

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532S): **Multimodal Learning with Vision, Language and Sound**

Lecture 12: Generative Models



Supervised Learning

Data: (x, y) x is data, y is label

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

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





This image is CC0 public domain

Cat



Supervised Learning

Data: (x, y) x is data, y is label

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

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



DOG, DOG, CAT

Object Detection

This image is CC0 public domain



Supervised Learning

Data: (x, y) x is data, y is label

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

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



GRASS, CAT, TREE, SKY

Semantic Segmentation

This image is CC0 public domain



Supervised Learning

Data: (x, y) x is data, y is label

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

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



A cat sitting on a suitcase on the floor

Image Captioning

This image is CC0 public domain



Unsupervised Learning

Data: X Just data, no labels!

Goal: Learn some underlying hidden structure of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



k-means clustering

This image is CC0 public domain



Unsupervised Learning

Data: X Just data, no labels!

Goal: Learn some underlying hidden structure of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



dimensionality reduction

This image is CC0 public domain



Unsupervised Learning

Data: X Just data, no labels!

Goal: Learn some underlying hidden structure of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-dim density estimation



2-dim density estimation

2-d density images left and right are CC0 public domain

Supervised Learning

Data: (x, y) x is data, y is label

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

Examples: Classification, regression, object detection, semantic segmentation, image captioning, *etc.*

Jasepervised Learning



Generative Models

Given training data, generate new samples from the same distribution



Training data ~ $p_{data}(\mathbf{x})$



Generated samples $\sim p_{\text{model}}(\mathbf{x})$

Want to learn $p_{model}(x)$ similar to $p_{data}(x)$



Generative Models

Given training data, generate new samples from the same distribution



Training data ~ $p_{data}(x)$

Want to learn $p_{\text{model}}(\mathbf{x})$ similar to $p_{\text{data}}(\mathbf{x})$

Addresses density estimation, a core problem in unsupervised learning

- **Explicit** density estimation: explicitly define and solve for $p_{model}(x)$
- Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it



Generated samples $\sim p_{\text{model}}(\mathbf{x})$



Taxonomy of Generative Models



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



Taxonomy of Generative Models



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



Why Generative Models?

- Realistic samples for artwork, super-resolution, colorization, etc.









Why **Generative** Models?

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for **simulation**, **predictions** and planning (reinforcement learning applications)
- Training generative models can also enable inference of latent representation that can be useful as general features
- **Dreaming** / hypothesis visualization



PixelRNN and PixelCNN

Explicit Density model

Use chain rule to decompose likelihood of an image x into product of (many) 1-d distributions



then maximize likelihood of training data

[van der Oord et al., 2016]

$$p(x_i | x_1, ..., x_{i-1})$$

Probability of i'th pixel value given all previous pixels



Explicit Density model

Use chain rule to decompose likelihood of an image x into product of (many) 1-d distributions



then maximize likelihood of training data

[van der Oord et al., 2016]

$$p(x_i|x_1,...,x_{i-1})$$

given all previous pixels

Complex distribution over pixel values, so lets model using neural network





Explicit Density model

Use chain rule to decompose likelihood of an image x into product of (many) 1-d distributions



then maximize likelihood of training data

[van der Oord et al., 2016]

$$p(x_i|x_1,...,x_{i-1})$$

Probability of i'th pixel value given all previous pixels

Complex distribution over pixel values, so lets model using neural network

Also requires defining ordering of "previous pixels"









Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

[van der Oord et al., 2016]





PixelRNN

Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

[van der Oord et al., 2016]





PixelRNN

Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

[van der Oord et al., 2016]









[van der Oord et al., 2016]



Generate image pixels starting from the corner

Dependency on previous pixels model using an RNN (LSTM)

Problem: sequential generation is slow

[van der Oord et al., 2016]





PixelCNN

Still generate image pixels starting from the corner

Dependency on previous pixels now modeled using a CNN over context region

[van der Oord et al., 2016]





PixelCNN

Still generate image pixels starting from the corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, \dots, x_{i-1})$$

[van der Oord et al., 2016]

Softmax loss at each pixel





PixeICNN

Still generate image pixels starting from the corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, \dots, x_{i-1})$$

[van der Oord et al., 2016]



Generation is still slow (sequential), but learning is faster







Generated Samples



32x32 CIFAR-10

[van der Oord et al., 2016]



32x32 ImageNet



PixelRNN and PixelCNN

Pros:

- Can explicitly compute likelihood p(x)
- Explicit likelihood of training data gives good evaluation metric
- Good samples

Con:

— Sequential generation => slow

Improving PixelCNN performance

- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...



Multi-scale PixelRNN

Take sub-sampled pixels as additional input pixels

Can capture better global information (more visually coherent)

[van der Oord et al., 2016]





Multi-scale PixelRNN



[van der Oord et al., 2016]





Conditional Image Generation

vector **h**

 $p(\mathbf{x}) = p(x_1, x_2, \dots, x_{n^2})$ $p(\mathbf{x}|\mathbf{h}) = p(x_1, x_2, ..., x_{n^2}|\mathbf{h})$



[van der Oord et al., 2016]

Similar to PixelRNN/CNN but conditioned on a high-level image description



Conditional Image Generation



African elephant



[van der Oord et al., 2016]

Sandbar



Attention RNN: Structured Spatial Attention Mechanism



Siddhesh Khandelwal



val Leonid Sigal



Motivation

Attention is widely used in vision: helps identify relevant regions of the image





Q: What color is a hydrant?



A: It is red





Motivation

Attention is widely used in vision: helps identify relevant regions of the image

Attention is **applied**:

On the output of CNN architecture



After each CNN layer



(Image taken from Seo et al. BMVC 2018)




Attention is widely used in vision: helps identify relevant regions of the image

Attention is **applied**:

On the **output of CNN** architecture



Existing attention mechanisms are either global or local

After each CNN layer



(Image taken from Seo et al. BMVC 2018)





Attention is widely used in vision: helps identify relevant regions of the image

Attention is **applied**:

On the **output of CNN** architecture



Existing attention mechanisms are either global or local



Global:

Pros:

- Can model arbitrary context

Cons:

Cannot be applied at high res.

After each CNN layer



(Image taken from Seo et al. BMVC 2018)





Attention is widely used in vision: helps identify relevant regions of the image

Attention is **applied**:

On the **output of CNN** architecture



Existing attention mechanisms are either **global** or **local**



Global:

Pros:

- Can model arbitrary context

Cons:

Cannot be applied at high res.

After each CNN layer



(Image taken from Seo et al. BMVC 2018)



Local:

Pros:

Can be applied at any resolution

Cons:

- Can only model local context





Attention is widely used in vision: helps identify relevant regions of the image

Attention is **applied**:

On the **output of CNN** architecture



Existing attention mechanisms are either global or local



Global:

Pros:

- Can model arbitrary context

Cons:

Cannot be applied at high res.

Neither can account (explicitly) for structure in the attention variables

After each CNN layer



(Image taken from Seo et al. BMVC 2018)



Local:

Pros:

Can be applied at any resolution

Cons:

Can only model local context





Novel autoregressive attention mechanism that can encode structural dependencies among attention values

- Inspired by diagonal Bi-LSTM architecture from PixelRNN
- Spatial attention values are generated sequentially
- Image is traversed diagonally from top-left to bottom-right



Each attention value depends on

- Local image context
- Previously generated attention values

Novel autoregressive attention mechanism that can encode structural dependencies among attention values

- Inspired by diagonal Bi-LSTM architecture from PixelRNN
- Spatial attention values are generated sequentially
- Image is traversed diagonally from top-left to bottom-right



Each attention value depends on

– Local image context

Previously generated attention values



Image

Each attention value depends on

- Local image context
- Previously generated attention values





Attention Mask

Each attention value depends on

- Local image context
- Previously generated attention values





Attention Mask

Each attention value depends on

- Local image context
- Previously generated attention values





Attention Mask

Each attention value depends on

- Local image context
- Previously generated attention values





Image

Attention Mask

Each attention value depends on

- Local image context
- Previously generated attention values





Attention Mask

Each attention value depends on

- Local image context
- Previously generated attention values



Block Attention RNN: Scalability

LSTMs don't work well over large sequences (say beyond 50 x 50)

We utilize symmetric down-sampling and up-sampling scheme to deal with larger resolution attention maps



Experiments: Visual Attribute Prediction

Task: Given an image, predict a color of the number specified by a query

	MREF	MDIST
CNN+SAN [23]	81.28	78.06
$CNN+\neg CTX$ [23]	94.95	90.06
CNN+CTX [23]	97.92	95.80
CNN+ARNN	99.11	97.04



(a) MREF

(b) MDIST

(c) MBG

MBG

- 52.36 Global (only top layer)
- 66.83 Local (all layers)
- 79.26 Local with context (all layers)
- **86.07** Local with context + structure (all layers)

Experiments: Visual Digit Prediction

Task: Given an image, predict a digit number specified by a query color





	SAN	¬ CTX	СТХ	ARNN
ectness	0.1258	0.2017	0.2835	0.3729

Experiments: Visual Question Answering

Task: Visual Question Answering

Accuracy
24.04
51.27
53.23

HCA — Hierarchical Co-attention

[Lu et al., NIPS 17]



Q: what is the man holding a snowboard on top of a snow covered? A: mountain



what is the man holding a snowboard on top of a snow covered



snowboard on top of a snow covered ?



what is the man holding a snowboard on top of a snow covered ?



Q: what is the color of the bird? **A**: white



Q: how many snowboarders in formation in the snow, four is sitting? A: 5



what is the color of the bird?

how many snowboarders in

formation in the snow, four is

sitting ?



what is the color of the bird?





how many snowboarders in formation in the snow, four is sitting?







how many snowboarders in formation in the snow, four is sitting?

Conclusions

Modeling structure among the latent attention variables is useful





PixelCNNs define tractable density function, optimize likelihood of training data:

n $p(x) = \prod p(x)$ i=1

$$p(x_i|x_1,...,x_{i-1})$$





$$p(x) = \prod_{i=1}^{n} p(x_i | x_1, \dots, x_{i-1})$$

VAEs define intractable density function with latent variables z (that we need to marginalize):

$$p_{\theta}(x) = \int f$$

cannot optimize directly, derive and optimize lower bound of likelihood instead

PixelCNNs define tractable density function, optimize likelihood of training data:

$p_{\theta}(z)p_{\theta}(x|z)dz$



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



Originally: Linear + nonlinearity (sigmoid) Later: Deep, fully-connected Later: ReLU CNN





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



Originally: Linear + nonlinearity (sigmoid) Later: Deep, fully-connected Later: ReLU CNN





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



Originally: Linear + nonlinearity (sigmoid) Later: Deep, fully-connected Later: ReLU CNN





Train such that features can reconstruct original data best they can







Train such that features can reconstruct original data best they can





Reconstructed data



Encoder: 4-layer conv Decoder: 4-layer upconv

Input data







Reconstructed data



Encoder: 4-layer conv Decoder: 4-layer upconv

Input data







Reconstructed data



Encoder: 4-layer conv Decoder: 4-layer upconv

Input data















Probabilistic spin on autoencoder - will let us sample from the model to generate Assume training data is generated from underlying unobserved (latent)

representation z



[Kingma and Welling, 2014]



Probabilistic spin on autoencoder - will let us sample from the model to generate Assume training data is generated from underlying unobserved (latent)

representation z



[Kingma and Welling, 2014]

Intuition: *x* is an image, *z* is latent factors used to generate x (e.g., attributes, orientation, etc.)



We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]



We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **represent** this model?



We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **represent** this model?

Choose prior p(z) to be simple, e.g., Gaussian Reasonable for latent attributes, e.g., pose, amount of smile



We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **represent** this model?

Choose prior p(z) to be simple, e.g., Gaussian Reasonable for latent attributes, e.g., pose, amount of smile

Conditional $p(\mathbf{x}|\mathbf{z})$ is complex (generates image) Represent with Neural Network





We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **train** this model?


We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **train** this model?

Remember the strategy from earlier — learn model parameters to maximize likelihood of training data $p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$

(now with latent z that we need to marginalize)



We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **train** this model?

Remember the strategy from earlier — learn model parameters to maximize likelihood of training data $p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$

(now with latent z that we need to marginalize)

What is the problem with this?



We want to estimate the true parameters θ^* of this generative model



[Kingma and Welling, 2014]

How do we **train** this model?

Remember the strategy from earlier — learn model parameters to maximize likelihood of training data $p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$

(now with latent z that we need to marginalize)

Intractable !



Data likelihood: $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$

[Kingma and Welling, 2014]





[Kingma and Welling, 2014]





[Kingma and Welling, 2014]

Decoder Neural Network

• •





[Kingma and Welling, 2014]

Decoder Neural Network

 $\int p_{ heta}(z) p_{ heta}(x|z) dz$

• •





Posterior density is also intractable: $p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$

[Kingma and Welling, 2014]

Decoder Neural Network





Posterior density is also intractable: $p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$

[Kingma and Welling, 2014]

Decoder Neural Network





Posterior density is also intractable: p

Solution: In addition to decoder network modeling $p_{\theta}(x|z)$, define additional encoder network $q_{\Phi}(z|x)$ that approximates $p_{\theta}(z|x)$ - Will see that this allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize

[Kingma and Welling, 2014]

Decoder Neural Network

$$p_{ heta}(x|z)dz$$

$$p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$$



Optional subtitle

Since we are modeling probabilistic generation of data, encoder and decoder networks are probabilistic (they model distributions)



[Kingma and Welling, 2014]



Since we are modeling probabilistic generation of data, encoder and decoder networks are probabilistic (they model distributions)



[Kingma and Welling, 2014]



Since we are modeling probabilistic generation of data, encoder and decoder networks are probabilistic (they model distributions)



[Kingma and Welling, 2014]





Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta})$$

Taking expectation with respect to z (using encoder network) will come in handy later

[Kingma and Welling, 2014]

 $(x^{(i)})$ Does not depend on z)



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] \qquad (p_{\theta})$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \qquad (Ba)$$

[Kingma and Welling, 2014]

 $(x^{(i)})$ Does not depend on z)

ayes' Rule)



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood: $(x^{(i)})$ Does not depend on z)

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta})$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (Ba)$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x)}{q_{\phi}(z \mid x)}\right]$$

[Kingma and Welling, 2014]

ayes' Rule)

 $\left| \frac{r^{(i)}}{r^{(i)}} \right|$ (Multiply by constant)



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood: $(x^{(i)})$ Does not depend on z)

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) = \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (Bay)$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} = \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] = \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] = \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] = \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] = \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] = \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] = \mathbf{E}_{z} \left[\log$$

[Kingma and Welling, 2014]

yes' Rule)

 $\left| \frac{r^{(i)}}{r^{(i)}} \right|$ (Multiply by constant) $\frac{\phi(z \mid x^{(i)})}{p_{\theta}(z)} + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood: $(x^{(i)})$ Does not depend on z)

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) = \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (Bay)$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} = \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}}{q_{\phi}(z \mid x^{(i)})}\right] = \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}))$$

Expectation with respect to z (using encoder network) leads to nice KL terms

[Kingma and Welling, 2014]

yes' Rule)

 $\left[\frac{i^{(i)}}{i^{(i)}}\right]$ (Multiply by constant) $\frac{p_{\theta}(z \mid x^{(i)})}{p_{\theta}(z)} + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$ $|x^{(i)}|| p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z \mid x^{(i)}))|$



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood: $(x^{(i)})$ Does not depend on z)

$$\begin{split} \log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] & (p_{\theta}(z)) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] & (Bay) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} \right] \\ &= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \\ &= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)})) \right] \end{split}$$

Decoder network gives $p_A(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through **reparam. trick**, see paper.)

[Kingma and Welling, 2014]

yes' Rule)

 $\left| \frac{r^{(i)}}{r^{(i)}} \right|$ (Multiply by constant) $\frac{\phi(z \mid x^{(i)})}{p_{\theta}(z)} + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$ $|x^{(i)}|| p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z \mid x^{(i)}))|$ p_A(z x) intractable (saw earlier), can't This KL term (between Gaussians) for encoder and z prior) has nice compute this KL term :(closed-form solution! But we know KL divergence always ≥ 0 .



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\begin{split} \log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) || p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + D_{KL}(q_{\phi}(z \mid x^{(i)}) || p_{\theta}(z \mid x^{(i)})) \end{split}$$

Tractable lower bound which we can take gradient of and optimize! ($p\theta(x|z)$ differentiable, KL term differentiable) [Kingma and Welling, 2014]



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

Variational lower bound ("**ELBO**")

[Kingma and Welling, 2014]

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$



Derivation of lower bound of the data likelihood

Now equipped with **encoder** and **decoder** networks, let's see (log) data likelihood:

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$
$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$
$$\text{Reconstruct} \qquad \text{Make approximate posterior}$$

Reconstruct
Input DataMake approximate posterior
close to the prior
$$= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z \mid x^{(i)}))$$
$$L(x^{(i)}, \theta, \phi)$$
$$\log p_{\theta}(x^{(i)}) \ge \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound ("**ELBO**")

[Kingma and Welling, 2014]

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$



Putting it all together:

maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Lets look at **computing the bound** (forward pass) for a given mini batch of input data





Putting it all together:

maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$





Putting it all together:

maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$
Make approximate posterior distribution close to prior





Putting it all together:

maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$
Make approximate posterior distribution close to prior





Putting it all together:

maximizing the likelihood lower bound

$$\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$
Make approximate posterior distribution close to prior





Putting it all together:

maximizing the likelihood lower bound

Maximize likelihood of original input being reconstructed

$$\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))$$

 $\mathcal{L}(x^{(i)}, \theta, \phi)$

Make approximate posterior distribution close to prior





Putting it all together: maximizing the likelihood lower

bound

Maximize likelihood of original input being reconstructed

$$\underbrace{\mathbf{E}_{z}\left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid\mid p_{\theta}(z))}_{\mathbf{V}}$$

 $\mathcal{L}(x^{(i)}, \theta, \phi)$

Make approximate posterior distribution close to prior

For every minibatch of input data: compute this forward pass, and then backprop!





Use decoder network and sample z from **prior**



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



Use decoder network and sample z from **prior**



Data manifold for 2-d z



Diagonal prior on z => independent latent variables

Different dimensions of z encode interpretable factors of variation

Data manifold for 2-d z



Vary z_1

(degree of smile)

(head pose)



Diagonal prior on z => independent latent variables

Different dimensions of z encode interpretable factors of variation

Also good feature representation that can be computed using $q_{\phi}(z|x)!$

Data manifold for 2-d z



Vary z_1

(degree of smile)

(head pose)





32x32 CIFAR-10



Labeled Faces in the Wild



Conditional VAEs







(a) Frame 1



(b) Frame 2 (ground truth)

[Xue et al., 2016]



(c) Frame 2 (Sample 1) (d) Frame 2 (Sample 2)
Variational Autoencoders

Probabilistic spin to traditional autoencoders => allows generating data Defines an intractable density = derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Allows inference of q(z|x), can be useful feature representation for other tasks

Cons:

- Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
- Samples blurrier and lower quality compared to state-of-the-art (GANs)

Active area of research:

- More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian
- Incorporating structure in latent variables (our submission to CVPR)