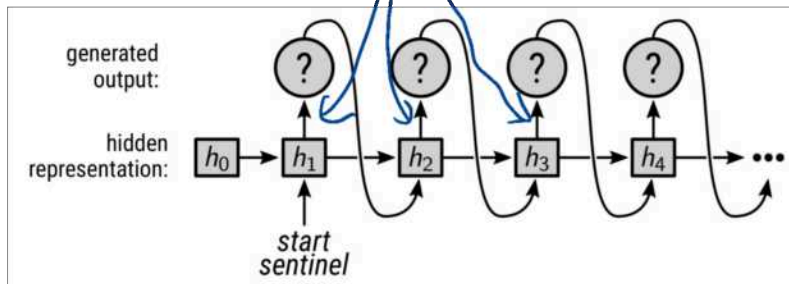


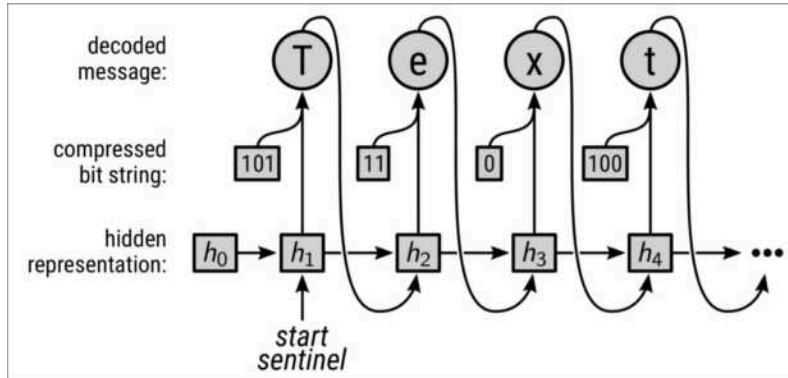
Data Compression with Deep Probabilistic Models

Reminder: Problem 3.2: compression with a learned autoregressive model

parameterizes a probability dist.
 \Rightarrow can be used for compression



\rightarrow when used for compression (here: decoder side):



autoregressive models: $P_{\theta}(X) = P_{\theta}(x_1) P_{\theta}(x_2|x_1) P_{\theta}(x_3|x_1, x_2)$

θ model parameters (neural network weights)
 \rightarrow optimize θ by minimizing an empirical estimate of cross entropy $H(P_{data}, P_{\theta})$

\rightarrow can we do the same thing with latent variable models

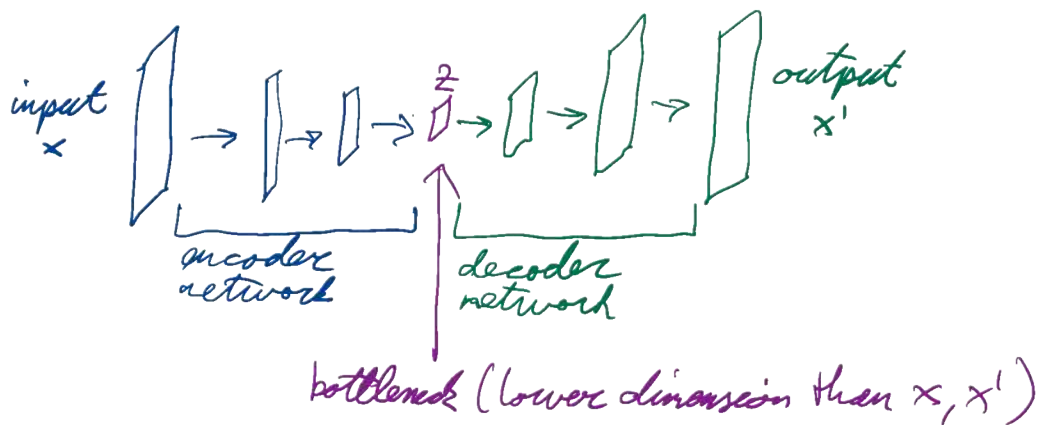
Deep Latent Variable Models & Scalable Approximate Bayesian Inference

Spoiler: variational autoencoders (VAEs)

→ a form of representation learning

→ often introduced with the following explanation:

"learn to map data to itself while squeezing it through a bottleneck"



use cases of VAEs for compression

↳ lossless compression

1) map x to z & encode z

2) map z to x' & encode residual $x - x'$

↳ lossy compression: leave out residual

⇒ 3 training objectives

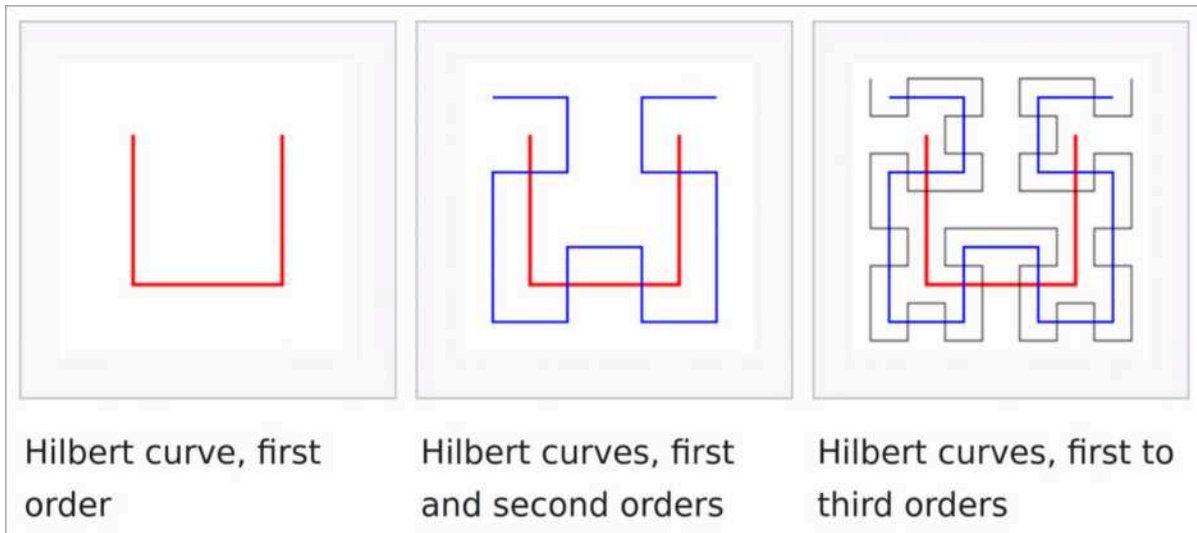
(i) decoder network should reconstruct the data well

(⇒ residual $x' - x$ small / low entropy)

(ii) encoder network decorrelates data

→ need probabilistic model (we want $P(z) = \prod_{i=1}^k P(z_i; i)$)

Note: just squeezing data through a lower-dimensional bottleneck does not in itself imply compression
→ think about information theoretical measures rather than dim.
(iii) keep $M(Z)$ low to enable effective compression



"Hilbert Curve"
(drawings from Wikipedia)

[ACM Transactions on Graphics (TOG) 35.4 (2016)]

A Compiler for 3D Machine Knitting

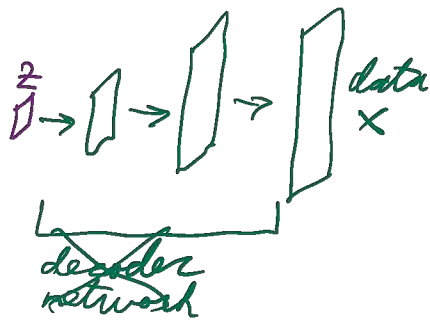
James McCann¹ Lea Albaugh¹ Vidya Narayanan¹ April Grow^{1,2}
Wojciech Matusik³ Jen Mankoff^{1,4} Jessica Hodgins¹

¹Disney Research ²UC Santa Cruz ³Massachusetts Institute of Technology ⁴Carnegie Mellon University



Deep Latent Variable Models

- look at decoder network only



→ interpret as a latent variable

model:

$$P_{\theta}(X, z) = P_{\theta}(z) P_{\theta}(X|z)$$

\uparrow learned model parameters (e.g. neural network weights)

common example:

→ prior is fully factorized, i.e., $P_{\theta}(z) = \prod_{i=1}^k P_{\theta}(z_i)$

→ likelihood: $P_{\theta}(x|z) = \mathcal{N}(x; f_{\theta}(z), \sigma^2 I)$

normal dist. (= Gaussian)

neural network

fixed or learned

lower case, density fct



Goal: minimize $H(P_{\text{data}}(X), P_{\theta}(X)) = \mathbb{E}_{P_{\text{data}}(X)} \left[\underbrace{-\log P_{\theta}(X)}_{\text{"evidence"}} \right]$

Problem $P_{\theta}(X=x) = \int P_{\theta}(x, z) dz$

\uparrow prohibitively expensive \uparrow high dimensional

We want to maximize evidence $P_{\theta}(X=x)$ when evaluated on data x from the training set.

Recall: bits - bad coding

$$R_{\text{net}}(x) = -\log P_{\theta}(X=x)$$

$$= -\log P_{\theta}(Z=z) - \log P_{\theta}(X=x|Z=z) + \log \underbrace{P_{\theta}(Z=z|X=x)}_{\text{posterior}}$$

problem: $P_{\theta}(Z|X) = \frac{P_{\theta}(X, Z)}{P_{\theta}(X)}$ ← intractable

• Idea: replace posterior with some other dist. $Q_{\lambda_x}(Z)$

(e.g.: $Q_{\lambda_x}(Z) = \prod_{i=1}^k \mathcal{N}(z_i; \mu_i, \sigma_i^2)$)

↑ ↑
make up λ_x

$$\Rightarrow \tilde{R}_{\text{net}}^{(z)}(x) = -\log P_{\theta}(X=x, Z=z) + \log Q_{\lambda_x}(Z=z)$$

$$\mathbb{E}_{Z \sim Q_{\lambda_x}(Z)} [\tilde{R}_{\text{net}}^{(z)}(x)] \geq R_{\text{net}}(x) = -\log P_{\theta}(X=x)$$

↑
equality if $Q_{\lambda_x}(Z) = P_{\theta}(Z|X=x)$

we want to minimize this

Notation & Naming Conventions

- $\log P_{\theta}(X=x)$ is called evidence (we want this to be high)
- $-\mathbb{E}_{Z \sim Q_{\lambda_x}(Z)} [\tilde{R}_{\text{net}}^{(z)}(x)] = \mathbb{E}_{Z \sim Q_{\lambda_x}(Z)} [\log P_{\theta}(X=x, Z=z) - \log Q_{\lambda_x}(Z=z)]$ is called the evidence lower bound (ELBO)

$$\Rightarrow \text{ELBO}(\theta, \lambda_x) \leq \underbrace{\log P_{\theta}(X=x)}_{\text{evidence}}$$

- parameters λ_x of the distribution $Q_{\lambda_x}(Z)$ are called "variational parameters"

- $Q_{\lambda_x}(Z)$ is called "variational distribution"

- Variational Inference (VI): approximate evidence $\log P_{\theta}(X=x)$ by $\text{ELBO}(\theta, \lambda_x^*)$ where

$$\lambda_x^* := \arg \max_{\lambda_x} \text{ELBO}(\theta, \lambda_x)$$

→ observation: this typically leads to a $Q_{\lambda_x^*}(z)$ which is "close" to true posterior $P_\theta(z|X=x)$.

(Reviews: Blei et al 2016, Zhang et al. 2018)

→ we now can approximate $\log P_\theta(X=x)$, but we still have to maximize it over θ .

→ idea: maximize our approximation $ELBO(\theta, \lambda_x^*)$ over θ .

Pseudocode:

for t in training-steps:

sample a minibatch B of training points

initialize λ_x randomly $\forall x \in B$

for t' in inner-training-steps:

perform gradient step for $\lambda_x \forall x \in B$

perform gradient step for θ on $ELBO(\theta, \lambda_x^*)$

nested loop
→ extremely expensive

$\left. \begin{array}{l} \text{VI,} \\ \text{finds} \\ \lambda_x^* \end{array} \right\}$

Remember: • model params θ are global (i.e., the same for all data points x)

• variational params λ_x parameterize an approximation of $P(z|X=x) \Rightarrow$ they are local (i.e., different for all data points x)

• we want to maximize $\mathbb{E}_{x \sim P_{\text{data}}} [\log P_\theta(X=x)]$
 \Rightarrow we have to sample a new minibatch in each iteration of outer loop
 \Rightarrow invalidates λ_x^* from previous iteration of outer loop.

→ "Variational Expectation Maximization"

(Dempster et al 1977, Beal & Ghahramani 2003)

Final additional trick: learn how to do variational inference

i.e., learn a function $g_\phi: x \mapsto \lambda_x$

set $\lambda_x = g_\phi(x)$ in the ELBO

notation: $Q_\phi(z|x) = Q_{\lambda_x}(z)$ with $\lambda_x = g_\phi(x)$

$$\text{ELBO}(\vartheta, \phi) = \mathbb{E}_{z \sim Q_\phi(z|x)} [\log P_\vartheta(X=x, z) - \log Q_\phi(z|x)]$$

$\uparrow \quad \uparrow$
 now both are
 global params

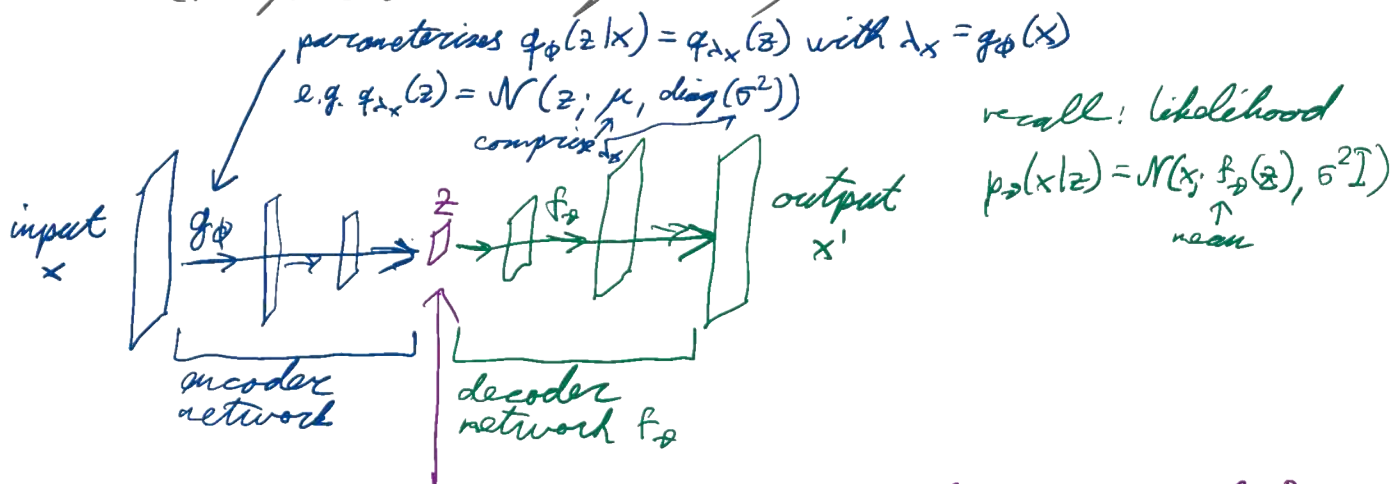
$$\leq \underbrace{\log P_\vartheta(X=x)}_{\text{evidence}}$$

→ maximize $\mathbb{E}_{x \sim P_{\text{data}}} [\text{ELBO}(\vartheta, \phi)]$ over both ϑ, ϕ
 often also just called "ELBO"

→ "Amortized Variational Expectation Maximization"

= "Variational Autoencoders" (VAEs)

(Kingma & Welling 2013)



- minimize entropy of this representation → more precisely D_{KL}
- we inject noise here: $z \sim Q_\phi(z|x)$

Interpretations of the ELBO (i.e. the objective function)

$$\begin{aligned} \text{ELBO}(\theta, \phi) &= \mathbb{E}_{z \sim Q_{\phi}(z|x)} \left[\log p_{\theta}(z) + \log p_{\theta}(x|z) - \log q_{\phi}(z|x) \right] \\ &\stackrel{\substack{\uparrow \\ \text{we maximize} \\ \text{this}}}{=} + \underbrace{\mathbb{E}_{z \sim Q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) \right]}_{\substack{\text{maximizing only this part} \\ \text{would be maximum likelihood} \\ \text{estimation (MLE)} \\ \rightarrow \text{it would make } Q_{\phi}(z|x) \\ \text{collapse to a } \delta\text{-function} \\ \text{peaked at the MLE} = \arg \max_z \log p_{\theta}(x|z)}} - \underbrace{D_{\text{KL}}(Q_{\phi}(z|x) \parallel P_{\theta}(z))}_{\substack{\text{think of this as a regularizer} \\ \rightarrow \text{try to make } Q_{\phi}(z|x) \text{ similar} \\ \text{to } P_{\theta}(z) \\ \rightarrow \text{at compression: want to encode} \\ z \text{ using } P_{\theta}(z); \text{ this term ensures} \\ \text{that } z \text{ is obtained from} \\ \text{encoder have high } P_{\theta}(z)}} \end{aligned}$$

$$= \underbrace{\log P_{\theta}(x=x)}_{\text{evidence}} - \underbrace{D_{\text{KL}}(Q_{\phi}(z|x) \parallel P_{\theta}(z|x=x))}_{\text{minimizing this makes the variational} \\ \text{dist. } Q_{\phi} \text{ similar to the true posterior} \\ \Rightarrow Q_{\phi} \text{ can be called the "approximate posterior"}}$$

\rightarrow maximizing this minimizes the info content of x under our model P_{θ} , i.e., the theoretical lower bound of the bitrate

\Rightarrow Q_{ϕ} can be called the "approximate posterior"

\rightarrow Goal: maximize ELBO over θ & ϕ

\cdot issue: $\text{ELBO}(\theta, \phi) = \mathbb{E}_{z \sim Q_{\phi}(z|x)} [\dots]$

distribution from which we have to sample depends on ϕ , by which we want to differentiate

\rightarrow see Problem set

(reparameterization grad: Kingma & Welling 2013
REINFORCE-gradients: Renzaglia et al. 2014)

Why all this fuss?

ongoing research on VI & related methods may be applicable to compression - or it may not be

\Rightarrow look into that literature & try out if it improves compression methods

- Examples:
- lots of research on tighter bounds of the evidence (tighter than the standard ELBO):
 - e.g. importance weighted VI, recently applied to compression by Thai's & Ho 2021
 - iterative amortized inference
 - Marino et al 2018
 - Campos et al 2019
 - other approximate Bayesian inference methods (alternatives to VI) exist (in particular: Markov Chain Monte Carlo = MCMC)
 - nontrivial how to use these for compression (pioneering work: Havasi et al., 2018)

References

- [Beal and Ghahramani, 2003] Beal, M. J. and Ghahramani, Z. (2003). The variational bayesian em algorithm for incomplete data: with application to scoring graphical model structures. *Bayesian statistics*, 7(453-464):210.
- [Blei et al., 2017] Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. (2017). Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518).
- [Campos et al., 2019] Campos, J., Meierhans, S., Djelouah, A., and Schroers, C. (2019). Content adaptive optimization for neural image compression. *arXiv preprint arXiv:1906.01223*.
- [Dempster et al., 1977] Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1).
- [Havasi et al., 2018] Havasi, M., Peharz, R., and Hernández-Lobato, J. M. (2018). Minimal random code learning: Getting bits back from compressed model parameters. *arXiv preprint arXiv:1810.00440*.
- [Kingma and Welling, 2014] Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *International Conference on Learning Representations*.
- [Marino et al., 2018] Marino, J., Yue, Y., and Mandt, S. (2018). Iterative amortized inference. In *International Conference on Machine Learning*.
- [Ranganath et al., 2014] Ranganath, R., Gerrish, S., and Blei, D. (2014). Black box variational inference. In *Artificial intelligence and statistics*, pages 814–822.
- [Theis and Ho, 2021] Theis, L. and Ho, J. (2021). Importance weighted compression. In *Neural Compression: From Information Theory to Applications–Workshop@ ICLR 2021*.
- [Zhang et al., 2018] Zhang, C., Bütepage, J., Kjellström, H., and Mandt, S. (2018). Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8).