

Cycle-Consistent Adversarial Learning as Approximate Bayesian Inference

Louis C. Tiao¹ Edwin V. Bonilla² Fabio Ramos¹

July 22, 2018

¹University of Sydney, ²University of New South Wales

Motivation: Unpaired Image-to-Image Translation

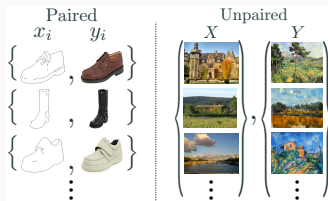
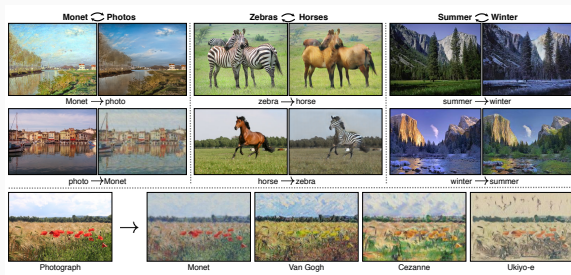


Figure 1: From Zhu et al. (2017)

Cycle-Consistent Adversarial Learning (CycleGAN)

- Introduced by Kim et al. (2017); Zhu et al. (2017)
- Forward and reverse mappings $\mathbf{m}_\phi : \mathbf{x} \mapsto \mathbf{z}$ and $\mu_\theta : \mathbf{z} \mapsto \mathbf{x}$
- Discriminators \mathbf{D}_α and \mathbf{D}_β

Distribution matching (GAN objectives)

Yield realistic outputs in the other domain.

$$\begin{aligned}\ell_{\text{GAN}}^{\text{reverse}}(\alpha; \phi) &= \mathbb{E}_{p^*(\mathbf{z})}[\log \mathbf{D}_\alpha(\mathbf{z})] + \mathbb{E}_{q^*(\mathbf{x})}[\log(1 - \mathbf{D}_\alpha(\mathbf{m}_\phi(\mathbf{x})))] , \\ \ell_{\text{GAN}}^{\text{forward}}(\beta; \theta) &= \mathbb{E}_{p^*(\mathbf{x})}[\log \mathbf{D}_\beta(\mathbf{x})] + \mathbb{E}_{p^*(\mathbf{z})}[\log(1 - \mathbf{D}_\beta(\mu_\theta(\mathbf{z})))] .\end{aligned}$$

Cycle-consistency losses

Encourage tighter correspondences—must be able to reconstruct output from input and vice versa. **May alleviate mode-collapse**

$$\begin{aligned}\ell_{\text{CONST}}^{\text{reverse}}(\theta, \phi) &= \mathbb{E}_{q^*(\mathbf{x})}[\|\mathbf{x} - \mu_\theta(\mathbf{m}_\phi(\mathbf{x}))\|_\rho^\rho] , \\ \ell_{\text{CONST}}^{\text{forward}}(\theta, \phi) &= \mathbb{E}_{p^*(\mathbf{z})}[\|\mathbf{z} - \mathbf{m}_\phi(\mu_\theta(\mathbf{z}))\|_\rho^\rho] .\end{aligned}$$

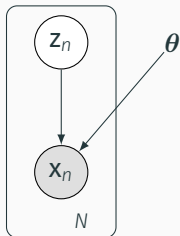
We cast the problem of learning **inter-domain correspondences without paired data** as **approximate Bayesian inference** in a **latent variable model (LVM)**.

1. We introduce **implicit latent variable models (ILVMS)**,
 - prior over latent variables specified flexibly as **implicit distribution**.
2. We develop a new variational inference (VI) algorithm based on
 - minimizing the **symmetric Kullback-Leibler (κ L)** divergence
 - between a variational and exact **joint distribution**.
3. We demonstrate that CYCLEGAN (Kim et al., 2017; Zhu et al., 2017) can be instantiated as a **special case** of our framework.

Implicit Latent Variable Models

Joint Distribution

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = \underbrace{p_{\theta}(\mathbf{x} | \mathbf{z})}_{\text{likelihood}} \underbrace{p^*(\mathbf{z})}_{\text{prior}}$$



Prescribed Likelihood

Likelihood $p_{\theta}(\mathbf{x}_n | \mathbf{z}_n)$ is **prescribed** (as usual)

Implicit Prior

Prior $p^*(\mathbf{z})$ over latent variables specified as **implicit** distribution

- Given only by a finite collection $\mathbf{Z}^* = \{\mathbf{z}_m^*\}_{m=1}^M$ of its samples,

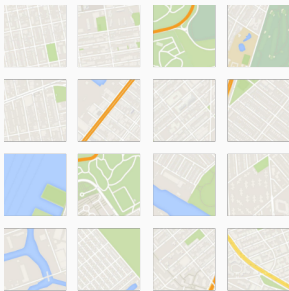
$$\mathbf{z}_m^* \sim p^*(\mathbf{z})$$

- Offers utmost degree of flexibility in treatment of prior information.

Implicit Latent Variable Models: Example

Unpaired Image-to-Image Translation

- Prior distribution $p^*(z)$ specified by images $Z^* = \{z_m^*\}_{m=1}^M$ from one domain.
- Empirical data distribution $q^*(x)$ specified by images $X^* = \{x_n\}_{n=1}^N$ from another domain.



(a) samples from $p^*(z)$



(b) a sample from $q^*(x)$

Inference in Implicit Latent Variable Models

Having specified the generative model, our aims are

- Optimize θ by maximizing marginal likelihood $p_{\theta}(\mathbf{x})$
- Infer hidden representations \mathbf{z} by computing posterior $p_{\theta}(\mathbf{z} | \mathbf{x})$

Both require **intractable** $p_{\theta}(\mathbf{x})$

- **must resort to approximate inference**

Classical Variational Inference

- Approximate **exact posterior** $p_{\theta}(\mathbf{z} | \mathbf{x})$ with **variational posterior** $q_{\phi}(\mathbf{z} | \mathbf{x})$
- Reduces **inference** problem to **optimization** problem

$$\min_{\phi} \text{KL} [q_{\phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\theta}(\mathbf{z} | \mathbf{x})]$$

Symmetric Joint-Matching Variational Inference

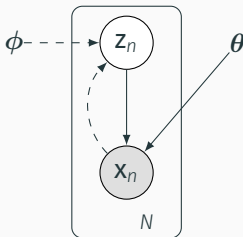
Joint-Matching Variational Inference

Variational Joint

- Consider instead directly approximating the **exact joint** with **variational joint**

$$q_{\phi}(\mathbf{x}, \mathbf{z}) = q_{\phi}(\mathbf{z} | \mathbf{x}) q^*(\mathbf{x})$$

- variational posterior $q_{\phi}(\mathbf{z} | \mathbf{x})$ also **prescribed**



Symmetric Joint-Matching Variational Inference

Minimize **symmetric KL divergence** between **joints**

$$\text{KL}_{\text{SYMM}} [p_{\theta}(\mathbf{x}, \mathbf{z}) \parallel q_{\phi}(\mathbf{x}, \mathbf{z})]$$

where

$$\text{KL}_{\text{SYMM}} [p \parallel q] = \underbrace{\text{KL} [p \parallel q]}_{\text{forward KL}} + \underbrace{\text{KL} [q \parallel p]}_{\text{reverse KL}}$$

Why?

1. Because we can:

- $\text{KL}_{\text{SYMM}} [p_{\theta}(\mathbf{x}, \mathbf{z}) \parallel q_{\phi}(\mathbf{x}, \mathbf{z})]$ **tractable**
- $\text{KL}_{\text{SYMM}} [p_{\theta}(\mathbf{z} \mid \mathbf{x}) \parallel q_{\phi}(\mathbf{z} \mid \mathbf{x})]$ **intractable**

2. Helps avoid under/over-dispersed approximations (see paper for details)

Reverse KL Variational Objective

- Minimizing **reverse** KL divergence between **joints** equivalent to maximizing usual **evidence lower bound (ELBO)**,

$$\begin{aligned}\text{KL}[q_\phi(\mathbf{x}, \mathbf{z}) \parallel p_\theta(\mathbf{x}, \mathbf{z})] &= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})} [\log q_\phi(\mathbf{x}, \mathbf{z}) - \log p_\theta(\mathbf{x}, \mathbf{z})] \\ &= \underbrace{\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})} [\log q_\phi(\mathbf{z} \mid \mathbf{x}) - \log p_\theta(\mathbf{x}, \mathbf{z})]}_{\mathcal{L}_{\text{NELBO}}(\theta, \phi)} - \underbrace{\mathbb{H}[q^*(\mathbf{x})]}_{\text{constant}}\end{aligned}$$

- Recall (negative) ELBO,

$$\mathcal{L}_{\text{NELBO}}(\theta, \phi) = \underbrace{\mathbb{E}_{q^*(\mathbf{x})q_\phi(\mathbf{z} \mid \mathbf{x})} [-\log p_\theta(\mathbf{x} \mid \mathbf{z})]}_{\mathcal{L}_{\text{NELL}}(\theta, \phi)} + \underbrace{\mathbb{E}_{q^*(\mathbf{x})} \text{KL}[q_\phi(\mathbf{z} \mid \mathbf{x}) \parallel p^*(\mathbf{z})]}_{\text{intractable}}$$

- KL term is intractable as prior $p^*(\mathbf{z})$ is unavailable—can only sample!

Forward KL Variational Objective

- Minimizing **forward** KL divergence between **joints**

$$\begin{aligned}\text{KL}[p_{\theta}(\mathbf{x}, \mathbf{z}) \parallel q_{\phi}(\mathbf{x}, \mathbf{z})] &= \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{x}, \mathbf{z})] \\ &= \underbrace{\mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{z})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z}) - \log q_{\phi}(\mathbf{x}, \mathbf{z})]}_{\mathcal{L}_{\text{NAPLBO}}(\theta, \phi)} - \underbrace{\mathbb{H}[p^*(\mathbf{z})]}_{\text{constant}}\end{aligned}$$

- New variational objective**, aggregate posterior lower bound (APLBO)

$$\mathcal{L}_{\text{NAPLBO}}(\theta, \phi) = \underbrace{\mathbb{E}_{p^*(\mathbf{z})p_{\theta}(\mathbf{x} \mid \mathbf{z})} [-\log q_{\phi}(\mathbf{z} \mid \mathbf{x})]}_{\mathcal{L}_{\text{NELP}}(\theta, \phi)} + \underbrace{\mathbb{E}_{p^*(\mathbf{z})} \text{KL}[p_{\theta}(\mathbf{x} \mid \mathbf{z}) \parallel q^*(\mathbf{x})]}_{\text{intractable}}$$

- KL term is intractable as empirical data distribution $q^*(\mathbf{x})$ is unavailable—can only sample!

Density Ratio Estimation and f -divergence Approximation

General f -divergence lower bound (Nguyen et al., 2010)

For convex lower-semicontinuous function $f : \mathbb{R}_+ \rightarrow \mathbb{R}$,

$$\underbrace{\mathbb{E}_{q^*(\mathbf{x})} \mathcal{D}_f[p^*(\mathbf{z}) \parallel q_\phi(\mathbf{z} | \mathbf{x})]}_{\text{intractable}} \geq \max_{\alpha} \underbrace{\mathcal{L}_f^{\text{latent}}(\alpha; \phi)}_{\text{tractable}},$$

where

$$\mathcal{L}_f^{\text{latent}}(\alpha; \phi) = \mathbb{E}_{q^*(\mathbf{x}) q_\phi(\mathbf{z} | \mathbf{x})} [f'(r_\alpha(\mathbf{z}; \mathbf{x}))] - \mathbb{E}_{q^*(\mathbf{x}) p^*(\mathbf{z})} [f^*(f'(r_\alpha(\mathbf{z}; \mathbf{x})))]$$

- Turns **divergence estimation** into an **optimization problem**
- Estimate divergence using a l.b. that just requires samples!
- r_α is a neural net with parameters α , with equality at

$$r_\alpha^*(\mathbf{z}; \mathbf{x}) = \frac{q_\phi(\mathbf{z} | \mathbf{x})}{p^*(\mathbf{z})}$$

KL divergence lower bound

Example: KL divergence lower bound

For $f(u) = u \log u$, we instantiate the KL lower bound

$$\underbrace{\mathbb{E}_{q^*(x)} \text{KL} [q_\phi(z|x) \parallel p^*(z)]}_{\text{intractable}} \geq \max_{\alpha} \underbrace{\mathcal{L}_{\text{KL}}^{\text{latent}}(\alpha; \phi)}_{\text{tractable}}$$

where

$$\mathcal{L}_{\text{KL}}^{\text{latent}}(\alpha; \phi) = \mathbb{E}_{q^*(x)q_\phi(z|x)}[\log r_\alpha(z; x)] - \mathbb{E}_{q^*(x)p^*(z)}[r_\alpha(z; x) - 1]$$

Yields estimate of the ELBO where all terms are tractable,

$$\begin{aligned} \mathcal{L}_{\text{NELBO}}(\theta, \phi) &= \underbrace{\mathcal{L}_{\text{NELL}}(\theta, \phi)}_{\text{tractable}} + \underbrace{\mathbb{E}_{q^*(x)} \text{KL} [q_\phi(z|x) \parallel p^*(z)]}_{\text{intractable}} \\ &\geq \max_{\alpha} \underbrace{\mathcal{L}_{\text{NELL}}(\theta, \phi)}_{\text{tractable}} + \underbrace{\mathcal{L}_{\text{KL}}^{\text{latent}}(\alpha; \phi)}_{\text{tractable}} \end{aligned}$$

CycleGAN as a Special Case

Cycle-consistency as Conditional Probability Maximization

For Gaussian likelihood and variational posterior

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{\theta}(\mathbf{z}), \tau^2 \mathbf{I}), \quad q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{m}_{\phi}(\mathbf{x}), t^2 \mathbf{I})$$

Can instantiate $\ell_{\text{CONST}}^{\text{reverse}}(\theta, \phi)$ from $\mathcal{L}_{\text{NELL}}(\theta, \phi)$

as posterior $q_{\phi}(\mathbf{z} | \mathbf{x})$ degenerates (as $t \rightarrow 0$)

Can instantiate $\ell_{\text{CONST}}^{\text{forward}}(\theta, \phi)$ from $\mathcal{L}_{\text{NELP}}(\theta, \phi)$

as likelihood $p_{\theta}(\mathbf{x} | \mathbf{z})$ degenerates (as $\tau \rightarrow 0$)

Cycle-consistency corresponds to maximizing conditional probabilities:

- **ELL.** forces $q_{\phi}(\mathbf{z} | \mathbf{x})$ to place mass on hidden representations that *recover* the **data**
- **ELP.** forces $p_{\theta}(\mathbf{x} | \mathbf{z})$ to generate observations that *recover* the **prior**

Distribution Matching as Regularization

For appropriate setting of f , and simplifying the mappings and discriminators,

- Can instantiate $\ell_{\text{GAN}}^{\text{reverse}}(\alpha; \phi)$ from $\mathcal{L}_f^{\text{latent}}(\alpha; \phi)$
- Can instantiate $\ell_{\text{GAN}}^{\text{forward}}(\beta; \theta)$ from $\mathcal{L}_f^{\text{observed}}(\beta; \theta)$

Approximately minimizes intractable divergences:

- $\mathcal{D}_f[p^*(z) \parallel q_\phi(z|x)]$ — forces $q_\phi(z|x)$ to match prior $p^*(z)$
- $\mathcal{D}_f[q^*(x) \parallel p_\theta(x|z)]$ — forces $p_\theta(x|z)$ to match data $q^*(x)$

Summary

$$\begin{aligned}\mathcal{L}_{\text{NELBO}}(\theta, \phi) &\geq \max_{\alpha} \underbrace{\mathcal{L}_{\text{NELL}}(\theta, \phi)}_{\ell_{\text{CONST}}^{\text{reverse}}(\theta, \phi)} + \underbrace{\mathcal{L}_{\text{KL}}^{\text{latent}}(\alpha; \phi)}_{\ell_{\text{GAN}}^{\text{reverse}}(\alpha; \phi)} \\ \mathcal{L}_{\text{NAPLBO}}(\theta, \phi) &\geq \max_{\beta} \underbrace{\mathcal{L}_{\text{NELP}}(\theta, \phi)}_{\ell_{\text{CONST}}^{\text{forward}}(\theta, \phi)} + \underbrace{\mathcal{L}_{\text{KL}}^{\text{observed}}(\beta; \theta)}_{\ell_{\text{GAN}}^{\text{forward}}(\beta; \theta)}\end{aligned}$$

Conclusion

- Formulated **implicit latent variable models**, which introduces **implicit prior** over latent variables
 - Offers utmost degree of flexibility in incorporating prior knowledge
- Developed new paradigm for variational inference
 - directly approximates exact **joint distribution**
 - minimizes the **symmetric KL divergence**
- Provided theoretical treatment of the links between **CycleGAN methods** and **Variational Bayes**

Poster Session

To find out more, come visit us at our poster!

Poster #14, Session 4 (17:10-18:00 Saturday, 14 July)

Questions?

References

- Kim, T., Cha, M., Kim, H., Lee, J. K., and Kim, J. (2017). Learning to Discover Cross-Domain Relations with Generative Adversarial Networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, volume 70, pages 1857–1865.
- Nguyen, X., Wainwright, M. J., and Jordan, M. I. (2010). Estimating Divergence Functionals and the Likelihood Ratio by Convex Risk Minimization. *IEEE Trans. Information Theory*, 56(11):5847–5861.
- Zhu, J.-Y., Park, T., Isola, P., and Efros, A. A. (2017). Unpaired image-to-image translation using cycle-consistent adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*.

Symmetric Joint-Matching KL Minimization i

- KL divergence is **asymmetric** $\text{KL}[p \parallel q] \neq \text{KL}[q \parallel p]$
- $\text{KL}[q_\phi(\mathbf{z} \mid \mathbf{x}) \parallel p_\theta(\mathbf{z} \mid \mathbf{x})]$ (reverse) **underestimates** support
- $\text{KL}[p_\theta(\mathbf{z} \mid \mathbf{x}) \parallel q_\phi(\mathbf{z} \mid \mathbf{x})]$ (forward) **overestimates** support
- Consider **symmetric** KL: $\text{KL}_{\text{SYMM}}[p \parallel q] = \text{KL}[p \parallel q] + \text{KL}[q \parallel p]$
- Forward KL involves expectation under intractable posterior $p_\theta(\mathbf{z} \mid \mathbf{x})$ —what we're trying to approximate in the first place

$$\text{KL}[p_\theta(\mathbf{z} \mid \mathbf{x}) \parallel q_\phi(\mathbf{z} \mid \mathbf{x})] = \mathbb{E}_{p_\theta(\mathbf{z} \mid \mathbf{x})} \left[\log \frac{p_\theta(\mathbf{z} \mid \mathbf{x})}{q_\phi(\mathbf{z} \mid \mathbf{x})} \right]$$

Symmetric Joint-Matching KL Minimization ii

- Can show

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z | x) \parallel p_{\theta}(z | x)] = \arg \min_{\phi} \text{KL} [q_{\phi}(x, z) \parallel p_{\theta}(x, z)]$$

$$\arg \min_{\phi} \text{KL} [p_{\theta}(z | x) \parallel q_{\phi}(z | x)] = \arg \min_{\phi} \text{KL} [p_{\theta}(x, z) \parallel q_{\phi}(x, z)]$$

- Already showed

$$\arg \max_{\phi} \mathcal{L}_{\text{ELBO}}(\theta, \phi) = \arg \min_{\phi} \text{KL} [q_{\phi}(x, z) \parallel p_{\theta}(x, z)]$$

- Can we find something similar for $\text{KL} [p_{\theta}(x, z) \parallel q_{\phi}(x, z)]$?