#### Introduction to variational autoencoders

Abstract

Variational autoencoders are interesting generative models, which combine ideas from deep learning with statistical inference. They can be used to learn a low dimensional representation Z of high dimensional data X such as images (of e.g. faces). In contrast to standard auto encoders, X and Z are random variables. It's therefore possible to sample X from the distribution P(X| Z), thus creating e.g. images of faces, MNIST Digits, or speech.

In this talk I will in some detail describe the paper of Kingma and Welling. *"Auto-Encoding Variational Bayes, International Conference on Learning Representations." ICLR,* 2014. <u>arXiv:1312.6114</u> [stat.ML].

I will also show some code. A TensorFlow notebook can be found at: <a href="https://github.com/oduerr/dl\_tutorial/blob/master/tensorflow/vae/vae\_demo.ipynb">https://github.com/oduerr/dl\_tutorial/blob/master/tensorflow/vae/vae\_demo.ipynb</a>

#### Introduction to variational autoencoders

#### Oliver Dürr

Datalab-Lunch Seminar Series Winterthur, 11 May, 2016

#### **Motivation: Generating Faces**



#### Other examples

- <u>random faces</u>
- <u>MNIST</u>
- Speech

just google vae...

These are not part of the trainingset!

https://www.youtube.com/watch?v=XNZIN7Jh3Sg

#### Motivation: Generating Hand Written Digits



http://www.dpkingma.com/sgvb\_mnist\_demo/demo.html



#### Recap: Auto Encoders ('classical')



A simple autoencoder more see **Beates talk on Autoencoders** 

#### **Recap: Linear Regression**

Most people think of linear regression as points and a straight line:



Strange axis names, to be compatible with later notation

#### **Recap: Linear Regression**

Benefits of having an error model:

- How likely is a data point
- Confidence bounds
- Compare models



Strange axis names, to be compatible with later notation

Statisticians additionally have  $P_{\theta}(X | Z)$ 

#### Recap: Linear Regression (as a graphical model)



See Beates talk on <u>Causal inference with graphical models</u>

#### Going from $R^1$ to $R^{10000}$



Is R<sup>2</sup> "big enough" to create images from R<sup>100000</sup>?...

#### Manifold hypothesis

- X high dimensional vector
- Data is concentrated around a low dimensional manifold



• Hope finding a representation Z of that manifold.

#### Variational auto encoders (idea of low dim manifold)

1D

Low Dimensional representation a line



2D

High Dimensional (number of pixels)



3D



credit: http://www.deeplearningbook.org/

#### Variational auto encoders (idea of low dim manifold)

Examples of successful unfolding (2D  $\rightarrow$  R<sup>28x28</sup>, R<sup>20x26</sup>)



Frey Face dataset

2000 pictures of Brendan Frey (20x26)

# How did they do that?

#### Variational Autoencoders ("history")

Simultaneously discovered by

- Kingma and Welling. "Auto-Encoding Variational Bayes, International Conference on Learning Representations." ICLR, 2014. arXiv:1312.6114 [stat.ML] (20 December 2013, Amsterdam University) Talk
- Rezende, Mohamed and Wierstra. "Stochastic back-propagation and variational inference in deep latent Gaussian models." ICML, 2014 <u>arXiv:1401.4082</u> [stat.ML] (16 January 2014, Google DeepMind)

Alternative approach (for binary distributions)

- Gregor, Danihelka et all. "Deep autoregressive networks." ICML 2014
  - Has a more information theoretic ansatz (codings length)
  - Lecture given at <u>Nando de Freitas ML Course (University of Oxford)</u> (a bit hand waving argument but with nice examples)
- We focus on the approach as in Kingma, Welling

#### Principle Idea decoder network (graphical model)

- We have a set of N-observations (e.g. images) {x<sup>(1)</sup>,x<sup>(2)</sup>,...,x<sup>(N)</sup>}
- Complex model parameterized with θ
- There is a latent space z with

 $z \sim p(z)$  multivariate Gaussian  $x|z \sim p_{\theta}(x|z)$ 



#### One Example

Wish to learn  $\theta$  from the N training observations  $x^{(i)}$  i=1,...,N

#### The model for the decoder network

- For illustration z one dimensional x 2D
- Want a complex model of distribution of x given z
- Idea: NN + Gaussian (or Bernoulli) here with diagonal covariance Σ



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

#### Training as an autoencoder



Training use maximum likelihood of p(x) given the training data

Problem:

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

Cannot be calculated:

Solution:

- MCMC (too costly)
- Approximate p(z|x) with q(z|x)

#### The model for the encoder network





The complete auto-encoder



Learning the parameters  $\phi$  and  $\theta$  via backpropagation



Determining the loss function

#### **Training: Loss Function**

• What is (one of the) most beautiful idea in statistics?

- Max-Likelihood, tune  $\Phi$ ,  $\theta$  to maximize the likelihood
- We maximize the (log) likelihood of a given "image" x<sup>(i)</sup> of the training set. Later we sum over all training data (using minibatches)

#### Lower bound of likelihood (mathematical sleight of hand)

Likelihood, for an image  $x^{(i)}$  from training set. Writing  $x=x^{(i)}$  for short.

$$L = \log (p(x))$$

$$= \sum_{z} q(z|x) \log (p(x)) \qquad \text{multiplied with } 1$$

$$= \sum_{z} q(z|x) \log \left(\frac{p(z,x)}{p(z|x)}\right)$$

$$= \sum_{z} q(z|x) \log \left(\frac{p(z,x)}{q(z|x)}\frac{q(z|x)}{p(z|x)}\right)$$

$$= \sum_{z} q(z|x) \log \left(\frac{p(z,x)}{q(z|x)}\right) + \sum_{z} q(z|x) \log \left(\frac{q(z|x)}{p(z|x)}\right)$$

$$= L^{\vee} + D_{\text{KL}} (q(z|x)||p(z|x))$$

 $D_{KL}$  KL-Divergence >= 0 depends on how good q(z|x) can approximate p(z|x) L<sup>v</sup> "lower variational bound of the (log) likelihood" L<sup>v</sup> =L for perfect approximation

#### Approximate Inference (rewriting L<sup>v</sup>)

$$L^{v} = \sum_{z} q(z|x) \log \left(\frac{p(z,x)}{q(z|x)}\right) \quad \text{with } p(z,x) = p(x|z) p(z)$$

$$= \sum_{z} q(z|x) \log \left(\frac{p(x|z)p(z)}{q(z|x)}\right)$$

$$= \sum_{z} q(z|x) \log \left(\frac{p(z)}{q(z|x)}\right) + \sum_{z} q(z|x) \log (p(x|z))$$

$$= -D_{\text{KL}} \left(q(z|x)||p(z)\right) + \mathbb{E}_{q(z|x)} \left(\log (p(x|z))\right) \quad \text{putting in } x^{(i)} \text{ for } x$$

$$= -D_{\text{KL}} \left(q(z|x^{(i)})||p(z)\right) + \mathbb{E}_{q(z|x^{(i)})} \left(\log (p(x^{(i)}|z)\right)\right)$$

$$= -D_{\text{KL}} \left(q(z|x^{(i)})||p(z)\right) + \mathbb{E}_{q(z|x^{(i)})} \left(\log (p(x^{(i)}|z)\right)$$

$$= -D_{\text{KL}} \left(q(z|x^{(i)})||p(z)\right) + \mathbb{E}_{q(z|x^{(i)})} \left(\log (p(x^{(i)}|z)\right)\right)$$

$$= -D_{\text{KL}} \left(q(z|x^{(i)})||p(z)\right) + \mathbb{E}_{q(z|x^{(i)})} \left(\log (p(x^{(i)}|z)\right)$$

$$= -D_{\text{KL}} \left(q(z|x^{(i)}||z\right)$$

$$= -D_{\text{K}} \left(q(z|x^{(i)}||z\right)$$

$$= -D_{\text{K}} \left(q(z|x^{(i)}||z\right)$$

$$= -D_{\text{K}} \left(q(z|x^{(i)}||z\right)$$

$$= -D_{\text{K}} \left(q(z|x^{(i)}|z\right)$$

$$= -D_{\text{K}} \left(q$$

Calculation the regularization  $-D_{KL}(q(z|x^{(i)})||p(z))$ 

Use N(0,1) as prior for p(z)

 $q(z|x^{(i)})$  is Gaussian with parameters ( $\mu^{(i)}, \sigma^{(i)}$ ) determined by NN

$$-D_{\text{KL}}\left(q(z|x^{(i)})||p(z)\right) = \frac{1}{2}\sum_{j=1}^{J}\left(1 + \log(\sigma_{z_j}^{(i)^2}) - \mu_{z_j}^{(i)^2} - \sigma_{z_j}^{(i)^2}\right)$$



## Sampling to calculate $\mathbb{E}_{q(z|x^{(i)})} \left( \log(p(x^{(i)}|z)) \right)$

Approximating  $\mathbb{E}_{q(z|x^{(i)})}$  with sampling form the distribution  $q(z|x^{(i)})$ 

With 
$$z^{(i,l)} \ l = 1, 2, ... L$$
 sampled from  $z^{(i,l)} \sim q(z|x^{(i)})$   
 $L^{\vee} = -D_{\mathrm{KL}} \left( q(z|x^{(i)})||p(z) \right) + \mathbb{E}_{q(z|x^{(i)})} \left( \log \left( p(x^{(i)}|z) \right) \right)$   
 $L^{\vee} \approx -D_{\mathrm{KL}} \left( q(z|x^{(i)})||p(z) \right) + \frac{1}{L} \sum_{i=1}^{L} \log \left( p(x^{(i)}|z^{(i,l)}) \right)$ 



L is often very small (often just L=1)

#### One last trick

Backpropagation not possible through random sampling!



Sampling (reparametrization trick)

Cannot back propagate through a random drawn number

 $z^{(i,l)} \sim N(\mu^{(i)}, \sigma^{2(i)})$  $z^{(i,l)} = \mu^{(i)} + \sigma^{(i)} \odot \varepsilon_i \quad \varepsilon_i \sim N(0,1)$ 

z has the same distribution, but now one can back propagate.

Writing z in this form, results in a deterministic part and noise.

Image from: NIPS Workshop 2015 (Kingma & Welling)

#### Putting it all together

Prior  $p(z) \sim N(0,1)$  and p, q Gaussian, extension to dim(z) > 1 trivial



Cost: Regularisation

$$-D_{\mathrm{KL}}\left(q(z|x^{(i)})||p(z)\right) = \frac{1}{2}\sum_{j=1}^{J}\left(1 + \log(\sigma_{z_j}^{(i)^2}) - \mu_{z_j}^{(i)^2} - \sigma_{z_j}^{(i)^2}\right)$$

**Cost: Reproduction** 

$$-\log(p(x^{(i)}|z^{(i)})) = \sum_{j=1}^{D} \frac{1}{2}\log(\sigma_{x_j}^2) + \frac{(x_j^{(i)} - \mu_{x_j})^2}{2\sigma_{x_j}^2}$$

We use mini batch gradient decent to optimize the cost function over all  $x^{(i)}$  in the mini batch

Least Square for constant variance

### Use the source Luke

#### Simple example 2-D distribution

https://github.com/oduerr/dl\_tutorial/blob/master/tensorflow/vae/vae\_demo-2D.ipynb

Simple MNIST Example https://github.com/oduerr/dl\_tutorial/blob/master/tensorflow/vae/vae\_demo.ipynb

## Recent developments of VAE

#### Recent developments in VAE / generative models (subjective overview)

- Authors of VAE Amsterdam University and Google DeepMind teamed up and wrote a paper on semi-supervised learning:
  - Diederik P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, Max Welling.
     <u>"Semi-supervised learning with deep generative models"</u> (2014)
- Karl Gregor et al. extended the (binary autoencoder) with attention
  - DRAW: A Recurrent Neural Network For Image Generation <u>https://arxiv.org/abs/1502.04623</u> (2015)
  - <u>https://www.youtube.com/watch?v=Zt-7MI9eKEo</u>
- Adversial networks as a non-statistical way to generate high dimensional data
  - Play a game:
    - Fist network invents some data  $\rightarrow P(X)$  to fool second network
    - Second network tells if first network is a liar.

#### Semisupervised learning



Slide: Kingma, Rezendem Nohamed, Welling

Semisupervised learning

#### VAEs are SOTA on semi-supervised learning on MNIST

	That's 1 per class	0 5!
	100 labels	
AtlasRBF (Pitelis et al., 2014)	8.10% (±0.95)	
Deep Generative Model (M1+M2) (Kingma et al., 2014)	3.33% (±0.14)	
Virtual Adversarial (Miyato et al., 2015)	2.12%	
Ladder (Rasmus et al., 2015)	$1.06\%$ ( $\pm 0.37$ )	)
Auxiliary Deep Generative Model (1 MC)	$2.25\% (\pm 0.08)$	
Auxiliary Deep Generative Model (10 MC)	$0.96\% \ (\pm \ 0.02)$	K

"Improving Semi-Supervised Learning with Auxiliary Deep Generative Models" [Maaløe, Sønderby, Sønderby and Winter, 2015]

# Thank you, questions?