

Semi-Supervised Variational Autoencoder for Survival Prediction

Sveinn Pálsson^{*1}, Stefano Cerri^{*1}, Andrea Dittadi^{*2}, and Koen Van Leemput^{1,3}

¹ Department of Health Technology, Technical University of Denmark, Denmark

² Department of Applied Mathematics and Computer Science, Technical University of Denmark, Denmark

³ Athinoula A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital, Harvard Medical School, USA

Abstract. In this paper we propose a semi-supervised variational autoencoder for classification of overall survival groups from tumor segmentation masks. The model can use the output of any tumor segmentation algorithm, removing all assumptions on the scanning platform and the specific type of pulse sequences used, thereby increasing its generalization properties. Due to its semi-supervised nature, the method can learn to classify survival time by using a relatively small number of labeled subjects. We validate our model on the publicly available dataset from the Multimodal Brain Tumor Segmentation Challenge (BraTS) 2019.

Keywords: Survival time · deep generative models · semi-supervised VAE.

1 Introduction

Brain tumor prognosis involves forecasting the future disease progression in a patient, which is of high potential value for planning the most appropriate treatment. Glioma is the most common primary brain tumor and patients suffering from its most aggressive form, glioblastoma, have generally very poor prognosis. Glioblastoma patients have a median overall survival (OS) of less than 15 months, and a 5-year OS rate of only 10% even when they receive treatment [1]. Automatic prediction of overall survival of glioblastoma patients is an important but unsolved problem, with no established method available in clinical practice.

The last few years have seen an increased interest in brain tumor survival time prediction from magnetic resonance (MR) images, often using discriminative methods that directly encode the relationship between image intensities and prediction labels [2]. However, due to the flexibility of MR imaging, such methods do not generalize well to images acquired at different centers and with different scanners, limiting their potential applicability in clinical settings. Furthermore, being supervised methods, they require “labeled” training data where for each

^{*} Authors contributed equally.

training subject both imaging data and ultimate survival time are available. Although public imaging databases with survival information have started to be collected [3,4,5,6], the requirement of such labeled data fundamentally limits the number of subjects available for training, severely restricting the prediction performance attainable with current methods.

In this paper, we explore whether the aforementioned issues with supervised intensity-based methods can be ameliorated by using a semi-supervised approach instead, using only segmentation masks as input. In particular, we adapt a semi-supervised variational autoencoder model [7] to predict overall survival from a small amount of labeled training subjects, augmented with *unlabeled* subjects in which only imaging data is available. The method only takes segmentation masks as input, thereby removing all assumptions on the image modalities and scanners used.

The Multimodal Brain Tumor Segmentation Challenge (BraTS) [3] has been held every year since 2012, and focuses on the task of segmenting three different brain tumors structures (“enhancing tumor”, “tumor core” and “whole tumor”) and “background” from multimodal MR images. Since 2017, BraTS has also included the task of OS prediction. In this paper we focus on the latter, classifying the scans into three prognosis groups: **long-survivors** (>15 months), **short-survivors** (<10 months), and **mid-survivors** (between 10 and 15 months), all relative to the time of diagnosis.

2 Model

We begin by formally describing the problem we aim to solve. The available training data consists of a set of N_l labeled pairs $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{N_l}, y_{N_l})\}$, possibly augmented with a set of N_u *unlabeled* data points $\{\mathbf{x}_{N_l+1}, \dots, \mathbf{x}_{N_l+N_u}\}$, where $\mathbf{x}_i \in \{1, \dots, M_x\}^D$ is the i -th subject’s image data in the form of a segmentation map with D voxels, and the target variable $y_i \in \{1, \dots, M_y\}$ denotes the survival group the subject belongs to. In our case we have the segmentation of $M_x = 4$ different tumor structures as input to the model, and $M_y = 3$ different survival groups. For convenience, we will omit the index i when possible in the remainder.

We assume that the data is generated by a random process, illustrated in Figure 1, that involves some latent variables $\mathbf{z} \in \mathcal{R}^L$, assumed to be independent of y , where $L \ll D$. These latent variables encode high-level tumor shape and location features shared across survival groups. Specifically, we assume a generative model of the form

$$p_{\boldsymbol{\theta}}(\mathbf{x}, y, \mathbf{z}) = p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{y}, \mathbf{z})p(\mathbf{z})p(y), \quad (1)$$

where $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ is a zero-mean isotropic multivariate Gaussian, $p(y) \propto 1$ is a flat categorical prior distribution over y , and $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{y}, \mathbf{z})$ is a conditional distribution parameterized by $\boldsymbol{\theta}$.

Our task is to find the maximum likelihood parameters, i.e., the parameter values $\boldsymbol{\theta}$ that maximize the probability of the training data under the model.

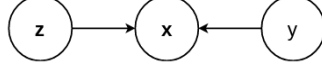


Fig. 1. Probabilistic graphical model of the generative process.

This is equivalent to maximizing

$$\sum_{i=1}^{N_l} \log p_{\theta}(\mathbf{x}_i, y_i) + \sum_{i=N_l+1}^{N_l+N_u} \log p_{\theta}(\mathbf{x}_i) \quad (2)$$

with respect to θ , where

$$p_{\theta}(\mathbf{x}, y) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}, y, \mathbf{z}) d\mathbf{z} \quad (3)$$

and

$$p_{\theta}(\mathbf{x}) = \sum_y p_{\theta}(\mathbf{x}, y). \quad (4)$$

Once suitable parameter values are found, the survival group of a new subject with image data \mathbf{x} can be predicted by assessing $p_{\theta}(y|\mathbf{x}) = p_{\theta}(\mathbf{x}, y)/p_{\theta}(\mathbf{x})$.

2.1 Semi-supervised variational autoencoder

Maximizing eq. (2) for θ directly is not feasible due to intractability of the integral over the latent variables in eq. (3). We therefore use an Expectation-Maximization (EM) [8] algorithm to exploit the fact that the optimization would be easier if the latent variables were known. The algorithm iteratively constructs and maximizes a lower bound to eq. (2) in a process that involves “filling in” the missing latent variables using their posterior distribution. Since this posterior distribution is intractable, we follow [7] and approximate $p_{\theta}(\mathbf{z}, y|\mathbf{x})$ using a specific functional form $q_{\phi}(\mathbf{z}|\mathbf{x}, y)$ with parameters ϕ :

$$q_{\phi}(\mathbf{z}, y|\mathbf{x}) = q_{\phi}(\mathbf{z}|\mathbf{x}, y)q_{\phi}(y|\mathbf{x}),$$

where $q_{\phi}(\mathbf{z}|\mathbf{x}, y)$ is a multivariate Gaussian distribution with diagonal covariance matrix, and $q_{\phi}(y|\mathbf{x})$ is a categorical distribution. This approximation can be used to obtain a lower bound to eq. (2) as follows. The probability of each *labeled* data point (first term in eq. (2)) can be rewritten as:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}, y) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} [\log p_{\theta}(\mathbf{x}, y)] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} \left[\log \left[\frac{p_{\theta}(\mathbf{x}, y, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}, y)} \right] \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} \left[\log \left[\frac{p_{\theta}(\mathbf{x}, y, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} \frac{q_{\phi}(\mathbf{z}|\mathbf{x}, y)}{p_{\theta}(\mathbf{z}|\mathbf{x}, y)} \right] \right] \\ &= \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} \left[\log \left[\frac{p_{\theta}(\mathbf{x}, y, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} \right] \right]}_{=\mathcal{L}_{\theta, \phi}(\mathbf{x}, y)} + \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, y)} \left[\log \left[\frac{q_{\phi}(\mathbf{z}|\mathbf{x}, y)}{p_{\theta}(\mathbf{z}|\mathbf{x}, y)} \right] \right]}_{=D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}, y) || p_{\theta}(\mathbf{z}|\mathbf{x}, y))} \end{aligned}$$

where D_{KL} denotes the Kullback-Leibler (KL) divergence. Since the KL divergence is always non-negative, we have that

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}, y) \geq \mathcal{L}_{\boldsymbol{\theta}, \phi}(\mathbf{x}, y). \quad (5)$$

Using a similar derivation, the probability of each *unlabeled* data point can be bounded as follows:

$$\begin{aligned} \log p_{\boldsymbol{\theta}}(\mathbf{x}) &\geq \mathbb{E}_{q_{\phi}(y, \mathbf{z}|\mathbf{x})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, y, \mathbf{z})}{q_{\phi}(\mathbf{z}|y, \mathbf{x})} - \log q_{\phi}(y|\mathbf{x}) \right] \\ &= \sum_y q_{\phi}(y|\mathbf{x}) (\mathcal{L}_{\boldsymbol{\theta}, \phi}(\mathbf{x}, y)) + \mathcal{H}(q_{\phi}(y|\mathbf{x})) = \mathcal{U}_{\boldsymbol{\theta}, \phi}(\mathbf{x}), \end{aligned} \quad (6)$$

where $\mathcal{H}(\cdot)$ denotes the entropy of a probability distribution.

By combining (5) and (6), a lower bound to eq. (2) is finally obtained as:

$$\mathcal{J}_{\boldsymbol{\theta}, \phi} = \sum_{i=1}^{N_l} \mathcal{L}_{\boldsymbol{\theta}, \phi}(\mathbf{x}_i, y_i) + \sum_{i=N_l+1}^{N_l+N_u} \mathcal{U}_{\boldsymbol{\theta}, \phi}(\mathbf{x}_i), \quad (7)$$

which we optimize with respect to both the variational parameters ϕ and the generative parameters $\boldsymbol{\theta}$. We use stochastic gradient ascent for the optimization, approximating gradients of the expectations in (7) as described in [9]. Implementation details are discussed in Section 4.

From an information theory point of view, the latent unobserved variables \mathbf{z} can be interpreted as a code. Therefore, we can refer to the distributions $q_{\phi}(\mathbf{z}|\mathbf{x}, y)$ and $p_{\boldsymbol{\theta}}(\mathbf{x}|y, \mathbf{z})$ as a probabilistic *encoder* and *decoder*, respectively [9]. The label predictive distribution $q_{\phi}(y|\mathbf{x})$ has the form of a discriminative *classifier*, and can be used as an approximation to $p_{\boldsymbol{\theta}}(y|\mathbf{x})$ for classifying new cases after training.

2.2 Model modifications

Here we describe a few model modifications for making the parameter learning process faster and less prone to overfitting.

Classification objective Note that in the objective function (7), the label predictive approximation $q_{\phi}(y|\mathbf{x})$ only appears in the bound for unlabeled data. To let $q_{\phi}(y|\mathbf{x})$ also learn from labeled data, we follow [7] and add a weak classification loss, resulting in the modified objective

$$\mathcal{J}_{\boldsymbol{\theta}, \phi}^{\alpha} = \mathcal{J}_{\boldsymbol{\theta}, \phi} + \alpha \sum_{i=1}^{N_l} \log q_{\phi}(y_i|\mathbf{x}_i) \quad (8)$$

where α controls the relative weight between generative and purely discriminative learning.

Gumbel-Softmax One of the issues of training a semi-supervised VAE is that the marginalization over $q_\phi(y|\mathbf{x})$ in eq. (6) can be computationally expensive. This marginalization can be avoided by using Gumbel-Softmax [10,11], a continuous distribution on the probability simplex that approximates a categorical sample and can be smoothly annealed (through a temperature parameter) to the categorical distribution. Gumbel-Softmax is reparameterizable so that the gradient of the loss function can be propagated back through the sampling step $y \sim q_\phi(y|\mathbf{x})$ for single-sample gradient estimation.

Regularization The lower bound for labeled data can be rewritten as

$$\begin{aligned}\mathcal{L}_{\theta,\phi}(\mathbf{x}, y) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, y)} \left[\log \frac{p_\theta(\mathbf{x}, y, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x}, y)} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, y)} \left[\log p_\theta(\mathbf{x}|\mathbf{z}, y) \right] + \log p(y) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}, y) || p(\mathbf{z}))\end{aligned}$$

where $\log p(y)$ is a constant, the first term can be interpreted as expected negative reconstruction error, and the last term is the negative KL divergence from the prior to the approximate posterior. Similarly, we can express the bound for unlabeled data as follows:

$$\mathcal{U}_{\theta,\phi}(\mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}, y|\mathbf{x})} \left[\log p_\theta(\mathbf{x}|\mathbf{z}, y) \right] - D_{KL}(q_\phi(\mathbf{z}, y|\mathbf{x}) || p(\mathbf{z}, y))$$

In both cases, the KL divergence acts as a regularization term that encourages the approximate posterior to be close to the prior, thereby constraining the amount of information encoded in the latent variables. The overall lower bound (7) thus trades off reconstruction error with this regularization term. When training a VAE, we can control such trade-off in order to favor more accurate reconstructions or more constrained latent space, by simply multiplying the KL term by a factor $\beta > 0$ as proposed in [12]. Similarly, we found it beneficial in practice to scale the entropy of $q_\phi(y|\mathbf{x})$ in eq. (6) by a factor $\gamma > 1$. Intuitively, the entropy term acts as a regularizer in the classifier by encouraging $q_\phi(y|\mathbf{x})$ to have high entropy: the amplification of this term helps to further reduce overfitting in the classifier.

3 Data and models

The BraTS 2019 challenge is composed of a training, a validation and a test set. The training set is composed of 335 delineated tumor images, in which 210 images have survival labels. The validation set is composed of 125 non-delineated images without survival labels, in which only 29 images with resection status of GTR (i.e., Gross Total Resection) are part of the online evaluation platform (CBICA’s Image Processing Portal). Finally, the test set will be made available to the challenge participants during a limited time window, and the results will be part of the BraTS 2019 workshop.

In all our experiments we performed 3-fold cross-validation by randomly splitting the BraTS 2019 training set with survival labels into a “private” training (75%) and validation set (25%) in each fold, in order to have an alternative to the online evaluation platform. This help us having a more informative indication of the model performance, since the online evaluation platform includes just 29 cases (vs. 53 cases in our private validation sets). With this set-up, which we call **S0** in the remainder, we effectively trained the model on a training set of $\mathbf{N}_l = 157$ and $\mathbf{N}_u = 125$ for each of the three cross-validation folds. These models were subsequently tested on their corresponding private validation sets of 53 subjects, as well as on the standard BraTS 2019 validation set of 29 subjects.

In order to evaluate just how much the proposed method is able to learn from *unlabeled* data (i.e., subjects with tumor delineations but no survival time information), we used three open-source methods [13,14,15] to automatically segment both the entire BraTS 2019 training and validation sets in order to have many more unlabeled training subjects available. We further augmented these unlabeled data sets by flipping the images in the coronal plane. With this new set-up, which we call **S1**, we then trained the model on an “augmented” private training set of $\mathbf{N}_l = 157$ and $\mathbf{N}_u = 2268$ for each of the three cross-validation folds. Ideally, dramatically increasing the set of unlabeled data points this way should help the model learn to better encode tumor representations, thereby increasing classification accuracy.

4 Implementation

We implemented the encoder $q_\phi(\mathbf{z}|\mathbf{x}, y)$, the decoder $p_\theta(\mathbf{x}|\mathbf{z}, y)$ and the classifier $q_\phi(y|\mathbf{x})$ all as deep convolutional networks using PyTorch [16]. The segmentation volumes provided in the BraTS challenge have size $240 \times 240 \times 155$, but since large parts of the volume are always zero, we cropped the volumes to $146 \times 188 \times 128$ without losing any tumor voxels. We further reduced the volume by a factor of 2 in all dimensions, resulting in a shape of $73 \times 94 \times 64$, roughly a 95% overall reduction in input image size. This leads to faster training and larger batches fitting in memory, while losing minimal information.

We optimized the model end-to-end with Adam optimizer [17], using a batch size of 32, learning rate $2 \cdot 10^{-5}$, latent space size 32, $\alpha = 10^{-5} \cdot D \approx 4.4$ with D the data dimensionality (number of voxels), β from 0 to $6 \cdot 10^3$ in $3 \cdot 10^4$ steps, $\gamma = 50$, and exponentially annealing the Gumbel-Softmax sampling temperature from 1.0 to 0.2 in $5 \cdot 10^4$ steps. Hyperparameters were found by grid search, although not fine-tuned because of the computational cost. The total number of parameters in the model is around 2.7×10^6 .

4.1 Network architecture

The three networks consist of 3D convolutional layers, with the exception of a few fully connected layers in the classifier. There are nonlinearities (Scaled Exponential Linear Units, [18]) and dropout [19] after each layer, except when

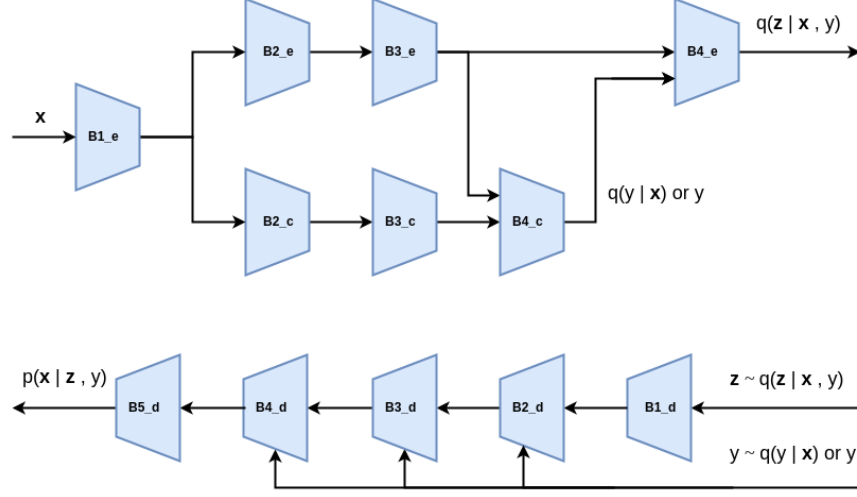


Fig. 2. Networks architectures: encoder, decoder and classifier architectures.

noted. What follows is a high-level description of the network architecture, represented in diagrams in Figure 2. For more details, the code is available at <https://github.com/sveinnpalsson/semivaebrats>.

The inference network consists of a convolutional layer ($B1_e$) with large kernel size and stride (7 and 4, respectively), followed by two residual blocks [20] ($B2_e$ and $B3_e$). The input to each block is processed in parallel in two branches, one consisting of two convolutional layers, the other of average pooling followed by a linear transformation (without nonlinearities). The results of the two branches are added together. The output of the first layer is also fed into the classifier network, which outputs the class scores (these will be used to compute the classification loss for labeled data). A categorical sample from $q_\phi(y|x)$ is drawn using the Gumbel-Softmax reparameterization given the class scores, and is embedded by a fully connected layer into a real vector space. Such embedding is then concatenated to the output of the two encoder blocks, so that the means and variances of the approximate posterior $q_\phi(z|x, y)$, that are computed by a final convolutional layer, are conditioned on the sampled label. The classifier consists of two residual blocks similar to the ones in the encoder ($B2_c$ and $B3_c$), followed by two fully connected layers ($B4_c$).

The decoder network consists of two convolutional layers ($B1_d$ and $B2_d$), two residual blocks similar to those in the encoder ($B3_d$ and $B4_d$), and a final convolution followed by a sigmoid nonlinearity ($B5_d$). In the decoder, most convolutions are replaced by transposed convolutions (for upsampling), and pooling in the residual connections is replaced by nearest neighbour interpolation. The input to the decoder network is a latent vector z sampled from the approximate posterior. The embedding of y , computed as in the final stage of the inference network, is also concatenated to the input of each layer (except the ones in

the middle of a block) to express the conditioning of the likelihood function on the label. Here, the label is either the ground truth (for labeled examples) or a sample from the inferred posterior (for unlabeled examples).

5 Results

5.1 Conditional generation

We visually tested whether the decoder $p_{\theta}(\mathbf{x}|y, \mathbf{z})$ is able to generate tumor-like images after training, and whether it can disentangle the classes. For this purpose we sampled \mathbf{z} from $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and varied y between the three classes, namely, short survivor, mid survivor and long survivor. Figure 3 shows the three shapes generated accordingly by one of the models trained in set-up **S0**. From the images we can see that the generated tumor for the short survivor class has an irregular shape with jagged edges while the long survivor generated tumor has a more compact shape with rounded edges.

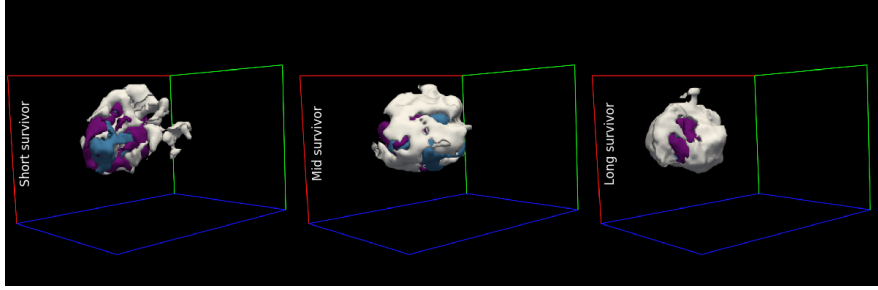


Fig. 3. Generated tumor from $p_{\theta}(\mathbf{x}|y, \mathbf{z})$ where we sampled \mathbf{z} from $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and we varied y between short survivor, mid survivor and long survivor.

5.2 Quantitative evaluation

All the classification accuracies are reported with binomial confidence interval with normal approximation [21], defined as

$$a \pm z^* \sqrt{\frac{a(1-a)}{n}}$$

where a is the classification accuracy, $z^* = 1.96$ is the critical value with confidence level at 95% and n is the number of subjects. In Table 1 we show the classification accuracy of the proposed method on the “private” validation set of 53 subjects for each of the three cross-validation folds, both for the set-up with fewer (**S0**) and more (**S1**) unlabeled training subjects. The corresponding

results based on the online evaluation platform (29 validation subjects) are summarized in Table 2, where we submitted the majority vote for survival group prediction across the three models trained in the cross-validation folds. The online evaluation platform takes the estimated number of days as input and returns the accuracy along with mean- and median squared error and Spearman’s rank correlation coefficient. To make these predictions we input the average survival from each class. Our scores on the challenge leaderboard for set-up **S0** are as follows: 37.9% accuracy, 111214.828 mean squared error, 51076.0 median squared error and a correlation of 0.36. When testing the models we found that they are insensitive to the segmentation method used to produce the input.

Table 1. Classification accuracies [%] for both set-ups on the “private” validation set for each of the three cross-validation folds.

Set-up	Fold 1	Fold 2	Fold 3	Avg
S0	42.18 ± 13.30	35.90 ± 12.91	39.53 ± 13.16	39.20 ± 7.59
S1	47.55 ± 13.45	41.13 ± 13.40	42.91 ± 13.32	43.86 ± 7.71

Table 2. Classification accuracies [%] for both set-ups on the BraTS 2019 online evaluation platform.

Set-up	Majority voting
S0	37.90 ± 17.57
S1	31.00 ± 16.83

The results show that in none of the experiments our model achieved a significant improvement over always predicting the largest class, which constitutes around 40% of the labeled cases.

6 Discussion and conclusions

In this paper we evaluated the potential of a semi-supervised deep generative model for classifying brain tumor patients into three overall survival groups, based only on tumor segmentation masks. The main potential advantages of this approach are (1) its in-built invariance to MR intensity variations when different scanners and protocols are used, enabling wide applicability across clinics; and (2) its ability to learn from unlabeled data, which is much more widely available than fully-labeled data.

We compared two different set-ups: one where fewer unlabeled subjects were available for training, and one where their number was (largely artificially) increased using automatic segmentation and data augmentation. Although the latter set-up increased classification performance in our “private” experiments,

this increase did not reach statistically significant levels and was not replicated on the small BraTS 2019 validation set. We demonstrated visually that the proposed model effectively learned class-specific information, but overall failed to achieve classification accuracies significantly higher than predicting always the largest class.

The results described here are only part of a preliminary analysis. More real unlabeled data, obtained from truly different subjects pooled across treatment centers, and more clinical covariates of the patients, such as age and resection status, may be necessary to reach better classification accuracies. Future work may also involve stacking hierarchical generative models to further increase the classification performance of the model [7].

7 Acknowledgements

This project was funded by the European Unions Horizon 2020 research and innovation program under the Marie Skłodowska-Curie project TRABIT (agreement No 765148).

References

1. Poulsen et al. The prognostic value of fet pet at radiotherapy planning in newly diagnosed glioblastoma. *European journal of nuclear medicine and molecular imaging*, 44(3):373–381, 2017.
2. Bakas S. et al. Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge. *CoRR*, abs/1811.02629, 2018.
3. B. H. Menze et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE Transactions on Medical Imaging*, 34(10):1993–2024, Oct 2015.
4. Sotiras A Bilello M Rozycki M Kirby JS Freymann JB Farahani K Davatzikos C. Bakas S, Akbari H. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. 2017.
5. Bakas S. et al. Segmentation labels and radiomic features for the pre-operative scans of the tcga-gbm collection, 07 2017.
6. Bakas S. et al. Segmentation labels and radiomic features for the pre-operative scans of the tcga-lgg collection, 07 2017.
7. Diederik P. Kingma, Danilo J. Rezende, Shakir Mohamed, and Max Welling. Semi-Supervised Learning with Deep Generative Models. *arXiv e-prints*, page arXiv:1406.5298, Jun 2014.
8. Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
9. Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv e-prints*, page arXiv:1312.6114, Dec 2013.
10. Eric Jang, Shixiang Gu, and Ben Poole. Categorical Reparameterization with Gumbel-Softmax. *arXiv e-prints*, page arXiv:1611.01144, Nov 2016.
11. Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.

12. Higgins et al. beta-VAE: Learning basic visual concepts with a constrained variational framework.
13. Ourselin S Vercauteren T Wang G, Li W. Automatic brain tumor segmentation using cascaded anisotropic convolutional neural networks. BrainLes 2017, Springer LNCS 10670 (2018) 178190.
14. Wick W. Bendszus M. Maier-Hein K.H. Isensee F., Kickingereder P. No new-net. BrainLes 2018, Springer LNCS 11384 (2019) 234244.
15. Mehta S. Nuechterlein N. 3d-espnet with pyramidal refinement for volumetric brain tumor image segmentation. BrainLes 2018, Springer LNCS 11384 (2019) 245253.
16. Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS-W*, 2017.
17. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
18. Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. In *Advances in neural information processing systems*, pages 971–980, 2017.
19. Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
20. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
21. Lawrence D. Brown, T. Tony Cai, and Anirban DasGupta. Interval estimation for a binomial proportion. *Statist. Sci.*, 16(2):101–133, 05 2001.