Supplementary Material: Improved Multimodal Deep Learning with Variation of Information

Kihyuk Sohn, Wenling Shang and Honglak Lee University of Michigan Ann Arbor, MI, USA {kihyuks,shangw,honglak}@umich.edu

S1 Derivation of Equation (4)

The NLL objective function can be written as

$$
2\mathcal{L}^{NLL}(\theta) = -2\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X, Y)]
$$

\n
$$
= -\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X|Y) + \log P_{\theta}(Y)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(Y|X) + \log P_{\theta}(X)]
$$

\n
$$
= -\mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X|Y) + \log P_{\theta}(Y|X)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X) + \log P_{\theta}(Y)]
$$

\n
$$
= \mathcal{L}^{VI}(\theta) - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(X)] - \mathbb{E}_{P_{\mathcal{D}}} [\log P_{\theta}(Y)]
$$
(S1)

$$
= \mathcal{L}^{VI}(\theta) + \underbrace{\mathbb{E}_{P_{\mathcal{D}}}\left[\log \frac{P_{\mathcal{D}}(X)}{P_{\theta}(X)}\right]}_{KL(P_{\mathcal{D}}(X)||P_{\theta}(X))} + \underbrace{\mathbb{E}_{P_{\mathcal{D}}}\left[\log \frac{P_{\mathcal{D}}(Y)}{P_{\theta}(Y)}\right]}_{KL(P_{\mathcal{D}}(Y)||P_{\theta}(Y))}
$$
(S2)

$$
-\mathbb{E}_{P_{\mathcal{D}}}[\log P_{\mathcal{D}}(X)] - \mathbb{E}_{P_{\mathcal{D}}}[\log P_{\mathcal{D}}(Y)]
$$

= $\mathcal{L}^{\text{VI}}(\theta) + KL(P_{\mathcal{D}}(X)||P_{\theta}(X)) + KL(P_{\mathcal{D}}(Y)||P_{\theta}(Y)) + C_1$ (S3)

where Equation [\(S1\)](#page-0-0) holds by the definition of $\mathcal{L}^{VI}(\theta)$. Note that C_1 is independent of θ . Similarly, we can rewrite the MinVI objective as

$$
\mathcal{L}^{VI}(\theta) = -\mathbb{E}_{P_{\mathcal{D}}} \left[\log P_{\theta}(X|Y) + \log P_{\theta}(Y|X) \right]
$$
(S4)

$$
= \mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(X|Y)}{P_{\theta}(X|Y)} \right] + \mathbb{E}_{P_{\mathcal{D}}} \left[\log \frac{P_{\mathcal{D}}(Y|X)}{P_{\theta}(Y|X)} \right]
$$
(S5)

$$
- \mathbb{E}_{P_{\mathcal{D}}} \left[\log P_{\mathcal{D}}(X|Y) \right] - \mathbb{E}_{P_{\mathcal{D}}} \left[\log P_{\mathcal{D}}(Y|X) \right]
$$

$$
\underbrace{-\mathbb{E}_{P_{\mathcal{D}}} \left[\log P_{\mathcal{D}}(X|Y) \right] - \mathbb{E}_{P_{\mathcal{D}}} \left[\log P_{\mathcal{D}}(Y|X) \right]}_{C_2}
$$

where in Equation [\(S5\)](#page-0-1), we have

$$
\mathbb{E}_{P_{\mathcal{D}}}\left[\log\frac{P_{\mathcal{D}}(X|Y)}{P_{\theta}(X|Y)}\right] = \sum_{y} P_{\mathcal{D}}(y)\mathbb{E}_{P_{\mathcal{D}}(X|y)}\left[\log\frac{P_{\mathcal{D}}(X|y)}{P_{\theta}(X|y)}\right]
$$
(S6)

$$
= \mathbb{E}_{P_{\mathcal{D}}(Y)} \left[KL \left(P_{\mathcal{D}}(X|Y) \| P_{\theta}(X|Y) \right) \right] \tag{S7}
$$

Finally, we have

$$
\mathcal{L}^{\text{VI}}(\theta) = \mathbb{E}_{P_{\mathcal{D}}(X)} \left[KL \left(P_{\mathcal{D}}(Y|X) \| P_{\theta}(Y|X) \right) \right] +
$$

$$
\mathbb{E}_{P_{\mathcal{D}}(Y)} \left[KL \left(P_{\mathcal{D}}(X|Y) \| P_{\theta}(X|Y) \right) \right] + C_2.
$$
 (S8)

 C_2 is independent of θ and by setting $C = C_1 + C_2$, we derive the Equation (4).

S2 Proof of Theorem 2.1

Proposition S2.1 ([\[1,](#page-5-0) [2\]](#page-5-1)). Let X be a finite state space. Let irreducible transition matrices T_n *and* T converge to $\pi_n(X)$ and $\pi(X)$, respectively, where $\pi(X) = P_{\mathcal{D}}(X)$ is a data-generating *distribution of* X. If T_n *converges to* T *in the induced matrix norm, which is denoted by* $\|\cdot\|$ *, then* $\pi_n(X)$ converges to $P_{\mathcal{D}}(X)$ in l^2 norm.

Proof. Let $|\mathcal{X}|$ be the number of states. For simplicity, we denote $\pi = \pi(X)$ and $\pi_n = \pi_n(X)$. Since π is a stationary distribution of irreducible transition matrix T, π is uniquely defined and it satisfies the following:

$$
T\pi = \pi, \ \mathbf{1}^\top \pi = 1. \tag{S9}
$$

Combining above two equations, we have

$$
\begin{bmatrix}\nT_{1,1} - 1 & T_{1,2} & \cdots & T_{1,|\mathcal{X}|} \\
T_{2,1} & T_{2,2} - 1 & \cdots & T_{2,|\mathcal{X}|} \\
\vdots & \cdots & \cdots & \vdots \\
T_{|\mathcal{X}|-1,1} & \cdots & \cdots & T_{|\mathcal{X}|-1,|\mathcal{X}|-1} - 1\n\end{bmatrix}\n\pi = \begin{bmatrix}\n0 \\
0 \\
\vdots \\
1\n\end{bmatrix}
$$
\n(S10)

Since π exists and unique, \widetilde{T} is invertible and we have

$$
\boldsymbol{\pi} = \tilde{T}^{-1} \begin{bmatrix} 0 & 0 & \cdots & 1 \end{bmatrix}^\top \tag{S11}
$$

and similarly,

$$
\pi_n = \widetilde{T}_n^{-1} \begin{bmatrix} 0 & 0 & \cdots & 1 \end{bmatrix}^\top
$$
 (S12)

Since T_n (entrywise) converges to T, T_n^{-1} also converges to T^{-1} . Therefore, we conclude π_n converges to $\pi = P_{\mathcal{D}}(X)$. \Box

Now, we provide a proof of Theorem 2.1.

Proof of Theorem 2.1. To prove the convergence of marginal distributions, it is sufficient to show the convergence of transition operators. Since $|\mathcal{X}|$ and $|\mathcal{Y}|$ are finite, for any $\epsilon > 0$, there exists N such that $\forall n \geq N$, with probability at least $1 - \epsilon$, $\forall x \in \mathcal{X}, \forall y \in \mathcal{Y}$,

$$
|P_{\theta_n}(y|x) - P_{\mathcal{D}}(y|x)| < \epsilon, \ |P_{\theta_n}(x|y) - P_{\mathcal{D}}(x|y)| < \epsilon
$$

The transition operators are defined as follows:

$$
T_n^{\mathcal{Y}}(y[t]|y[t-1]) = \sum_{x \in \mathcal{X}} P_{\theta_n}(y[t]|x) P_{\theta_n}(x|y[t-1]),
$$

$$
T^{\mathcal{Y}}(y[t]|y[t-1]) = \sum_{x \in \mathcal{X}} P_{\mathcal{D}}(y[t]|x) P_{\mathcal{D}}(x|y[t-1])
$$

where $P_{\theta_n}(x|y)$ and $P_{\theta_n}(y|x)$ are derived from the joint distribution $P_{\theta_n}(x,y)$ and similarly for data-generating distribution, $P_D(x|y)$ and $P_D(y|x)$ are derived from $P_D(x, y)$. Then, for $n \geq N$, we have, for any $y_t, y_{t-1} \in \mathcal{Y}$, with probability at least $1 - \epsilon$,

$$
\left| T_n^{\mathcal{Y}} \left(y_t | y_{t-1} \right) - T^{\mathcal{Y}} \left(y_t | y_{t-1} \right) \right|
$$
\n
$$
\leq \left| \sum_{x \in \mathcal{X}} P_{\theta_n} \left(y_t | x \right) P_{\theta_n} \left(x | y_{t-1} \right) - P_{\mathcal{D}} \left(y_t | x \right) P_{\mathcal{D}} \left(x | y_{t-1} \right) \right|
$$
\n
$$
\leq \left| \mathcal{X} \right| \max_{x \in \mathcal{X}} \left| P_{\theta_n} \left(y_t | x \right) P_{\theta_n} \left(x | y_{t-1} \right) - P_{\mathcal{D}} \left(y_t | x \right) P_{\mathcal{D}} \left(x | y_{t-1} \right) \right| \tag{S13}
$$
\n
$$
\leq \left| \mathcal{X} \right| \left(2\epsilon \right)
$$

As we assume finite sets X and Y , this proves the convergence (in probability) of transition operator $T_n^{\mathcal{Y}}$ to $T^{\mathcal{Y}}$. The same argument holds for the convergence of transition operator $T_n^{\mathcal{X}}$ to $T^{\mathcal{X}}$. With

Proposition [S2.1,](#page-1-0) we proved the convergence of asymptotic marginal distribution $\pi_n(X)$ and $\pi_n(Y)$ to data-generating marginal distributions $P_D(X)$ and $P_D(Y)$, respectively.

Now, let's look at the joint probability distributions $P_{\theta_n}(x, y) = P_{\theta_n}(x|y)P_{\theta_n}(y)$ and similarly, $P_{\mathcal{D}}(x,y) = P_{\mathcal{D}}(x|y)P_{\mathcal{D}}(y)$. As we proved above, the following inequalities hold $\forall n \geq N'$:

$$
\left| P_{\theta_n}(y) - P_{\mathcal{D}}(y) \right| < \epsilon, \left| P_{\theta_n}(x|y) - P_{\mathcal{D}}(x|y) \right| < \epsilon \tag{S14}
$$

Therefore, using the similar argument in Equation [\(S13\)](#page-1-1), we have

$$
\left| P_{\theta_n}(x, y) - P_{\mathcal{D}}(x, y) \right| < 2\epsilon \tag{S15}
$$

and this completes the proof.

 \Box

S3 Retrieval Task

We provide more results of retrieval task with multimodal queries on MIR-Flickr database.

4

home, modern, chair chair chair

portrait, blackandwhite, nikon, bw, por

girl, newyork, best

blackandwhite, 365d

self, friends, d50, hair

california, h

graffiti, streetart, graffiti, portugal, streetart, graffiti, 2007, graf, graffiti, 2007, graf, graffiti, nyc, streetart chile, rio lisboa, lisbon tags, graff tags, graff tags, graff tags, graff tags, graff tags, graff tag

bw, portrait, japan,

bw, portrait, nikon40, hands bw, chile, mujer light, window, blackandwhite

Figure S1: Retrieval results with multimodal queries on MIR-Flickr database. The leftmost imagetext pairs are multimodal queries and those in the right side of the bar are retrieved samples with the highest similarities to the query.

day, interior, rainbow, chair, home home desk
books, library, apartment, decor

References

- [1] Y. Bengio, L. Yao, G. Alain, and P. Vincent. Generalized denoising auto-encoders as generative models. In *NIPS*, 2013.
- [2] Y. Bengio, E. Thibodeau-Laufer, G. Alain, and J. Yosinski. Deep generative stochastic networks trainable by backprop. In *ICML*, 2014.