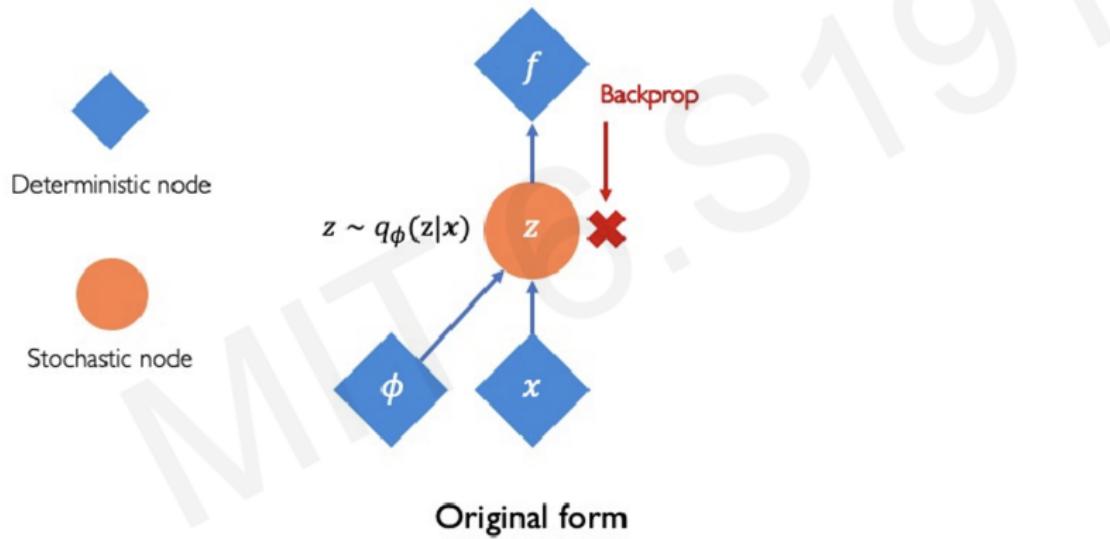
Consider the sampled latent vector z as a sum of

- a fixed μ vector;
- and fixed σ vector, scaled by random constants drawn from the prior distribution

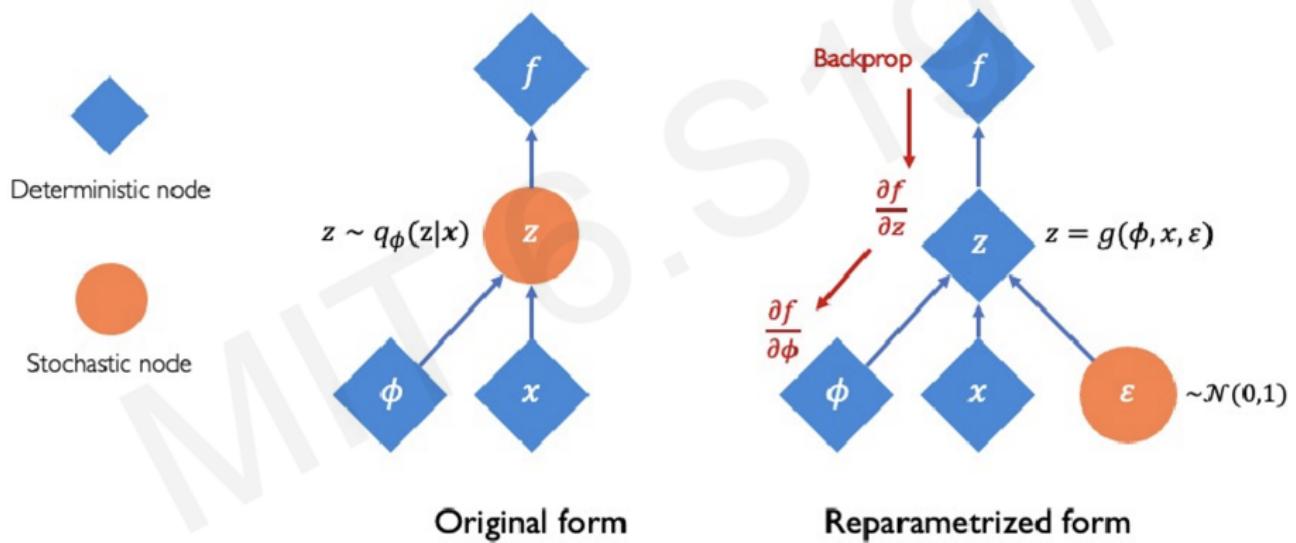
$$\Rightarrow z = \mu + \sigma \odot \varepsilon$$

where $\varepsilon \sim \mathcal{N}(0, 1)$

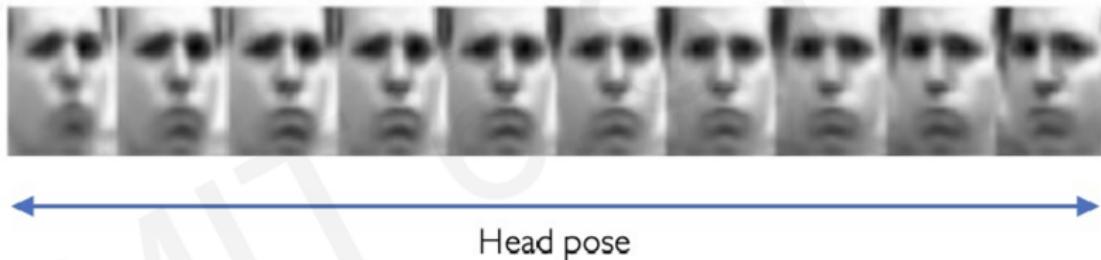
Reparametrizing the sampling layer



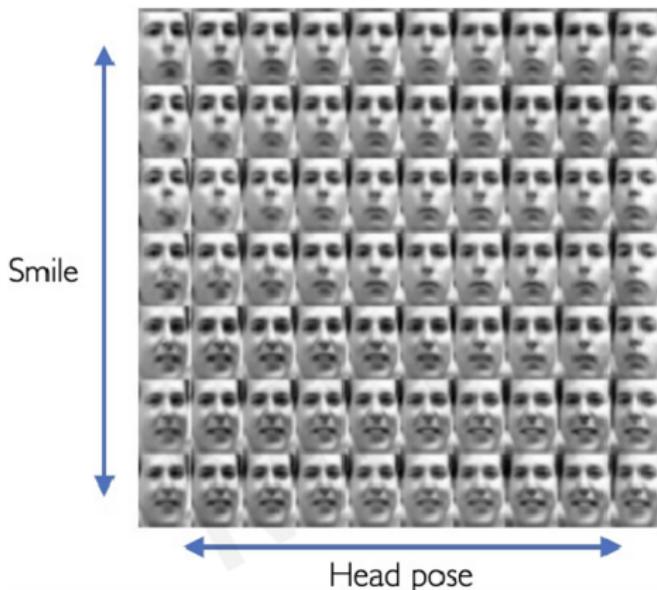
Reparametrizing the sampling layer



Slowly increase or decrease a **single latent variable**
Keep all other variables fixed



Different dimensions of z encodes **different interpretable latent features**



Ideally, we want latent variables that are uncorrelated with each other

Enforce diagonal prior on the latent variables to encourage independence

Disentanglement

Latent space distenglement with β – VAEs

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$$

Reconstruction term

Regularization term

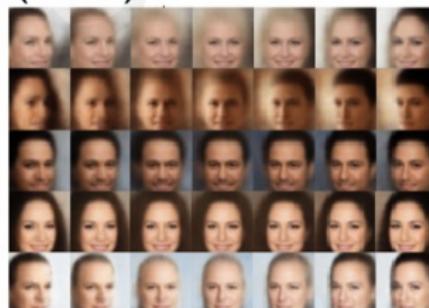
$\beta > 1$: constrain latent bottleneck, encourage efficient latent encoding \rightarrow disentanglement

Head rotation (azimuth)

Smile also changing!



Standard VAE ($\beta = 1$)



β -VAE ($\beta = 250$)

Smile relatively constant!

Why latent variable models? Debiasing

Capable of uncovering **underlying features** in a dataset



Homogeneous skin color, pose

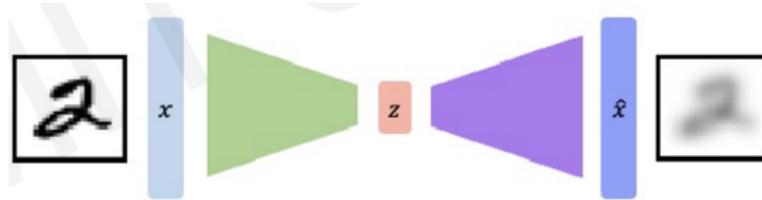
vs



Diverse skin color, pose, illumination

How can we use this information to create fair and representative datasets?

- ➊ Compress representation of world to something we can use to learn
- ➋ Reconstruction allows for unsupervised learning (no labels!)
- ➌ Reparameterization trick to train end-to-end
- ➍ Interpret hidden latent variables using perturbation
- ➎ Generation of new samples



Trilemma for VAE

- **Low-fidelity samples** : due to the overlap of several distributions in the latent space
- **High diversity samples** : Likelihood maximization forces to cover all modes
- **Fast Sampling** : require just one pass of the decoder network to turn a latent representations into a realistic data sample

