

---

# Semi-supervised Multiple Instance Learning using Variational Auto-Encoders

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We consider the multiple-instance learning (MIL) paradigm, which is a special case  
2 of supervised learning where training instances are grouped into bags. In MIL, the  
3 hidden instance labels do not have to be the same as the label of the comprising bag.  
4 On the other hand, the hybrid modelling approach is known to possess advantages  
5 basically due to the smooth consolidation of both discriminative and generative  
6 components. In this paper, we investigate whether we can get the best of both  
7 worlds (MIL and hybrid modelling), especially in a semi-supervised learning (SSL)  
8 setting. We first integrate a variational autoencoder (VAE), which is a powerful  
9 deep generative model, with an attention-based MIL classifier, then evaluate the  
10 performance of the resulting model in SSL. We assess the proposed approach on  
11 an established benchmark as well as a real-world medical dataset.

## 12 1 Introduction

13 In the standard form of supervised learning, it is assumed that the learner encounters training data  
14 in a flat form where each instance, e.g., an image, belongs to a class (category). However, another  
15 setting which can be more practical in representing many real-world applications is multiple-instance  
16 learning (MIL), where training instances are grouped together into bags. In MIL, both bags and  
17 instances have labels, but an instance within a bag may have a different label from that of the bag.  
18 Only the bag label is available for learning since instance labels are not observed. Several applications  
19 can be cast as MIL problems, e.g., in medical imaging [Quellec et al., 2017] and computational  
20 biology [Dietterich et al., 1997].

21 The principal goal of MIL is to learn a model which can predict the bag label. This corresponds to the  
22 molecule binding property in the above example or to the all-important medical diagnosis in medical  
23 imaging applications. Nonetheless, inferring which instances are the most influential in predicting  
24 the bag label is of major importance due to several reasons including interpretability of the obtained  
25 prediction (especially in medical diagnosis) and related issues like GDPR (General Data Protection  
26 Regulation) which forces the right to understand in sensitive applications like self-driving cars and  
27 medical applications.

28 In this work, we investigate how the MIL framework fares in the semi-supervised learning paradigm  
29 (SSL, Zhu et al., 2003, Chapelle et al., 2006, Kingma et al., 2014, Siddharth et al., 2017). In SSL,  
30 the data presented to the learner typically consists of a few labeled examples as well as numerous  
31 unlabeled examples. The main goal of a semi-supervised learner is to utilize the unlabeled data in  
32 order to improve the model's performance on the supervised subset of the data. In case of the SSL  
33 MIL setting, the supervision is at the bag level. This means that the learner encounters both labeled  
34 and unlabeled bags.

35 To deal with both the labeled and unlabeled data, we propose to learn a joint distribution over  
 36 instances and a bag label within the hybrid modeling framework. Hybrid models are known to  
 37 combine the advantages of (standard supervised) discriminative models with those of generative  
 38 models [Jaakkola and Haussler, 1999, Tulyakov et al., 2017, Nalisnick et al., 2019]. Hybrid models  
 39 have also been exploited in other frameworks including semi-supervised learning [Ilse et al., 2020,  
 40 Nalisnick et al., 2019] and anomaly detection [Maaloe et al., 2019, Liu and Abbeel, 2020]. In this  
 41 work, we propose an MIL framework which leverages the prowess of hybrid models so that they can  
 42 excel in problems and applications possessing the bag-instance nature modelled by MIL. We build  
 43 our modelling on top of the seminal attention-based deep MIL classifier [Ilse et al., 2018], mainly  
 44 due to its permutation-invariant characteristics and its ability to give instance weights which can  
 45 be interpreted as the contributions of each instance to the bag label. As a result, we formulate a  
 46 latent variable model that could be seen as a Variational Auto-Encoder (VAE, Kingma and Welling,  
 47 2014, Rezende et al., 2014) for instances and a classifier that is fed with the outputs of the VAE’s  
 48 encoder. We evaluate the SSL performance of the proposed framework on a common benchmark and  
 49 a real-world medical data.

50 As such, our main contributions can be summarized as follows: (1) Integrating an attention-based  
 51 Deep MIL classifier with a deep generative model in the form of a VAE. (2) Developing an SSL  
 52 framework based on the proposed hybrid MIL approach. (3) Evaluating the proposed hybrid approach  
 53 on the semi-supervised MIL scenario and comparing it with baselines on two datasets (MNIST-BAGS,  
 54 COLON-CANCER).

## 55 2 Methodology

### 56 2.1 Multiple-Instance Learning

57 In standard binary classification, the main goal is to establish a model which predicts the target variable  
 58  $y \in \{0, 1\}$  for a data instance  $\mathbf{x} \in \mathbb{R}^D$ . On the other hand, each data sample in an MIL paradigm  
 59 comes in the form of a bag of unordered and independent<sup>1</sup> instances  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K\}$ , where  
 60 the number of instances, referred to as  $K$  can differ for different bags. An MIL model must learn to  
 61 predict the bag label  $Y$ , which is observed for the training data instances. In addition, there are also  
 62 instance labels  $y_1, y_2, \dots, y_K$  which are all hidden even for the training data. The standard MIL rule  
 63 on how to infer the bag label  $Y$  given its instance labels  $y_1, y_2, \dots, y_K$  can be expressed as follows:

$$Y = \begin{cases} 0, & \text{iff } \sum_k y_k = 0, \\ 1, & \text{otherwise.} \end{cases} \quad (1)$$

64 The MIL model we develop is trained by optimizing the log-likelihood (LL) function where the bag  
 65 label is distributed according to a Bernoulli distribution  $\theta(X) \in [0, 1]$ , which depicts the probability  
 66  $Y = 1$  given a bag  $X$  of instances. Also note that, since we assume bags of unordered and independent  
 67 instances, the bag probability  $\theta(X)$  must be permutation-invariant.

68 We pursue a three-step approach to predict bag labels, in which: (1) instances  $\mathbf{x}_k$  are first transformed  
 69 into a low-dimensional representation  $\mathbf{z}_k = f_\psi(\mathbf{x}_k)$ , (2) a combination of the transformed instances  
 70 is formed via a permutation-invariant function (referred to as the MIL pooling), and (3) in order  
 71 to form a bag representation, another transformation is applied over the combined instances, after  
 72 which a classifier  $\theta(X)$  is used for the resulting bag representation. We adopt a deep neural network  
 73 to parameterize all the transformations. Thus, the whole model can be optimized in an end-to-end  
 74 fashion via backpropagation.

### 75 2.2 Hybrid MIL

76 **Joint distribution** As mentioned earlier, we assume that instances within a bag  $X$  are identically  
 77 and independently distribution. This assumption is crucial in our methodology. Further, we are  
 78 interested in calculating the joint distribution over  $X$  and  $Y$  given the number of points in the bag  $X$ ,

<sup>1</sup>We refer to the standard MIL case which assumes independence among instances within a bag. Nonetheless, there are a few works which study MIL settings where instances within a bag do not follow the IID assumption, e.g. [Zhou et al., 2009, Zhang, 2021]

79  $p(X, Y|K)$ . Moreover, we consider the following generative model with shared latent variables:

$$p(X, Y|K) = \int p(Y, Z, X|K) dZ \quad (2)$$

$$= \int p(Y|Z, X)p(X, Z|K) dZ \quad (3)$$

$$= \int p(Y|Z)p(X|Z, K)p(Z|K) dZ \quad (4)$$

$$\stackrel{iid}{=} \int p(Y|Z) \left( \prod_{k=1}^K p(\mathbf{x}_k|\mathbf{z}_k)p(\mathbf{z}_k) \right) dZ, \quad (5)$$

80 where  $Z = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K\}$ .

81 **Variational inference** We parameterize these distributions using neural networks, thus, calculating  
 82 the integral becomes analytically intractable. In order to overcome this issue, we propose to use varia-  
 83 tional inference which allows calculating the lower bound to the logarithm of the joint distribution (the  
 84 ELBO). Considering the following family of variational posteriors  $q_\phi(Z|X, K) = \prod_{k=1}^K q_\phi(\mathbf{z}_k|\mathbf{x}_k)$   
 85 yields:

$$\log p_\vartheta(X, Y|K) = \log \int p_\vartheta(X, Y, Z|K) \frac{q_\phi(Z|X, K)}{q_\phi(Z|X, K)} dZ \quad (6)$$

$$\geq \mathbb{E}_{q_\phi(Z|X)} \left[ \log p_\vartheta(Y|Z) + \sum_{k=1}^K \left( \log p_\vartheta(\mathbf{x}_k|\mathbf{z}_k) + \log p_\vartheta(\mathbf{z}_k) - \log q_\phi(\mathbf{z}_k|\mathbf{x}_k) \right) \right] \quad (7)$$

$$\stackrel{df}{=} -\mathcal{L}(X, Y, K|\vartheta, \phi) \quad (8)$$

86 Notice that in the ELBO we have a component for the classification of a bag,  $\log p(Y|Z)$ , and a  
 87 sum of objectives for each object in the bag  $X$  that coincide with the formulation of Variational  
 88 Auto-Encoders [Kingma and Welling, 2014, Rezende et al., 2014].

89 **Semi-supervised learning** Since the ELBO consists of a sum of two objectives, namely, one  
 90 for the classifier and one for the marginal over objects, the proposed approach is well-suited for  
 91 semi-supervised learning. Let us denote the part with  $X$  as follows:

$$\mathcal{U}(X, K|\vartheta, \phi) \stackrel{df}{=} -\mathbb{E}_{q_\phi(Z|X)} \left[ \sum_{k=1}^K \left( \log p_\vartheta(\mathbf{x}_k|\mathbf{z}_k) + \log p_\vartheta(\mathbf{z}_k) - \log q_\phi(\mathbf{z}_k|\mathbf{x}_k) \right) \right]. \quad (9)$$

92 For two given sources of data, namely, labelled data  $(X, Y) \sim p_l(X, Y)$ , and unlabelled data  $X \sim$   
 93  $p_u(X)$ , we can formulate a joint learning objective by minimizing the combination of  $\mathcal{L}(X, Y, K|\vartheta, \phi)$   
 94 and  $\mathcal{U}(X, K|\vartheta, \phi)$ . However, typically we have more unlabelled data, therefore we consider a  
 95 weighted objective:

$$\mathcal{J}(\vartheta, \phi) = \alpha \cdot \sum_{(X, Y) \sim p_l} \mathcal{L}(X, Y, K|\vartheta, \phi) + \sum_{X \sim p_u} \mathcal{U}(X, K|\vartheta, \phi), \quad (10)$$

96 where  $\alpha > 0$ . This approach is known as *hybrid modeling* [Lasserre et al., 2006].

97 **Modeling  $p(Y|Z)$**  In this paper, we pursue an attention-based MIL pooling approach for modeling  
 98  $p(Y|Z)$  due to several reasons: Attention-based MIL pooling is more flexible, adaptive, and more  
 99 trainable than the max and mean pooling operators. It is also more interpretable due to the data-driven  
 100 adjustment of instance weights according to the task and data at hand, which can potentially provide  
 101 instance scores signifying the most relevant instances w.r.t. the bag label prediction. Attention-based  
 102 pooling is depicted in the form of a weighted averaging with learnable parameters. To ensure  
 103 invariance to the size (i.e. number of instances) of a bag, the weights are constrained to sum up to 1.

104 Assuming a bag of  $K$  instance representation embeddings  $Z = \{\mathbf{z}_1, \dots, \mathbf{z}_K\}$ , the MIL pooling is  
 105 expressed as:

$$\mathbf{h} = \sum_{k=1}^K a_k \mathbf{z}_k, \quad (11)$$

106 where:

$$a_k = \frac{\exp\{\mathbf{w}^\top (\tanh(\mathbf{V}\mathbf{z}_k^\top) \odot \text{sigm}(\mathbf{U}\mathbf{z}_k^\top))\}}{\sum_{j=1}^K \exp\{\mathbf{w}^\top (\tanh(\mathbf{V}\mathbf{z}_j^\top) \odot \text{sigm}(\mathbf{U}\mathbf{z}_j^\top))\}}, \quad (12)$$

107 where  $\mathbf{w} \in \mathbb{R}^{L \times 1}$ ,  $\mathbf{V} \in \mathbb{R}^{L \times M}$  and  $\mathbf{U} \in \mathbb{R}^{L \times M}$  are parameters, and  $\tanh(\cdot)$  is an element-wise  
 108 hyperbolic tangent nonlinearity. Element-wise multiplication is depicted by  $\odot$ , and  $\text{sigm}(\cdot)$  refers to  
 109 the sigmoid nonlinearity which grants the adoption of a gating mechanism, potentially avoiding some  
 110 troublesome linearity issues associated with  $\tanh(\cdot)$  [Ilse et al., 2018].

111 Eventually, the classifier works as follows:

- 112 1.  $X$  is transformed to  $Z$  through a shared stochastic encoder  $q_\phi(Z|X, K)$ , i.e., we calculate a  
 113 sample  $Z \sim q_\phi(Z|X, K)$ .
- 114 2. An embedding  $\mathbf{h}$  is calculated through the attention-based MIL pooling operator (see Eq. 11)  
 115 for given  $Z$ .
- 116 3. A neural network is used to calculate probabilities of class labels,  $\theta(\mathbf{h})$ .

### 117 3 Experiments

118 We quantitatively and qualitatively evaluate the proposed framework, which we refer to as semi-  
 119 supervised multiple-instance learning variational autoencoder (ssMILVAE). The conducted experi-  
 120 ments mainly address the following issues: (i) To assess the (accuracy) performance of the proposed  
 121 ssMILVAE in the SSL paradigm, and (ii) to gauge the degree of interpretability granted by ssMILVAE  
 122 and whether the learned instance weights can provide information on the contributions of each  
 123 instance to the bag label prediction.

124 We assess ssMILVAE on two datasets, MNIST-BAGS which is an MNIST-based image dataset, and  
 125 COLON CANCER which is a real-world histopathology dataset. We use 10-fold cross-validation  
 126 and repeat each experiment five times. To compare on common ground, we follow most of the  
 127 settings and modelling choices pursued by Ilse et al. [2018]. We refer to the latter method here as  
 128 AD-MIL. The MIL pooling layers are located right below the top layer of the model. In addition to  
 129 the classification accuracy, we compare the bag level performance based on: recall (true positive rate),  
 130 the area under the receiver operating characteristic curve (AUC) and (bag) classification accuracy. All  
 131 the experiments have been run for 100 epochs. Adam [Kingma and Ba, 2015] is the optimizer used,  
 132 with values of  $\beta_1$  and  $\beta_2$  set equal to 0.9 and 0.999, respectively. Weights are initialized according to  
 133 [He et al., 2015]. The hyperparameter  $\alpha$  (i.e., the weighting between the labelled objective and the  
 134 unlabelled objective) was determined through the model selection on the validation set.

#### 135 3.1 MNIST-BAGS

136 MNIST-BAGS is based on the well-known MNIST image data. We sample images from the MNIST  
 137 training (test) set to form training (test) bags, respectively. Each bag consists of a random number  
 138 of  $28 \times 28$  greyscale handwritten MNIST images. Number of images within a bag is Gaussian  
 139 distributed where the closest integer value is the chosen bag size. Since the number ‘9’ can possibly  
 140 be confused with ‘7’ and ‘4’, we rate a bag as positive if it contains at least one image of the digit ‘9’.

141 The ROC and accuracy results are displayed in Figures 1. The results demonstrate the supremacy  
 142 of the proposed ssMILVAE when the learner encounters a small number of labeled bags. The  
 143 performance of ssMILVAE is nearly equalled by AD-MIL with a larger number of labeled bags.

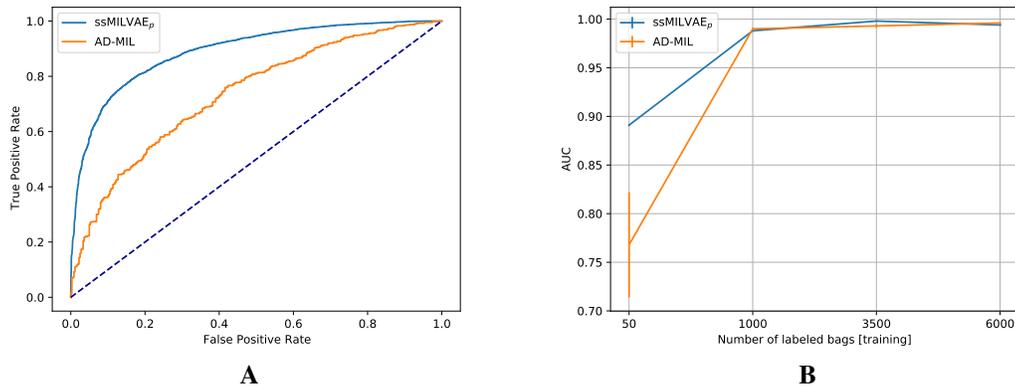


Figure 1: A comparison between ssMILVAE and AD-MIL. **A:** The ROC curve results for a bag size of 10 instances on the MNIST-BAGS dataset. **B:** The bag AUC results for 10-instance bags on the MNIST-BAGS dataset.

144 We next evaluate the attention mechanism of the proposed ssMILVAE algorithm on the MNIST-BAGS  
 145 dataset, and compare it with the seminal AD-MIL approach. We compare the two algorithms based on  
 146 a rather limited number of labeled bags, which is 50 bags. The bags displayed in Figure 4 have been  
 147 correctly classified by both algorithms and not cherry-picked. The proposed ssMILVAE is capable of  
 148 assigning higher weights to the positive instances than AD-MIL. This suggests that ssMILVAE may  
 149 provide more *interpretable* bag label predictions than AD-MIL, when trained on a limited number of  
 150 labeled bags, since the instance weights convey the relevance of the respective instances for the bag  
 151 labeling decision.

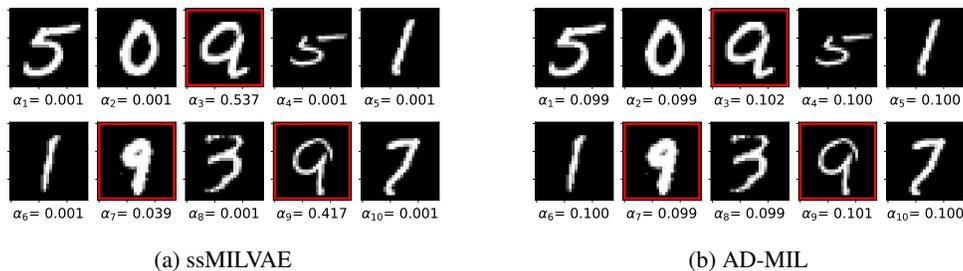


Figure 2: Evaluation of the attention mechanism of the proposed ssMILVAE algorithm compared to that of AD-MIL, tested on bags containing multiple positive ('9') instances from the MNIST-BAGS dataset.

### 152 3.2 COLON CANCER

153 The COLON CANCER dataset consists of real-world histopathology data [Sirinukunwattana et al.,  
 154 2016]. The data contains cancerous regions in hematoxylin and eosin (H&E) stained whole-slide  
 155 images. There are a total of 22,444 nuclei labeled as epithelial, inflammatory, fibroblast or miscella-  
 156 neous. It consists of 100 H&E images originating from a variety of tissue appearances from healthy  
 157 and malignant regions [Ilse et al., 2018]. Each bag consists of  $27 \times 27$  patches. A bag is labeled as  
 158 positive if it contains at least one epithelial nuclei. Colon cancer clinically originates from epithelial  
 159 cells, and this is why epithelial nuclei are very informative about the diagnosis here.

160 The accuracy results for experiments on the COLON CANCER dataset are displayed in Figure 3.  
 161 We experiment with the following number of labeled training bags: 22, 92 and 162. The proposed  
 162 ssMILVAE algorithm is more accurate when trained on a small number of training bags. When the  
 163 number of available labeled training bags increases, AD-MIL begins to outperform ssMILVAE.

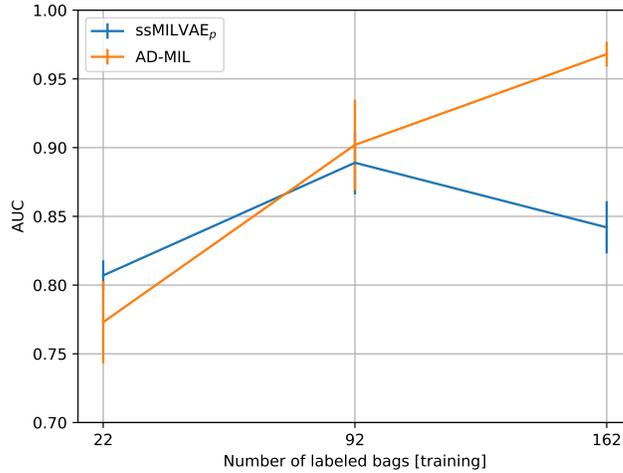


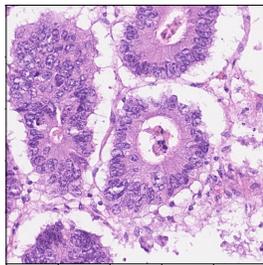
Figure 3: The bag AUC results on the COLON CANCER dataset for the proposed ssMILVAE and AD-MIL given a small number of labeled training bags.

164 Regarding the attention mechanism, we compare the proposed ssMILVAE with AD-MIL in terms of  
 165 the resulting regions of interest (ROIs), which are of paramount importance in medical diagnosis.  
 166 The raw histopathological image is displayed in Figure 4a. The histopathological image is split into  
 167 smaller patches containing single cells. A heatmap is generated by multiplying cell images by their  
 168 respective attention weights. The attention weights are then rescaled using  $a' = \frac{a_k - \min(a)}{\max(a) - \min(a)}$ .  
 169 As can be noticed in Figure 4d, the proposed attention mechanism by ssMILVAE achieves a much  
 170 better outcome in spotting the relevant cells compared to AD-MIL. As such, the attention mechanism  
 171 of the proposed ssMILVAE provides more interpretable predictions by identifying the key patches  
 172 responsible for the diagnosis.

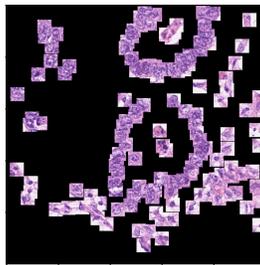
## 173 4 Conclusion

174 In this paper, we have presented an extension of the MIL classification problem to learning a joint  
 175 distribution in the semi-supervised setting. We have proposed a latent variable model for the MIL  
 176 generative model with a shared parameterization between the classifier and the unsupervised part.  
 177 In the experiments, we have shown that the proposed approach is beneficial in cases with a limited  
 178 number of labeled data.

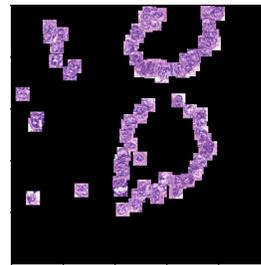
179 In many applications, (especially in the medical domain), it is difficult to obtain huge sizes of labeled  
 180 cases, and in such cases ssMILVAE seems to represent a recommended choice due to its ability to  
 181 learn from limited numbers of labeled bags (medical cases). Moreover, the attention mechanism  
 182 allows assisting a human expert (e.g., a physician) in interpreting results.



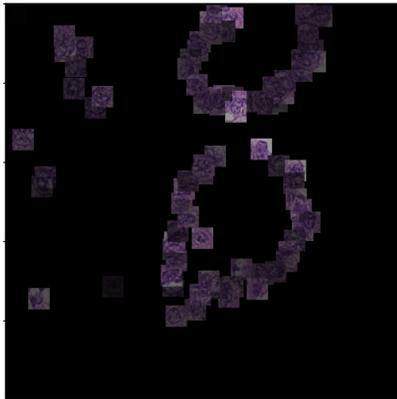
(a) Raw image



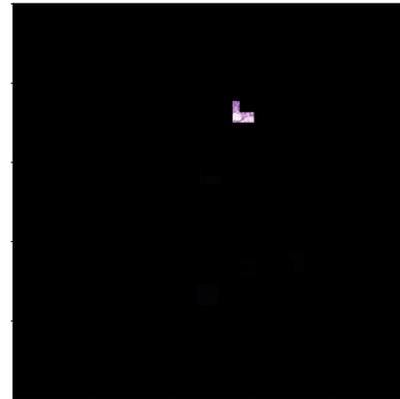
(b) All cells



(c) The ground-truth relevant cells



(d) Attention mechanism of the proposed ssMILVAE.



(e) Attention mechanism of AD-MIL.

Figure 4: Evaluation of the attention mechanism of the proposed ssMILVAE algorithm compared to that of AD-MIL, tested on the COLON CANCER dataset. Compared to AD-MIL, ssMILVAE assigns significantly higher weights to most of the relevant cells.

## 183 **References**

- 184 O. Chapelle, B. Scholkopf, and A. Zien. *Semi-supervised learning*. MIT Press, 2006.
- 185 T. Dietterich, R. Lathrop, and T. Lozano-Perez. Solving the multiple instance problem with axis-  
186 parallel rectangles. *Artificial Intelligence*, 89:31–71, 1997.
- 187 K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level perfor-  
188 mance on Imagenet classification. *IEEE International Conference on Computer Vision (ICCV)*,  
189 2015.
- 190 M. Ilse, J. Tomczak, and M. Welling. Attention-based deep multiple instance learning. *International  
191 Conference on Machine Learning (ICML)*, 2018.
- 192 M. Ilse, J. Tomczak, C. Louizos, and M. Welling. DIVA: Domain Invariant Variational Autoencoders.  
193 *Medical Imaging with Deep Learning*, 2020.
- 194 T. Jaakkola and D. Haussler. Exploiting generative models in discriminative classifiers. *Advances in  
195 Neural Information Processing Systems (NIPS)*, 1999.
- 196 D. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *International Conference on  
197 Learning Representations (ICLR)*, 2015.

- 198 D. Kingma and M. Welling. Auto-encoding variational Bayes. *International Conference on Learning*  
199 *Representations (ICLR)*, 2014.
- 200 D. Kingma, D. Rezende, S. Mohamed, and M. Welling. Semi-supervised learning with deep generative  
201 models. *Advances in neural information processing systems (NIPS)*, 28:3581–3589, 2014.
- 202 Julia A Lasserre, Christopher M Bishop, and Thomas P Minka. Principled hybrids of generative  
203 and discriminative models. In *2006 IEEE Computer Society Conference on Computer Vision and*  
204 *Pattern Recognition (CVPR'06)*, volume 1, pages 87–94. IEEE, 2006.
- 205 H. Liu and P. Abbeel. Hybrid discriminative-generative training via contrastive learning. *arXiv*  
206 *preprint arXiv:2007.09070*, 2020.
- 207 L. Maaloe, M. Fraccaro, V. Lievin, and O. Winther. BIVA: A very deep hierarchy of latent variables  
208 for generative modeling. *Advances in neural information processing systems (NeurIPS)*, 2019.
- 209 E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan. Hybrid models with  
210 deep and invertible features. *International Conference on Machine Learning (ICML)*, 2019.
- 211 G. Quellec, G. Cazuguel, B. Cochener, and M. Lamard. Multiple-instance learning for medical image  
212 and video analysis. *IEEE Reviews in Biomedical Engineering*, 2017.
- 213 D. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference  
214 in deep generative models. *International Conference on Machine Learning (ICML)*, 31, 2014.
- 215 N. Siddharth, B. Paige, J. van den Meent, A. Demaison, N. Goodman, P. Kohli, F. Wood, and P. Torr.  
216 Learning disentangled representations with semi-supervised deep generative models. *Advances in*  
217 *neural information processing systems (NIPS)*, 2017.
- 218 K. Sirinukunwattana, S. Raza, Y. Tsang, D. Snead, I. Cree, and N. Rajpoot. Locality sensitive deep  
219 learning for detection and classification of nuclei in routine colon cancer histology images. *IEEE*  
220 *Transactions on Medical Imaging*, 35(5):1196–1206, 2016.
- 221 S. Tulyakov, A. Fitzgibbon, and S. Nowozin. Hybrid VAE: Improving deep generative models using  
222 partial observations. *arXiv preprint arXiv:1711.11566*, 2017.
- 223 W. Zhang. Non-I.I.D. multi-instance learning for predicting instance and bag labels using variational  
224 auto-encoder. *arXiv preprint arXiv:2105.01276*, 2021.
- 225 Z. Zhou, Y. Sun, and Y. Li. Multi-instance learning by treating instances as non-I.I.D. samples.  
226 *International Conference on Machine Learning (ICML)*, 2009.
- 227 X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using Gaussian fields and harmonic  
228 functions. *International Conference on Machine Learning (ICML)*, 2003.