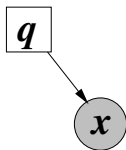


4F10: Variational AutoEncoder

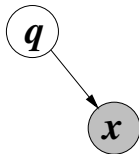
Mark Gales

Michaelmas 2018

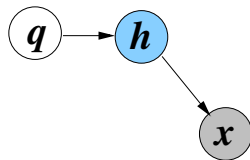
Latent Variable Examples



Gaussian Mixture Model



Factor Analysis



VAE

- Latent variables enables rich models: indicator variable q
 - introduce **discrete** latent variable c_m
 - introduce **continuous** latent variable z
- For discrete case possible to write

$$p(\mathbf{x}) = \sum_{m=1}^M P(c_m) p(\mathbf{x}|c_m)$$

- Generative latent variable models have form

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

- \mathbf{x} is the d -dimensional observation
 - \mathbf{z} is the p -dimensional latent variable
- Focus on **Gaussian** distributions of the form

$$\begin{aligned}p(\mathbf{x}|\mathbf{z}) &= \mathcal{N}(\mathbf{x}; \mathbf{f}(\mathbf{z}), \sigma^2 \mathbf{I}) \\ p(\mathbf{z}) &= \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})\end{aligned}$$

- Consider linear relationship

$$\mathbf{f}(\mathbf{z}) = \mathbf{A}\mathbf{z}$$

- this is a restricted form of **factor analysis**
- also called **probabilistic PCA**
- The overall distribution can be written as

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{A}\mathbf{A}^T + \sigma^2\mathbf{I})$$

- parameters can be estimated using **Maximum Likelihood** (ML)
- EM can be used: auxiliary function in this case has the form

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\lambda}, \tilde{\boldsymbol{\lambda}}) &= \int p(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}}) \log(p(\mathbf{x}|\mathbf{z}; \boldsymbol{\lambda})) d\mathbf{z} \\ p(\mathbf{z}|\mathbf{x}) &\propto p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \quad \text{Gaussian distribution} \end{aligned}$$

- Rather than restricting relationship to be linear
 - consider general mappings from \mathbf{z} to the mean
 - use a (deep) neural network with parameters λ

$$p(\mathbf{x}|\mathbf{z}; \lambda) = \mathcal{N}(\mathbf{x}; \mathbf{f}(\mathbf{z}; \lambda), \sigma^2 \mathbf{I})$$

- In general no simple closed-form solutions to integral
 - cannot compute the log-likelihood for inference
 - cannot compute gradient to estimate λ and σ^2
- To address this [variational](#) approaches are adopted
 - first revision of [Kullback-Leibler Divergence](#)
 - then EM and variational EM discussed

Kullback-Leibler Divergence

- Consider two PDFs, $p(\mathbf{x})$ and $q(\mathbf{x})$.
 - Kullback-Leibler divergence, $\mathcal{KL}(p(\mathbf{x})\|q(\mathbf{x}))$, is

$$\mathcal{KL}(p(\mathbf{x})\|q(\mathbf{x})) = \int p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) d\mathbf{x} = - \int p(\mathbf{x}) \log \left(\frac{q(\mathbf{x})}{p(\mathbf{x})} \right) d\mathbf{x}$$

Using $\log(y) \leq y - 1$, we can write

$$\begin{aligned} \int p(\mathbf{x}) \log \left(\frac{q(\mathbf{x})}{p(\mathbf{x})} \right) d\mathbf{x} &\leq \int p(\mathbf{x}) \left(\frac{q(\mathbf{x})}{p(\mathbf{x})} - 1 \right) d\mathbf{x} \\ &= \int (q(\mathbf{x}) - p(\mathbf{x})) d\mathbf{x} = 0 \end{aligned}$$

- This gives the following inequality

$$\mathcal{KL}(p(\mathbf{x})\|q(\mathbf{x})) = \int p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) d\mathbf{x} \geq 0$$

KL Divergence for Gaussians

- Consider two Gaussian distributions

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1); \quad q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

- The KL divergence between the two is given by

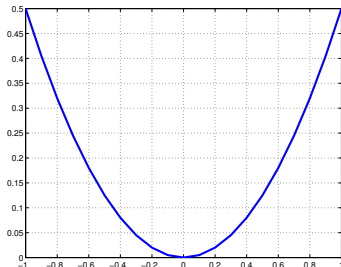
$$\mathcal{KL}(p(\mathbf{x})||q(\mathbf{x})) = \frac{1}{2} \left(\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 - \mathbf{I}) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \log \left(\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) \right)$$

- Take simple example:

$$p(x) = \mathcal{N}(x; 0, 1);$$

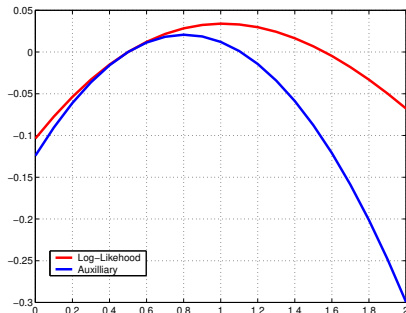
$$q(x) = \mathcal{N}(x; \mu, 1)$$

- Vary μ - plot KL



(Variational) EM

Expectation Maximisation (EM)



- Maximise an **auxiliary function** rather than log-likelihood
 - iterative optimisation for iteration $k + 1$, parameters $\lambda^{(k+1)}$

$$\lambda^{(k+1)} = \arg \max_{\lambda} \left\{ \mathcal{Q}(\lambda^{(k)}, \lambda) \right\} = \arg \max_{\lambda} \left\{ \int p(\mathbf{z}|\mathbf{x}; \lambda^{(k)}) \log (p(\mathbf{x}, \mathbf{z}; \lambda)) d\mathbf{z} \right\}$$

- Maximising auxiliary function simpler than log-likelihood
 - expected value of \log joint distribution, $\log(p(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}))$
 - consider factor analysis (expectation over \mathbf{z})

$$\begin{aligned}\log(p(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda})) &= \log(p(\mathbf{x}|\mathbf{z}; \boldsymbol{\lambda})) + \log(p(\mathbf{z})) \\ &= \log(\mathcal{N}(\mathbf{x}; \mathbf{A}\mathbf{z}, \boldsymbol{\Sigma}_{\text{diag}})) + \log(\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})) \\ &= -\frac{1}{2}(\mathbf{z}^T \mathbf{A}^T \boldsymbol{\Sigma}_{\text{diag}}^{-1} \mathbf{A} \mathbf{z} - 2\mathbf{z}^T \mathbf{A}^T \boldsymbol{\Sigma}_{\text{diag}}^{-1} \mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}_{\text{diag}}^{-1} \mathbf{x}) \\ &\quad - \frac{1}{2} \log(|\boldsymbol{\Sigma}_{\text{diag}}|) + C\end{aligned}$$

- C is a constant not a function of $\boldsymbol{\lambda}$
- a function of the first and second moments of \mathbf{z}

- Estimate parameters, λ , based on **log-likelihood**, $\mathcal{L}(\lambda)$
 - for simplicity only consider a single training observation \mathbf{x}

$$\begin{aligned}\mathcal{L}(\lambda) &= \log(p(\mathbf{x}; \lambda)) \\&= \int p(\mathbf{z}|\mathbf{x}; \lambda) \log(p(\mathbf{x}; \lambda)) d\mathbf{z} \\&= \int p(\mathbf{z}|\mathbf{x}; \lambda) \log\left(\frac{p(\mathbf{x}; \lambda)p(\mathbf{z}|\mathbf{x}; \lambda)}{p(\mathbf{z}|\mathbf{x}; \lambda)}\right) d\mathbf{z} \\&= \left\langle \log\left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{p(\mathbf{z}|\mathbf{x}; \lambda)}\right) \right\rangle_{p(\mathbf{z}|\mathbf{x}; \lambda)}\end{aligned}$$

- Need posterior distribution of the latent variables $p(\mathbf{z}|\mathbf{x}; \lambda)$

Log-Likelihood Expression (cont)

- BUT what happens if not possible to compute $p(\mathbf{z}|\mathbf{x}; \lambda)$
- Consider **any** valid distribution $q(\mathbf{z}; \tilde{\lambda})$

$$\begin{aligned}\mathcal{L}(\lambda) &= \log(p(\mathbf{x}; \lambda)) \\ &= \int q(\mathbf{z}; \tilde{\lambda}) \log(p(\mathbf{x}; \lambda)) d\mathbf{z} \\ &= \int q(\mathbf{z}; \tilde{\lambda}) \log\left(\frac{p(\mathbf{x}; \lambda)p(\mathbf{z}|\mathbf{x}; \lambda)}{p(\mathbf{z}|\mathbf{x}; \lambda)}\right) d\mathbf{z} \\ &= \left\langle \log\left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{p(\mathbf{z}|\mathbf{x}; \lambda)}\right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})}\end{aligned}$$

- For any parameter values, e.g. $\tilde{\lambda}$, and distribution $q(\mathbf{z}; \tilde{\lambda})$,

$$\begin{aligned}\mathcal{L}(\lambda) &= \left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{p(\mathbf{z}|\mathbf{x}; \lambda)} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})} \\ &= \left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{q(\mathbf{z}; \tilde{\lambda})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})} + \left\langle \log \left(\frac{q(\mathbf{z}; \tilde{\lambda})}{p(\mathbf{z}|\mathbf{x}; \lambda)} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})} \\ &\geq \left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{q(\mathbf{z}; \tilde{\lambda})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})} = \mathcal{F}(q(\mathbf{z}; \tilde{\lambda}), \lambda)\end{aligned}$$

- uses KL-divergence to yield a lower-bound
- equality when $q(\mathbf{z}; \tilde{\lambda}) = p(\mathbf{z}|\mathbf{x}; \lambda)$

- EM expressed based in this form - at iteration $k + 1$
 - initially $q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}^{(k)})$

$$\mathcal{L}(\boldsymbol{\lambda}^{(k)}) = \mathcal{F}(q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}), \boldsymbol{\lambda}^{(k)})$$

- Maximise $\mathcal{F}(q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}), \boldsymbol{\lambda})$ to find parameters $\boldsymbol{\lambda}^{(k+1)}$

$$\boldsymbol{\lambda}^{(k+1)} = \arg \max_{\boldsymbol{\lambda}} \left\{ \left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda})}{q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)})} \right\}$$

- this maximises a lower-bound on $\mathcal{L}(\boldsymbol{\lambda}^{(k)})$
- Minimise KL divergence for “variational” approximation

$$q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k+1)}) = \arg \min_{q(\mathbf{z}; \tilde{\boldsymbol{\lambda}})} \left\{ \left\langle \log \left(\frac{q(\mathbf{z}; \tilde{\boldsymbol{\lambda}})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}^{(k+1)})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\boldsymbol{\lambda}})} \right\}$$

- occurs when $q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}^{(k+1)})$

- Yields the following sequence of inequalities for iteration $k + 1$

$$\mathcal{L}(\boldsymbol{\lambda}^{(k)}) = \mathcal{F}\left(q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}), \boldsymbol{\lambda}^{(k)}\right) \leq \mathcal{F}\left(q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}), \boldsymbol{\lambda}^{(k+1)}\right) \leq \mathcal{L}(\boldsymbol{\lambda}^{(k+1)})$$

- **provided** $q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}) = p(\mathbf{z}|\mathbf{x}; \boldsymbol{\lambda}^{(k)})$
- Iterate until convergence:
 - each iteration **guaranteed not to decrease the likelihood**
 - finds a **local** maximum of the likelihood
 - final solution depends on initial parameters $\boldsymbol{\lambda}^{(0)}$

- Start from the same basic equation as before at iteration k

$$\mathcal{L}(\lambda^{(k)}) = \left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda^{(k)})}{q(\mathbf{z}; \tilde{\lambda}^{(k)})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda}^{(k)})} + \left\langle \log \left(\frac{q(\mathbf{z}; \tilde{\lambda}^{(k)})}{p(\mathbf{z}|\mathbf{x}; \lambda^{(k)})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda}^{(k)})}$$

- General form of $q(\mathbf{z}; \tilde{\lambda}^{(k)})$
 - second-term is **greater than or equal to zero** initially
 - yields an initial **inequality**

$$\mathcal{L}(\lambda^{(k)}) \geq \mathcal{F}\left(q(\mathbf{z}; \tilde{\lambda}^{(k)}), \lambda^{(k)}\right)$$

- Repeat same process as previous slide
 - estimate $\lambda^{(k+1)}$ using $\mathcal{F}\left(q(\mathbf{z}; \tilde{\lambda}^{(k)}), \lambda\right)$
 - estimate $q(\mathbf{z}; \tilde{\lambda}^{(k+1)})$ by minimising KL-divergence

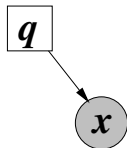
- Yields the following inequalities

$$\mathcal{L}(\boldsymbol{\lambda}^{(k)}) \geq \mathcal{F}\left(q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}), \boldsymbol{\lambda}^{(k)}\right) \leq \mathcal{F}\left(q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}^{(k)}), \boldsymbol{\lambda}^{(k+1)}\right) \leq \mathcal{L}(\boldsymbol{\lambda}^{(k+1)})$$

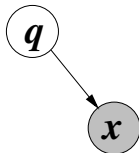
- no guarantees of not decreasing log-likelihood
 - **but** allows any form of expression for $q(\mathbf{z}; \tilde{\boldsymbol{\lambda}})$
- One standard form is the **mean-field approximation** where

$$q(\mathbf{z}; \tilde{\boldsymbol{\lambda}}) = \prod_{i=1}^n q_i(z_i; \tilde{\boldsymbol{\lambda}})$$

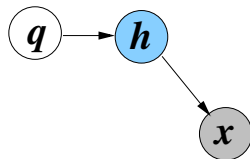
Variational AutoEncoder



Gaussian Mixture Model



Factor Analysis



VAE

- Consider a neural network with parameters λ

$$p(\mathbf{x}; \lambda) = \int p(\mathbf{x}|\mathbf{z}; \lambda)p(\mathbf{z})d\mathbf{z} = \int \mathcal{N}(\mathbf{x}; \mathbf{f}(\mathbf{z}; \lambda), \sigma^2\mathbf{I})p(\mathbf{z})d\mathbf{z}$$

- $p(\mathbf{z})$ known - standard Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I})$
- Cannot compute integral - $\mathbf{f}(\mathbf{z}; \lambda)$ non-linear

- Iterative maximisation involves

$$\begin{aligned}\lambda^{(k+1)} &= \arg \max_{\lambda} \left\{ \mathcal{F}(q(\mathbf{z}; \tilde{\lambda}^{(k)}), \lambda) \right\} \\ &= \arg \max_{\lambda} \left\{ \left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{q(\mathbf{z}; \tilde{\lambda}^{(k)})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda}^{(k)})} \right\}\end{aligned}$$

- need to consider form of $q(\mathbf{z}; \tilde{\lambda})$
 - needs to approximate $p(\mathbf{z}|\mathbf{x}; \lambda)$
 - influences how tight the bound is for optimisation
- VEM allows EM to be applied for a wider range of models
 - useful when maximising auxiliary function is easy
 - **not** (normally) possible for deep learning

(Stochastic) Gradient Descent

- Rather than variational EM, use gradient descent
 - approximate gradient by lower-bound gradient

$$\begin{aligned}\nabla \mathcal{L}(\lambda) &\approx \nabla \left(\left\langle \log \left(\frac{p(\mathbf{x}, \mathbf{z}; \lambda)}{q(\mathbf{z}; \tilde{\lambda})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})} \right) \\ &= \nabla \left(\langle \log(p(\mathbf{x}|\mathbf{z}; \lambda)) \rangle_{q(\mathbf{z}; \tilde{\lambda})} + \left\langle \log \left(\frac{p(\mathbf{z})}{q(\mathbf{z}; \tilde{\lambda})} \right) \right\rangle_{q(\mathbf{z}; \tilde{\lambda})} \right)\end{aligned}$$

- $p(\mathbf{x}|\mathbf{z}; \lambda)$ is Gaussian $\mathcal{N}(\mathbf{x}; \mathbf{f}(\mathbf{z}; \lambda), \sigma^2 \mathbf{I})$
 - $p(\mathbf{z})$ is Gaussian (and known)
- Iterative optimisation
 - optimise λ (model parameters)
 - optimise $\tilde{\lambda}$ (variational approximation)

- Use a density network as the variational approximation
 - introduce a dependence on \mathbf{x} for variational approximation

$$q(\mathbf{z}; \tilde{\lambda}) \rightarrow q(\mathbf{z}|\mathbf{x}; \tilde{\lambda}) = \mathcal{N}(\mathbf{z}; \mathbf{f}_{\mu}(\mathbf{x}; \tilde{\lambda}), \mathbf{f}_{\Sigma}(\mathbf{x}; \tilde{\lambda}))$$

- everything is Gaussian!
- Likelihood can now be written as:

$$\mathcal{L}(\lambda) = \left\langle \log \left(\frac{q(\mathbf{z}|\mathbf{x}; \tilde{\lambda})}{p(\mathbf{z}|\mathbf{x}; \lambda)} \right) + \log (p(\mathbf{x}|\mathbf{z}; \lambda)) + \log \left(\frac{p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x}; \tilde{\lambda})} \right) \right\rangle_{q(\mathbf{z}|\mathbf{x}; \tilde{\lambda})}$$

- Consider the three terms (left-to-right)
 1. **error**: need variational approximation to be close to \mathbf{z} posterior
 2. **decoding**: compute probability given \mathbf{z}
 3. **encoding**: encode information about \mathbf{x} into \mathbf{z}

- Gradient lower-bound for optimisation

$$\nabla \left(\left\langle \log(p(\mathbf{x}|\mathbf{z}; \boldsymbol{\lambda})) + \log\left(\frac{p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})}\right) \right\rangle_{q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})} \right)$$

- First term is tricky - remember

$$\langle \log(p(\mathbf{x}|\mathbf{z}; \boldsymbol{\lambda})) \rangle_{q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})} = \langle \log(\mathcal{N}(\mathbf{x}; \mathbf{f}(\mathbf{z}; \boldsymbol{\lambda}), \sigma^2 \mathbf{I})) \rangle_{q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})}$$

- need to compute (for example) $\langle \mathbf{f}(\mathbf{z}) \rangle_{q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})}$
 - difficult if $\mathbf{f}(\mathbf{z})$ highly non-linear
- Second term on RHS (-) KL-Divergence between Gaussians
 - simple closed-form solution - see earlier slide

- Use standard integration approximation

$$\begin{aligned}\langle \log(p(\mathbf{x}|\mathbf{z}; \boldsymbol{\lambda})) \rangle_{q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})} &\approx \frac{1}{K} \sum_{i=1}^K \log(p(\mathbf{x}|\mathbf{z}^{(i)}; \boldsymbol{\lambda})) \\ \mathbf{z}^{(i)} &\sim \mathcal{N}(\mathbf{f}_{\mu}(\mathbf{x}; \tilde{\boldsymbol{\lambda}}), \mathbf{f}_{\Sigma}(\mathbf{x}; \tilde{\boldsymbol{\lambda}})) \\ q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}}) &= \mathcal{N}(\mathbf{z}; \mathbf{f}_{\mu}(\mathbf{x}; \tilde{\boldsymbol{\lambda}}), \mathbf{f}_{\Sigma}(\mathbf{x}; \tilde{\boldsymbol{\lambda}}))\end{aligned}$$

- using SGD will approximate gradient using a subset of \mathcal{D}
 - can use a single sample of \mathbf{z} ($K = 1$)
- This approximation allows $\boldsymbol{\lambda}$ gradients to be estimated
 - but this is using samples from $q(\mathbf{z}|\mathbf{x}; \tilde{\boldsymbol{\lambda}})$
 - not possible to get gradients for $\tilde{\boldsymbol{\lambda}}$

Reparameterisation Trick

- To introduce the dependence on $\tilde{\lambda}$ rewrite

$$\begin{aligned}\mathbf{z}^{(i)} &= \mathbf{f}_{\mu}(\mathbf{x}; \tilde{\lambda}) + \mathbf{f}_{\Sigma}(\mathbf{x}; \tilde{\lambda})^{1/2} \boldsymbol{\epsilon}^{(i)} \\ \boldsymbol{\epsilon}^{(i)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I})\end{aligned}$$

- Now the expression becomes

$$\begin{aligned}\mathcal{L}(\lambda) &\geq \left\langle \log(p(\mathbf{x} | \mathbf{f}_{\mu}(\mathbf{x}; \tilde{\lambda}) + \mathbf{f}_{\Sigma}(\mathbf{x}; \tilde{\lambda})^{1/2} \boldsymbol{\epsilon}; \lambda)) \right\rangle_{\mathcal{N}(\mathbf{0}, \mathbf{I})} \\ &\quad + \left\langle \log \left(\frac{p(\mathbf{z})}{q(\mathbf{z}; \tilde{\lambda})} \right) \right\rangle_{q(\mathbf{z} | \mathbf{x}; \tilde{\lambda})}\end{aligned}$$

- the first term is a function of λ and $\tilde{\lambda}$
- second term - closed-form solution ((-)KL-divergence) - function of $\tilde{\lambda}$

Generated Examples



- “Fictional Celebrities” - Alec Radford (search YouTube)

Conclusions

Conclusions [4, 3, 7]

- Deep generative models are popular at the moment
- Other variants developed, for example
 - (restricted) Boltzmann machines (RBMs)
 - generative adversarial networks (GANs)
 - WaveNet (speech synthesis)
- GAN-generated art ... expensive (©Obvious)



- [1] L. Baum and J. Eagon, "An Inequality with Applications to Statistical Estimation for Probabilistic Functions of Markov Processes and to a Model for Ecology," *Bull Amer Math Soc*, vol. 73, pp. 360–363, 1967.
- [2] J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, M. W. (eds, M. J. Beal, and Z. Ghahramani, "The variational bayesian em algorithm for incomplete data: with application to scoring graphical model structures," 2003.
- [3] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 2672–2680. [Online]. Available: <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>
- [4] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural Computation*, vol. 14, no. 8, pp. 1771–1800, 2002. [Online]. Available: <https://doi.org/10.1162/089976602760128018>
- [5] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, no. 2014, 2013.
- [6] M. E. Tipping and C. M. Bishop, "Probabilistic principal component analysis," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 61, no. 3, pp. 611–622, 1999.
- [7] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," *CoRR*, 2016. [Online]. Available: <https://arxiv.org/pdf/1609.03499.pdf>