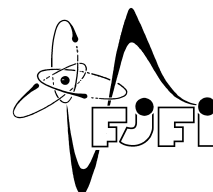




CZECH TECHNICAL UNIVERSITY IN PRAGUE
Faculty of Nuclear Sciences and Physical
Engineering



Hybrid Discriminative-Generative Training for Set data

Hybridní diskriminativně-generativní modely pro množinová data

Master Thesis

Author: **Bc. Jakub Bureš**
Supervisor: **doc. Ing. Václav Šmídl, Ph.D.**
Academic year: 2021/2022

Acknowledgment:

I would like to thank doc. Ing. Václav Šmídl, Ph.D. for his expert guidance and patience during this academic year.

Author's declaration:

I declare that this Master Thesis is entirely my own work and I have listed all the used sources in the bibliography.

Prague, May 2, 2022

Bc. Jakub Bureš

Název práce:

Hybridní diskriminativně-generativní modely pro množinová data

Autor: Bc. Jakub Bureš

Obor: Matematické inženýrství

Zaměření: Aplikované matematicko-stochastické metody

Druh práce: Diplomová práce

Vedoucí práce: doc. Ing. Václav Šmídl, Ph.D.

Abstrakt: Tato diplomová práce se zabývá hybridními diskriminativními a generativními modely a jejich možným využitím v multi–instančním učení, kde je jeden vzorek tvořen množinou vektorů. Nejprve se seznámíme s technikáliemi tohoto přístupu, po té provedeme jednoduchý experiment a nakonec se budeme věnovat samotnému multi–instančnímu učení. Hlavním cílem je pak využít strukturu HMill s knihovnou Mill.jl a rozšířit je o kontrastivní učení na množinová data, kde ukážeme výhody oproti diskriminativnímu učení.

Klíčová slova: hybridní diskriminativní a generativní modely, multi–instanční učení

Title:

Hybrid Discriminative-Generative Training for Set data

Author: Bc. Jakub Bureš

Abstract: This Master Thesis deals with hybrid discriminative and generative modeling and its possible utilization in multiple–instance learning. At first, technicalities of this approach are introduced; consequently a simple experiment is performed and lastly multi–instance learning itself is taken care of. The main aim of this work is to use the HMill framework with the Mill.jl library and involve contrastive learning into it, where the benefits of this approach are shown.

Key words: hybrid discriminative and generative modeling, multiple instance learning

Contents

Introduction	10
1 Theoretical introduction	11
1.1 Mathematical notation	11
1.2 Probability theory	11
1.3 Supervised Learning	12
1.3.1 Prediction	13
1.4 Unsupervised learning	14
1.5 Bayesian Inference	15
1.5.1 Choice of prior distribution	16
1.5.2 Prediction	17
2 Discriminative vs. Generative Models	18
2.1 Overview	18
2.2 Discriminative modeling	18
2.2.1 One-hot encoding	20
2.3 Generative modeling	21
2.3.1 Variational autoencoder	21
2.3.2 Noise-contrastive estimation	26
3 Hybrid Generative and Discriminative models	30
3.1 Energy-Based Models	30
3.2 Contrastive learning	31
3.3 Hybrid Discriminative and Generative Models, HDGM	31
3.4 Toy problem - Polynomial Regression	33
3.4.1 Experiment setup and results	34
4 Multiple Instance Learning	38
4.1 Fundamentals	38
4.2 Embedded-space paradigm	38
4.3 Training	39
4.4 Cross-validation on MIL datasets	39
4.4.1 Setup	40
4.4.2 Results	41

4.5	MIL to HDGM problem	42
4.5.1	Setup	42
4.5.2	Results	44
	Conclusion	47
	Bibliography	48
A	Computational formulas	50
A.1	Solution of $D_{\text{KL}}(\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \sigma^2 \mathbb{I}_P) \parallel \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbb{I}_P))$	50

List of Figures

1.1	Splitting the data into $K = 5$ roughly equal-sized parts.	13
1.2	Evaluation of prediction error as a function of model complexity.	15
2.1	Discriminative approach.	19
2.2	Generative approach.	19
2.3	VAE diagram.	22
2.4	Reparametrization trick.	24
2.5	Results of the NCE experiment for one-dimensional Gaussian case.	28
2.6	Results of the NCE experiment for two-dimensional Gaussian case.	29
3.1	Sensitivity of the polynomial model to parameter τ with parameter s held fixed for six different cases.	36
3.2	Sensitivity of the polynomial model to the order of polynomial $s - 1$ for six different cases.	37
4.1	The difference between standard ML and MIL [2]. Standard ML is special case of MIL with $ b = 1$	39
4.2	NN example for $z = 5$, where input and output layer are only illustrative.	40
4.3	Evaluation of prediction error with the use of training data and testing data on MIL datasets Musk1, Musk2, Fox and Tiger.	41
4.4	Evaluation of the prediction error $\widehat{\text{Err}}(\alpha)$ with the use of training data and testing data on MIL datasets.	45
4.5	Comparison of the prediction error $\widehat{\text{Err}}(z)$ for HDGM $\alpha = 0.5$ and the discriminative part only.	46

List of Tables

2.1	Transformation of a label encoding (left) to the one-hot encoding (right).	21
4.1	Results of CV evaluated on the testing data.	42
4.2	Prediction error statistics for HDGM in case of $z = 10$.	44
4.3	Comparison of prediction error statistics for HDGM $\alpha = 0.5$ and discriminative part only. Pay attention especially to the last column in each approach, $\widehat{\text{Err}}(z = 10)$.	44

List of Symbols

The next list describes several symbols that will be later used within the body of the document

$\mathbb{E} [.]$ expected value

\mathbb{R} the set of real numbers

$p(.)$ probability distribution function

List of Acronyms

PDF	probability density function
KL	Kullback-Leibler (divergence)
MLE	maximum likelihood estimation
ML	machine learning
CE	cross-entropy
CV	cross-validation
SL	supervised learning
UL	unsupervised learning
i.i.d.	independent and identically distributed
EBM	energy-based model
JEM	joint energy-based model
HDGM	hybrid discriminative-generative model
MIL	multiple instance learning

Introduction

In the field of supervised learning [18] has been achieved a tremendous progress and success in recent years. Examples of such successes include speech recognition [20] or anomaly detection [3]. A classification task is typically addressed minimizing a cross entropy loss, which is defined as an expected value of logarithm of Softmax function.

Contrastive learning [10, 11] is a machine learning method often used in representation learning for image classification or video understanding. For training such models is, most of the time, minimized the contrastive loss, which reduces the 'distance' between representations of different augmented views of the same image and increases the distance between representations of augmented views of different images.

In this research project, these two objectives are brought together and utilized in the form of a hybrid combination [12], which is used instead of the aforementioned cross entropy loss. Consequently, this approach is applied on multiple instance learning problems, where is taken advantage of the unified framework HMill and Mill.jl library [2] implemented in Julia programming language [19].

This work is arranged into 3 chapters in a logical sequence. In the first chapter is written theoretical introduction needed for a better understanding of the whole work. The second chapter consists of discriminative and generative modeling and its hybrid combination, where the simple polynomial regression experiment is performed. In the last, third, chapter is introduced the multiple instance learning with following experiments. The primary goal of this work is to test out the hybrid approach on the real data and, eventually, show its benefits in comparison to discriminative learning.

Chapter 1

Theoretical introduction

1.1 Mathematical notation

It will be most appropriate to begin by introducing the basic notation that will be used throughout this thesis. This will avoid any confusion, even though the notation is quite standard.

Notation for random variables using upper case letters of the Latin alphabet is widely used. Typically, the letters used are from the end of the alphabet, i.e. X, Y or Z . The realization of a random variable, also known as an observed value or simply an observation, will be denoted by the appropriate lower case letters. Thus, the realization x corresponds to the random variable X , which holds by analogy for other random variables. Bold symbols, for instance, \mathbf{x} or \mathbf{y} , will be used to distinguish vectors and scalars. All vectors are assumed to be column vectors, so $\mathbf{x} = (x_1, x_2, \dots, x_D)^\top$ and hence \mathbf{x}^\top is a row vector. Of course, random variables can also be vectors, but this thesis will not work with them; the primary interest is in realizations, since realizations of random variables are all data.

Any matrices will be denoted by blackboard bold Latin letters, for example, if one has N values of D -dimensional vector of observations $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, it can be simply combined into a $D \times N$ data matrix \mathbb{X} in which the j^{th} row of \mathbb{X} corresponds to the row vector \mathbf{x}_j^\top . The symbol \mathbb{I}_N denotes the square $N \times N$ identity matrix, i.e., matrix with ones on the main diagonal and zeros elsewhere.

1.2 Probability theory

Mathematical models are very well described by probability and for this reason, this section will look at some of the basic concepts of probability theory that we will need. The most important such concept is probability density (PDF). The symbol $p(x)$ will be used predominantly for the PDF, which is a function of x . In addition, this will be used for both discrete and continuous x . In this way can be achieved significant simplification and unification of all formulas and equations. Any PDF is a non-negative function and its integration over the entire space is equal to 1. This applies to multivariate case as well as it applies to univariate

case, therefore integration of joint PDFs $p(\mathbf{x}) = p(x_1, x_2, \dots, x_D)$ over the entire space is also equal to 1. In mathematical terms one can express it as follows

$$\int_{\mathbb{R}^D} p(\mathbf{x}) d\mathbf{x} = 1. \quad (1.1)$$

In other parts of this thesis we will use conditional PDFs such as $p_{\boldsymbol{\theta}}(\mathbf{x}) \equiv p(\mathbf{x}|\boldsymbol{\theta})$ that are conditioned by known parameters $\boldsymbol{\theta}$. The constraint (1.1) can be always fulfilled by redefining the PDF as

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{p_{\boldsymbol{\theta}}^0(\mathbf{x})}{Z(\boldsymbol{\theta})}, \quad Z(\boldsymbol{\theta}) = \int_{\mathbb{R}^D} p_{\boldsymbol{\theta}}^0(\mathbf{x}) d\mathbf{x}, \quad (1.2)$$

where $p_{\boldsymbol{\theta}}^0(\mathbf{x})$ specifies the functional form of the $p_{\boldsymbol{\theta}}(\mathbf{x})$ and does not need to integrate to 1. The normalization constant $Z(\boldsymbol{\theta})$ is often called the partition function.

The average value of some function $g(x)$ under a probability distribution $p(x)$ is typically denoted by $\mathbb{E}[g(x)]$ and it is called expected value or mean [1]. For a continuous variable, expected value are expressed in terms of an integration with respect to the corresponding probability density

$$\mathbb{E}[g(x)] = \int_{\mathbb{R}} p(x)g(x)dx. \quad (1.3)$$

In the case of a discrete variable, one has to keep in mind that an integration turns into a sum over all x . To specify over which PDF the expectation is calculated, the notation $\mathbb{E}_{p(x)}[g(x)]$ can be used.

1.3 Supervised Learning

Supervised learning (SL), in less academic terms called "learning with a teacher", is one of the machine learning tasks [18]. The goal of this approach is to make a good prediction of the output y (sometimes also called target variable), denoted by the symbol \hat{y} , with given input \mathbf{x} . This prediction is obtained through learning a model $f_{\boldsymbol{\theta}}(\mathbf{x}) \equiv f(\mathbf{x}; \boldsymbol{\theta})$ that minimizes a loss function $\mathcal{L}(f_{\boldsymbol{\theta}}(\mathbf{x}), y)$ (also known as the error function), where $\boldsymbol{\theta}$ are the parameters of the model.

To construct this prediction one needs data, hence it is supposed that we have available set of independent and identically distributed (i.i.d.) observations, input–output paired samples denoted by $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, eventually, this may in fact be

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N, \quad \mathbf{x}_i \in \mathbb{R}^D, \quad y_i \in \mathbb{R}, \quad \forall i = 1, \dots, N. \quad (1.4)$$

Such setting is usually known as training data and its applications are regression problems. As an example, we can mention a couple of typically used loss functions for such problems. They are squared error and absolute error

$$\mathcal{L}(\hat{y}, y) = \begin{cases} (y - \hat{y})^2 \\ |y - \hat{y}| \end{cases} \quad (1.5)$$

where $\hat{y} = \hat{f}_{\hat{\theta}}(\mathbf{x}) = f(\mathbf{x}; \hat{\theta})$. As we are not quite interested in regression problems in this thesis, we will mainly deal with the second approach. That is classification problems, i.e. when $y \in \mathcal{C}$ is qualitative output and where \mathcal{C} is a finite set. A typical example is binary classification, where $\mathcal{C} = \{0, 1\}$. However, classification will be object of interest later in Section 2.

Here it is clear why the term "learning with a teacher" is used. This metaphor means that the student presents output $\hat{y} = \hat{f}_{\hat{\theta}}(\mathbf{x}) = f(\mathbf{x}; \hat{\theta})$ and the teacher provides either a correct answer and/or an error that corresponds to the student's answer.

1.3.1 Prediction

The generalization performance, i.e. the performance on out-of-sample data of the models learned by the algorithm relates to its prediction capability on independent test data \mathcal{T} . Assessment of this performance is extremely important in practice, since it guides the choice of learning method or model, and gives us a measure of the quality of the ultimately chosen model. There are in fact two separate goals to achieve:

1. *Model selection* - estimating the performance of different models in order to choose the best one.
2. *Model assessment* - having chosen a final model, estimating its prediction error (generalization error) on new data.

1.3.1.1 Cross-Validation

The simplest and most widely used method for estimating prediction error of the model $\hat{f}_{\hat{\theta}}$ is called *cross-validation* (CV) [17]. It is used for direct estimating of the expected extra-sample error

$$\text{err} = \mathbb{E} [\mathcal{L}(y, \hat{y})], \quad (1.6)$$

the measure how accurately is the model able to predict output values for previously unseen data - independent test sample. In an ideal case, if we have sufficient number of data, we can set aside a test set and use it to assess the performance of our prediction model. Since data are often scarce, this is usually not possible.

Very elegant solution to this problem is via *K-fold cross-validation* [17]. It uses part of the available data for fitting the model, and a different part for testing. We split the data into K roughly equal-sized parts, for example, when $K = 5$, the scenario is shown in Figure 1.1. For

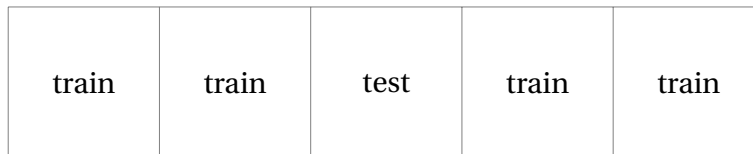


Figure 1.1: Splitting the data into $K = 5$ roughly equal-sized parts.

the j^{th} part (third in Figure 1.1), we train the model to the other $K - 1$ parts of the data, and calculate the prediction error of the fitted model when predicting the j^{th} part of the data. We repeat this process for $j = \{1, 2, \dots, K\}$ and combine the K estimates of prediction error. Let $\gamma : \{1, \dots, N\} \rightarrow \{1, \dots, K\}$ be an indexing function that indicates the partition to which observation j is allocated by the randomization. Symbol $\hat{f}_{\theta}^{-j}(\mathbf{x})$ denotes the fitted model, computed with the j^{th} part of the data removed. Then the cross-validation estimate of prediction error is defined by

$$\text{CV}(\hat{f}_{\theta}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, \hat{f}_{\theta}^{-\gamma(i)}(\mathbf{x}_i)). \quad (1.7)$$

Typical choices of K are 5 or 10 and even case $K = N$ that is known as *leave-one-out* cross-validation. Generally, there is not an universal way of choosing K , since it strongly depends on the available number of data.

The biggest problem of this method is a fact that it is computationally very expensive. For extremely complicated and complex models that are trained for hours or days, is cross-validation inconvenient approach of estimating the prediction error.

1.3.1.2 Model complexity

We shall also discuss the problem of the model complexity. Consider a polynomial regression problem, where the model is defined by

$$f_{\theta}(x) = \sum_{i=0}^{s-1} \theta_i x^i. \quad (1.8)$$

Here, overfitting occurs very frequently. The complexity of the model of this case is very intuitive, as it is just the order of the polynomial, $s - 1$. Smaller orders of the polynomial (may) give rather poor fits to the data in contrast to a much higher order polynomial giving an excellent fit. However, such a polynomial passes exactly through each data point, oscillates wildly, and gives a poor prediction for the new input variable x_0 .

To obtain some quantitative insight into the dependence of the generalization performance on model complexity, consider a separate test set of data (testing data) used to assess the performance of the model. In general, the prediction error evaluated on the training data for increasing the complexity of the model approaches zero. On the other hand, the prediction error evaluated on the testing data for increasing model complexity is (from a certain point) increasing as well. The typical scenario is illustrated in Figure 1.2.

1.3.1.3 Akaike information criterion

1.4 Unsupervised learning

The previous section dealt with input–output paired samples D . The second approach is a logical modification of SL, based on data without labels. Such setting is called unsupervised

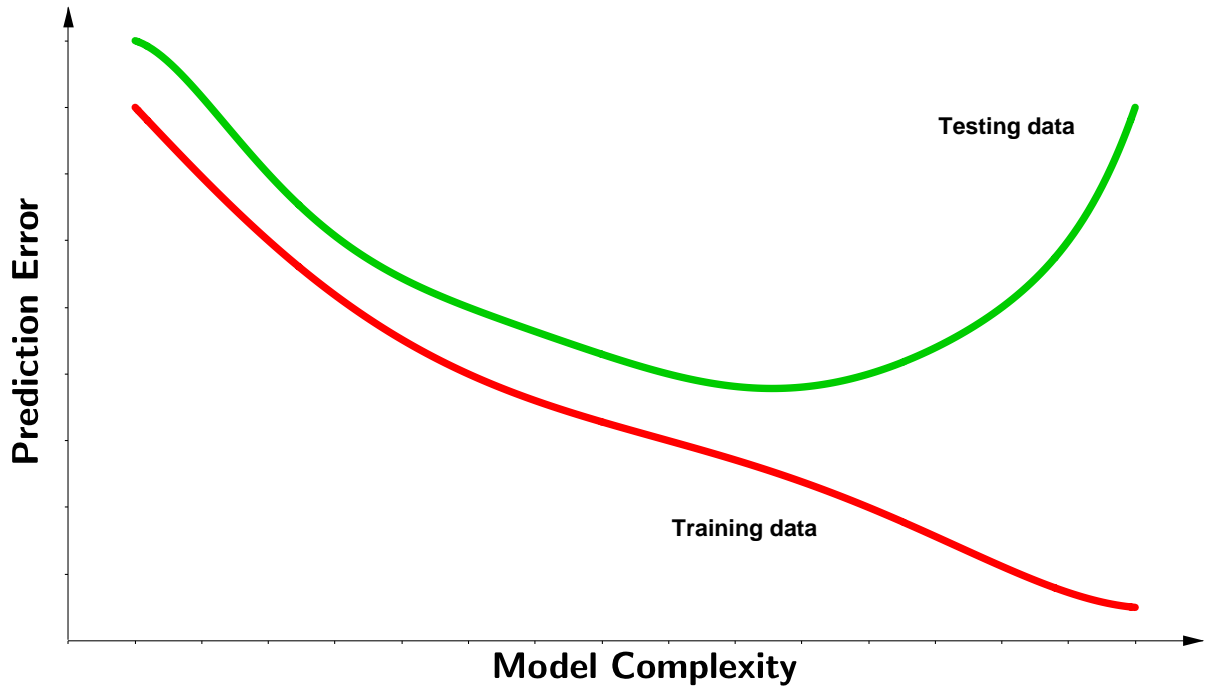


Figure 1.2: Evaluation of prediction error as a function of model complexity.

learning (UL) or "learning without a teacher". Unlike SL, one has a set of N observations in the form of $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$ and nothing more. In this case, the student learns without any feedback from a supervisor or teacher providing correct answers. The goal is to directly infer the properties of $p(\mathbf{x})$.

1.5 Bayesian Inference

The Bayesian methodology is a well established approach to statistical inference and became very important technique in statistics and data analysis. As its name suggests, Bayesian statistics is based on application of Bayes' rule. In this chapter, we briefly review basic concept of this approach, which was suggested here [9].

Let the measured data be denoted by D , defined according to previous section 1.1. A parametric probabilistic model of the data D is given by the probability density function $p(D|\boldsymbol{\theta})$, where $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^s$ denotes parameters of the model. The main idea behind Bayesian theory is the treatment of the unknown parameters $\boldsymbol{\theta}$ as a random variable. Bayes' rule is applied to infer model parameters $\boldsymbol{\theta}$, therefore

$$p(\boldsymbol{\theta}|D) = \frac{p(\boldsymbol{\theta}, D)}{p(D)} = \frac{p(D|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int_{\Theta} p(D|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}}. \quad (1.9)$$

Since $p(D)$ is just the normalization constant, Equation (1.9) is often simplified to

$$p(\boldsymbol{\theta}|D) \propto p(D|\boldsymbol{\theta})p(\boldsymbol{\theta}). \quad (1.10)$$

Symbol \propto means equal up to the normalization constant. The term $p(\boldsymbol{\theta}|D)$ is known as the *posterior* distribution, $p(D|\boldsymbol{\theta})$ as the *observation model*, and $p(\boldsymbol{\theta})$ is called the *prior* distribution of the $\boldsymbol{\theta}$. Note that evaluation of the normalization constant can be computationally expensive, in higher dimension even intractable.

Popular choices for an optimal value of the point estimate are:

1. Maximum A posteriori estimate (MAP)

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}|D) \quad (1.11)$$

This method estimates $\boldsymbol{\theta}$ as the mode of the posterior distribution. It appears to be computationally attractive, as it is not necessary to evaluate the normalization constant.

2. Mean or expected value

$$\hat{\boldsymbol{\theta}}_{\text{B}} = \int_{\Theta} \boldsymbol{\theta} p(\boldsymbol{\theta}|D) d\boldsymbol{\theta} = \mathbb{E}_{p(\boldsymbol{\theta}|D)} [\boldsymbol{\theta}] \quad (1.12)$$

Mean value, unlike MAP estimate, may be very expensive to compute because of the required integration. This may lead to further approximations such as EM algorithm [13].

1.5.1 Choice of prior distribution

For the posterior computation, it is necessary to specify the prior distribution $p(\boldsymbol{\theta})$, unfortunately, this might not be easily determined. This can be achieved through knowledge of previous models, expert knowledge, their combination, or even uncertainty about $\boldsymbol{\theta}$ being a viable option.

There are also many practical aspects of priors:

- *Regularization* - supplementing the data if there are scarce, insufficient data, or poorly defined model.
- *Restrictive conditions* - imposing various restrictions on the parameters $\boldsymbol{\theta}$ reflecting physical constraints. The choice of a prior distribution with bounded support will also result in a posterior distribution with bounded support.
- *Non-informative prior* - if the data are informative enough to make a prediction, it is proposed to choose a prior with minimal impact on the posterior distribution, such as uniform distribution. However, typical choices of non-informative priors are the so-called *Jeffreys priors* [5].

1.5.2 Prediction

We are not usually interested in the value of $\hat{\boldsymbol{\theta}}$ itself, but rather, once the model is estimated, we are interested in making a prediction of the output variable y_0 for the new input variable \mathbf{x}_0 . Note that the symbol D contains all previously given data \mathbf{x} and y . The posterior predictive distribution is then determined by the distribution of y_0 , marginalized over the posterior

$$p(y_0|\mathbf{x}_0, D) = \int_{\Theta} p(y_0|\mathbf{x}_0, \boldsymbol{\theta}) p(\boldsymbol{\theta}|D) d\boldsymbol{\theta}. \quad (1.13)$$

When the distribution $p(\boldsymbol{\theta}|D)$ is not available, we have to approximate leveraging the Dirac delta function $\delta(x)$

$$p(y_0|\mathbf{x}_0, D) = \int_{\Theta} p(y_0|\mathbf{x}_0, \boldsymbol{\theta}) \delta(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) d\boldsymbol{\theta} = p(y_0|\mathbf{x}_0, \hat{\boldsymbol{\theta}}), \quad (1.14)$$

causing an error. In typical MAP, this is known as *over-fitting*.

Chapter 2

Discriminative vs. Generative Models

2.1 Overview

Machine learning models can be classified into two main categories, discriminative and generative models. Simply put, a discriminative model makes predictions based on conditional probability $p(y|\mathbf{x})$ and is used for classification or regression problems. In other words, discriminative models distinguish the decision boundary between the classes. It corresponds to learning parameters that maximize the conditional probability distribution $p(y|\mathbf{x})$. On the contrary, a generative model revolves around the distribution of a data set to return a probability for a given example. Rather than looking at classes and trying to find something to separate them, it focuses only on the one class at the time and builds a model what that certain class looks like, then turns attention to the other class. To express it more formally, generative models learn parameters that maximize $p(\mathbf{x}|y)$ and $p(y)$. Since

$$p(\mathbf{x}, y) = p(\mathbf{x}|y) \cdot p(y), \quad (2.1)$$

with joint PDF it is possible to generate new (\mathbf{x}', y') pairs. In some cases, the use of the second decomposition $p(\mathbf{x}, y) = p(y|\mathbf{x}) \cdot p(\mathbf{x})$ is also an option. Note that in an unsupervised setting, the task is reduced to inferring only $p(\mathbf{x})$.

2.2 Discriminative modeling

In this section, we review the basics of discriminative modeling proposed in [12]. Given a data distribution through the probability density $p(\mathbf{x})$ and a label distribution with probability density $p(y|\mathbf{x})$ containing C categories. In this thesis, we focus on classification problems, where the label y is now a qualitative variable, taking on C possible values, and comes from a finite set \mathcal{C} . A classification problem is typically solved using a parametric function $f_{\theta} : \mathbb{R}^D \rightarrow \mathcal{C}$, where θ denotes the parameters of the model. In practice, the function f_{θ} is often used in the form of $\mathbb{R}^D \rightarrow \mathbb{R}^C$. This function maps each data point $\mathbf{x} \in \mathbb{R}^D$ to C real-valued numbers known as logits. It should be noted that \mathbb{R}^C is allowed here due to the utilization of

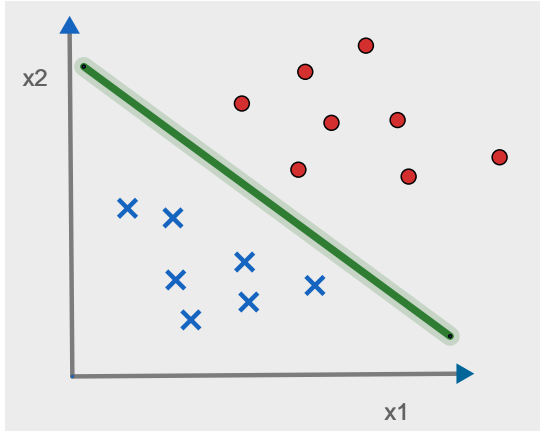


Figure 2.1: Discriminative approach.

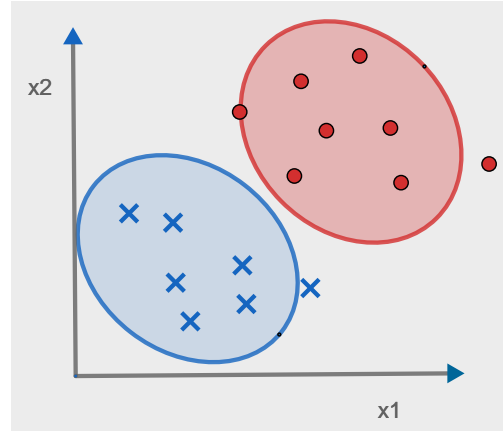


Figure 2.2: Generative approach.

one-hot encoding, which will be explained in Section 2.2.1. Logits are used to parameterize a categorical distribution through the function

$$q_{\theta}(y|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{y \in \mathcal{C}} \exp(f_{\theta}(\mathbf{x})[y])}, \quad (2.2)$$

which is known as the Softmax. In other words, the data density $p(y|\mathbf{x})$ is modeled by a parameterized family of functions $\{q_{\theta}(y|\mathbf{x}) | \theta \in \Theta\}$ and thus $p(y|\mathbf{x})$ is assumed to belong to this family. Note that the convention $f_{\theta}(\mathbf{x})[y]$ means the element y^{th} of $f_{\theta}(\mathbf{x})$. For learning f_{θ} is usually minimized cross-entropy loss

$$\text{CE}(\theta) = -\mathbb{E}_{p_{\text{data}}(y, \mathbf{x})} [\log q_{\theta}(y|\mathbf{x})] \approx -\frac{1}{N} \sum_{i=1}^N \log q_{\theta}(y_i | \mathbf{x}_i). \quad (2.3)$$

The rationale for this objective comes from minimizing the Kullback-Leibler (KL) divergence with a target distribution $p(y|\mathbf{x})$ [7]. In general, the KL divergence (or KL distance) from $p(y|\mathbf{x})$ to $q_{\theta}(y|\mathbf{x})$ is defined as

$$D_{\text{KL}}(p(y|\mathbf{x}) \| q_{\theta}(y|\mathbf{x})) = \int p(y|\mathbf{x}) \log \frac{p(y|\mathbf{x})}{q_{\theta}(y|\mathbf{x})} dy = \mathbb{E}_{p(y|\mathbf{x})} \left[\log \frac{p(y|\mathbf{x})}{q_{\theta}(y|\mathbf{x})} \right] \quad (2.4)$$

and has the following properties:

1. $D_{\text{KL}}(p(y|\mathbf{x}) \| q_{\theta}(y|\mathbf{x})) \geq 0$,
2. $D_{\text{KL}}(p(y|\mathbf{x}) \| q_{\theta}(y|\mathbf{x})) = 0$ iff $p(y|\mathbf{x}) = q_{\theta}(y|\mathbf{x})$ almost everywhere,
3. $D_{\text{KL}}(p(y|\mathbf{x}) \| q_{\theta}(y|\mathbf{x})) \neq D_{\text{KL}}(q_{\theta}(y|\mathbf{x}) \| p(y|\mathbf{x}))$ and KL divergence does not obey the triangle inequality.

The third property indicates that care is needed in the syntax describing KL divergence. We say that (2.4) is from $p(y|\mathbf{x})$ to $q_{\theta}(y|\mathbf{x})$. Using the logarithmic property, (2.4) can be further rewritten in the form

$$\mathbb{E}_{p(y|\mathbf{x})} \left[\log \frac{p(y|\mathbf{x})}{q_{\theta}(y|\mathbf{x})} \right] = \mathbb{E}_{p(y|\mathbf{x})} [\log p(y|\mathbf{x})] - \mathbb{E}_{p(y|\mathbf{x})} [\log q_{\theta}(y|\mathbf{x})], \quad (2.5)$$

where subscript θ emphasizes that $q_{\theta}(y|\mathbf{x})$ is the approximative density we get to control. Note that the first term does not depend on θ and therefore minimizing either CE or KL divergence is equivalent. Finally, by minimizing with respect to θ we obtain

$$\min_{\theta} D_{\text{KL}}(p(y|\mathbf{x}) \| q_{\theta}(y|\mathbf{x})) = \min_{\theta} -\mathbb{E}_{p(y|\mathbf{x})} [\log q_{\theta}(y|\mathbf{x})]. \quad (2.6)$$

For the sake of clarity, the expected value will be discussed. In practice, it is dealt with discrete data, so the term $\mathbb{E}_{p(y|\mathbf{x})} [\log q_{\theta}(y|\mathbf{x})]$ takes the form of

$$\mathbb{E}_{p(y|\mathbf{x})} [\log q_{\theta}(y|\mathbf{x})] \approx \sum_{k=1}^C p(y_k|\mathbf{x}) \log q_{\theta}(y_k|\mathbf{x}). \quad (2.7)$$

This part deserves further discussion for a few reasons:

- Maximum likelihood estimation (MLE) of θ is equivalent to minimizing the KL distance.
- One may encounter the concepts of minimization or maximization of CE.

To address these reasons, it is necessary to briefly review the MLE. The principle of MLE assumes that the most reasonable values for θ are those for which the probability of the observed sample is highest. Since $q_{\theta}(y|\mathbf{x})$ is model PDF, we have to follow the objective function

$$\mathcal{L}(\theta) = \sum_{i=1}^N \log q_{\theta}(y_i|\mathbf{x}_i). \quad (2.8)$$

Further optimization of $\mathcal{L}(\theta)$ gives

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log q_{\theta}(y_i|\mathbf{x}_i) \quad (2.9)$$

$$= \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^N \log q_{\theta}(y_i|\mathbf{x}_i). \quad (2.10)$$

2.2.1 One-hot encoding

Machine learning (ML) algorithms can misinterpret the numeric values of labels if there exists a hierarchy between them. One-hot encoding is a very common approach for dealing with this issue, in order to improve the algorithm performance. Each unique category value is

transformed into a new column, and these dummy variables are then filled with 0 or 1 (0 for FALSE and 1 for TRUE). For the sake of clarity, the transformation of a label encoding into a one-hot encoding is illustrated in the following table 2.1.

However, this method has its own downsides. For example, it creates new variables and if there exist many unique category values, the models have to deal with a large number of predictors, which leads to the so-called *Big-p problem* [8]. Also, one-hot encoding causes multicollinearity between the individual variables, which may lead to reducing the model's accuracy.

Food Name	Categorical #	Calories	\Rightarrow	Pizza	Hamburger	Caviar	Calories
Pizza	1	266		1	0	0	266
Hamburger	2	295		0	1	0	295
Caviar	3	264		0	0	1	264

Table 2.1: Transformation of a label encoding (left) to the one-hot encoding (right).

2.3 Generative modeling

2.3.1 Variational autoencoder

The first generative modeling approach that will be discussed is the variational autoencoder (VAE). In this section, motivation will be addressed and individual mathematical aspects will be discussed in detail.

2.3.1.1 Problem scenario

Assume that the data $\mathbf{X} = \{x_i\}_{i=1}^N$ are generated by some random process involving an unobserved continuous variable \mathbf{z} , which will be referenced as a latent variable or code. The objective is again to find the PDF of the given data in parametric form $p_{\theta}(\mathbf{x})$. One can choose an approximative distribution in the form of

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}, \quad (2.11)$$

But such an approximation is usually very expensive to compute or can even be intractable. Intractability of the $p_{\theta}(\mathbf{x})$ makes posterior PDF $p_{\theta}(\mathbf{z}|\mathbf{x})$ also intractable.

2.3.1.2 Naive approach

One of the simplest ways to solve this problem may seem to be to build a model depending on the latent variable $f_{\theta}(\mathbf{z})$ and try to train its parameters. For simplicity, let $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbb{I}_P)$, where P denotes the dimension of the latent space \mathbf{z} , and also let

$$\mathbf{x} = f_{\theta}(\mathbf{z}) + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \cdot \mathbb{I}_D) \quad (2.12)$$

which actually gives

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; f_{\theta}(\mathbf{z}), \sigma^2 \cdot \mathbb{I}_D). \quad (2.13)$$

It may be noted that the subscript D represents the dimension of the data point \mathbf{x} . The true PDF of the given data can be cleverly written using the empirical PDF, i.e., in the form of $p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}_i)$, which can be exploited by finding the θ parameters by minimizing the $D_{\text{KL}}(p_{\text{emp}}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x}))$. Since minimizing the KL distance is equivalent to ML estimation and using the approximative form (2.11), the following holds

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) \quad (2.14)$$

$$= \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^N \log \int \mathcal{N}(\mathbf{x}_i; f_{\theta}(\mathbf{z}), \sigma^2 \cdot \mathbb{I}_D) \cdot \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbb{I}_P) d\mathbf{z} \quad (2.15)$$

$$= \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^N \log \sum_{j=1}^P \exp\left(-\frac{1}{2\sigma^2} (\mathbf{x}_i - f_{\theta}(\mathbf{z}_j))^{\top} (\mathbf{x}_i - f_{\theta}(\mathbf{z}_j))\right). \quad (2.16)$$

Integration over \mathbf{z} is represented by sampling. In iterations, for an incorrect value of θ , all the generated samples may be away from the samples of \mathbf{x} , and the gradient is poor.

2.3.1.3 Variational Bayes approach

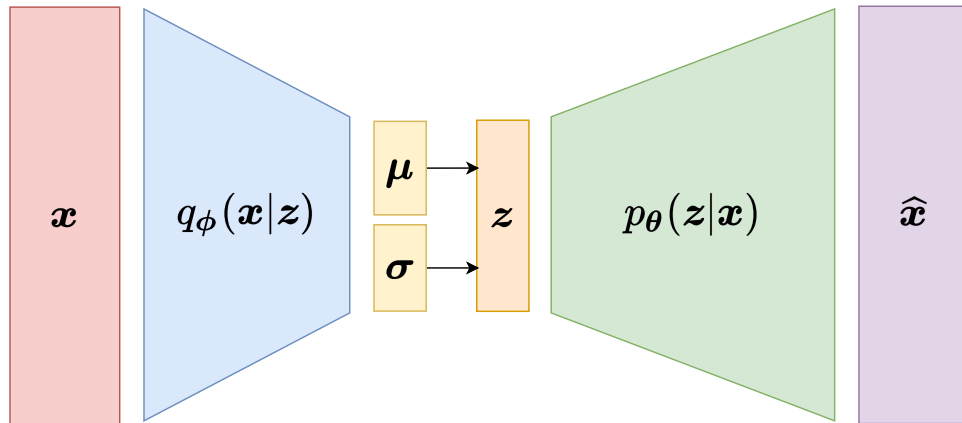


Figure 2.3: VAE diagram.

To solve this problem, it is necessary to introduce a further approximative posterior distribution $q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p_{\theta}(\mathbf{z}|\mathbf{x})$ with parameters ϕ , preferably Gaussian. The standard terminology refers to the model $q_{\phi}(\mathbf{z}|\mathbf{x})$ as a probabilistic *encoder* and $p_{\theta}(\mathbf{x}|\mathbf{z})$ is called a probabilistic

decoder. For VAE, the idea is to use the KL distance from $q_\phi(\mathbf{z}|\mathbf{x})$ to $p_\theta(\mathbf{z}|\mathbf{x})$, which produces

$$\begin{aligned}
 D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}|\mathbf{x})) &= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\
 &= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x}) p_\theta(\mathbf{x})}{p_\theta(\mathbf{x}|\mathbf{z}) p_\theta(\mathbf{z})} d\mathbf{z} \\
 &= \log p_\theta(\mathbf{x}) + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{x}|\mathbf{z}) p_\theta(\mathbf{z})} d\mathbf{z} \\
 &= \log p_\theta(\mathbf{x}) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} - \log p(\mathbf{x}|\mathbf{z}) \right] \\
 &= \log p_\theta(\mathbf{x}) + D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z})) - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})].
 \end{aligned} \tag{2.17}$$

Using the last equality of (2.22), it is possible to rewrite the equation in its typical form

$$\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}|\mathbf{x})) = \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z}))}_{= L(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})}, \tag{2.18}$$

where the right-hand side is called *variational lower bound*. There is no uniformity in terminology, and thus one can also encounter the name evidence lower bound (ELBO). The first term on the right-hand side is known as reconstruction loss, and the second term is often called a regularization term. As a KL distance is always non-negative, it holds

$$\log p_\theta(\mathbf{x}) \geq L(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}). \tag{2.19}$$

The objective is to maximize the log-likelihood $\log p_\theta(\mathbf{x})$ which is equivalent to minimizing the negative log-likelihood and that is what will be used here. At this point, we have a lower bound for one data point \mathbf{x} , but we need to include all observations in the lower bound. The joint log-likelihood can be rewritten as a sum over the marginal log-likelihoods of individual observations $\log p_\theta(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \sum_{i=1}^N \log p_\theta(\mathbf{x}_i)$ that completes all the building blocks needed to determine the optimization equation. This formulation provides one major advantage, which is that it is now possible to jointly optimize both the generative parameters $\boldsymbol{\theta}$ and the variational parameters $\boldsymbol{\phi}$ as follows

$$\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\theta}, \boldsymbol{\phi}}{\operatorname{argmin}} - \sum_{i=1}^N \log p_\theta(\mathbf{x}_i) \tag{2.20}$$

$$= \underset{\boldsymbol{\theta}, \boldsymbol{\phi}}{\operatorname{argmin}} - \sum_{i=1}^N L(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}_i) \tag{2.21}$$

$$= \underset{\boldsymbol{\theta}, \boldsymbol{\phi}}{\operatorname{argmin}} - \sum_{i=1}^N \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}_i)} [\log p_\theta(\mathbf{x}_i|\mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}_i) \| p_\theta(\mathbf{z})). \tag{2.22}$$

For a better understanding of the problem, a VAE diagram is shown in Figure 2.3. Note that the latent space is usually much smaller than the input space, and for this reason, it is also sometimes called the bottleneck.

2.3.1.4 Reparameterization trick

The key success of VAE lies in the fact that (2.22) can be efficiently computed using *reparameterization trick*, where we express z as a deterministic variable

$$z = g_{\phi}(\epsilon, x), \quad (2.23)$$

where ϵ stands for an auxiliary variable with independent marginal $p(\epsilon)$ and $g_{\phi}(\cdot)$ is a function parameterized by ϕ .

A common explanation for this trick is that during the optimization the gradient cannot back-propagate through a random node. So, in the case of VAE, the reparameterization trick shifts the source of randomness to another variable different from z and allows differentiation with respect to z . However, this explanation may not be sufficient, and for this reason we will state a more formal justification. Consider taking the gradient with respect to θ of $\mathbb{E}_{p(z)}[f_{\theta}(z)]$. It can be easily computed as

$$\nabla_{\theta} \mathbb{E}_{p(z)}[f_{\theta}(z)] = \nabla_{\theta} \int p(z) f_{\theta}(z) dz \quad (2.24)$$

$$= \int p(z) \nabla_{\theta} f_{\theta}(z) dz \quad (2.25)$$

$$= \mathbb{E}_{p(z)}[\nabla_{\theta} f_{\theta}(z)]. \quad (2.26)$$

The result is obvious; the gradient of the expectation is equal to the expectation of the gradient. However, the gradient of the expectation becomes much more interesting if the PDF $p_{\theta}(z)$ is also parameterized by θ , resulting in

$$\nabla_{\theta} \mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)] = \nabla_{\theta} \int p_{\theta}(z) f_{\theta}(z) dz \quad (2.27)$$

$$= \int p_{\theta}(z) \nabla_{\theta} f_{\theta}(z) dz + \int f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz \quad (2.28)$$

$$= \mathbb{E}_{p_{\theta}(z)}[\nabla_{\theta} f_{\theta}(z)] + \int f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz. \quad (2.29)$$

The second term of (2.29) is not guaranteed to be an expectation and this very fact indicates that backpropagation would not compute an estimate of $\nabla_{\theta} \mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)]$. That being the case, if we apply the reparameterization trick $z = g_{\theta}(\epsilon, x)$ to this simple example, we get

$$\mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)] = \mathbb{E}_{p(\epsilon)}[f(g_{\theta}(\epsilon, x))]. \quad (2.30)$$

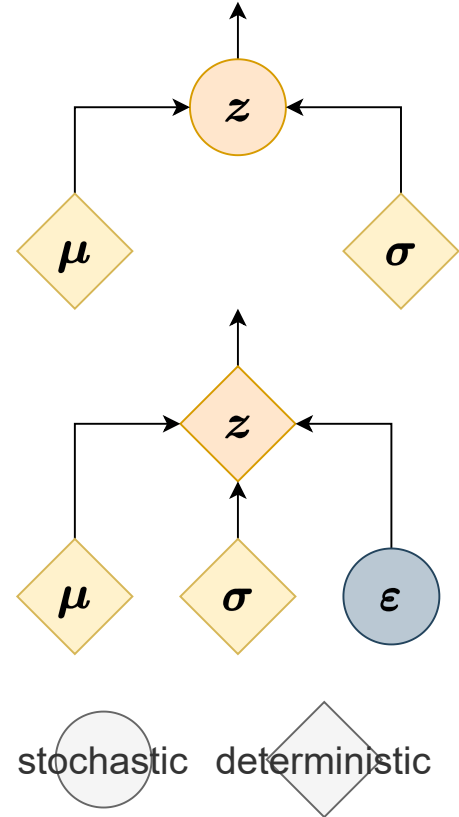


Figure 2.4: Reparametrization trick.

At this point, it is possible to take the gradient $\nabla_{\theta} \mathbb{E}_{p(\epsilon)} [f(g_{\theta}(\epsilon, \mathbf{x}))]$ analogously to that in (2.26). To be perfectly clear, the authors of [] proposed an easy exercise. Take the univariate Gaussian case $p(z|x) = \mathcal{N}(\mu, \sigma^2)$. In such a case, proper reparametrization takes the shape of

$$z = \mu + \sigma \epsilon, \quad (2.31)$$

where $\epsilon \sim \mathcal{N}(0, 1)$ and, therefore, the expectation

$$\mathbb{E}_{\mathcal{N}(z; \mu, \sigma^2)} [f(z)] = \mathbb{E}_{\mathcal{N}(\epsilon; 0, 1)} [f(\mu + \sigma \epsilon)] \approx \frac{1}{M} \sum_{j=1}^M f(\mu + \sigma \epsilon_j). \quad (2.32)$$

Note that this is nothing more than a transformation of a random variable. And this is exactly the problem with optimizing ELBO (2.22). We need to rewrite the expectation $\mathbb{E}_{q_{\phi}(z|x_i)}$ so that the Monte Carlo estimate of the expected value is differentiable with respect to ϕ .

2.3.1.5 Variational autoencoder

So far we have only dealt with VAE in general. In this section, we put everything together and specify the individual parts of the ELBO (2.22). Let the probabilistic encoder be a multivariate Gaussian with a diagonal covariance matrix

$$q_{\phi}(z|x) = \mathcal{N}(z; \mu, \sigma^2 \mathbb{I}_P) \quad (2.33)$$

and let the probabilistic decoder takes the form depending on the type of given data and model. This is typically either multivariate Gaussian or Bernoulli. Finally, let the prior $p_{\theta}(z)$ be the centered isotropic multivariate Gaussian, i.e.

$$p_{\theta}(z) = p(z) = \mathcal{N}(z; \mathbf{0}, \mathbb{I}_P), \quad (2.34)$$

where the generative parameters θ are omitted, since the chosen prior distribution lacks parameters. When using (2.33), μ and σ are non-linear functions of the data point \mathbf{x} and the variational parameters ϕ . This setting actually allows us to take the reparameterization trick in a form similar to that of Equation (2.31), which means that

$$z_{i,j} = \mu_i + \sigma_i \odot \epsilon_j, \quad (2.35)$$

where the symbol \odot denotes the Hadamard product, i.e., the element-wise product and the auxiliary variable $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$. Another important fact is that the KL distance from a Gaussian distribution to a Gaussian distribution has an analytical solution, so $D_{\text{KL}}(q_{\phi}(z|x) \| p_{\theta}(z))$ can be expressed in closed form:

$$\begin{aligned} D_{\text{KL}}(q_{\phi}(z|x) \| p_{\theta}(z)) &= D_{\text{KL}}(\mathcal{N}(z; \mu, \sigma^2 \mathbb{I}_P) \| \mathcal{N}(z; \mathbf{0}, \mathbb{I}_P)) \\ &= \frac{1}{2} \sum_{j=1}^P \left(1 + \log \sigma_j^2 - \mu_j^2 - \sigma_j^2 \right). \end{aligned} \quad (2.36)$$

Now all that is left is to plug everything into equation (2.22), which leads to the final form

$$\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\theta}, \boldsymbol{\phi}}{\operatorname{argmin}} - \sum_{i=1}^N \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_i|\mathbf{z})] - D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}_i) \| p_{\boldsymbol{\theta}}(\mathbf{z})) \quad (2.37)$$

$$= \underset{\boldsymbol{\theta}, \boldsymbol{\phi}}{\operatorname{argmin}} - \sum_{i=1}^N \left(\frac{1}{P} \sum_{j=1}^P \log p_{\boldsymbol{\theta}}(\mathbf{x}_i|\mathbf{z}_{i,j}) + \frac{1}{2} \sum_{j=1}^P \left(1 + \log \sigma_{i,j}^2 - \mu_{i,j}^2 - \sigma_{i,j}^2 \right) \right). \quad (2.38)$$

2.3.2 Noise-contrastive estimation

Suppose one has to estimate a model that is specified by a non-normalized probability density function $q_{\boldsymbol{\theta}}^0(\mathbf{x})$. In such a case, one can utilize noise-contrastive estimation (NCE). The first step is to introduce another parameter c among the estimated parameters $\boldsymbol{\theta}$. For clarity, the symbol $\boldsymbol{\theta}^* = \{\boldsymbol{\theta}, c\}$ is introduced for the set of estimated parameters, including c . Using this notation, the following equality hold

$$\log q_{\boldsymbol{\theta}^*}(\mathbf{x}) = \log q_{\boldsymbol{\theta}^*}^0(\mathbf{x}) + c, \quad (2.39)$$

which means that the newly introduced parameter c is an estimate of the negative logarithm of the normalization constant $Z(\boldsymbol{\theta})$. As the name suggests, we use noise to estimate. By our convention, let $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be the observations and $\boldsymbol{\Xi} = \{\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2, \dots, \boldsymbol{\epsilon}_N\}$ be the artificially generated noise data with known distribution $\psi(\boldsymbol{\epsilon})$. The estimate $\hat{\boldsymbol{\theta}}^*$ is then defined as

$$\hat{\boldsymbol{\theta}}^* = \underset{\boldsymbol{\theta}^*}{\operatorname{argmax}} \mathcal{L}^{\text{NC}}(\boldsymbol{\theta}^*) \quad (2.40)$$

$$= \underset{\boldsymbol{\theta}^*}{\operatorname{argmax}} \frac{1}{2N} \sum_{i=1}^N \log S_{\boldsymbol{\theta}^*}(\mathbf{x}_i) + \log(1 - S_{\boldsymbol{\theta}^*}(\boldsymbol{\epsilon}_i)) \quad (2.41)$$

$$= \underset{\boldsymbol{\theta}^*}{\operatorname{argmin}} - \frac{1}{2N} \sum_{i=1}^N \log S_{\boldsymbol{\theta}^*}(\mathbf{x}_i) + \log(1 - S_{\boldsymbol{\theta}^*}(\boldsymbol{\epsilon}_i)) \quad (2.42)$$

where $S_{\boldsymbol{\theta}^*}$ stands for a logistic function,

$$S_{\boldsymbol{\theta}^*}(\mathbf{x}) = \frac{1}{1 + \exp(-G_{\boldsymbol{\theta}^*}(\mathbf{x}))} \quad (2.43)$$

and finally, the function $G_{\boldsymbol{\theta}^*}$ represents the difference of the log-likelihoods of $q_{\boldsymbol{\theta}^*}$ and ψ , hence

$$G_{\boldsymbol{\theta}^*}(\mathbf{x}) = \log q_{\boldsymbol{\theta}^*}(\mathbf{x}) - \log \psi(\mathbf{x}). \quad (2.44)$$

It may be noted that equation (2.41) also appears in SL tasks and is called binary CE loss. It is actually a special case of CE itself. Thus, it is used for the classification of two classes. This gives an intuitive insight into how noise-contrastive estimation really works. When data and

noise are compared, the model is learned, so this method can be called learning by comparison. To make the connection with SL more explicit, denote $U = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{2N}\}$ the union of two sets \mathbf{X} and $\mathbf{\Xi}$. Then each data point \mathbf{u}_i is assigned a binary class label y_i , where $y_i = 1$ if $\mathbf{u}_i \in \mathbf{X}$ and $y_i = 0$ if $\mathbf{u}_i \in \mathbf{\Xi}$. The aim is to estimate the posterior probabilities of the classes given the data \mathbf{u}_i . To do this, one needs the class-conditional PDFs that are given by

$$p(\mathbf{u}|y=1) = q_{\theta^*}(\mathbf{u}) \quad p(\mathbf{u}|y=0) = \psi(\mathbf{u}). \quad (2.45)$$

Class labels are equally likely, so that $P(y=1) = P(y=0) = \frac{1}{2}$ and the posteriors are determined as follows

$$P(y=1|\mathbf{u}) = \frac{q_{\theta^*}(\mathbf{u})}{q_{\theta^*}(\mathbf{u}) + \psi(\mathbf{u})} = S_{\theta^*}(\mathbf{u}), \quad (2.46)$$

$$P(y=0|\mathbf{u}) = 1 - S_{\theta^*}(\mathbf{u}). \quad (2.47)$$

The class labels y_i are Bernoulli-distributed so that for the log-likelihood we get

$$\mathcal{L}^{\text{NC}}(\theta) = \sum_{i=1}^{2N} y_i \log P(y=1|\mathbf{u}_i) + (1 - y_i) \log P(y=0|\mathbf{u}_i) \quad (2.48)$$

$$= \sum_{i=1}^N \log S_{\theta^*}(\mathbf{x}_i) + \log(1 - S_{\theta^*}(\mathbf{\epsilon}_i)), \quad (2.49)$$

which is the equation that is optimized in (2.41) or (2.42).

2.3.2.1 Choice of the contrastive noise PDF

The noise distribution $\psi(\epsilon)$ can be considered as a design parameter. But this choice is not completely arbitrary because in practice the noise distribution should meet certain conditions. These are:

1. It is easy to sample from, because NCE approach relies on artificially generated noise data $\epsilon_1, \epsilon_2, \dots, \epsilon_N$.
2. Closed form for $\log \psi(\cdot)$ in order to smoothly evaluate (2.44).
3. It leads to a small mean squared error $\mathbb{E} \left[\left(\widehat{\theta} - \theta^* \right)^2 \right]$.

The authors of [] suggest using Gaussian or uniform distribution, eventually a Gaussian mixture.

EXAMPLE 2.1 (One-dimensional Gaussian distribution).

To test this approach, we performed a simple experiment. There are a total of $N = 100$ i.i.d. and one-dimensional observations x_1, x_2, \dots, x_N from an unknown distribution that is assumed to be non-normalized and Gaussian. Therefore, it is of the form

$$q_{\theta^*}(x) = \exp \left(-\frac{1}{2} \cdot \frac{(x - \mu)^2}{\sigma^2} + c \right), \quad (2.50)$$

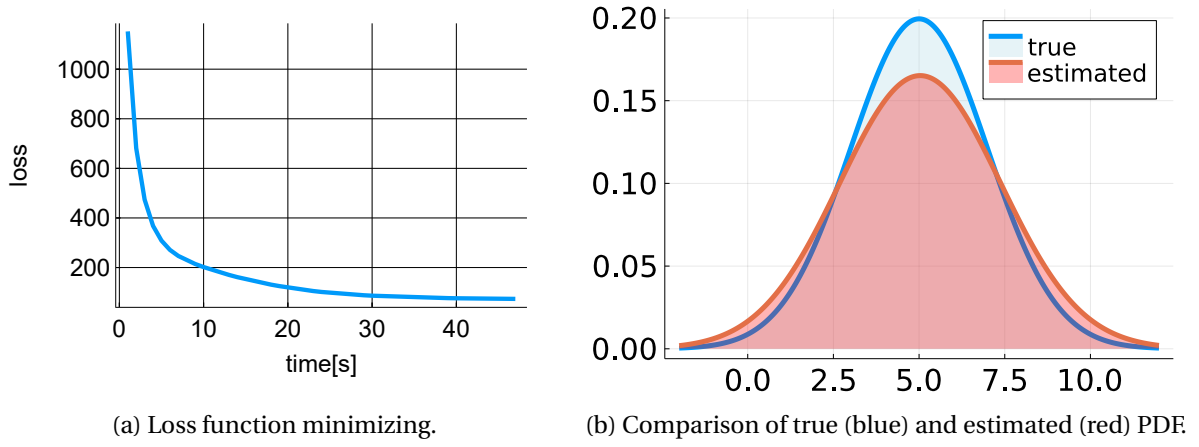


Figure 2.5: Results of the NCE experiment for one-dimensional Gaussian case.

where $\boldsymbol{\theta}^* = \{\mu, \sigma^2, c\}$. Next, we artificially generate noise data e_1, e_2, \dots, e_N , which is again easier to do using a Gaussian distribution. This means that it can be chosen, for example,

$$\psi(e) = \frac{1}{\sqrt{2\pi}10} \exp\left(-\frac{1}{2} \cdot \frac{e^2}{10}\right). \quad (2.51)$$

We choose the noise PDF intentionally so widely spread from its mean value because these two PDFs, i.e. (2.50) and (2.51), should at least partially overlap. At this point, we have all components available and it is possible to construct a function $-\mathcal{L}^{\text{NC}}(\boldsymbol{\theta}^*)$ that is minimized by using stochastic gradient descent. The following figure shows the training process and the comparison between the estimated distribution and the true one. As can be seen in Figure 2.5, this approach works quite well and for more observations, the results would be even better. In addition, the minimization of $-\mathcal{L}^{\text{NC}}(\boldsymbol{\theta}^*)$ is very fast.

EXAMPLE 2.2 (Two-dimensional Gaussian distribution).

The one-dimensional case may seem too simple, and therefore an example with a two-dimensional Gaussian distribution was performed. The experimental setup remains nearly the same; only the dimensionality of the problem differs.

Recall that the non-normalized multivariate Gaussian distribution in \mathbb{R}^2 can be written as

$$q_{\boldsymbol{\theta}^*}(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) + c\right), \quad (2.52)$$

where $\boldsymbol{\mu} \in \mathbb{R}^2$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{2 \times 2}$ is a symmetric and positive semidefinite covariance matrix. As noise PDF is chosen $\psi(\mathbf{e}) = \mathcal{N}(\mathbf{e}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$, where $\boldsymbol{\mu}_1 = (2, 2)^\top$ and $\boldsymbol{\Sigma}_1 = 10 \cdot \mathbb{I}_2$. Figure 2.6 shows the results in a similar vein to the previous case.

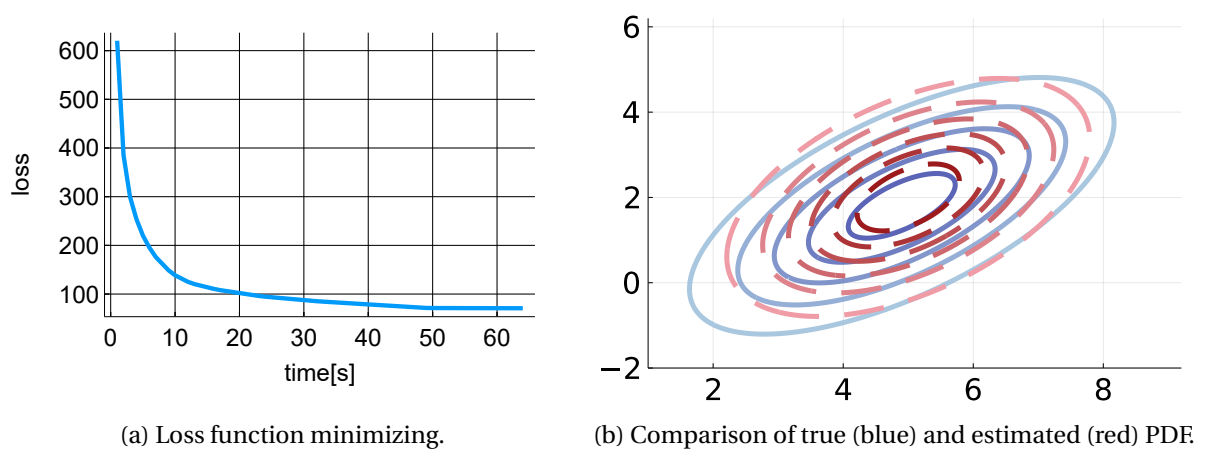


Figure 2.6: Results of the NCE experiment for two-dimensional Gaussian case.

Chapter 3

Hybrid Generative and Discriminative models

3.1 Energy-Based Models

Energy-based models (EBM) were first introduced here [14], but the following nomenclature was established here [15]. The probability densities $p(\mathbf{x})$ for $\mathbf{x} \in \mathbb{R}^D$ for EBM are assumed to be expressed in the form

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{\exp(-E_{\boldsymbol{\theta}}(\mathbf{x}))}{Z(\boldsymbol{\theta})}, \quad (3.1)$$

where the function $E_{\boldsymbol{\theta}} : \mathbb{R}^D \rightarrow \mathbb{R}$ is called energy and maps each data point \mathbf{x} to a scalar. The denominator $Z(\boldsymbol{\theta})$ is a normalization constant (also known as a partition function); thus,

$$Z(\boldsymbol{\theta}) = \sum_{i=1}^N \exp(-E_{\boldsymbol{\theta}}(\mathbf{x}_i)), \quad (3.2)$$

where the summation is over all available data points \mathbf{x} . The sum becomes an integral for a continuous \mathbf{x} . Very important observations made by the authors in [15], where they show that classifiers in supervised learning are secretly energy-based models on $p(\mathbf{x}, y)$ and can be expressed as

$$p(\mathbf{x}, y) = \frac{\exp(f_{\boldsymbol{\theta}}(\mathbf{x})[y])}{Z(\boldsymbol{\theta})}. \quad (3.3)$$

The above objective is called a joint energy model and it is obvious that $f_{\boldsymbol{\theta}}(\mathbf{x})[y] = -E_{\boldsymbol{\theta}}(\mathbf{x}, y)$. It may be useful to have only a density model of data points $p(\mathbf{x})$ without labels. This could be achieved by marginalizing $p(\mathbf{x}, y)$ over y

$$p(\mathbf{x}) = \frac{\sum_{i=1}^C \exp(f_{\boldsymbol{\theta}}(\mathbf{x})[y_i])}{Z(\boldsymbol{\theta})}, \quad (3.4)$$

where the energy is given by $E_{\boldsymbol{\theta}}(\mathbf{x}) = -\log \sum_{i=1}^C \exp(f_{\boldsymbol{\theta}}(\mathbf{x})[y_i])$. A very useful property appears when computing $p(y|\mathbf{x})$. We can take advantage of the definition of a conditional distribution

$p(y|\mathbf{x}) = \frac{p(\mathbf{x},y)}{p(\mathbf{x})}$, resulting in

$$p_{\theta}(y|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{i=1}^C \exp(f_{\theta}(\mathbf{x})[y_i])}. \quad (3.5)$$

Note that the normalization constant $Z(\theta)$ canceled out and we ended up with the same function, which was introduced in (2.2).

3.2 Contrastive learning

Contrastive learning [10, 11] is an ML technique used to learn the so-called general features of a data set by teaching the model which data points are similar or different. All this happens without labels; therefore, contrastive learning is often called the *self-supervised* technique of ML. In contrastive learning problems, it is very common to optimize an objective often called contrastive loss, which can be written in the form as follows:

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\log \frac{\exp(h_{\theta}(\mathbf{x}) \cdot h_{\theta}(\mathbf{x}'))}{\sum_{i=1}^N \exp(h_{\theta}(\mathbf{x}) \cdot h_{\theta}(\mathbf{x}_i))} \right]. \quad (3.6)$$

The function $h_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}^H$ maps each data point to a representation space of dimension H , while \mathbf{x} and \mathbf{x}' are two different augmented views of the same data point. If \mathbf{x} is an image, then an augmented view of \mathbf{x} can be obtained, for example, by rotating or colorizing that image. Note that the inner product between two vectors can be replaced with any distance metric, for instance the Euclidean distance.

What this objective does is that it tries to maximally distinguish an input \mathbf{x}_i from an alternative input \mathbf{x}'_i . In other words, (3.6) reduces the distance between the representations of different augmented views of the same image \mathbf{x}, \mathbf{x}' (positive pairs) and increases the distance between the representations of augmented views of different images (negative pairs). This means that the model should be able to distinguish between different types of image without even knowing what these images really are.

3.3 Hybrid Discriminative and Generative Models, HDGM

In this section, we will discuss an approach to combine both of the already mentioned models. The authors of the article [12] proposed a solution; however, the rationale for this objective originates from [6], where the authors show that hybrid models can outperform their purely generative or purely discriminative counterparts.

The primary goal is to train a model that can classify \mathbf{x} into classes y . Secondarily, learned models should be capable of out-of-distribution detection and serve as a generative model. To achieve these goals, a hybrid model consists of a discriminative conditional and a generative conditional by maximizing the sum of both conditional log-likelihoods

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x},y)} [\log q_{\theta}(y|\mathbf{x}) + \log q_{\theta}(\mathbf{x}|y)], \quad (3.7)$$

where the first term

$$q_{\theta}(y|\mathbf{x}) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{i=1}^C \exp(f_{\theta}(\mathbf{x})[y_i])} \quad (3.8)$$

is a standard Softmax neural net classifier (as mentioned in Equation (2.2)) and the second term

$$q_{\theta}(\mathbf{x}|y) = \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{j=1}^N \exp(f_{\theta}(\mathbf{x}_j)[y])}, \quad (3.9)$$

is a contrastive loss (3.6) attributed a label. This objective often struggles with the unknown partition function $\sum_{j=1}^N \exp(f_{\theta}(\mathbf{x}_j)[y])$, which is often intractable, specifically if the number of datapoints is very high. This obstacle is typically addressed using an approximation

$$\mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} [\log q_{\theta}(\mathbf{x}|y)] \quad (3.10)$$

$$= \mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} \left[\log \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{j=1}^N \exp(f_{\theta}(\mathbf{x}_j)[y])} \right] \quad (3.11)$$

$$\approx \mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} \left[\log \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{j=1}^M \exp(f_{\theta}(\mathbf{x}_j)[y])} \right], \quad (3.12)$$

where $M < N$ denotes the number of normalization samples. To have an adequate approximation, M must be sufficiently large, becoming exact in the limit $M \rightarrow N$. In practice, increasing M is not simple as it requires a larger memory. However, this does not apply to our experiments.

Now it is possible to substitute the approximation (3.10) with Equation (3.7), which gives a hybrid combination of supervised learning and contrastive learning in the form of

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} [\alpha \log q_{\theta}(y|\mathbf{x}) + (1 - \alpha) \log q_{\theta}(\mathbf{x}|y)] \quad (3.13)$$

$$\approx \min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} \left[\alpha \log \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{i=1}^C \exp(f_{\theta}(\mathbf{x})[y_i])} + (1 - \alpha) \log \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{i=1}^M \exp(f_{\theta}(\mathbf{x}_i)[y])} \right]. \quad (3.14)$$

The parameter α is a weight between $[0, 1]$. It is obvious that in the case of $\alpha = 1$, the objective reduces to the standard cross-entropy loss, while in $\alpha = 0$, the objective is reduced to a case called *end-to-end supervised version of contrastive learning*. The choice of parameter α is a decision of the experiment designer; however, the authors of [12] evaluated many possible variants in the experiments performed and found that the choice of $\alpha = 0.5$ produces the highest classification accuracy performance. Unfortunately, these experiments involved only image classification.

The hybrid combination of supervised learning and contrastive learning (3.14) is absolutely crucial for us as we extend this approach to the multi-instance learning problem, but this is discussed more in further sections.

3.4 Toy problem - Polynomial Regression

At first, we would like to try the hybrid discriminative and generative approach on a simple example before moving on to more difficult cases. The goal is to train a model of the form (3.14) that was derived in the previous section.

Assume data $\{x_i, y_i\}_{i=1}^N$, where $x_i, y_i \in \mathbb{R}$, therefore, this is only a two-dimensional problem. According to the energy-based models 3.1, we know that for the joint distribution, it holds

$$p(x, y) = \frac{\exp(f_{\theta}(x)[y])}{Z(\theta)}, \quad (3.15)$$

where the model is given by

$$f_{\theta}(x)[y] = -E_{\theta}(x, y). \quad (3.16)$$

At this point, we should transform our problem into a polynomial regression. We need to be aware of the discriminative term in Equation (3.7), because we do not want to classify, but we would like to find the best fit to the given data. For this reason, we replace that with the typical regression loss

$$S = S(\theta) = \sum_{k=1}^N \left(y_k - \sum_{i=0}^{s-1} \theta_i x_k^i \right)^2. \quad (3.17)$$

Any joint probability distribution can be broken down into parts by the chain rule.

$$p(x, y) = p(y, x) = p(y|x) \cdot p(x), \quad (3.18)$$

Therefore, we need to find $p(y|x)$ and $p(x)$. From polynomial regression we can obtain the conditional probability density

$$p(y|x, \theta) = \mathcal{N} \left(y; \sum_{i=0}^{s-1} \theta_i x^i, \sigma^2 \right), \quad (3.19)$$

where the symbol $\mathcal{N}(\cdot)$ denotes the probability density function of the Normal distribution and σ^2 is the known variance. In this case, we also need to determine the prior distribution of x . To keep this example simple, let the PDF takes the form

$$p(x|\tau) = \mathcal{N}(x; 0, \tau^2), \quad (3.20)$$

where the choice of parameter τ is based on the fact that we would like to have a non-informative prior, thus τ should be adequately high. If the value of τ is high, the data are spread very far from their expected value. Substituting equations (3.20) and (3.19) into (3.18) results

$$p(x, y) = \mathcal{N}(x; 0, \tau^2) \cdot \mathcal{N} \left(y; \sum_{i=0}^{s-1} \theta_i x^i, \sigma^2 \right) = \frac{1}{2\pi\sigma\tau} \exp \left(-\frac{\left(y - \sum_{i=0}^{s-1} \theta_i x^i \right)^2}{2\sigma^2} - \frac{x^2}{2\tau^2} \right), \quad (3.21)$$

whereas our desirable model is given by

$$f_{\theta}(x)[y] = -\frac{\left(y - \sum_{i=0}^{s-1} \theta_i x^i\right)^2}{2\sigma^2} - \frac{x^2}{2\tau^2}. \quad (3.22)$$

We can now substitute (3.22) and (3.17) in the equation (3.7), resulting in

$$\min_{\theta} \left\{ \alpha S(\theta) - \mathbb{E}_{p_{\text{data}}(x,y)} \left[(1 - \alpha) \log q_{\theta}(x|y) \right] \right\} = \quad (3.23)$$

$$\min_{\theta} \left\{ \alpha S(\theta) - \mathbb{E}_{p_{\text{data}}(x,y)} \left[(1 - \alpha) \log \frac{\exp(f_{\theta}(x)[y])}{\sum_{i=1}^N \exp(f_{\theta}(x_i)[y])} \right] \right\} = \quad (3.24)$$

$$\min_{\theta} \left\{ \alpha S(\theta) - \mathbb{E}_{p_{\text{data}}(x,y)} \left[(1 - \alpha) \log \frac{\exp\left(-\frac{\left(y - \sum_{i=0}^{s-1} \theta_i x^i\right)^2}{2\sigma^2} - \frac{x^2}{2\tau^2}\right)}{\sum_{k=1}^N \exp\left(-\frac{\left(y - \sum_{i=0}^{s-1} \theta_i x_k^i\right)^2}{2\sigma^2} - \frac{x_k^2}{2\tau^2}\right)} \right] \right\}. \quad (3.25)$$

Finally, we simplify the generative term $\log q_{\theta}(x|y)$ into

$$\log q_{\theta}(x|y) = \left(-\frac{\left(y - \sum_{i=0}^{s-1} \theta_i x^i\right)^2}{2\sigma^2} - \frac{x^2}{2\tau^2} \right) - \log \sum_{k=1}^N \exp \left(-\frac{\left(y - \sum_{i=0}^{s-1} \theta_i x_k^i\right)^2}{2\sigma^2} - \frac{x_k^2}{2\tau^2} \right). \quad (3.26)$$

Note that for $\alpha = 1$ we get purely polynomial regression, and for $\alpha = 0$, the term $S(\theta)$ is not involved at all. Now, we have everything we need to carry out the experiment.

3.4.1 Experiment setup and results

We would like to test the sensitivity of this approach to the unknown parameter τ and the order of the polynomial $s - 1$. In addition, we would like to observe how the estimated model behaves in relation to α .

We generate synthetic data, two clusters consisting of 5 data points each, then the model was fitted for different weights $\alpha \in \{0.0, 0.1, 0.2, \dots, 1.0\}$, giving 11 different models in total. The models estimated for $\alpha \in \{0.0, 0.5, 1.0\}$ are highlighted because they are more important to us than the other models. At first, we trained the models mentioned for the fixed order of the polynomial, but for 6 different values of the parameter τ . The results obtained (Figure 3.1) for small τ barely vary from those for high τ , which is exactly what we hoped for, as the previous distribution should be non-informative (3.20). This is a very exciting discovery because there is no need to know much about the data distribution. Second, we trained our polynomial models for the fixed value of τ with 6 different values of the order of the polynomial. The goal of this part of the experiment is to observe how the contrastive part of Equation (3.23) affects the loss of polynomial regression for different values of $s - 1$. As can be seen in Figure 3.2, for small $s - 1$, such as $s - 1 = 2$, the term $\log q_{\theta}(x|y)$ does not affect the polynomial regression

too much. However, we get a considerable difference for higher orders of the polynomial. Furthermore, it seems that the curve $\alpha = 0$ prefers not to oscillate. This could also be very interesting because combining discriminative and generative models could result in better model predictions.

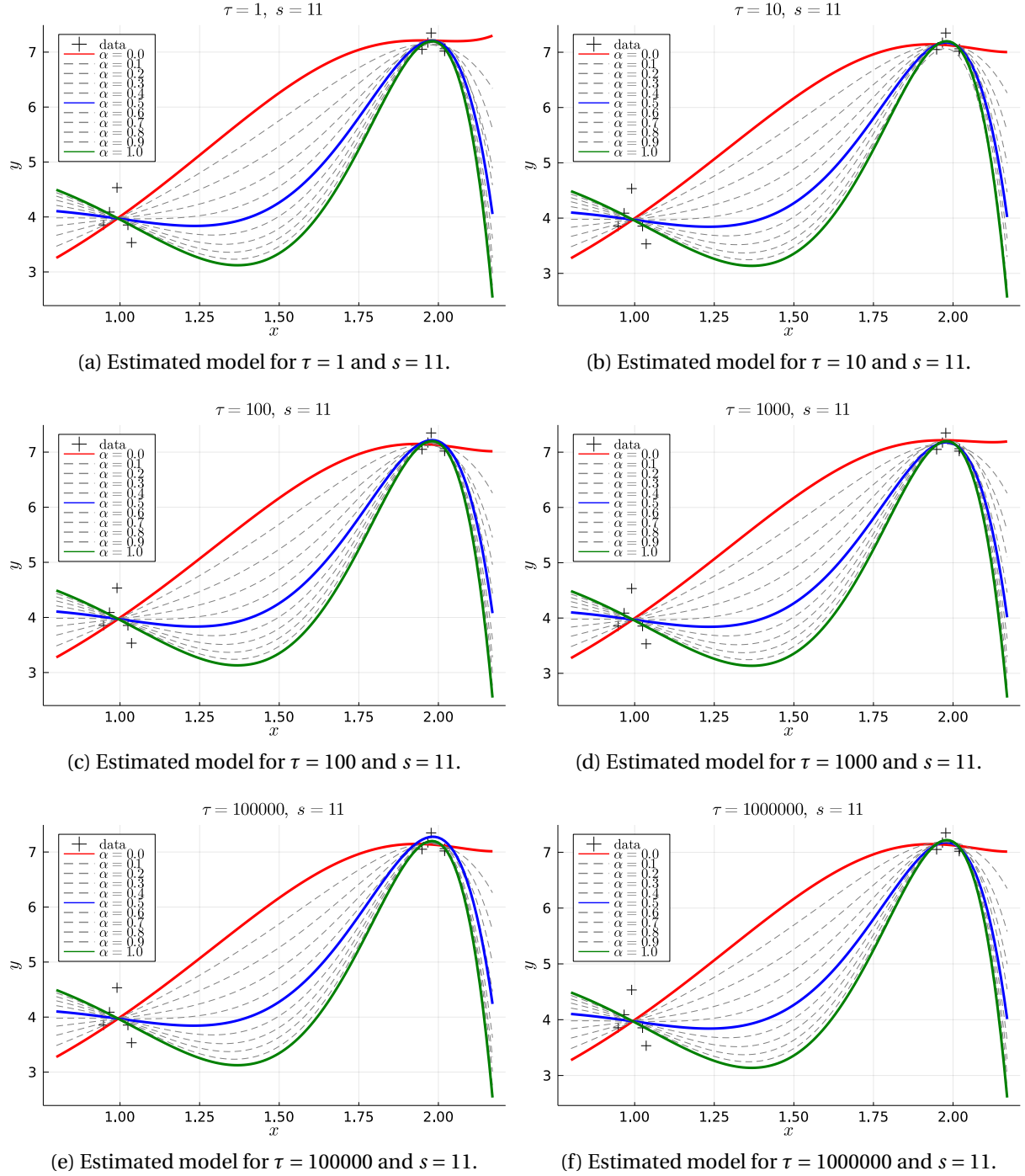


Figure 3.1: Sensitivity of the polynomial model to parameter τ with parameter s held fixed for six different cases.

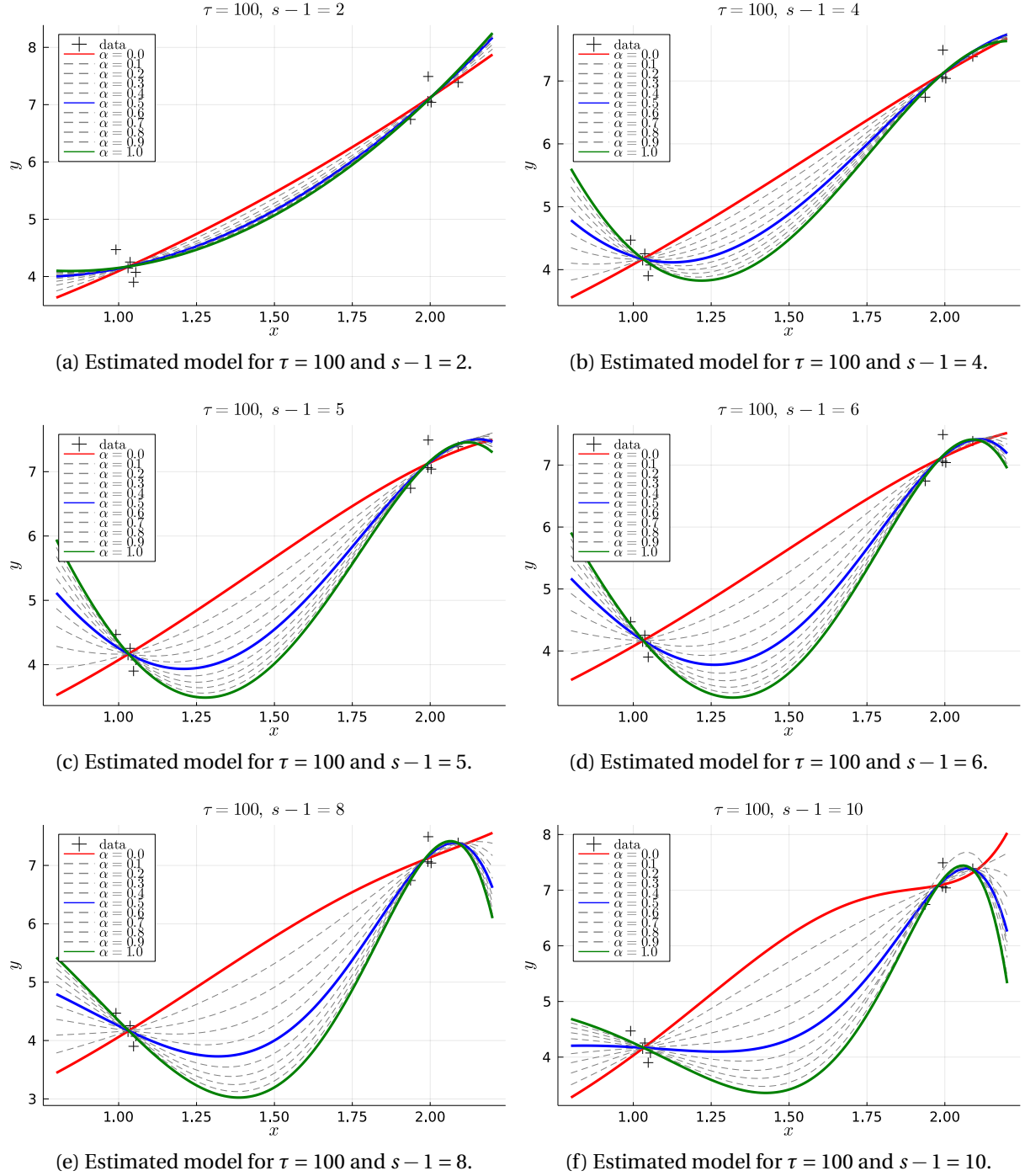


Figure 3.2: Sensitivity of the polynomial model to the order of polynomial $s - 1$ for six different cases.

Chapter 4

Multiple Instance Learning

4.1 Fundamentals

The term multiple instance learning originates from [16] and in [2], the authors proposed the following nomenclature for MIL, which will be reviewed and will be used gladly in our work.

In standard machine learning problems, each sample is represented by a fixed vector \mathbf{x} of observations, however, in multiple instance learning (MIL) it is dealt with samples which are represented by a set of vectors. These vectors are called *instances* and come from an instance space \mathcal{X} , for example, \mathbb{R}^n . The sets of these instances are called *bags* and come from the bag space $\mathcal{B} = \mathcal{P}_F(\mathcal{X})$, where $\mathcal{P}_F(\mathcal{X})$ denotes all finite subsets of \mathcal{X} . With this in mind, we can easily write any bag as $b = \{\mathbf{x} \in \mathcal{X}\}_{\mathbf{x} \in b}$. Each bag b can be arbitrarily large or empty, thus the size of the bag is defined in the form $|b| \in \mathbb{N}_0$. There may exist intrinsic labeling of instances, but we are only interested in labeling at the bag levels. Bag labels come from a finite set \mathcal{C} , and what we want in MIL is to learn a predictor in the form $f_{\theta} : \mathcal{B}(\mathcal{X}) \rightarrow \mathcal{C}$ that can also be rewritten in the form $f_{\theta}(\{\mathbf{x}\}_{\mathbf{x} \in b})$. In contrast to ML, where a predictor is learned in the form $f_{\theta} : \mathbb{R}^n \rightarrow \mathcal{C}$. We consider the supervised setting, in which each sample of the data set is attributed a label. We can denote the available data by the notation $\mathcal{D} = \{(b_i, y_i) \in \mathcal{B} \times \mathcal{C} \mid i \in \{1, 2, \dots, |\mathcal{D}|\}\}$, where $|\mathcal{D}|$ apparently denotes the size of \mathcal{D} .

4.2 Embedded-space paradigm

Embedded space paradigm [2] defines a vector space for bag representation and specifies a mapping from each bag $b \in \mathcal{B}$ to this space. Assume that the target vector space is \mathbb{R}^d , then partial mapping $\phi_i : \mathcal{B} \rightarrow \mathbb{R}$, $\forall i = 1, 2, \dots, d$ is defined and overall embedding $\phi : \mathcal{B} \rightarrow \mathbb{R}^d$ is given by

$$\phi(b) = (\phi_1(b), \phi_2(b), \dots, \phi_d(b)). \quad (4.1)$$

Mappings ϕ_i are instrumental in obtaining and aggregating the appropriate information on the level of instances. They can be defined by some instance transformation $k : \mathcal{X} \rightarrow \mathbb{R}^d$ and

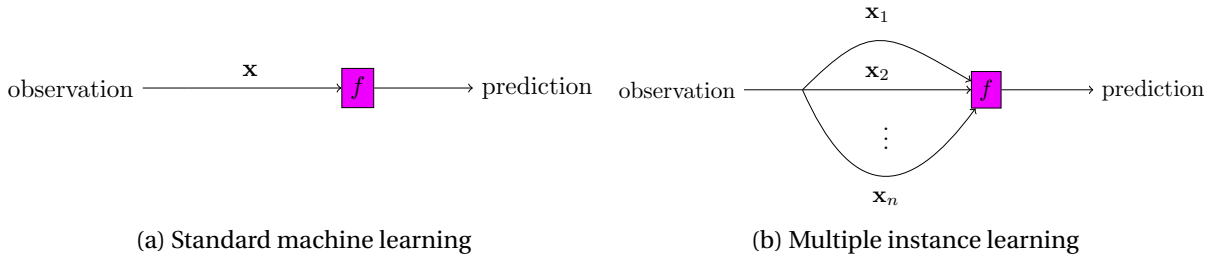


Figure 4.1: The difference between standard ML and MIL [2]. Standard ML is special case of MIL with $|b| = 1$.

an aggregation function $g : \mathcal{P}_F(\mathbb{R}^d) \rightarrow \mathbb{R}^s$ of the form

$$\phi_i(b) = g(k\{\mathbf{x}\}_{\mathbf{x} \in b}). \quad (4.2)$$

On the resulting embedded representation of bag samples, any standard machine learning algorithm can be applied, that is, training a bag-level classifier $f_{\theta}^B : \mathbb{R}^s \rightarrow \mathcal{C}$ using adjusted dataset $D_{ES} = \{(\phi(b_i), y_i) \in \mathbb{R}^s \times \mathcal{C} \mid i \in \{1, 2, \dots, |D|\}\}$. The most used aggregation functions are minimum, maximum, or mean value.

4.3 Training

Authors of [2] proposed a versatile unified framework called HMill (Hierarchical multi-instance learning library) for the definition and training of the models and even implemented this functionality in the *Julia* programming language. Furthermore, the framework was published as an open source project entitled *Mill.jl* under the MIT license. The aim of this work is not to rigorously derive the MIL model since it is fairly complicated and requires a considerable amount of work. Taking into account that we settled for the MIL model being a neural network (NN) that utilizes aggregation functions on the level of instances and refer to [2] for more details about the definition of the model and its composition. Learning of the MIL model f_{θ} is also supervised, a specific case called binary classification, and therefore accomplished by minimizing standard cross-entropy

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(\mathbf{x}, y)} \left[\log \frac{\exp(f_{\theta}(\mathbf{x})[y])}{\sum_{i=1}^C \exp(f_{\theta}(\mathbf{x})[y_i])} \right], \quad (4.3)$$

already mentioned in section 2. This is very important for us and we will take advantage of that in our experiments.

4.4 Cross-validation on MIL datasets

For the MIL testing, we have four datasets available, namely Musk1, Musk2, Tiger, and Fox. All of them will be used to assess the performance of the MIL model.

4.4.1 Setup

In this experiment, data sets are 100 times randomly split into 2 sets in advance, train and test sets, with 80% of observations being in train set and 20% of observations belonging to the test set. For future simplification, let a number of random splits be denoted by r , and thus $r = 100$. We fit the model on train set, then we evaluate the prediction error on the train data via cross-entropy 2.3. The objective here is to plot the dependence of the prediction error on the complexity of the model. A smaller number of random splits r were tested, but the results obtained were too noisy. For this reason, such a high r was selected, although the experiment became noticeably expensive to compute.

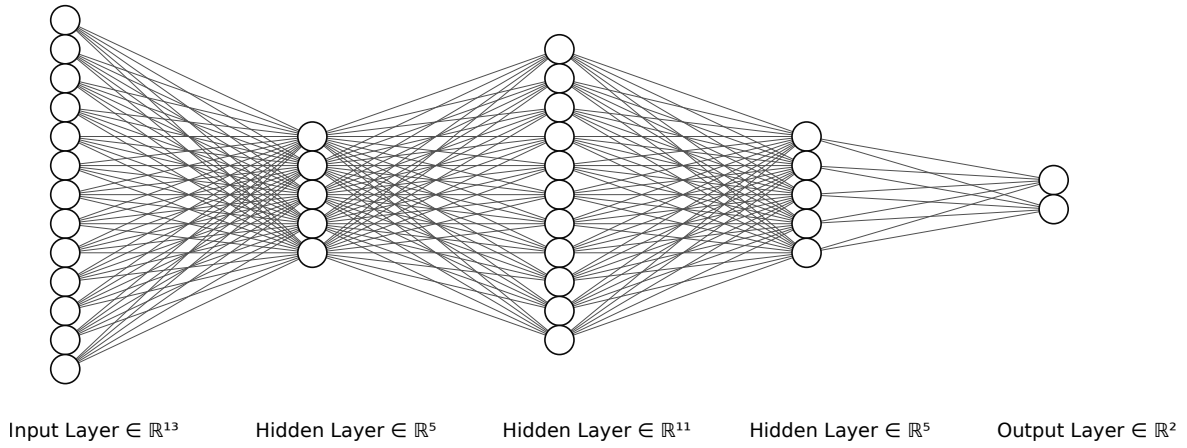


Figure 4.2: NN example for $z = 5$, where input and output layer are only illustrative.

Model Complexity As was mentioned in section 4.3, defining such model for the MIL problem is very complex task, therefore choosing a right model complexity metric is not trivial. Consider a neural network consisting of input layer, 3 hidden layers $h_1 \in \mathbb{R}^z$, $h_2 \in \mathbb{R}^{2z+1}$, $h_3 \in \mathbb{R}^z$ and output layer. Then $z \in \{1, 2, 3, \dots, 20\}$ was selected as model complexity metrics, because this is one of the easiest ways, how to control complexity of the defined MIL model. To gain a better insight, example of such NN is illustrated in Figure 4.2, however this is not the exact NN used in HMill.

Prediction Error For prediction error metrics, we simply used the standard cross-entropy loss as outlined at the beginning of this section. There is no need for any trickier objective. Let the total cross-entropy loss evaluated in the k^{th} random split for fixed z be denoted by $\mathcal{L}_k(z)$, then the estimated prediction error is given by

$$\widehat{\text{Err}}(z) = \frac{1}{r} \sum_{k=1}^r \mathcal{L}_k(z). \quad (4.4)$$

To summarize, we fit 100 models for selected z on train data and evaluate the prediction error for all of them on test data, then we take the mean value. This process is repeated for each $z \in \{1, 2, 3, \dots, 20\}$, giving 2000 models in total.

4.4.2 Results

As can be seen in Figure 4.3, the results obtained are totally expected. The prediction error evaluated on the training data for the higher complexity of the model approaches zero. However, testing data give oscillating curves with an increasing trend (with a little exception of Musk1), therefore, model selection is necessary. This applies to each data set. Furthermore, Table 4.1 numerically summarizes the results evaluated in the test data.

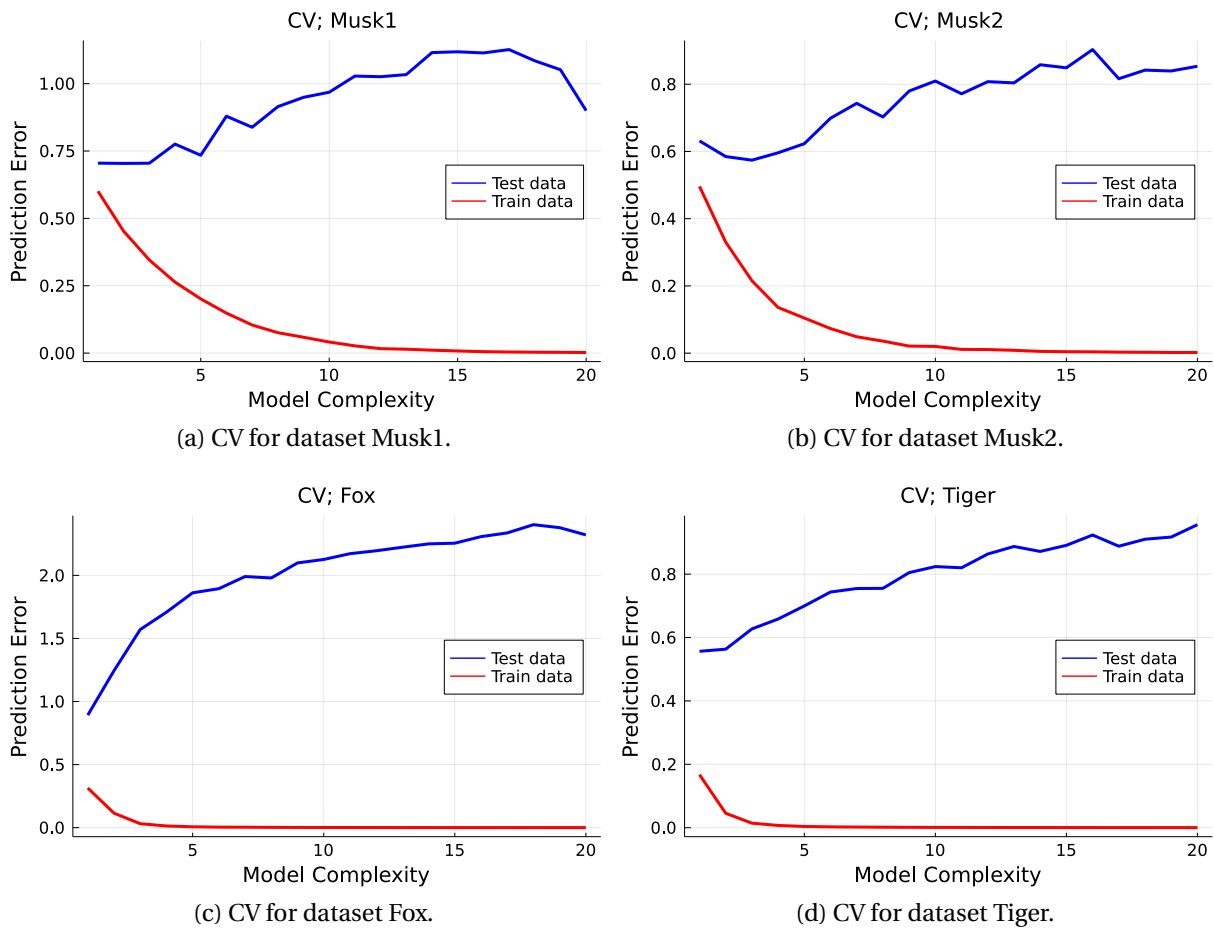


Figure 4.3: Evaluation of prediction error with the use of training data and testing data on MIL datasets Musk1, Musk2, Fox and Tiger.

Dataset	$\text{argmin} \widehat{\text{Err}}(z)$	$\min \widehat{\text{Err}}(z)$	$\widehat{\text{Err}}(z = 10)$
Musk1	2	0.70	0.97
Musk2	3	0.57	0.81
Fox	1	0.89	2.13
Tiger	1	0.56	0.82

Table 4.1: Results of CV evaluated on the testing data.

4.5 MIL to HDGM problem

In the previous part, the CV experiment was performed with expected results. However, the estimated prediction error seems to be rather high. Logically, the question of whether the prediction error can be reduced has been raised and also whether it is possible to make r smaller has been raised.

4.5.1 Setup

On the initiative to reduce prediction error, it is proposed to train a MIL model that is obtained by minimizing the hybrid loss function

$$\min_{\theta} -\mathbb{E}_{p_{\text{data}}(x,y)} \left[\alpha \log \frac{\exp(f_{\theta}(x)[y])}{\sum_{i=1}^C \exp(f_{\theta}(x)[y_i])} + (1 - \alpha) \log \frac{\exp(f_{\theta}(x_i)[y])}{\sum_{j=1}^N \exp(f_{\theta}(x_j)[y])} \right]. \quad (4.5)$$

Since the discriminative part is already used in HMill framework, the only task is to add the generative part into it. At this point are available all normalization samples, thus $M = N$. This modification should lead to a reduced prediction error evaluated on training data.

In the first part of this experiment, we would like to train models in relation to parameter α . For this setup, we need to choose fixed z . Since authors of [2] usually use $z = 10$, we use this value as well. For the prediction error evaluation is used standard cross-entropy as in the previous experiment, again with $r = 100$ and $\alpha \in \{0.0, 0.1, 0.2, \dots, 1.0\}$. Therefore estimated prediction error can be written in the form

$$\widehat{\text{Err}}(\alpha) = \frac{1}{r} \sum_{k=1}^r \mathcal{L}_k(z = 10, \alpha). \quad (4.6)$$

We hope to see a curve in the shape of a bowl, having its global minimum at a point $\alpha = 0.5$ or somewhere near.

In the second part, evaluating the prediction error is approached in the other way. The fixed $\alpha = 0.5$ is selected and the dependency on $z \in \{1, 2, 3, \dots, 20\}$ is evaluated as in Section 4.4. Finally, the estimated prediction error in this case is defined by

$$\widehat{\text{Err}}(z) = \frac{1}{r} \sum_{k=1}^r \mathcal{L}_k(z, \alpha = 0.5). \quad (4.7)$$

In other words, the setup is the same as in 4.4, therefore, the results obtained will be added to the figure 4.3 and the table 4.1 for a convenient comparison. Note that curves for train data from this experiment will be omitted, since they are not important at this point. We are only interested in predictions for test data.

4.5.2 Results

Results of the first part are represented in Figure 4.4 and in Table 4.2, where can be seen that adding the generative term into the MIL loss function decreased the prediction error evaluated on all datasets. This improvement is considerable on datasets Musk1 and Musk2, where a can be seen a nice bowl. On datasets Fox and Tiger such improvement does not occur. This means that HDGM approach works, thus a regularization in the form of very simple generative term may bring improvement in predictions. Unfortunately, the choice of $\alpha = 0.5$ was not confirmed as the best in our experiment, see Table 4.2.

In the second part of the experiment results are shown in Figure 4.5 and in Table 4.3. Here is the situation very similar to the previous part of this experiment. Improvement of the prediction error is quite noticeable on first two datasets Musk1 and Musk2, while Fox and Tiger does not look so convincingly. In addition, number of random splits $r = 100$ is still needed to get rid of the noise on the prediction error.

Overall, it can be said that HDGM approach leads to the decreased prediction error in a small way.

Dataset	$\operatorname{argmin} \widehat{\text{Err}}(\alpha)$	$\min \widehat{\text{Err}}(\alpha)$
Musk1	0.4	0.68
Musk2	0.2	0.54
Fox	0.7	1.89
Tiger	0.4	0.74

Table 4.2: Prediction error statistics for HDGM in case of $z = 10$.

Dataset	Discriminative part only			HDGM; $\alpha = 0.5$		
	$\operatorname{argmin} \widehat{\text{Err}}(z)$	$\min \widehat{\text{Err}}(z)$	$\widehat{\text{Err}}(z = 10)$	$\operatorname{argmin} \widehat{\text{Err}}(z)$	$\min \widehat{\text{Err}}(z)$	$\widehat{\text{Err}}(z = 10)$
Musk1	2	0.70	0.97	6	0.64	0.69
Musk2	3	0.57	0.81	6	0.55	0.66
Fox	1	0.89	2.13	1	0.95	1.96
Tiger	1	0.56	0.82	2	0.58	0.80

Table 4.3: Comparison of prediction error statistics for HDGM $\alpha = 0.5$ and discriminative part only. Pay attention especially to the last column in each approach, $\widehat{\text{Err}}(z = 10)$.

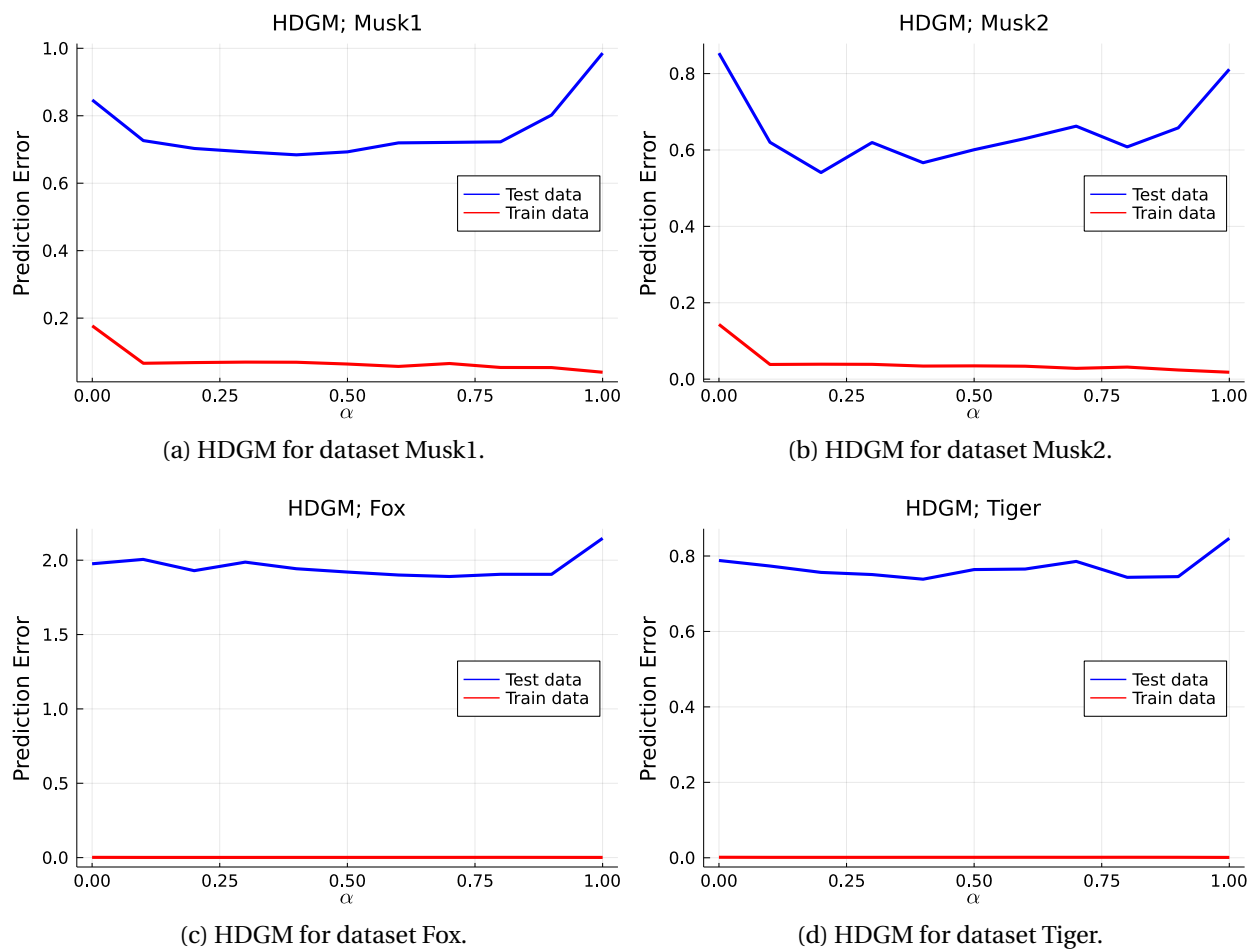


Figure 4.4: Evaluation of the prediction error $\widehat{\text{Err}}(\alpha)$ with the use of training data and testing data on MIL datasets.

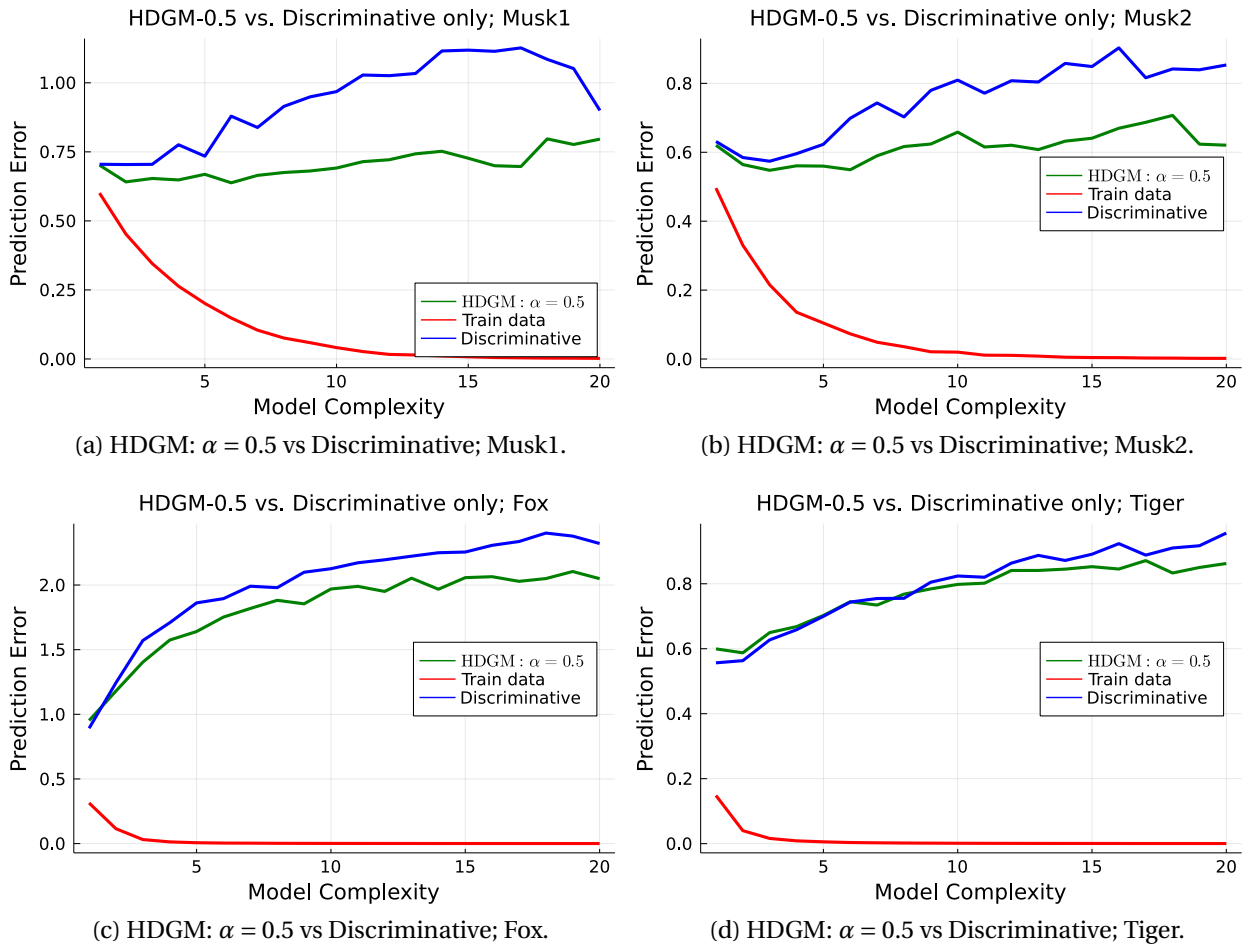


Figure 4.5: Comparison of the prediction error $\widehat{\text{Err}}(z)$ for HDGM $\alpha = 0.5$ and the discriminative part only.

Conclusion

At the beginning of this work, supervised learning and energy-based models were introduced. The following passage consists of contrastive learning, which was briefly reviewed, and thereafter supervised learning and contrastive learning was merged into hybrid discriminative and generative models. Subsequently, this approach was applied and tested on a simple example consisting of polynomial regression. After this example, we briefly introduced multiple instance learning with a short description of embedded space and its way of training. As a first MIL experiment, cross-validation was performed on four MIL datasets, where enough space was given to a proper definition of the model complexity metrics. The results obtained were totally expected; they included decreasing prediction errors for train data and increasing for test data. In the next step, a solution was searched on how to decrease the prediction error evaluated on the test data. In this initiative, HDGM was trained instead of a standard discriminative model. In the consecutive result, it was found that HDGM leads to a decrease in prediction error, however, nothing significant. In addition, a high number of random splits of the data set was still necessary to eliminate the noise from the prediction error. In conclusion, this approach was proven to be functional, although a proposed generative regularization is very simple and can be replaced by much complicated models. From this point of view, this approach has great potential that has not yet been fully exploited.

Bibliography

- [1] Christopher M. Bishop: *Pattern recognition and machine learning*. [New York]: Springer, 2006. Information science and statistics. ISBN 0-387-31073-8.
- [2] S. Mandlik.: *Mapping the Internet: Modelling Entity Interactions in Complex Heterogeneous Networks*. arXiv preprint arXiv:2104.09650.
- [3] T. Pevny and P. Somol: *Discriminative models for multi-instance problems with tree structure*. In Proceedings of the 2016 ACM Workshop on Artificial Intelligence and Security, 2016, 83-91.
- [4] J. Wu, S. Pan, X. Zhu, C. Zhang, X. Wu: *Multi-instance learning with discriminative bag mapping*. IEEE Transactions on Knowledge and Data Engineering, 30(6), 2018, 1065-1080.
- [5] H. Jeffrey: *An invariant form for the prior probability in estimation problems*. Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences. 1946, 1-9.
- [6] M. I. Jordan and A. Y. Ng. : *On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes*. Advances in neural information processing systems. 2002.
- [7] D. Commenges: *Information Theory and Statistics: an overview*. ArXiv preprint arXiv:1511.00860, 2015, 1-22.
- [8] J. Brownlee: *How to Handle Big-p, Little-n ($p \gg n$) in Machine Learning*. (2020). [on-line]. Available from: <https://machinelearningmastery.com/how-to-handle-big-p-little-n-p-n-in-machine-learning/>.
- [9] V. Smidl: *The Variational Bayes Approach in Signal processing*. PhD Thesis. Trinity College Dublin. 2004.
- [10] M. Zhuang and M. Collins: *Ma, Zhuang, and Michael Collins. "Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency*. arXiv preprint arXiv:1809.01812, 2018.
- [11] M. Gutmann and A. Hyvärinen: *Noise-contrastive estimation: A new estimation principle for unnormalized statistical models*. Proceedings of the thirteenth international conference on artificial intelligence and statistics. JMLR Workshop and Conference Proceedings, 2010.

- [12] H. Liu and P. Abbeel: *Hybrid discriminative-generative training via contrastive learning*. arXiv preprint arXiv:2007.09070, 2020.
- [13] S.K. Ng, T. Krishnan and G.J.McLachlan: *Handbook of computational statistics*. The EM algorithm. Springer, Berlin, Heidelberg. 2012, 139-172.
- [14] W. Gratwohl, K.C. Wang and J.H. Jacobsen: *Your classifier is secretly an energy based model and you should treat it like one*. GRATHWOHL, Will, et al. Your classifier is secretly an energy based model and you should treat it like one. arXiv preprint arXiv:1912.03263, 2019.
- [15] Y. LeCun, S. Chopra, R.Hadsell, M. Ranzato and F. Huang: *A tutorial on energy-based learning*. Predicting structured data, 2006, 1.0.
- [16] T.G. Dietterich, R.H. Lathrop and T. Lozano-Pérez: *Solving the multiple instance problem with axis-parallel rectangles*. Artificial intelligence, 1997, 89.1-2: 31-71.
- [17] J. Friedman, T. Hastie and R. Tibshirani: *The elements of statistical learning*. [New York]: Springer, 2001. Series in statistics.
- [18] M. Talabis, R. McPherson, I. Miyamoto, J. Martin and D.Kaye: *Information Security Analytics*. Syngress, 2015.
- [19] J.Bezanson, S. Karpinski, V.B. Skah: *A fast dynamic language for technical computing*. arXiv preprint arXiv:1209.5145, 2012.
- [20] D. Yu and L. Deng: *Automatic Speech Recognition*. Springer london limited, 2016.

Appendix A

Computational formulas

A.1 Solution of $D_{\text{KL}}(\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \sigma^2 \mathbb{I}_P) \parallel \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbb{I}_P))$

In the VAE section, the ELBO (2.22) is derived and subsequently optimized. One of the ELBO expressions is the KL distance mentioned above. For two multivariate Gaussian distributions we have a KL distance analytical solution. The complete calculation is given here. First, recall that PDF for a multivariate Gaussian distribution in \mathbb{R}^P with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ is defined as

$$\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^P \det \boldsymbol{\Sigma}}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu})\right). \quad (\text{A.1})$$

$$\text{KL} = D_{\text{KL}}(\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \parallel \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} [\log \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) - \log \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)] \quad (\text{A.2})$$

$$= \frac{1}{2} \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} \left[-\log \det \boldsymbol{\Sigma}_1 - (\mathbf{z} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1}(\mathbf{z} - \boldsymbol{\mu}_1) + \log \det \boldsymbol{\Sigma}_2 + (\mathbf{z} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\mathbf{z} - \boldsymbol{\mu}_2) \right] \quad (\text{A.3})$$

$$= \frac{1}{2} \log \frac{\det \boldsymbol{\Sigma}_2}{\det \boldsymbol{\Sigma}_1} + \frac{1}{2} \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} \left[-(\mathbf{z} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1}(\mathbf{z} - \boldsymbol{\mu}_1) + (\mathbf{z} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\mathbf{z} - \boldsymbol{\mu}_2) \right] \quad (\text{A.4})$$

$$= \frac{1}{2} \log \frac{\det \boldsymbol{\Sigma}_2}{\det \boldsymbol{\Sigma}_1} + \frac{1}{2} \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} \left[-\text{Tr}(\boldsymbol{\Sigma}_1^{-1}(\mathbf{z} - \boldsymbol{\mu}_1)(\mathbf{z} - \boldsymbol{\mu}_1)^\top) + \text{Tr}(\boldsymbol{\Sigma}_2^{-1}(\mathbf{z} - \boldsymbol{\mu}_2)(\mathbf{z} - \boldsymbol{\mu}_2)^\top) \right] \quad (\text{A.5})$$

$$= \frac{1}{2} \log \frac{\det \boldsymbol{\Sigma}_2}{\det \boldsymbol{\Sigma}_1} + \frac{1}{2} \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} \left[-\text{Tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_1) + \text{Tr}(\boldsymbol{\Sigma}_2^{-1}(\mathbf{z} \mathbf{z}^\top - 2\mathbf{z} \boldsymbol{\mu}_2^\top + \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^\top)) \right] \quad (\text{A.6})$$

$$= \frac{1}{2} \log \frac{\det \boldsymbol{\Sigma}_2}{\det \boldsymbol{\Sigma}_1} - \frac{1}{2} n + \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_2^{-1}(\boldsymbol{\Sigma}_1 + \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^\top - 2\boldsymbol{\mu}_2 \boldsymbol{\mu}_1^\top + \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^\top)) \quad (\text{A.7})$$

$$= \frac{1}{2} \left(\log \frac{\det \boldsymbol{\Sigma}_2}{\det \boldsymbol{\Sigma}_1} - n + \text{Tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + \text{Tr}(\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_1 - 2\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 + \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) \right) \quad (\text{A.8})$$

$$= \frac{1}{2} \left(\log \frac{\det \boldsymbol{\Sigma}_2}{\det \boldsymbol{\Sigma}_1} - n + \text{Tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_2^{-1}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \right) \quad (\text{A.9})$$