

Variational Interaction Information Maximization for Cross-domain Disentanglement

Reporter: 陈思玉

2024.3.23



Author



HyeongJoo Hwang

KAIST

在 ai.kaist.ac.kr 的电子邮件经过验证

关注

创建我的个人资料

标题	引用次数	年份
Demodice: Offline imitation learning with supplementary imperfect demonstrations CCF none	60	2021
GH Kim, S Seo, J Lee, W Jeon, HJ Hwang, H Yang, KE Kim International Conference on Learning Representations		
Variational interaction information maximization for cross-domain disentanglement CCF A	40	2020
HJ Hwang, GH Kim, S Hong, KE Kim Advances in Neural Information Processing Systems 33, 22479-22491		
Multi-view representation learning via total correlation objective CCF A	29	2021
HJ Hwang, GH Kim, S Hong, KE Kim Advances in Neural Information Processing Systems 34, 12194-12207		
Regularized Behavior Cloning for Blocking the Leakage of Past Action Information CCF A		2024
S Seo, HJ Hwang, H Yang, KE Kim Advances in Neural Information Processing Systems 36		
Information-theoretic state space model for multi-view reinforcement learning CCF A	2023	
HJ Hwang, S Seo, Y Jang, S Kim, GH Kim, S Hong, KE Kim		
사전학습을 간선을 수행하는 교차 엔트로피 계획법 CCF none	2020	
황형주, 장영수, 박재영, 김기웅 정보과학회논문지 47 (1), 88-94		

引用次数



Cross-Domain Disentanglement

定义

跨域解耦 (Cross-Domain Disentanglement) 从两个域中将表征划分为：

- 域不变 (Domain-Invariant) 表征
- 域特定 (Domain-Specific) 表征

举例

MNIST-CDCB 数据集

- 域不变表征：数字
- 域特定表征：前景、背景



Interaction Information Auto-Encoder

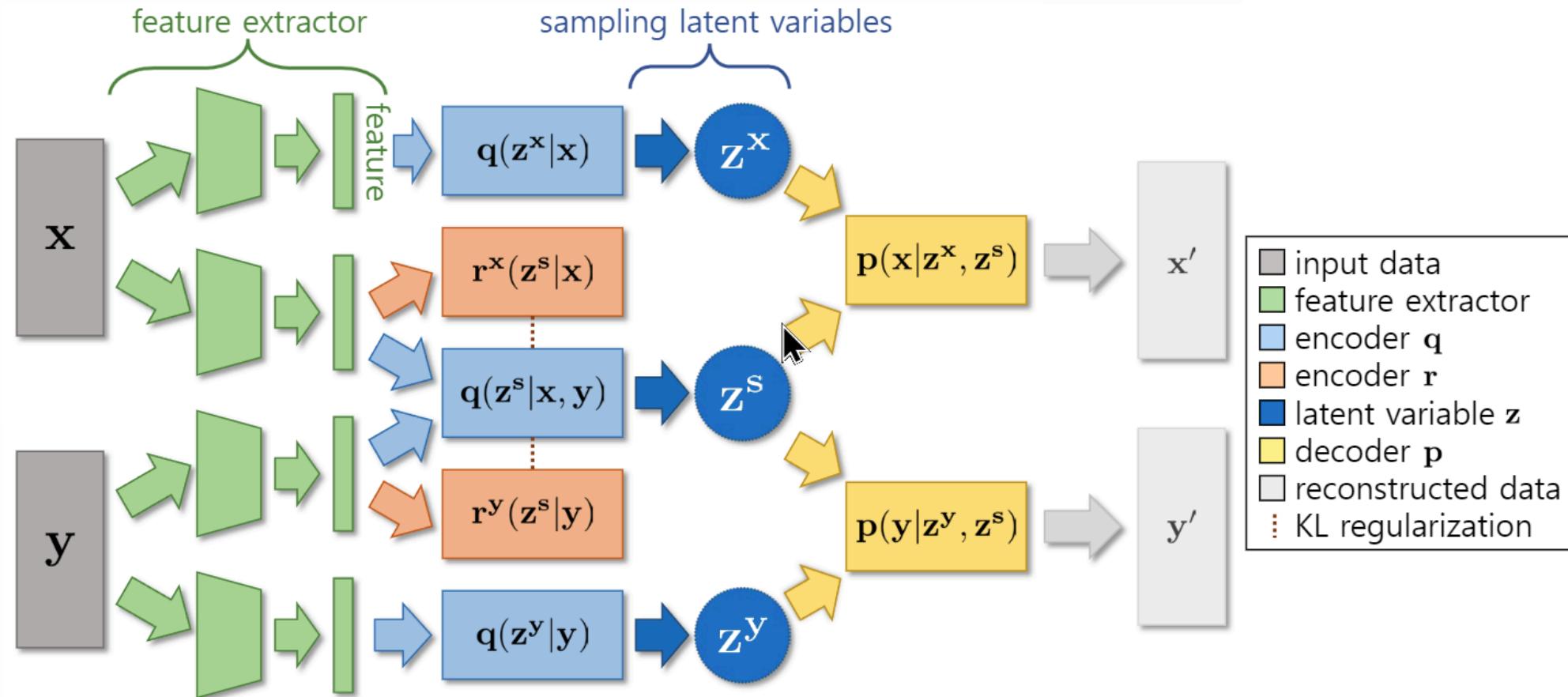


Figure 2: The architecture of Interaction Information Auto-Encoder.

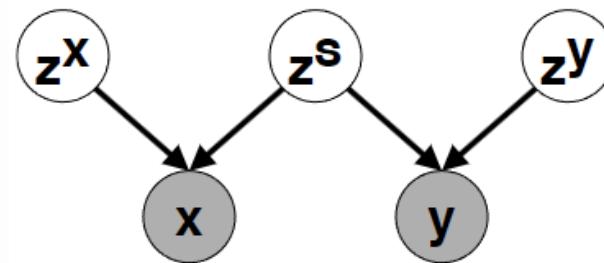
Formulation

未知联合分布 $(x, y) \sim p_D(x, y)$, 其中 x 和 y 来自不同的域 X 和 Y 。

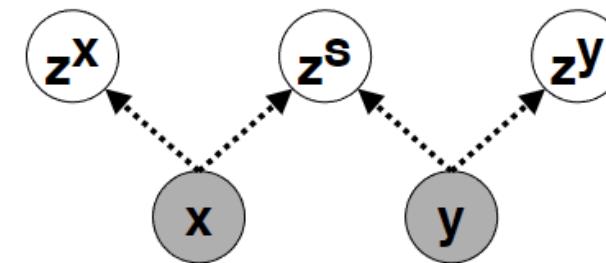
两个域存在特定于域的因素，同时存在一些共同的因素。

例如， x 和 y 可以是相同内容但风格不同的（例如草图和照片）的图像

问题：从给定的数据中解耦出三个部分：域特定特征 Z^X, Z^Y 和域共享特征 Z^S



(a) Generative model p_θ



(b) Approximate inference model q_ϕ

Figure 1: Graphical models for cross-domain disentanglement.

Formulation

Goal

$$\begin{aligned}
 p_{\theta}(x, y) &= \int dz^x dz^s dz^y p_{\theta_X}(x|z^x, z^s) p_{\theta_Y}(y|z^y, z^s) p(z^x) p(z^s) p(z^y) \\
 \Rightarrow q_{\phi}(z^x, z^s, z^y|x, y) &\underset{\text{近似}}{\approx} p_{\theta}(z^x, z^s, z^y|x, y) \\
 \Rightarrow \log p(x, y) &\geq \mathbb{E}_{q(z^x, z^s, z^y|x, y)} \left[\log \frac{p(x, y, z^x, z^s, z^y)}{q(z^x, z^s, z^y|x, y)} \right] \\
 &= \mathbb{E}_{q(z^x|x)q(z^s|x, y)} [\log p(x|z^x, z^s)] + \mathbb{E}_{q(z^y|y)q(z^s|x, y)} [\log p(y|z^y, z^s)] \\
 &\quad - D_{KL}[q(z^x|x)\|p(z^x)] - D_{KL}[q(z^y|y)\|p(z^y)] \\
 &\quad - D_{KL}[q(z^s|x, y)\|p(z^s)]
 \end{aligned}$$

不能满足：

- Z^S 与 Z^X 、 Z^Y 互斥
- Z^X, Z^Y 只包含域特定特征， Z^S 只包含域共享特征

Formulation

保证 Z^S 与 Z^X 、 Z^Y 互斥

$$\min I(Z^X; Z^S) = -I(X; Z^X, Z^S) + I(X; Z^X) + I(X; Z^S)$$

保证 Z^X, Z^Y 只包含域特定特征， Z^S 只包含域共享特征

$$\begin{aligned} \max I(X; Y; Z^S) &= I(X; Z^S) - I(X; Z^S|Y) \\ &= I(Y; Z^S) - I(Y; Z^S|X) \end{aligned}$$

联合优化

$$\begin{aligned} &\max_q I(X; Y; Z^S) - I(Z^X; Z^S) \\ &= \underbrace{I(X; Z^S) - I(X; Z^S|Y)}_{I(X; Y; Z^S)} + \underbrace{I(X; Z^X, Z^S) - I(X; Z^X) - I(X; Z^S)}_{-I(Z^X; Z^S)} \\ &= I(X; Z^X, Z^S) - I(X; Z^X) - I(X; Z^S|Y) \end{aligned} \tag{8}$$

Formulation

$$I(X; Z^X, Z^S)$$

$$q(x|z^x, z^s) = \frac{q(z^x, z^s|x) \underbrace{p_D(x)}_{\int p_D(x,y) q(z^x, z^s|x,y) dx dy}}{\int \underbrace{p_D(x,y)}_{q(x|z^x, z^s)} q(z^x, z^s|x,y) dx dy} \Rightarrow \text{难以计算}$$

$$\Rightarrow p(x|z^x, z^s) \xrightarrow{\text{计算}} q(x|z^x, z^s)$$

$$\begin{aligned} & I(X; Z^X, Z^S) \\ &= \mathbb{E}_{q(z^x, z^s|x)p_D(x)} \left[\log \frac{q(x|z^x, z^s)}{p_D(x)} \right] \\ &= H(X) + \mathbb{E}_{q(z^x, z^s|x)p_D(x)} [\log p(x|z^x, z^s)] + \mathbb{E}_{q(z^x, z^s)} [D_{KL}[q(x|z^x, z^s) \| p(x|z^x, z^s)]] \\ &\geq H(X) + \mathbb{E}_{q(z^x, z^s|x)p_D(x)} [\log p(x|z^x, z^s)] \\ &= H(X) + \mathbb{E}_{p_D(x,y)q(z^x|x)q(z^s|x,y)} [\log p(x|z^x, z^s)] \end{aligned}$$

Formulation

$$-I(X; Z^X)$$

$$\begin{aligned} q(z^x) &= \int \underbrace{p_D(x)}_{q(z^s|x)} q(z^s|x) dx \quad \Rightarrow \quad \text{难以计算} \\ \Rightarrow -\mathbb{E}_{p_D(x)}[D_{KL}[q(z^x|x) \| p(z^x)]] \end{aligned}$$

Formulation

$$-I(X; Z^S | Y)$$

$$\begin{aligned} q(z^s|y) &= \int \underbrace{p_D(x|y)}_{\Rightarrow \text{计算}} q(z^s|x, y) dx \quad \Rightarrow \quad \text{难以计算} \\ \Rightarrow r^y(z^s|y) &\quad \text{计算} \quad q(z^s|y) \end{aligned}$$

$$\begin{aligned} &-I(X; Z^S | Y) \\ &= -\mathbb{E}_{p_D(x,y)q(z^s|x,y)} \left[\log \frac{q(z^s|x,y)}{q(z^s|y)} \right] \\ &= -\mathbb{E}_{p_D(x,y)q(z^s|x,y)} \left[\log \frac{q(z^s|x,y)r^y(z^s|y)}{r^y(z^s|y)q(z^s|y)} \right] \\ &= -\mathbb{E}_{p_D(x,y)} [D_{KL}[q(z^s|x,y) \| r^y(z^s|y)]] + \mathbb{E}_{p_D(y)} [D_{KL}[q(z^s|y) \| r^y(z^s|y)]] \\ &\geq -\mathbb{E}_{p_D(x,y)} [D_{KL}[q(z^s|x,y) \| r^y(z^s|y)]] \end{aligned}$$

Formulation

最终目标

$$\begin{aligned}
 & (I(X; Y; Z^S) - I(Z^X; Z^S)) + (I(X; Y; Z^S) - I(Z^Y; Z^S)) \\
 = & 2 \cdot I(X; Y; Z^S) - I(Z^X; Z^S) - I(Z^Y; Z^S) \\
 = & I(X; Z^X, Z^S) + I(Y; Z^Y, Z^S) - I(X; Z^X) - I(Y; Z^Y) - I(X; Z^S|Y) - I(Y; Z^S|X) \\
 \geq & \mathbb{E}_{p_D(x,y)} \left[\mathbb{E}_{q(z^s|x,y)q(x^x|x)} [\log p(x|z^x, z^s)] + \mathbb{E}_{q(z^s|x,y)q(z^y|y)} [\log p(y|z^y, z^s)] \right] \\
 & - \mathbb{E}_{p_D(x,y)} [D_{KL}[q(z^x|x) \| p(z^x)] + D_{KL}[q(z^y|y) \| p(z^y)]] \\
 & - \mathbb{E}_{p_D(x,y)} [D_{KL}[q(z^s|x,y) \| r^y(z^s|y)] + D_{KL}[q(z^s|x,y) \| r^x(z^s|x)]] \\
 & + H(X) + H(Y)
 \end{aligned}$$

Interaction Information Auto-Encoder

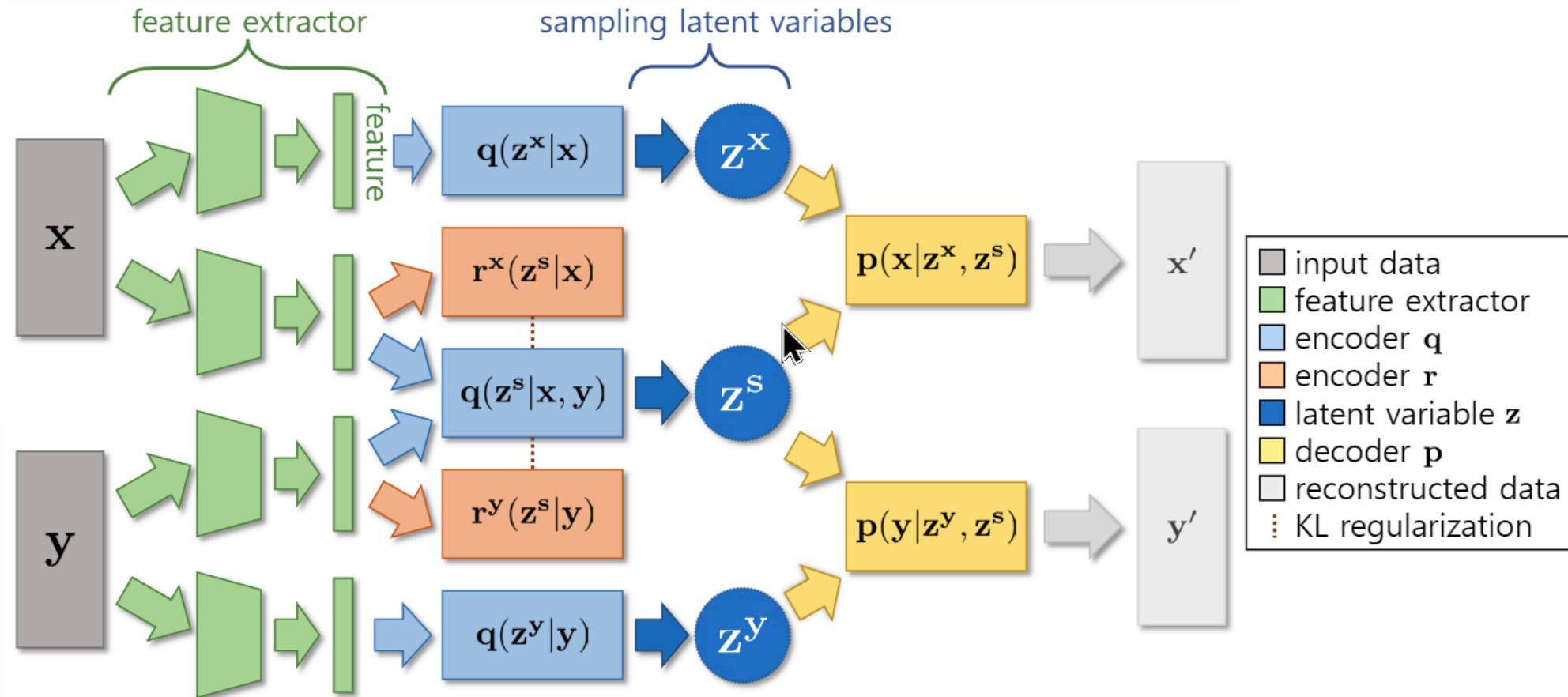


Figure 2: The architecture of Interaction Information Auto-Encoder.

LOSS

$$\begin{aligned}
 & \max_{p,q} \mathbb{E}_{q(z^x, z^s, z^y, x, y)} \left[\log \frac{p(x, y, z^x, z^s, z^y)}{q(z^x, z^s, z^y | x, y)} \right] + \lambda(2 \cdot I(X; Y; Z^S) - I(Z^X; Z^S) - I(Z^Y; Z^S)) \\
 \geq & \max_{p,q,r} (1 + \lambda) \cdot \mathbb{E}_{p_D(x,y)} [ELBO(p, q)] \\
 & + \lambda \cdot \mathbb{E}_{p_D(x,y)} [D_{KL}[q(z^s|x, y) \| p(z^s)]] \\
 & - \lambda \cdot \mathbb{E}_{p_D(x,y)} [D_{KL}[q(z^s|x, y) \| r^y(z^s|y)] + D_{KL}[q(z^s|x, y) \| r^x(z^s|x)]]
 \end{aligned}$$

Experiments

Cross-domain Image Translation

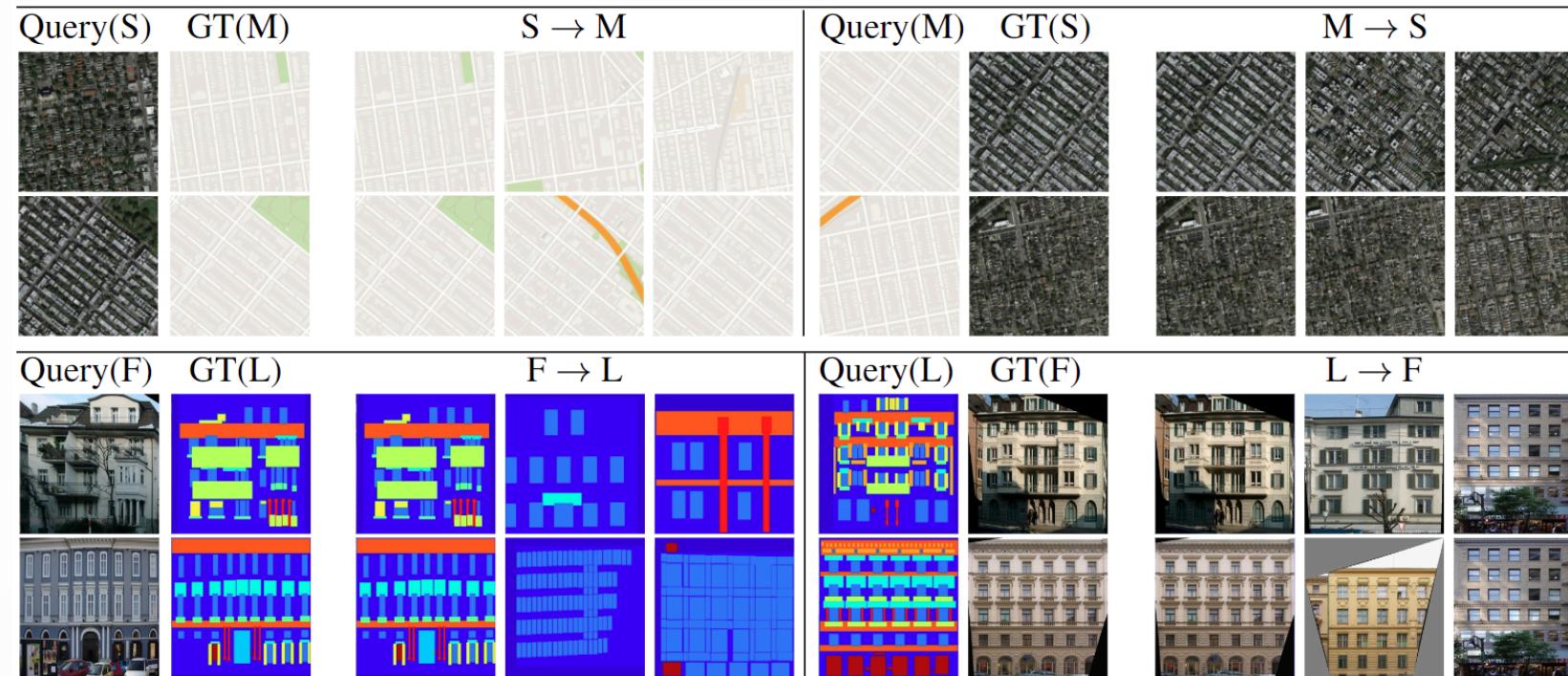
$X \rightarrow Y$						$Y \rightarrow X$						
Input	Outputs w/ different z^y					μ^y	Input	Outputs w/ different z^x				
	x	$z_1^y, z_2^y, z_3^y \sim p(z^y)$					y	$z_1^x, z_2^x, z_3^x \sim p(z^x)$				
4	4	4	4	4	4	4	4	4	4	4	4	
3	3	3	3	3	3	3	3	3	3	3	3	
8	8	8	8	8	8	8	8	8	8	8	8	
5	5	5	5	5	5	5	5	5	5	5	5	
6	6	6	6	6	6	6	6	6	6	6	6	



Cross-Domain retrieval

Table 2: Shared (exclusive) representation based retrieval on MNIST-CDCB [12], Maps [17], and Facades [41] dataset. CD/CB stand for colored digit/background, S/M stand for satellite/map, and F/L stand for facade/label respectively.

Dataset	MNIST-CDCB		Maps		Facades	
	Models	CD → CB	CB → CD	S → M	M → S	F → L
DRIT [26]	-	-	33.8 (0.09)	37.3 (0.09)	31.1 (0.94)	44.3 (0.94)
CdDN [12]	99.6 (0.0)	99.6 (0.0)	91.4 (0.18)	96.9 (0.09)	84.9 (0.94)	89.6 (0.0)
IIAE	99.7 (0.01)	99.7 (0.01)	96.6 (0.09)	97.3 (0.0)	96.2 (0.94)	99.1 (0.94)



Experiments

Zero-shot Sketch based Image Retrieval

Table 3: Evaluation on the Sketchy Extended dataset [29, 37]. WordEmb stands for word embedding.

Models	Feature Dimension	Evaluation metric		External knowledge		
		mAP	P@100	Attribute	WordEmb.	WordNet [33]
SAE [23]	300	0.216	0.293	✓	✓	-
FRWGAN [9]	512	0.127	0.169	✓	-	-
ZSIH [38]	64	0.258	0.342	-	✓	-
CAAE [22]	4096	0.196	0.284	-	-	-
SEM-PCYC [6]	64	0.349	0.463	-	✓	✓
LCALE [27]	64	0.476	0.583	-	✓	-
IIAE	64	0.573	0.659	-	-	-



Figure 4: Top-5 ZS-SBIR samples from IIAE on the Sketchy Extended dataset. Sketches in the first and seventh columns are queries and rest are retrieved candidates (Top-1 to 5 from the left to the right). Green checkmark indicates correct retrieval, whereas red crossmark indicates wrong retrieval.