

Lecture Notes for **Neural Networks and Machine Learning**



Generative Networks
and
Auto-Encoding Generators



Logistics and Agenda

- Logistics
 - Lab dates pushed back (see schedule)
 - Next Week: Student paper presentation
- Agenda
 - A historical perspective of generative Neural Networks
 - Variational Auto-Encoding
 - VAE in Keras Demo (if time)
 - Adversarial Auto-Encoders (if time)



Last Time



State of the Art in Audio Transfer

- FAIR results are compelling...



$$\hat{\mathbf{R}} = (1 - \alpha) \left(\mathbf{I} - \alpha \mathbf{D}^{-\frac{1}{2}} \mathbf{W} \mathbf{D}^{-\frac{1}{2}} \right)^{-1} \mathbf{Y}$$

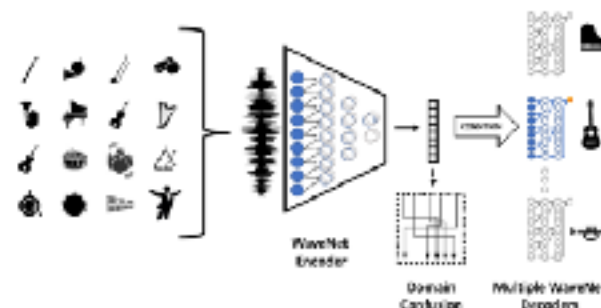
Labels for the equation components: \mathbf{I} (Identity), $\mathbf{D}^{-\frac{1}{2}}$ (NPNP), \mathbf{W} (NPNP), $\mathbf{D}^{-\frac{1}{2}}$ (NPNP), \mathbf{Y} (NPNP).

Diagram illustrating the Laplacian of graph \mathbf{D}_{ii} and the weight matrix \mathbf{W}_{ij} over each row.

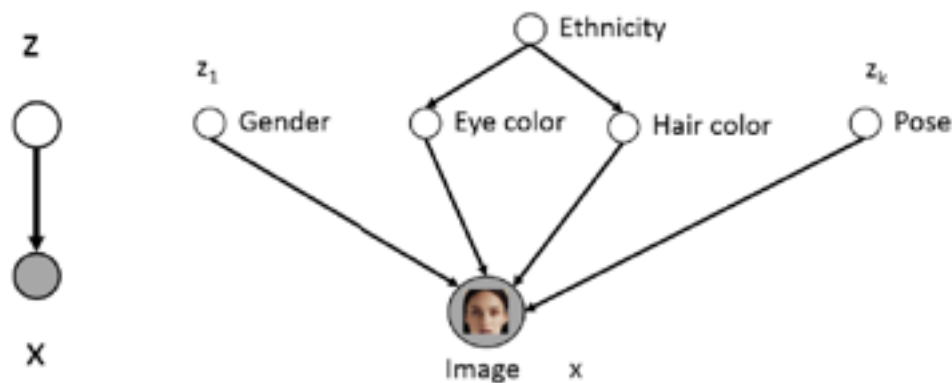
Labels for the diagram: \mathbf{D}_{ii} (Laplacian of graph), \mathbf{W}_{ij} (weight matrix).

- \mathbf{D} is diagonal and easily invertible
- \mathbf{W} is sparse and efficiently inverted after multiplications
- \mathbf{Y} is the stylized image pixels on a diagonal matrix
- \mathbf{R} can be converted to an image by returning the diagonal

- WaveNet is an autoencoder for speech and music, capable of capturing many aspects of music from time domain samples
- FAIR Paper: Train single encoder, multiple decoders

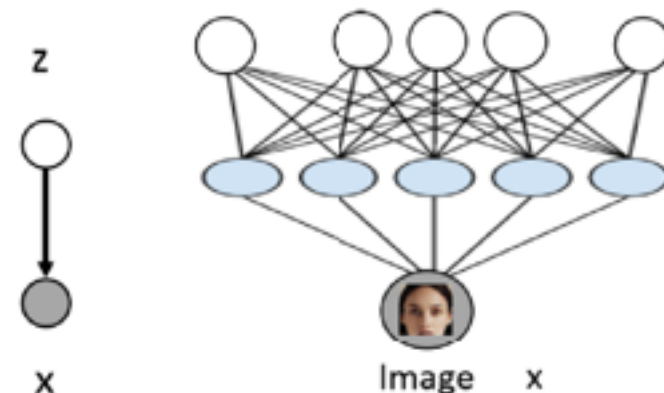


Motivations: Generative Latent Variables



$$p(\mathbf{x} | \mathbf{z})$$

Hard: \mathbf{z} is expertly chosen



$$p(\mathbf{x} | \mathbf{z})$$

Not as Hard: \mathbf{z} is trained,
latent variables are uncontrolled

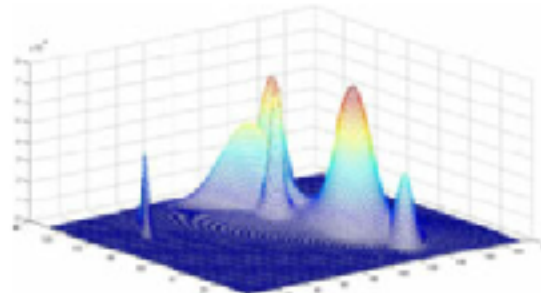
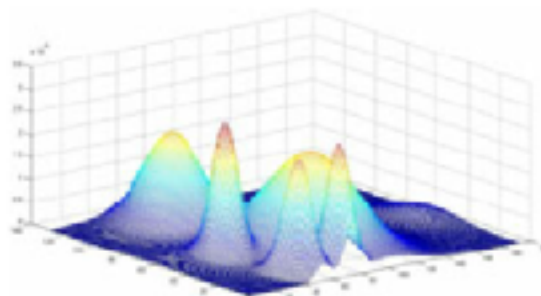
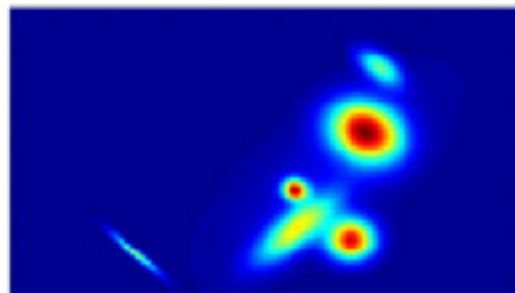
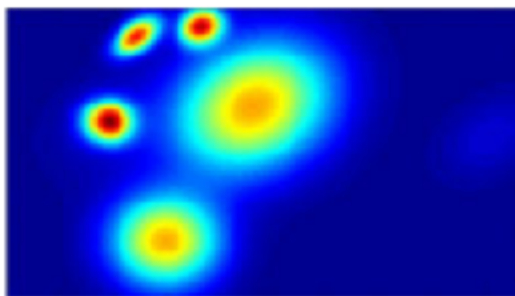
Want:
$$p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$$



Motivation: Mixtures for Simplicity

Want:
$$p(\mathbf{x}) \approx \sum_{\mathbf{z}} p(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z})$$

- Each latent variable is mostly independent of other latent variables
- The sum of various mixtures can approximate most any distribution
- Good choice for conditional is Normal Distribution
- Can parameterize $p(\mathbf{x} | \mathbf{z})$ to be a Neural Network



$$p_{\theta}(\mathbf{x} | \mathbf{z} = k) = \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

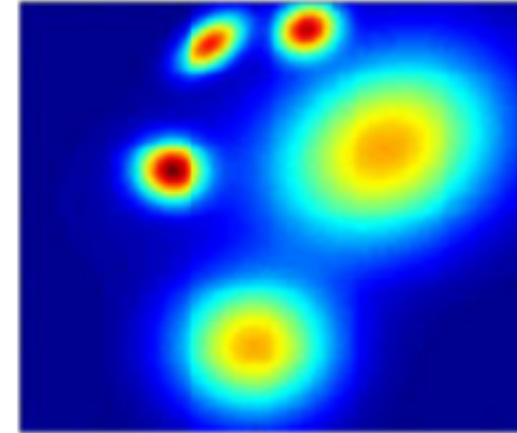
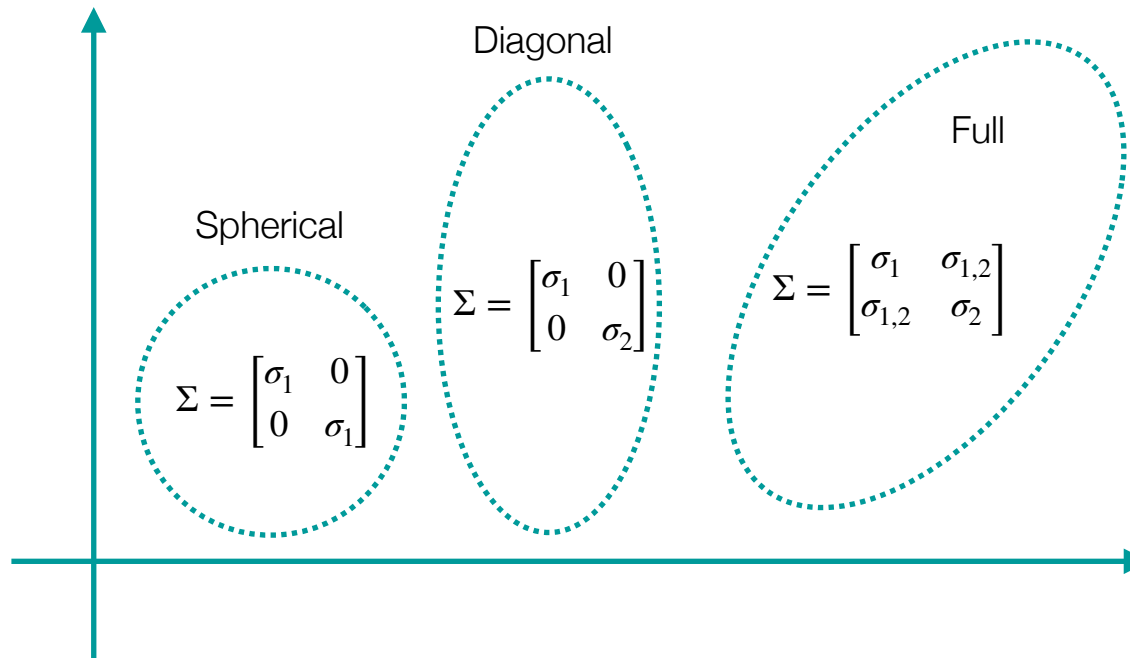
mean and covariance learned



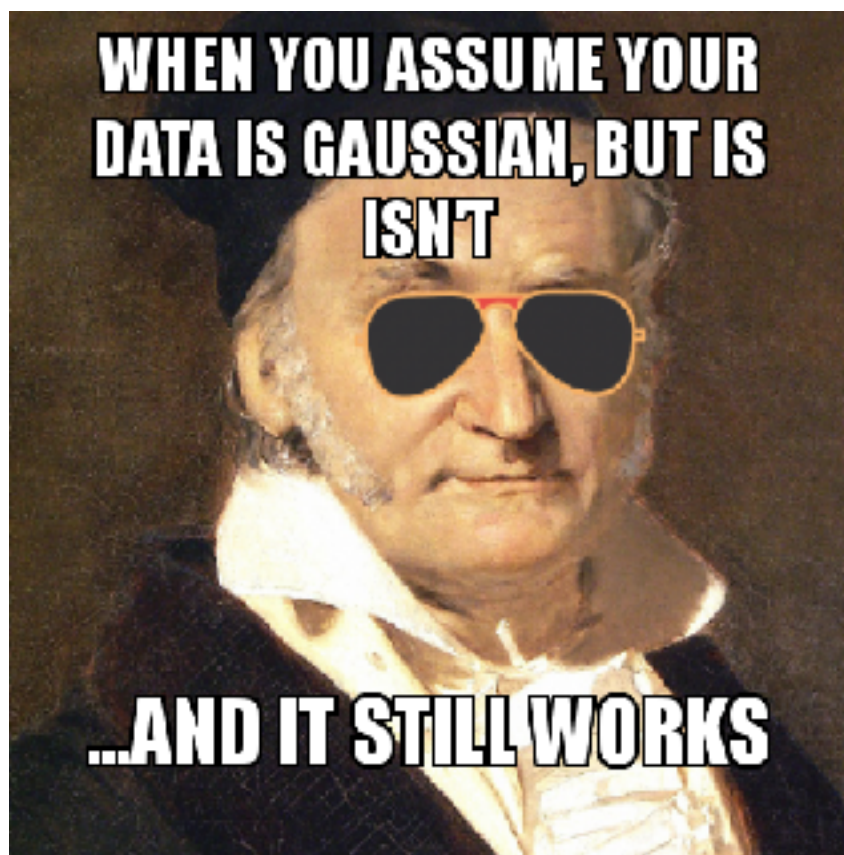
Motivation: Mixtures for Simplicity

$$= \mathcal{N}(\mathbf{x}, \mu_k, \Sigma_k)$$

mean and covariance learned



A History of Generative Networks



Aside: Notation

$$\overset{\text{some function}}{\mathbf{E}_{s \leftarrow q(s|x)}}[f(\cdot)] = \int q(s|x) \cdot f(x) dx \approx \overset{\text{could be neural networks}}{\sum_{\forall i} q(s|x^{(i)}) \cdot f(x^{(i)})}$$

Expected value of f under conditional distribution, q
 s is latent variable, $x^{(i)}$ is an observation

$$\mathbf{E}_{s \leftarrow q(s|x)}[\log f(\cdot)] = \sum_{\forall i} q(s|x^{(i)}) \cdot \log(f(x^{(i)}))$$

If function is a probability, this is just the negative of cross entropy of distributions:

$$H(q, p) = - \sum_x q(x) \cdot \log(p(x))$$

Recall that KL divergence is a measure of difference in two distribution, and is just:

$$D(p||q) = \sum_x p(x) \cdot \log \left(\frac{p(x)}{q(x)} \right) = \mathbf{E}_p \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$



Taxonomy of Generative Models

Taxonomy of Generative Models

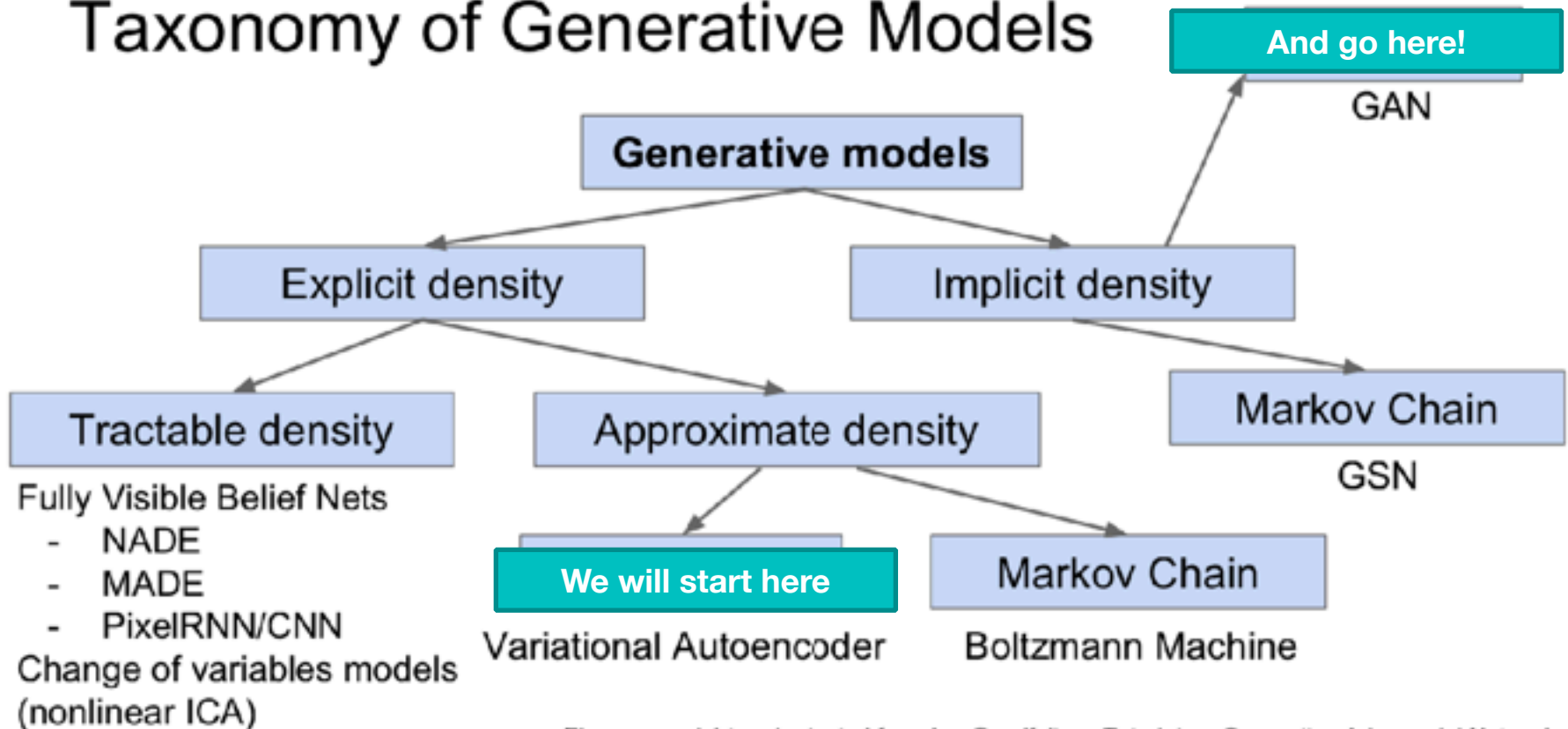


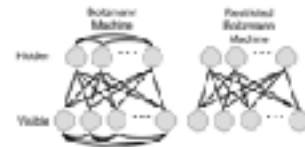
Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.



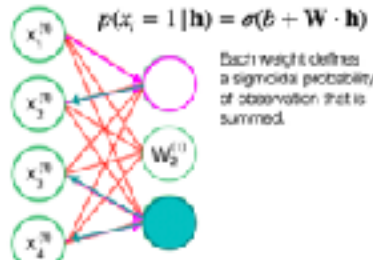
Abridged History of Generative Networks

• Restricted Boltzmann Machine

- Forward pass (visible to latent)
- Backward pass (latent to visible)
- Similar to an auto-encoder



AUTOENCODERS



RBM

<https://www.coursera.org/lecture/deep-learning-machine-learning/>

2006 Restricted Boltzmann Machine

• Deep Boltzmann Machine

$$P(v, h^{(1)}, h^{(2)}) = \frac{1}{2^N} \exp(-\sum_i v_i \lambda_i^{(1)} - \sum_{i,j} h_j^{(1)} W_{ij}^{(1)} - \sum_{i,j} h_j^{(2)} W_{ij}^{(2)} - \sum_i h_i^{(2)} \lambda_i^{(2)}). \quad (20.24)$$

To simplify our presentation, we omit the bias parameters below. The DBM energy function is then defined as follows:

$$E(v, h^{(1)}, h^{(2)}) = -v^T \lambda^{(1)} - h^{(1)T} W^{(1)} h^{(2)} - h^{(2)T} \lambda^{(2)}. \quad (20.25)$$

We now develop the mean field approach for the example with two hidden layers. Let $Q(h^{(1)}, h^{(2)} | v)$ be the approximation of $P(h^{(1)}, h^{(2)} | v)$. The mean field assumption implies that

$$Q(h^{(1)}, h^{(2)} | v) = \prod_i Q(h_i^{(1)} | v) \prod_j Q(h_j^{(2)} | v). \quad (20.26)$$

Not tractable: Can only optimize the Evidence lower bound, ELBO

One can consider off many ways of measuring how well $Q(h | v)$ fits $P(h | v)$. The mean field approach is to minimize:

$$\text{KL}(Q||P) = \sum_i Q(h_i^{(1)} | v) \log \left(\frac{Q(h_i^{(1)} | v)}{P(h_i^{(1)} | v)} \right) + \sum_j Q(h_j^{(2)} | v) \log \left(\frac{Q(h_j^{(2)} | v)}{P(h_j^{(2)} | v)} \right). \quad (20.27)$$

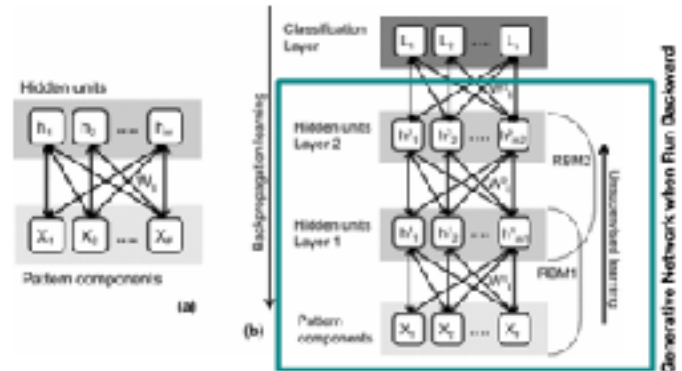
Approximate via MCMC via Gibbs Sampling

Goodfellow, Bengio, Courville: Deep learning, MIT press, 2016.

2009 Deep Boltzmann Machine Goodfellow, Bengio, Courville

• Deep Belief Network

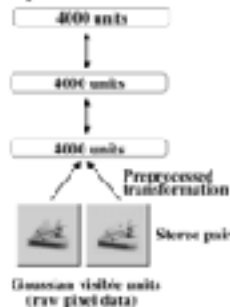
- Many RBM blocks together!



https://paperswithcode.com/abstract/10/ICCV/978-3-319-63132-7_23

2007, Deep Belief Networks RBMs with many layers

Deep Boltzmann Machine



Training Samples



Generated Samples



Figure 5: Left: The architecture of deep Boltzmann machine used for NCRB. Right: Random samples from the training set, and samples generated from the deep Boltzmann machines by running the Gibbs sampler for 10,000 steps.

2009, Practical Examples Salakhutdinov and Hinton



Contemporary Modeling

- DBNs and DBMs did not become very popular
 - Mathematics detracts from popular understanding
 - Often methods using sampling are not scalable
 - Cannot directly use Gradients (no Back Prop) 😓
- Popular method for calculating generative networks with Evidence Lower Bound (ELBO) approximation:
 - Variational Auto Encoding
 - ◆ No guarantees about global minimum
 - ◆ But scalable and will converge in finite time



Variational Auto Encoding

**“Mathematics is the
Khaleesi of sciences.”**

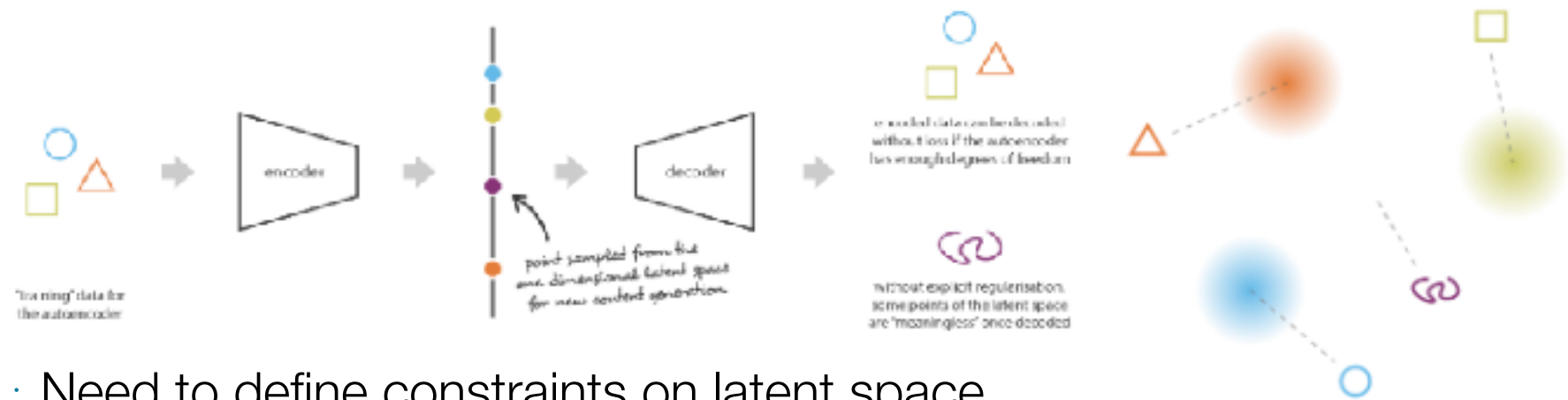
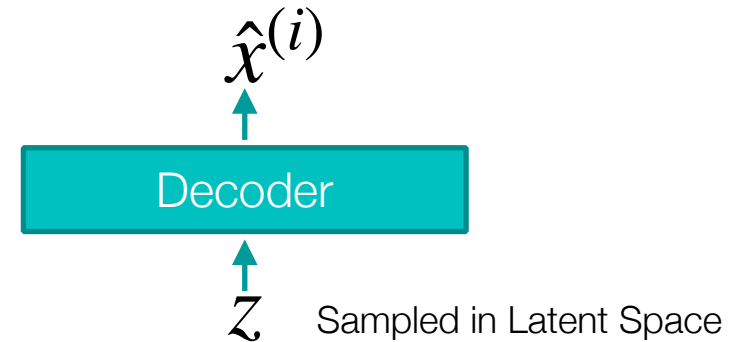
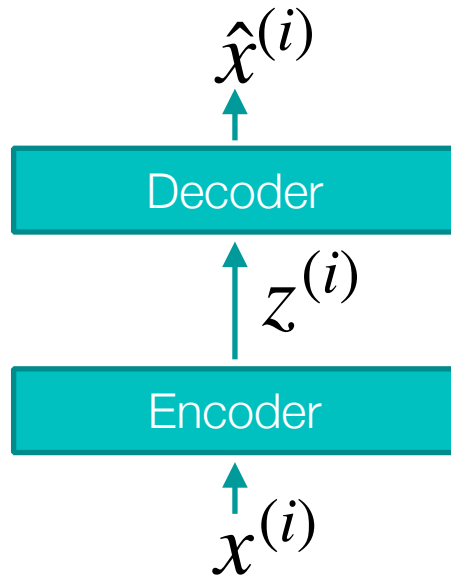


– Khal Friedrich Gauss



Can Auto Encoding Generate Samples?

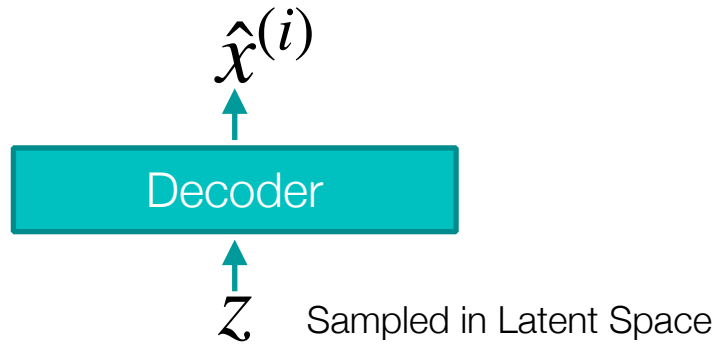
Once trained, is it possible to generate data?



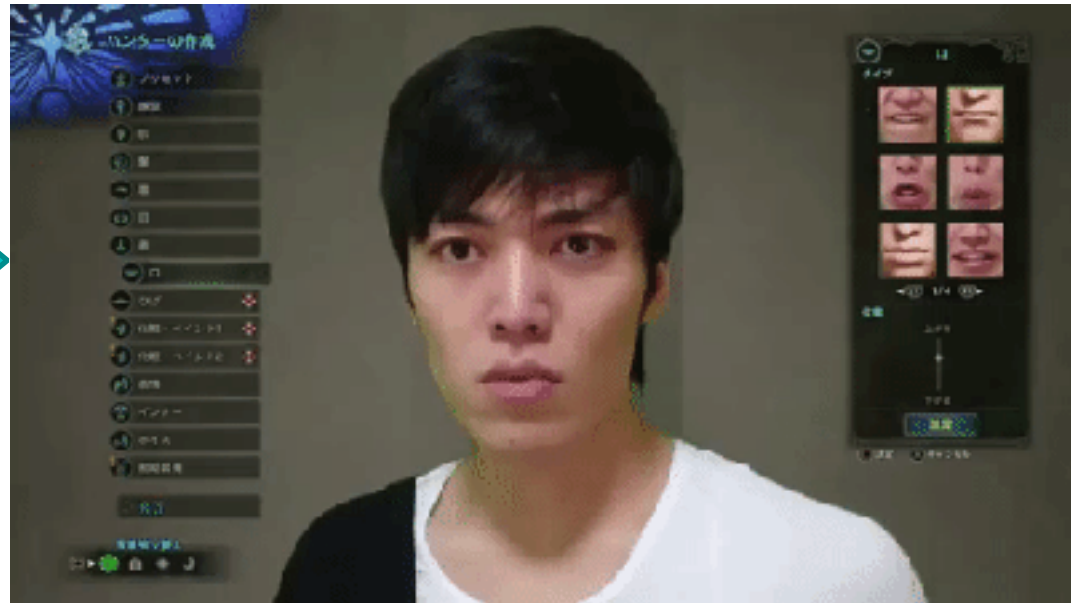
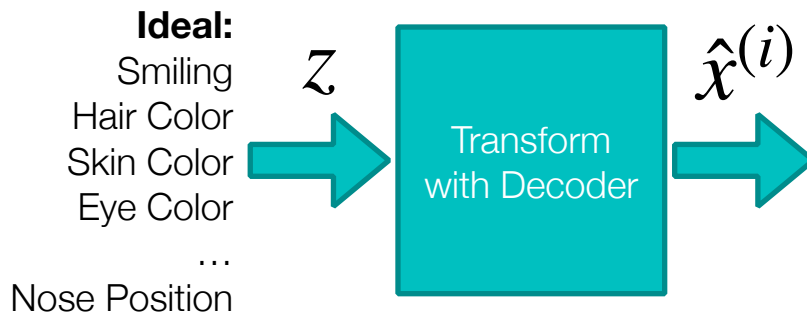
· Need to define constraints on latent space...



Reasonable constraints for $p(z)$?



- Should be simple, easy to sample from: **Normal**
- Each component should be i.i.d.:
Diag. Covariance
 - Encourages features that may be semantic, like expert might select



Optimizing

$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

We need this inference in order to compute latent variable

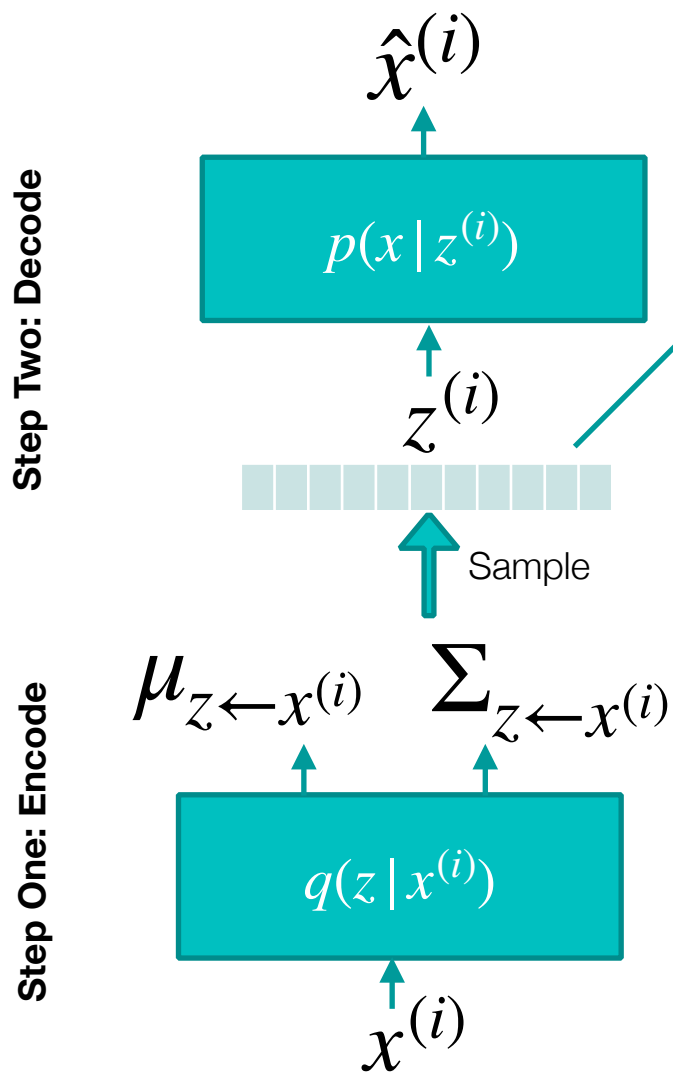
$$p(x) = \int p(x | z)p(z)dz$$

Denominator is of this form

- We can't compute! **Intractable computation** for all “ z ”
- So let's define this with **variational inference**:
 - AKA: Find the best approximation of desired distribution using a parametrized set of distributions (usually normal distributions)
 - Only needs to work for z **with observed** $x^{(i)}$
 - 1. **Encode** observed $x^{(i)}$ as Gaussian distribution via network $q(z | x^{(i)})$,
 - 2. Use $q(z | x^{(i)})$ to sample z appropriately, then **decode** with another neural network, $p(x^{(i)} | z^{(i)})$
 - 3. Make $q(z | x^{(i)})$ largest probability possible via Gaussian Distributions



Need a new formulation



Step Three: Make conditional p and q Similar

$$D_{KL} [q(z | x^{(i)}) || p(z | x^{(i)})] = \mathbf{E}_{q(z|x)} \left[\log \left(\frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right) \right]$$

Step Four: Use Variational Inference

Assume that a family of distributions can maximize likelihood of observing $x^{(i)}$:

$$\log p(x^{(i)}) \approx \mathbf{E}_{z \leftarrow q(z|x^{(i)})} [\log p(x^{(i)})]$$

Max Log Lik: maximize probability of observed $x^{(i)}$
given family of distributions q
hope this is a good approximation

Output of network, q , are the mean and covariance for sampling a variable z



Need a new formulation

$$\log p(x^{(i)}) \approx \mathbf{E}_{z \leftarrow q(z|x)} [\log p(x^{(i)})] \quad \text{Maximize!}$$

$$= \mathbf{E}_q \left[\log \frac{p(x^{(i)} | z) p(z)}{p(z | x^{(i)})} \frac{q(z | x^{(i)})}{q(z | x^{(i)})} \right] \quad \begin{array}{l} \text{Variational + multiply by one} \\ p(z | x^{(i)}) \text{ this is still a problem} \end{array}$$

$$= \mathbf{E}_q [\log p(x^{(i)} | z)] + \mathbf{E}_q \left[\log \frac{p(z)}{q(z | x^{(i)})} \right] + \mathbf{E}_q \left[\log \frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right]$$

$$= \mathbf{E}_q [\log p(x^{(i)} | z)] - \mathbf{E}_q \left[\log \frac{q(z | x^{(i)})}{p(z)} \right] + \mathbf{E}_q \left[\log \frac{q(z | x^{(i)})}{p(z | x^{(i)})} \right]$$

$$= \mathbf{E}_q [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)] + D_{KL} [q(z | x^{(i)}) || p(z | x^{(i)})]$$

always non-negative

$$\log p(x^{(i)}) \geq \mathbf{E}_q [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)] \quad \text{Will Maximize Lower Bound}$$

Can we motivate this in a different way?



The Loss Function

Maximize through
Error of Reconstruction
Same as minimizing cross entropy

want $p(z)$ to be $\mathcal{N}(\mu = 0, \Sigma = I)$
because it makes nice latent space
 $q(z|x^{(i)}) \rightarrow (\mu_{z|x}, \Sigma_{z|x}) \quad p(z) \rightarrow \mathcal{N}(0, 1)$

$$\begin{aligned}
 D_{KL}((\mu, \Sigma) \parallel \mathcal{N}(0, 1)) &= \frac{1}{2} \left(\text{tr}(\Sigma) + \mu \cdot \mu^T - \underbrace{k}_{|z|} - \log(\det(\Sigma)) \right) \begin{array}{l} \text{Determinant of diagonal} \\ \text{matrix is simple.} \\ \text{Motivates diagonal} \\ \text{covariance...} \end{array} \\
 \text{Can get this by manipulating} \\
 \text{the KL for normal distribution} \\
 &= \frac{1}{2} \left(\sum_k \Sigma_{k,k} + \sum_k \mu_k^2 - \sum_k 1 - \log \left(\prod_k \Sigma_{k,k} \right) \right) \\
 &\geq \mathbf{E}_{q(z|x^{(i)})} \left[\log p(x^{(i)} | z) - D_{KL}[q(z|x^{(i)}) \parallel p(z)] \right] \\
 &= \frac{1}{2} \sum_k (\Sigma_{k,k} + \mu_k^2 - 1 - \log \Sigma_{k,k})
 \end{aligned}$$



The Covariance Output

$$\geq \mathbf{E}_{q(z|x^{(i)})} [\log p(x^{(i)} | z)] - D_{KL} [q(z | x^{(i)}) || p(z)]$$

Maximize through
Error of Reconstruction
Same as minimizing cross entropy

want $p(z)$ to be $\mathcal{N}(\mu = 0, \Sigma = I)$
because it makes nice latent space
 $q(z | x^{(i)}) \rightarrow (\mu_{z|x}, \Sigma_{z|x}) \quad p(z) \rightarrow \mathcal{N}(0, 1)$

$$= \frac{1}{2} \sum_k (\Sigma_{k,k} + \mu_k^2 - 1 - \log \Sigma_{k,k})$$

raw covariance is not numerically stable because of underflow

$$\log \Sigma_{k,k} = \widehat{\Sigma_{k,k}}$$

predicted by
 $q(z | x^{(i)})$

$$= \frac{1}{2} \sum_k \left(\exp \left(\widehat{\Sigma_{k,k}} \right) + \mu_k^2 - 1 - \widehat{\Sigma_{k,k}} \right)$$

so we will have the neural network output log variance

Also, remember we assume **diagonal covariance**, so z 's are not correlated

This means covariance is only a vector of variances (the diagonal of Σ)

