



Hybrid Discriminative-Generative Training for Set Data

Jakub Bureš

Department of Mathematics

June 23, 2022

Content

- 1 Motivation
- 2 Discriminative vs. Generative modeling
- 3 Hybrid Discriminative–Generative Modeling
- 4 Variational Autoencoder
- 5 Multiple Instance Learning
- 6 Conclusion

Motivation

In machine learning, one can encounter two standard techniques:

- discriminative modeling,
- generative modeling.

Hybrid combination of these two approaches can improve a performance of a model. We focus on two hybrid variants:

- 1 HDGM [Abeel et. al 2020],
- 2 Semi-supervised VAE [Kingma et. al 2014],

which we apply to more complex cases, specifically, set data.

Discriminative modeling

- Data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where the input variable is $\mathbf{x} \in \mathbb{R}^D$ and the output variable $y \in \mathcal{C}$.
- *Classification problems* - the goal is to classify a new sample \mathbf{x} into some category y from the finite set \mathcal{C} , so we are looking for distribution $p(y|\mathbf{x})$.
- We train the model $f_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathcal{C}$, where θ denotes the parameters of this model.
- In practice, *one-hot encoding* is often used, thus the model is in the form of $f_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^C$.

Discriminative modeling

- Let the symbol $f_{\theta}(\mathbf{x})[y]$ denotes y^{th} element of $f_{\theta}(\mathbf{x})$.
- For distribution modeling is typically used Softmax

$$q_{\theta}(y|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{y \in \mathcal{C}} \exp(f_{\theta}(\mathbf{x})[y])}, \quad (1)$$

which is the basis for the definition of *cross-entropy*.

- Optimization of the model is performed by minimizing total cross-entropy (corresponding to the MLE), therefore

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} [\log q_{\theta}(y|\mathbf{x})]. \quad (2)$$

Generative modeling - Contrastive learning

- Generative models capture the joint probability $p(\mathbf{x}, y)$, or just $p(\mathbf{x})$ if there are no labels.
- Contrastive learning [Abeel et. al 2020] is a ML method typically utilized in the image classification.
- Here, we usually optimize

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\log \frac{\exp(m_{\theta}(\mathbf{x}) \cdot m_{\theta}(\mathbf{x}'))}{\sum_{i=1}^M \exp(m_{\theta}(\mathbf{x}) \cdot m_{\theta}(\mathbf{x}_i))} \right] \quad (3)$$

- We define a function $m_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^H$ that maps each sample to the representation space of dimension H .
- The sample \mathbf{x}' is called augmented view.

Generative modeling - Contrastive learning

Modification to our case:

- there is no need for augmented views x'
- instead of mapping m we simply use $f_{\theta}(\mathbf{x})[y]$ with labels y , yielding generative term

$$q_{\theta}(\mathbf{x}|y) \approx \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{j=1}^M \exp(f_{\theta}(\mathbf{x}_j)[y])}, \quad (4)$$

where $M < N$ is a number of normalization samples.

Hybrid Discriminative–Generative Modeling

- At this point, we have defined discriminative component $q_{\theta}(y|\mathbf{x})$ and generative component $q_{\theta}(\mathbf{x}|y)$.
- We can finally minimize the hybrid, convex combination of these two components

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x},y)} [\alpha \log q_{\theta}(y|\mathbf{x}) + (1 - \alpha) \log q_{\theta}(\mathbf{x}|y)]. \quad (5)$$

- Parameter $\alpha \in [0, 1]$ weighs generative and discriminative counterparts.

Variational Autoencoder

- A generative modeling method, the goal is to find $p(\mathbf{x})$ using the latent variable z , specifically, encoder $q_\phi(z|\mathbf{x})$ and decoder $p_\theta(\mathbf{x}|z)$.
- It leads to the ELBO optimization

$$\mathbb{E}_{q_\phi(z|\mathbf{x})} [\log p_\theta(\mathbf{x}|z)] - D_{\text{KL}}(q_\phi(z|\mathbf{x}) \| p_\theta(z)) \quad (6)$$

- No labels!
- However, there is no combination of discriminative and generative models -> **semi-supervised VAE** [Kingma et al. 2014].

VAE, semi-supervised version

It is necessary to consider two cases.

- First, consider observation x that has its class label y :

$$\mathbb{E}_{q_{\phi}(z|x,y)} [\log p_{\theta}(x|z, y) + \log p_{\theta}(y)] - D_{\text{KL}}(q_{\phi}(z|x, y) \| p_{\theta}(z)) \quad (7)$$

- Secondly, observation x is lacking its class label y :

$$\mathbb{E}_{q_{\phi}(y,z|x)} [\log p_{\theta}(x|z, y) + \log p_{\theta}(y)] - D_{\text{KL}}(q_{\phi}(y, z|x) \| p_{\theta}(z)) \quad (8)$$

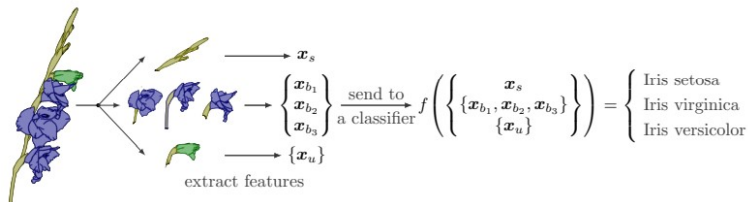
- We need a modification for a fully supervised dataset. We obtain

$$\tilde{J}_{HM}^{\nu}(\theta, \phi) = \tilde{J}_{HM}(\theta, \phi) + \nu \cdot \mathbb{E}_{\tilde{p}_l(x,y)} [-\log q_{\phi}(y|x)], \quad (9)$$

where parameter ν is a weight.

Multiple Instance Learning, MIL

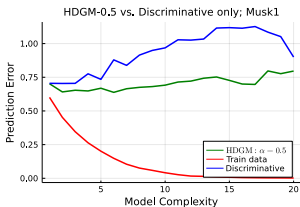
- In MIL, one sample is a set of vectors, these vectors are called instances and the sets are called bags.
- If the space of all instances is \mathcal{X} and the space of all bags is $\mathcal{B}(\mathcal{X})$, then in MIL the model is defined as $f_{\theta} : \mathcal{B}(\mathcal{X}) \rightarrow \mathcal{C}$.
- However, to properly define such a model, we need an embedded space and aggregation functions.
- We use the HMill framework and the Mill.jl package.



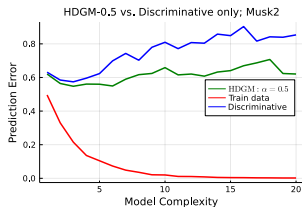
Obrázek: A general representation of set data. Credit [2].

Results - HDGM

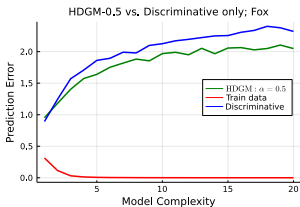
Cross-Validation of HDGM:



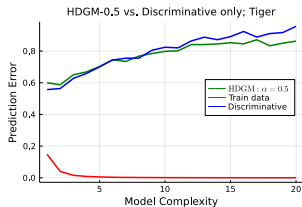
(a) Musk1.



(b) Musk2.



(c) Fox.



(d) Tiger.

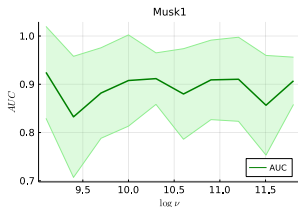
Results - HDGM

Average AUCs and st. deviations of individual data sets for both approaches:

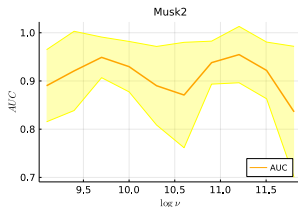
	D	HDGM $\alpha = 0.5$
Musk1	81.40 \pm 12.11 %	83.89 \pm 11.13%
Musk2	79.77 \pm 8.96%	82.78 \pm 7.38%
Tiger	91.22 \pm 2.34%	90.25 \pm 2.80%
Fox	54.43 \pm 3.77%	56.22 \pm 6.18%

Results - Hybrid VAE

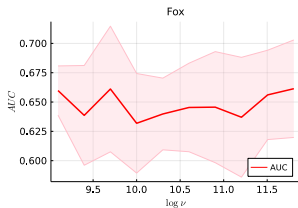
- Dependence of the average AUC on the logarithm of the parameter ν



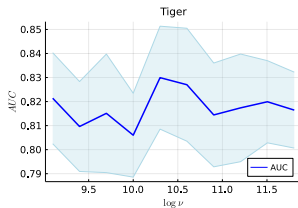
(e) Musk1.



(f) Musk2.



(g) Fox.






(h) Tiger.

Conclusion

- Slight reduction of prediction error for test data.
- Overall increase of AUC.
- High variance.

Thank you for your attention.

-  T. Pevny, S. Mandlik, *Mill.jl framework: a flexible library for (hierarchical) multi-instance learning* [on-line] Available from: <https://github.com/CTUAvastLab/Mill.jl>. Accessed 30 April 2020.
-  S.Mandlik.: *Mapping the Internet: Modelling Entity Interactions in Complex Heterogeneous Networks*. arXiv preprint arXiv:2104.09650.
-  D. P. Kingma and M. Welling, *Auto-Encoding Variational Bayes*. arXiv preprint arXiv:1312.6114, 2013.