

Deep Information Bottleneck

论文列表

- Diederik P Kingma and Max Welling. **Auto-encoding variational Bayes**. In ICLR, 2014.
- Fischer, I., Dillon, J.V., and Murphy, K. "**Deep Variational Information Bottleneck** [arXiv]." ArXiv (2016): 13 Pp. Web.
- Peng X B , Kanazawa A , Toyer S , et al. **Variational Discriminator Bottleneck: Improving Imitation Learning, Inverse RL, and GANs by Constraining Information Flow**[J]. 2018.on Bottleneck [arXiv]." ArXiv (2016): 13 Pp. Web.

目录

- 基本概念
 - 互信息
 - *KL*散度
 - 信息瓶颈
- 深度变分自编码器 (varitional auto-encoder)
- 深度变分信息瓶颈 (deep varitional information bottleneck)
- 信息瓶颈判别器 (varitional discriminator bottleneck)

互信息

- 设两个随机变量 (X, Y) 的联合分布为 $p(x, y)$, 边缘分布为 $p(x)$, $p(y)$, 互信息 $I(X; Y)$ 是联合分布 $p(x, y)$ 与边缘分布 $p(x)p(y)$ 的相对熵。

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

- 互信息可以看成是一个随机变量中包含的关于另一个随机变量的信息量, 或者说是一个随机变量由于已知另一个随机变量而减少的不肯定性。

$$I(X; Y) = H(X) - H(X | Y)$$

KL 散度与相对熵

- 设 $P(x), Q(x)$ 是随机变量 X 上的两个概率分布，则在离散和连续随机变量的情形下，相对熵的定义分别为：

$$\text{KL}(P\|Q) = \sum P(x) \log \frac{P(x)}{Q(x)}$$

$$\text{KL}(P\|Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx$$

- 非负性，非对称性

$$\text{KL}(P\|Q) \geq 0$$

$$\text{KL}(P\|Q) \neq \text{KL}(Q\|P)$$

信息瓶颈

- 相关信息

- 输入信号 $x \in X$ 的相关信息指信号 $x \in X$ 提供的和另一个信号 $y \in Y$ 有关的信息 \bar{x} 。

$$x \rightarrow \bar{x} \quad \bar{x} \rightarrow y$$

- 一方面，算法应该尽可能压缩原始信息，以互信息来衡量，即 $I(x, \bar{x})$ 应该尽可能小；另一方面，算法应该使得码本保留尽可能多的相关信息，即 $I(\bar{x}, y)$ 应该尽可能大。

$$I(x, \bar{x}) - \beta I(\bar{x}, y)$$

变分自编码器 (VAE)

- 。VAE 假设隐向量 z 从先验 分布 $p(z)$ 中采样得到，输入数据 x 从分布 $p(x|z)$ 中产生，通过变分方法最大化似然 函数从而实现输入数据的拟合。似然函数如下：

$$\log p(x_1, \dots, x_n) = \sum_{i=1}^n \log p(x_i)$$

变分自编码器 (VAE)

- 单个样本 x

$$\begin{aligned}\log p(x) &= \log p(x, z) - \log p(z | x) \\ &= \log \frac{p(x, z)}{q(z)} - \log \frac{p(z | x)}{q(z)} \\ &= \int q(z) \log \frac{p(x, z)}{q(z)} dz - \int q(z) \log \frac{p(z | x)}{q(z)} dz \\ &= L(q, x) + D_{KL}(q(z) \| p(z | x))\end{aligned}$$

- 在 VAE 中, $q(z)$ 被构造为一个包含编码参数 φ 的后验概率, 并服从高斯分布:

$$q(z) = q_{\varphi}(z | x) = N(z; \mu_z(x, \varphi), \sigma_z^2(x, \varphi))$$

变分自编码器 (VAE)

$$\log p(x) \geq L(q, x)$$

$$\begin{aligned} & \int q(z) \log \frac{p(x, z)}{q(z)} dz \\ &= E_{q(z)} [\log p(x \mid z)] - D_{KL}(q(z) \parallel p(z)) \\ &= -D_{KL}(q_{\varphi}(z \mid x) \parallel p(z)) + E_{q_{\varphi}(Z \mid x)} [\log p(x \mid z)] \end{aligned}$$

深度变分信息瓶颈 (DIB)

- 将神经网络中间层的输出 Z 当作输入 X 的编码，以 θ 表示神经网络的参数，则这一编码过程可以表示为 $p(z|x;\theta)$ 。以 Y 表示神经网络的输出，则根据信息瓶颈理论，神经网络的训练过程就可以最大化如下函数

$$R_{IB}(\theta) = I(Z, Y; \theta) - \beta I(Z, X; \theta)$$

- 假设： X, Y, Z 的联合概率分布

$$p(X, Y, Z) = p(Z | X, Y)p(Y | X)p(X) = p(Z | X)p(Y | X)p(X)$$

深度变分信息瓶颈 (DIB)

(1) $I(Z, Y; \theta)$

$$I(Z, Y) = \int dydz p(y, z) \log \frac{p(y, z)}{p(y)p(z)} = \int dydz p(y, z) \log \frac{p(y | z)}{p(y)}$$

- $p(y|z)$ 即是神经网络的解码器。通常情况下, $p(y|z)$ 是难以直接计算的, 使用变分方法来估计 $I(Z, Y)$ 的下界。使用 $q(y|z)$ 作为 $p(y|z)$ 的变分估计, 则根据 KL 散度的非负性, 有

$$\int dy p(y | z) \log p(y | z) \geq \int dy p(y | z) \log q(y | z)$$

深度变分信息瓶颈 (DIB)

- 得到 $I(Z, Y)$ 的一个下界:

$$\begin{aligned} I(Z, Y) &\geq \int dydz p(y, z) \log \frac{q(y | z)}{p(y)} \\ &= \int dydz p(y, z) \log q(y | z) - \int dy p(y) \log p(y) \\ &= \int dydz p(y, z) \log q(y | z) + H(Y) \\ p(y, z) &= \int dx p(x, y, z) = \int dx p(x) p(y | x) p(z | x) \\ I(Z, Y) &\geq \int dx dy dz p(x) p(y | x) p(z | x) \log q(y | z) \end{aligned}$$

深度变分信息瓶颈 (DIB)

(2) $I(Z, X; \theta)$

$$I(Z, X) = \int dz dx p(x, z) \log \frac{p(z|x)}{p(z)}$$

$$= \int dz dx p(x, z) \log p(z | x) - \int dz p(z) \log p(z)$$

- 通常情况下, $p(z)$ 是比较难以计算的, 因此使用 $r(z)$ 作为 $p(z)$ 的变分估计, 根据 KL 散度的非负性, 有:

$$\int dz p(z) \log p(z) \geq \int dz p(z) \log r(z)$$

- 得到 $I(Z, X)$ 的一个上界: $I(Z, X) \leq \int dx dz p(x) p(z | x) \log \frac{p(z | x)}{r(z)}$

深度变分信息瓶颈 (DIB)

- 结合 $I(Z, Y)$ 的下界和 $I(Z, X)$ 的上界可以得到 $RIB(\theta)$ 的一个下界

$$R_{IB}(\theta) = I(Z, Y; \theta) - \beta I(Z, X; \theta)$$

$$\geq \int dx dy dz p(x) p(y | x) p(z | x) \log q(y | z)$$

$$- \beta \int dx dz p(x) p(z | x) \log \frac{p(z | x)}{r(z)} = L$$

- 在神经网络的训练过程中，输入 X 和输出 Y 都是已知的。假设训练数据表示为 $\{X, Y\}$ ，每个样本表示为 $\{x_i, y_i\}$ ，则 X, Y 的联合概率密度

$$p(x, y) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(x) \delta_{y_i}(y)$$

$$\delta_{x_i}(x) = \begin{cases} 1 & \text{if } x_i = x \\ 0 & \text{else} \end{cases}$$

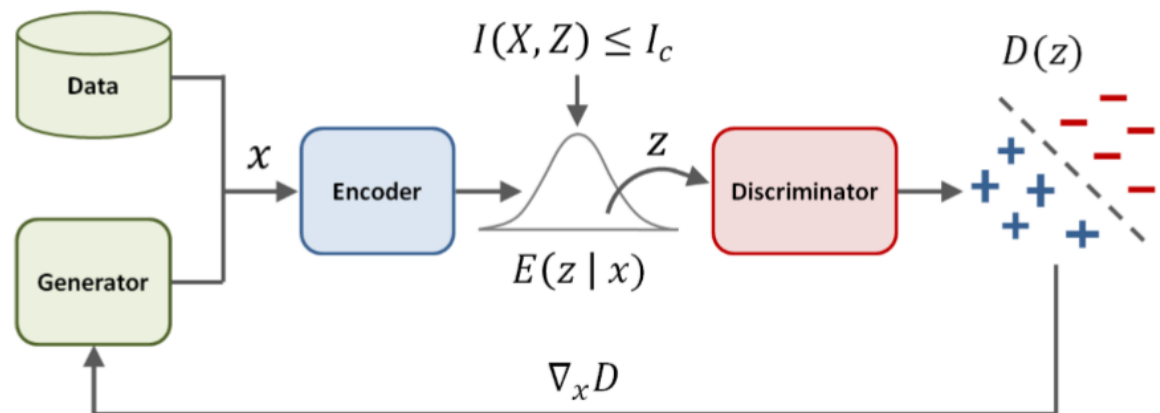
深度变分信息瓶颈 (DIB)

$$\begin{aligned} L &= \frac{1}{N} \sum_{i=1}^N \int dz p(z | x_i) \log q(y_i | z) - \beta p(z | x_i) \log \frac{p(z|x_i)}{r(z)} \\ &= \frac{1}{N} \sum_{i=1}^N E_{z \sim p(z|x_i)} \log q(y_i | z) - \beta D_{KL}(p(z | x_i) || r(z)) \end{aligned}$$

- 和VAE相比

$$L(q, x, \varphi) = -D_{KL}(q_{\varphi}(z | x) || p(z)) + E_{q_{\varphi}(z|x)}[\log p(x | z)]$$

信息瓶颈判别器 (VDB)



- 带约束的变分判别器

$$\begin{aligned} \max_G \min_D \quad & \mathbb{E}_{\mathbf{x} \sim p^*(\mathbf{x})} [-\log(D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{x})} [-\log(1 - D(\mathbf{x}))] \\ J(D, E) = \min_{D, E} \quad & \mathbb{E}_{\mathbf{x} \sim p^*(\mathbf{x})} [\mathbb{E}_{\mathbf{z} \sim E(\mathbf{z}|\mathbf{x})} [-\log(D(\mathbf{z}))]] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{x})} [\mathbb{E}_{\mathbf{z} \sim E(\mathbf{z}|\mathbf{x})} [-\log(1 - D(\mathbf{z}))]] \\ \text{s.t.} \quad & \mathbb{E}_{\mathbf{x} \sim \tilde{p}(\mathbf{x})} [\text{KL}[E(\mathbf{z} | \mathbf{x}) \| r(\mathbf{z})]] \leq I_c \end{aligned}$$

- 引入拉格朗日算子 β ,

$$\begin{aligned} J(D, E) = \min_{D, E} \max_{\beta \geq 0} \quad & \mathbb{E}_{\mathbf{x} \sim p^*(\mathbf{x})} [\mathbb{E}_{\mathbf{z} \sim E(\mathbf{z}|\mathbf{x})} [-\log(D(\mathbf{z}))]] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{x})} [\mathbb{E}_{\mathbf{z} \sim E(\mathbf{z}|\mathbf{x})} [-\log(1 - D(\mathbf{z}))]] \\ & + \beta (\mathbb{E}_{\mathbf{x} \sim \tilde{p}(\mathbf{x})} [\text{KL}[E(\mathbf{z} | \mathbf{x}) \| r(\mathbf{z})]] - I_c) \end{aligned}$$

信息瓶颈判别器 (VDB)

- 自适应更新 β

$$D, E \leftarrow \arg \min_{D, E} \mathcal{L}(D, E, \beta)$$

$$\beta \leftarrow \max(0, \beta + \alpha_{\beta} (\mathbb{E}_{\mathbf{x} \sim \tilde{p}(\mathbf{x})} [\text{KL}[E(\mathbf{z} \mid \mathbf{x}) \parallel r(\mathbf{z})]] - I_c))$$

- 优化生成器

$$\max_G \mathbb{E}_{\mathbf{x} \sim G(\mathbf{x})} [-\log(1 - D(\mu_E(\mathbf{x})))]$$

DVB用于图像生成 (VGAN)

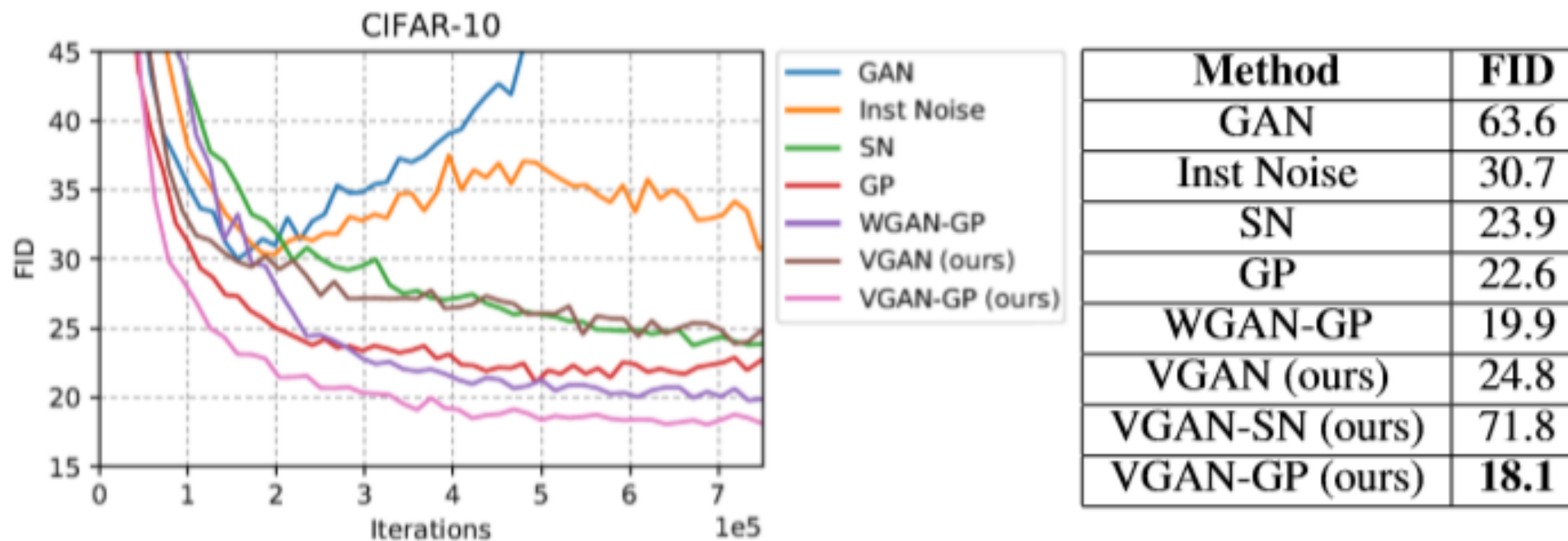


Figure 8: Comparison of VGAN and other methods on CIFAR-10, with performance evaluated using the Fréchet Inception Distance (FID).

DVB用于图像生成 (VGAN)



Figure 9: VGAN samples on CIFAR-10, CelebA 128×128 , and CelebAHQ 1024×1024 .

稳定性分析

| Method | Local convergence (a.c. case) | Local convergence (general case) |
|---|-------------------------------------|--|
| unregularized (Goodfellow et al., 2014) | ✓ | ✗ |
| WGAN (Arjovsky et al., 2017) | ✗ | ✗ |
| WGAN-GP (Gulrajani et al., 2017) | ✗ | ✗ |
| DRAGAN (Kodali et al., 2017) | ✓ | ✗ |
| Instance noise (Sønderby et al., 2016) | ✓ | ✓ |
| ConOpt (Mescheder et al., 2017) | ✓ | ✓ |
| Gradient penalties (Roth et al., 2017) | ✓ | ✓ |
| Gradient penalty on real data only | ✓ | ✓ |
| Gradient penalty on fake data only | ✓ | ✓ |

Table 1. Convergence properties of different GAN training algorithms for general GAN-architectures. Here, we distinguish between the case where both the data and generator distributions are absolutely continuous (a.c.) and the general case where they may lie on lower dimensional manifolds.

Mescheder L , Geiger A , Nowozin S . Which Training Methods for GANs do actually Converge?[J]. 2018.