

Supplements: Counterfactual Vision and Language Learning

Ehsan Abbasnejad, Damien Teney, Amin Parvaneh, Javen Shi, Anton van den Hengel

{ehsan.abbasnejad,damien.teney,amin.parvaneh,javen.shi,anton.vandenhengel}@adelaide.edu.au

Australian Institute for Machine Learning & The University of Adelaide, Australia

Theorem 1 (Theorem 4). Denote $u^i(\boldsymbol{\theta}) \equiv \ell(f_{\boldsymbol{\theta}}(\mathbf{q}_i, \mathbf{v}_i), a_i)\omega_i(\boldsymbol{\theta})$, $\bar{u} \equiv \sum_{i=1}^n u^i(\boldsymbol{\theta})/n$, $\hat{\mathbb{V}}(u) \equiv \sum_{i=1}^n (u^i(\boldsymbol{\theta}) - \bar{u})^2 / (n-1)$ and $\mathcal{Q}_\gamma \equiv \log(10 \cdot \epsilon/\gamma)$ for $0 < \gamma < 1$ and ϵ the ϵ -cover for the function class that predicts the answer. With probability at least $1 - \gamma$ for $n \geq 16$ we have

$$R(\boldsymbol{\theta}) \leq \hat{R}^M(\boldsymbol{\theta}) + \sqrt{18\hat{\mathbb{V}}(u)\mathcal{Q}_\gamma/n + 15M\mathcal{Q}_\gamma/(n-1)}$$

Proof. Follows the proof in Theorem 6 of [2]. \square

The density of the counterfactuals based on the observations, i.e.

$$p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) = \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p(\mathbf{q}, \mathbf{v})} \left[p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) \right] \quad (5)$$

Proof of Equation 5. We have:

$$\begin{aligned} p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) &= \int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})p^{\text{do}(I)}(\mathbf{u})d\mathbf{u} \\ &= \int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})p(\mathbf{u})d\mathbf{u} \\ &= \int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) \left(\int p(\mathbf{q}, \mathbf{v}, \mathbf{u}) dp(\mathbf{q}, \mathbf{v}) \right) d\mathbf{u} \\ &= \iint p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})p(\mathbf{u}|\mathbf{q}, \mathbf{v}) dp(\mathbf{q}, \mathbf{v}) d\mathbf{u} \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})p(\mathbf{u}|\mathbf{q}, \mathbf{v}) d\mathbf{u} \right] \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) \right] \end{aligned} \quad \square$$

Proof of Lemma (5). We have, \square

$$\begin{aligned} p^{\text{do}(I)}(a, \tilde{\mathbf{q}}, \tilde{\mathbf{v}}) &= p^{\text{do}(I)}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) \right]. \end{aligned}$$

Then using Jensen's inequality we have,

$$\begin{aligned} \log(\mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) \right]) \\ \geq \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right], \end{aligned}$$

We have:

$$\begin{aligned} &\mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right] \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\ &\quad \left. + \log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right] \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\ &\quad \left. + \log\left(\int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})p(\mathbf{u}|\mathbf{q}, \mathbf{v}) d\mathbf{u}\right) \right] \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\ &\quad \left. + \log\left(\int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})p(\mathbf{u}|\mathbf{q}, \mathbf{v}) d\mathbf{u}\right) \right] \end{aligned}$$

which is then lower-bounded as

$$\begin{aligned} &\geq \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\ &\quad \left. + \int \log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))p(\mathbf{u}|\mathbf{q}, \mathbf{v}) d\mathbf{u} \right] \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right] \\ &\quad + \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\int \log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))p(\mathbf{u}|\mathbf{q}, \mathbf{v}) d\mathbf{u} \right] \\ &= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right] \\ &\quad + \int \log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))p(\mathbf{u}) d\mathbf{u} \end{aligned}$$

Log-density of the joint for the question, image and answer as Eq. (8) in the paper:

$$\begin{aligned} \log(p^{\text{do}(I)}(a, \tilde{\mathbf{q}}, \tilde{\mathbf{v}})) &\geq \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p(\mathbf{q}, \mathbf{v})} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right] \\ &\quad + \mathbb{E}_{\mathbf{q}}[\log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))] \quad (8) \\ &\quad + H(\mathbf{q}) - H_{\mathbf{q}}(p). \end{aligned}$$

Proof of Equation 8. If we want to use an alternative distri-

bution

$$\begin{aligned}
& \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}) \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \log\left(\int p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) p(\mathbf{u}|\mathbf{q}, \mathbf{v}) d\mathbf{u}\right) \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \log\left(\int \frac{p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) p(\mathbf{u}|\mathbf{q}, \mathbf{v})}{q(\mathbf{u})} q(\mathbf{u}) d\mathbf{u}\right) \right]
\end{aligned}$$

which is then lower-bounded as

$$\begin{aligned}
& \geq \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \int \log\left(\frac{p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) p(\mathbf{u}|\mathbf{q}, \mathbf{v})}{q(\mathbf{u})}\right) q(\mathbf{u}) d\mathbf{u} \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \int \log\left(\frac{p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) p(\mathbf{u}|\mathbf{q}, \mathbf{v})}{q(\mathbf{u})}\right) q(\mathbf{u}) d\mathbf{u} \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \int (\log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) p(\mathbf{u}|\mathbf{q}, \mathbf{v})) \right. \\
&\quad \quad \left. - \log(q(\mathbf{u}))) q(\mathbf{u}) d\mathbf{u} \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \int (\log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}) p(\mathbf{u}|\mathbf{q}, \mathbf{v})) q(\mathbf{u}) d\mathbf{u} \right] + H(q) \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) \right. \\
&\quad \left. + \int \log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u})) q(\mathbf{u}) d\mathbf{u} \right. \\
&\quad \left. + \int \log(p(\mathbf{u}|\mathbf{q}, \mathbf{v})) q(\mathbf{u}) d\mathbf{u} \right] + H(q) \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) + H(q) \right. \\
&\quad \left. + \mathbb{E}_q[\log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))] + \mathbb{E}_q[\log(p(\mathbf{u}|\mathbf{q}, \mathbf{v}))] \right] \\
&= \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) + \mathbb{E}_q[\log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))] \right. \\
&\quad \left. + H(q) + \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \mathbb{E}_q[\log(p(\mathbf{u}|\mathbf{q}, \mathbf{v}))] \right] \\
&\geq \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} \left[\log(p^{\text{do}(I)|\mathbf{q}, \mathbf{v}}(a|\tilde{\mathbf{q}}, \tilde{\mathbf{v}})) + \mathbb{E}_q[\log(p^{\text{do}(I)}(\tilde{\mathbf{q}}, \tilde{\mathbf{v}}|\mathbf{u}))] \right. \\
&\quad \left. + H(q) - H_q(p) \right]
\end{aligned}$$

We have

$$\begin{aligned}
-\mathbb{E}_q \mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} [\log(p(\mathbf{u}|\mathbf{q}, \mathbf{v}))] &\geq -\mathbb{E}_q \log(\mathbb{E}_{(\mathbf{q}, \mathbf{v}) \sim p} [p(\mathbf{u}|\mathbf{q}, \mathbf{v})]) \\
&= -\mathbb{E}_q \log([p(\mathbf{u})]) = H_q(p)
\end{aligned}$$

□

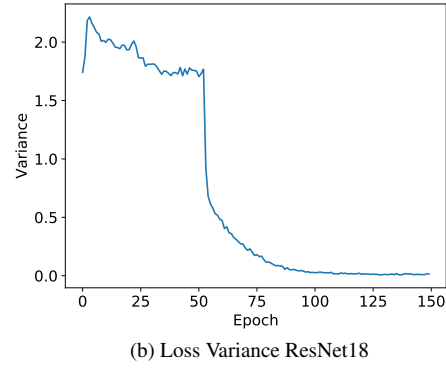
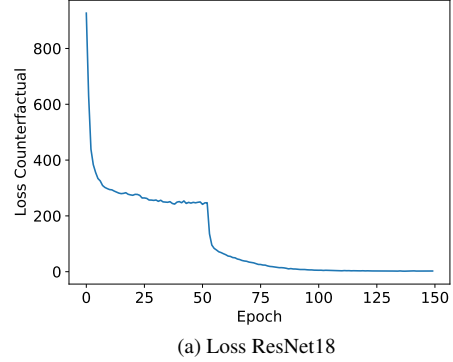
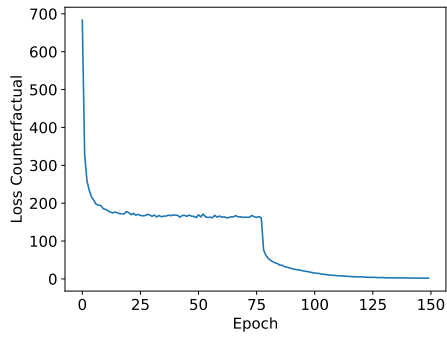


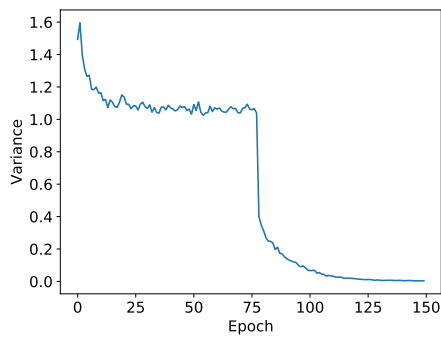
Figure 1: CIFAR10 results

1. Implementation details

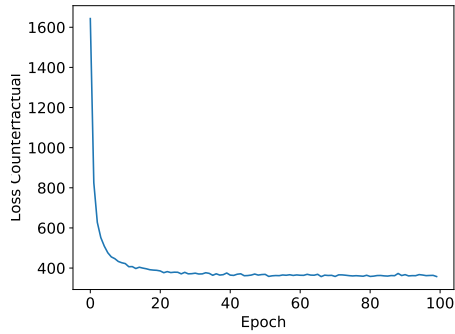
We implemented our approach on top of the original UpDn system [1]. The base system utilizes a Faster R-CNN head in conjunction with a ResNet-101 base network as the object detection module. For the VQA v2 experiment we utilize the ResNet-152 for detection. The detection head is pre-trained on the Visual Genome dataset. UpDn takes the final detection outputs and performs non-maximum suppression (NMS) for each object category using an IoU threshold of 0.7. Then, the convolutional features for the top 36 objects are extracted for each image as the visual features. For question embedding, we perform standard text pre-processing and tokenization. In particular, questions are first converted to lower case and then trimmed to a maximum of 14 words, and the words that appear less than 5 times are replaced with an “<unk>” token. We use GloVe embeddings and subsequently GRU for VQA-CP and LSTM for VQA v2A to sequentially process the word vectors and produce a sentential representation for the pre-processed question.



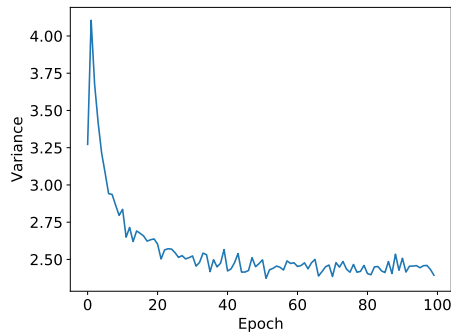
(a) Loss PreResNet18-CIFAR10



(b) Loss Variance PreResNet18-CIFAR10



(c) Loss PreResNet18-CIFAR100



(d) Loss Variance PreResNet18-CIFAR100

Figure 2: CIFAR results

References

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, 2018.
- [2] Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*, 2009.