TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2019

Expectation Maximization (EM)

The Evidence Lower Bound (the ELBO)

Variational Autoencoders (VAEs)

Latent Variable Models

We are often interested in models of the form

$$P_{\Phi}(y) = \sum_{z} P_{\Phi}(z) P_{\Phi}(y|z).$$

$$P_{\Phi}(y|x) = \sum_{z} P_{\Phi}(z|x) P_{\Phi}(y|z).$$

For example, CTC and probabilistic grammar models.

Expectation Maximization (EM) Mixture of Gaussian Modeling

$$\Phi = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_k, \mu_k, \Sigma_k)$$

$$p_{\Phi}(y) = \sum_{i} P(i)p(y|i)$$

$$= \sum_{i} \pi_{i} \frac{1}{Z_{i}} \exp\left(-\frac{1}{2}(y - \mu_{i})^{\top} \Sigma_{i}^{-1}(y - \mu_{i})\right)$$

i is the latent variable.

Expectation Maximization (EM) Mixture of Gaussian Modeling

$$\Phi = (\pi_1, \mu_1, \Sigma_1, \dots, \pi_k, \mu_k, \Sigma_k)$$

Train = $\{y_1, \dots, y_N\}$

Until Convergence:

$$P_{\Phi}(i|y_j) = \frac{\pi_i P(y_j|i)}{\sum_i \pi_i P(y_j|i)} \text{ Inference (E step)}$$

General EM

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} E_{y \sim \operatorname{Train}} - \ln P_{\Phi}(y)$$

$$P_{\Phi}(y) = \sum_{z} P_{\Phi}(z) P_{\Phi}(y|z).$$

$$\Phi^{t+1} = \underset{\Phi}{\operatorname{argmin}} \ E_{y \sim \operatorname{Train}} \ E_{z \sim P_{\Phi}t(z|y)} - \ln P_{\Phi}(z,y)$$
 Update Inference (M Step) (E Step)



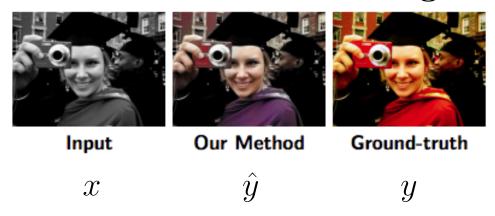
x is a black and white image.

y is a color image drawn from Pop(y|x).

 \hat{y} is an arbitrary color image.

 $P_{\Phi}(\hat{y}|x)$ is the probability that model Φ assigns to the color image \hat{y} given black and white image x.

Colorization with Latent Semantic Segmentation (TZ)



$$P_{\Phi}(\hat{y}|x) = \sum_{z} P_{\Phi}(z|x) P_{\Phi}(\hat{y}|z,x).$$

input x

$$P_{\Phi}(z|x) = \dots$$
 semantic segmentation

$$P_{\Phi}(\hat{y}|z,x) = \dots$$
 segment colorization

Maybe EM?

$$P_{\Phi}(y) = \sum_{z} P_{\Phi}(z) P_{\Phi}(y|z).$$

$$\Phi^{t+1} = \underset{\Phi}{\operatorname{argmin}} \ E_{y \sim \operatorname{Train}} \ E_{z \sim P_{\Phi}t(z|y)} \ - \ln P_{\Phi}(z,y)$$
 Update Inference

In most cases the inference is intractible!

Variational Inference:

The Evidence Lower Bound (The ELBO)

We introduce a friendly model $P_{\Psi}(z|y)$ to approximate $P_{\Phi}(z|y)$.

$$\ln P_{\Phi}(y) = E_{z \sim P_{\Psi}(z|y)} \ln P_{\Phi}(y)$$

$$= E_{z \sim P_{\Psi}(z|y)} \left(\ln P_{\Phi}(y) \frac{P_{\Phi}(z|y)}{P_{\Psi}(z|y)} + \ln \frac{P_{\Psi}(z|y)}{P_{\Phi}(z|y)} \right)$$

$$= \left(E_{z \sim P_{\Psi}(z|y)} \ln \frac{P_{\Phi}(z,y)}{P_{\Psi}(z|y)} \right) + KL(P_{\Psi}(z|y), P_{\Phi}(z|y))$$

$$= \text{ELBO} + KL(P_{\Psi}(z|y), P_{\Phi}(z|y))$$

EM is Alternating Maximization of the ELBO

ELBO =
$$E_{z \sim P_{\Psi}(z|y)} \ln \frac{P_{\Phi}(z,y)}{P_{\Psi}(z|y)}$$
 (1)
= $\ln P_{\Phi}(y) - KL(P_{\Psi}(z|y), P_{\Phi}(z|y))$ (2)

by (2)
$$\Psi^{t+1} = \underset{\Psi}{\operatorname{argmin}} E_{y \sim \operatorname{Train}} KL(P_{\Psi}(z|y), P_{\Phi^t}(z|y)) = \Phi^t$$

by (1)
$$\Phi^{t+1} = \underset{\Phi}{\operatorname{argmax}} E_{y \sim \operatorname{Train}} E_{z \sim P_{\Phi^t}(z|y)} \ln P_{\Phi}(z, y)$$

Different Ways of Writing the ELBO

ELBO =
$$E_{z \sim P_{\Psi}(z|y)} \ln \frac{P_{\Phi}(z,y)}{P_{\Psi}(z|y)}$$

= $\ln P_{\Phi}(y) - KL(P_{\Psi}(z|y), P_{\Phi}(z|y))$
= $\left(E_{z \sim P_{\Psi}(z|y)} \ln P(y|z)\right) - KL(P_{\Psi}(z|x), P_{\Phi}(z))$
= $\left(E_{z \sim P_{\Psi}(z|y)} P_{\Phi}(z,y)\right) + H(P_{\Psi}(z|y))$

Hard ELBO

Hard ELBO is to ELBO as hard EM is to EM.

$$\mathcal{L}_{\text{ELBO}}(y, \Phi, \Psi) = KL(P_{\Psi}(z|y), P_{\Phi}(z|y)) - \ln P_{\Phi}(y)$$

$$\mathcal{L}_{\text{ELBO}}(y, \Phi, \Psi) = E_{z \sim P_{\Psi}(z|y)} - \ln P_{\Phi}(z, y) + \ln P_{\Psi}(z|y)$$

$$\mathcal{L}_{\text{HELBO}}(y, \Phi, \Psi) = E_{z \sim P_{\Psi}(z|y)} - \ln P_{\Phi}(z, y)$$

Measuring the ELBO

ELBO =
$$E_{z \sim P_{\Psi}(z|y)} \ln \frac{P_{\Phi}(z,y)}{P_{\Psi}(z|y)}$$

If $P_{\Phi}(z)$, $P_{\Phi}(y|z)$, and $P_{\Psi}(z|y)$ are friendly (even when $P_{\Phi}(y)$ is not friendly) we can measure ELBO loss through sampling.

If we can measure it, we can do gradient descent on it (but perhaps with difficulty).

We want Ψ to adapt to Φ

$$\mathcal{L}_{\text{ELBO}}(y, \Phi, \Psi) = KL(P_{\Psi}(z|y), P_{\Phi}(z|y)) - \ln P_{\Phi}(y)$$

$$Q^*(z|y) = P_{\Phi}(z|y)$$

$$E_{y \sim \text{Pop}} \mathcal{L}_{\text{ELBO}}(y, \Phi, Q^*) = H(\text{Pop}, P_{\Phi})$$

However, Φ can ignore Ψ

$$\mathcal{L}_{\text{ELBO}}(y, \Phi, \Psi) = KL(P_{\Psi}(z|y), P_{\Phi}(z|y)) - \ln P_{\Phi}(y)$$

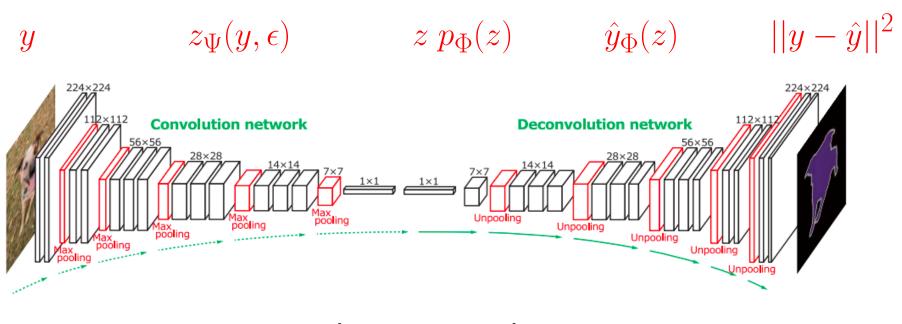
$$P^*(z) = P_{\Psi}(z)$$
$$P^*(y|z) = P_{\Phi}(y)$$

$$E_{y \sim \text{Pop}} \mathcal{L}_{\text{ELBO}}(y, P^*, \Psi) = H(\text{Pop}, P_{\Phi})$$

It seems important that $P_{\Phi}(y|z)$ have limited expressive power.

A VAE for Images

Auto-Encoding Variational Bayes, Diederik P Kingma, Max Welling, 2013.



Gaussian Distributions

$$p_{\Phi}(z) \propto \exp\left(\sum_{i} (z[i] - \mu[i])^{2} / (2\sigma[i]^{2})\right)$$

$$p_{\Phi}(y|z) \propto \exp\left(\sum_{j} (y[j] - y_{\Phi}(z)[j])^{2} / (2\gamma[j]^{2})\right)$$

$$p_{\Psi}(z|y) \propto \exp\left(\sum_{i} (z[i] - z_{\Psi}(y)[i])^{2} / (2\sigma_{\Psi}(y)[i]^{2})\right)$$

KL-Divergence Form for the ELBO

$$E_{z \in p_{\Psi}(z|y)} \ln p_{\Psi}(z|y) - \ln p_{\Phi}(z)p_{\Phi}(y|z)$$
 $\mathcal{L}_{\text{ELBO}}$

$$= KL(p_{\Psi}(z|y), p_{\Phi}(z)) + E_{z \in P_{\Psi}(z|y)} - \ln p_{\Phi}(y|z)$$

The ELBO is a KL-divergence + a cross entropy

Continuous KL-divergence is ok.

Continuous cross-entropy has issues — we will come back to that later.

Closed Form KL-Divergence

$$KL(p_{\Psi}(z|y), p_{\Phi}(z))$$

$$= \sum_{i} \frac{\sigma_{\Psi}(y)[i]^{2} + (z_{\Psi}(y)[i] - \mu[i])^{2}}{2\sigma[i]^{2}} + \ln \frac{\sigma[i]}{\sigma_{\Psi}(y)[i]} - \frac{1}{2}$$

Standardizing $p_{\Phi}(z)$

The KL-divergence term is

$$\sum_{i} \frac{\sigma_{\Psi}(y)[i]^{2} + (\boldsymbol{z}_{\Psi}(y)[i] - \boldsymbol{\mu}[i])^{2}}{2\boldsymbol{\sigma}[i]^{2}} + \ln \frac{\boldsymbol{\sigma}[i]}{\boldsymbol{\sigma}_{\Psi}(y)[i]} - \frac{1}{2}$$

We can adjust Ψ to Ψ' such that

$$z_{\Psi'}(y)[i] = z_{\Psi}(y)[i]/\sigma[i] + \mu[i]$$

$$\sigma_{\Psi'}(y)[i] = \sigma_{\Psi}(y)/\sigma[i]$$

We then get $KL(p_{\Psi}(z|y), p_{\Phi}(z)) = KL(p_{\Psi'}(z|y), \mathcal{N}(0, I)).$

Standardizing $p_{\Phi}(z)$

Without loss of generality the VAE becomes.

$$\min_{\Phi, \Psi} E_y KL(P_{\Psi}(z|y), \mathcal{N}(0, I)) + E_{z \in P_{\Psi}(z|y)} - \ln p_{\Phi}(y|z)$$

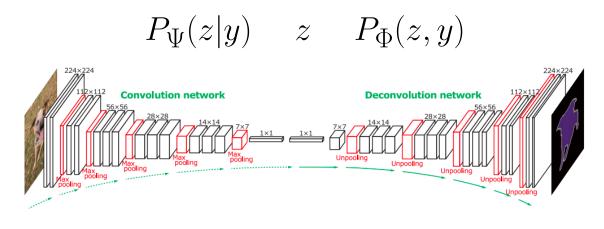
Reparameterization Trick for the Cross-Entropy

$$p_{\Psi}(z|y) \propto \exp\left(\sum_{i} (z[i] - z_{\Psi}(y)[i])^2 / (2\sigma_{\Psi}(y)[i]^2)\right)$$

$$E_{z \in p_{\Psi}(z|y)} \ln p_{\Phi}(y|z)$$

$$= E_{\epsilon \sim \mathcal{N}(0,I)} z[i] = z_{\Psi}(y)[i] + \sigma_{\Psi}(y)[i]\epsilon[i]; \quad \ln p_{\Phi}(y|z)$$

Sampling



[Hyeonwoo Noh et al.]

Sampling uses just the second half $P_{\Phi}(z, y)$.

Sampling



[Alec Radford]

Why Blurry?

A common explanation for the blurryness of images generated from VAEs is the use of L_2 as the distortion measure.

It does seem that L_1 works better.

However, training on L_2 distortion can produce sharp images in rate-distortion autoencoders.

Summary: ELBO and VAE

ELBO: $P_{\Phi}(z)$, $P_{\Phi}(y|z)$, $P_{\Psi}(z|y)$ friendly graphical models:

$$\Phi^*, \Psi^* = \underset{\Phi, \Psi}{\operatorname{argmin}} \ E_{y \sim \operatorname{Pop}, \ z \sim P_{\Psi}(z|y)} \ \ln P_{\Psi}(z|y) - \ln P_{\Phi}(z) P_{\Phi}(y|z)$$

VAE: $p_{\Phi}(z|y)$, $p_{\Phi}(y|z)$ Gaussian:

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} E_{y \sim \operatorname{Pop}} KL(p_{\Phi}(z|y), \mathcal{N}(0, I)) - E_{z \sim p_{\Phi}(z|y)} \ln p_{\Phi}(y|z)$$

\mathbf{END}