# Supplementary Material for
# Bayesian Triplet Loss: Uncertainty Quantification in Image Retrieval

Frederik Warburg[†], Martin Jørgensen[‡], Javier Civera[§], and Søren Hauberg[†]

[†]Technical University of Denmark, [‡]University of Oxford, [§]University of Zaragoza

[†]{frwa,sohau}@dtu.dk, [‡]martinj@robots.ox.ac.uk, [§]jcivera@unizar.es

## 1. Introduction

In this supplementary material we give additional details into the derivations of the mean and the variance for the Gaussian approximation of the proposed likelihood. We show through simulations that the Gaussian approximation of the likelihood is very accurate, even for a small number of dimensions. We provide the full derivation of the Expected Lower Bound (ELBO). Finally, we provide sample code for the proposed Bayesian Triplet loss, which highlights that the proposed method is simple to implement. We will make the entire code available upon acceptance.

## 2. Gaussian Likelihood

First, we provide the details into the derivations of the first two moments of the proposed Gaussian approximation of the likelihood.

### 2.1. The Mean

The expression

$$\mathbb{E}[\tau] = \mathbb{E}[p(p-2a)] - \mathbb{E}[n(n-2a)] \tag{1}$$

consists of two means of the same form. We first evaluate

$$\mathbb{E}[p(p-2a)] = \mathbb{E}[p^2] - 2\mathbb{E}[ap] = \mathbb{E}[p^2] - 2\mathbb{E}[a]\mathbb{E}[p], \tag{2}$$

where we have used that the expectation of the product of two independent variables is the product of the expectations. Using identical arguments for the second term of $\mathbb{E}[\tau]$ we get

$$\mathbb{E}[\tau] = \mathbb{E}[p^2] - 2\mathbb{E}[a]\mathbb{E}[p] - \mathbb{E}[n^2] + 2\mathbb{E}[a]\mathbb{E}[n] \tag{3}$$

$$= \mathbb{E}[p^2] - \mathbb{E}[n^2] - 2\mathbb{E}[a](\mathbb{E}[p] - \mathbb{E}[n]) \tag{4}$$

$$= \mu_p^2 + \sigma_p^2 - \mu_n^2 - \sigma_n^2 - 2\mu_a(\mu_p - \mu_n). \tag{5}$$

### 2.2. The variance

The expression

$$\text{var}[\tau] = \text{var}[p(p-2a) - n(n-2a)] \tag{6}$$

$$= \text{var}[p(p-2a)] + \text{var}[n(n-2a)] - 2\text{cov}[p(p-2a), n(n-2a)] \tag{7}$$

consists of three terms. We treat the terms one by one

$$\text{var}[p(p-2a)] = \text{var}[p^2 - 2ap] \tag{8}$$

$$= \text{var}[p^2] + 4\text{var}[ap] - 4\text{cov}[p^2, ap] \tag{9}$$

$$\begin{aligned} &= 2\sigma_p^4 + 4\mu_p^2\sigma_p^2 + 4\left(\text{var}[a] + \mathbb{E}[a]^2\right)\left(\text{var}[p] + \mathbb{E}[p]^2\right) \\ &\quad - 4\mathbb{E}[a]^2\mathbb{E}[p]^2 - 4\text{cov}[p^2, ap] \end{aligned} \tag{10}$$

$$\begin{aligned} &= 2\sigma_p^4 + 4\mu_p^2\sigma_p^2 + 4\left(\sigma_a^2 + \mu_a^2\right)\left(\sigma_p^2 + \mu_p^2\right) - 4\mu_a^2\mu_p^2 \\ &\quad - 4\text{cov}[p^2, ap] \end{aligned} \tag{11}$$

The last covariance term is

$$\text{cov}[p^2,ap]=\mathbb{E}[ap^3]-\mathbb{E}[ap]E[p^2] \tag{12}$$
$$=\mathbb{E}[a]\mathbb{E}[p^3]-\mathbb{E}[a]\mathbb{E}[p]E[p^2] \tag{13}$$
$$=\mu_a(\mu_p^3+3\mu_p\sigma_p^2)-\mu_a\mu_p(\mu_p^2+\sigma_p^2) \tag{14}$$
$$=\mu_a\left[\mu_p^3+3\mu_p\sigma_p^2-\mu_p^3-\mu_p\sigma_p^2)\right] \tag{15}$$
$$=2\mu_a\mu_p\sigma_p^2, \tag{16}$$

such that

$$\text{var}[p(p-2a)]=2\left[\sigma_p^4+2\mu_p^2\sigma_p^2+2\left(\sigma_a^2+\mu_a^2\right)\left(\sigma_p^2+\mu_p^2\right)-2\mu_a^2\mu_p^2-4\mu_a\mu_p\sigma_p^2\right]. \tag{17}$$

The second term of the variance follows the same structure,

$$\text{var}[n(n-2a)]=2\left[\sigma_n^4+2\mu_n^2\sigma_n^2+2\left(\sigma_a^2+\mu_a^2\right)\left(\sigma_n^2+\mu_n^2\right)-2\mu_a^2\mu_n^2-4\mu_a\mu_n\sigma_n^2\right]. \tag{18}$$

The last term of the variance is

$$\text{cov}[p(p-2a),n(n-2a)]=\text{cov}[p^2-2ap,n^2-2an] \tag{19}$$
$$=\text{cov}[p^2,n^2]-2\text{cov}[p^2,an]-2\text{cov}[ap,n^2]+4\text{cov}[ap,an] \tag{20}$$
$$=4\text{cov}[ap,an] \tag{21}$$
$$=4\mathbb{E}[a^2pn]-4\mathbb{E}[ap]\mathbb{E}[an] \tag{22}$$
$$=4\mathbb{E}[a^2]\mathbb{E}[p]\mathbb{E}[n]-4\mathbb{E}[a]^2\mathbb{E}[p]\mathbb{E}[n] \tag{23}$$
$$=4\mathbb{E}[p]\mathbb{E}[n]\left(\mathbb{E}[a^2]-\mathbb{E}[a]^2\right) \tag{24}$$
$$=4\mathbb{E}[p]\mathbb{E}[n]\text{var}[a] \tag{25}$$
$$=4\mu_p\mu_n\sigma_a^2. \tag{26}$$

Thus, given a mean and variance estimate for each embedding, we can analytically find the mean and variance of the Gaussian that approximate our proposed likelihood. In the next section, we show experimentally that this is good approximation.

### 2.3. Simulation Approximation of Gaussian Likelihood

In this section, we show that the Gaussian approximation of our proposed likelihood is accurate even for low dimensions. We sample $\mu_a$, $\mu_p$ and $\mu_n$ from a $D$-dimensional standard Gaussian, and $\sigma_a^2$, $\sigma_p^2$ and $\sigma_n^2$ from a 1-dimensional standard Gaussian, where we ensure that $\sigma_a^2$, $\sigma_p^2$ and $\sigma_n^2$ are positive by taking the absolute value of the sample. We then generate $N$ samples from $a \sim N(\mu_a,\sigma_a^2)$, $p \sim N(\mu_p,\sigma_p^2)$ and $n \sim N(\mu_n,\sigma_n^2)$, and calculate $\tau = \|a-p\|^2 - \|a-n\|^2$. We plot a histogram over the samples of $\tau$ as a Monte Carlo approximation of the true CDF (histograms in the figures below). We use the means and variances of $a$, $p$ and $n$ to analytically calculate the CDF of our Gaussian approximation (black line in the figures below). The figures show that the approximation is very accurate even for a very low number of dimensions (the deviation between the analytical CDF and the Monte Carlo simulation is only significant for $D<4$). This means that the approximation is accurate for any practical retrieval system.
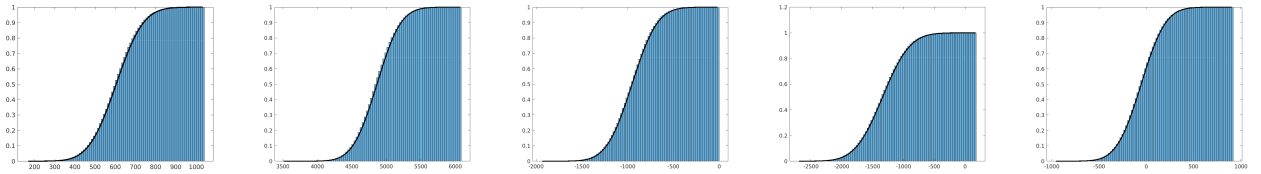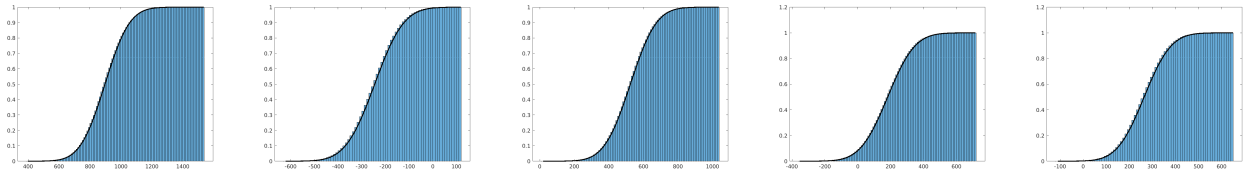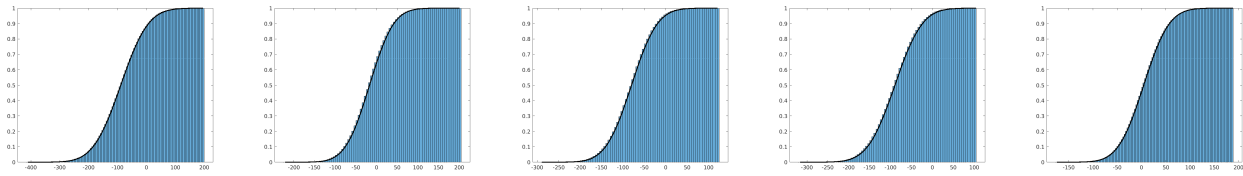


Figure 1: $D=2048$

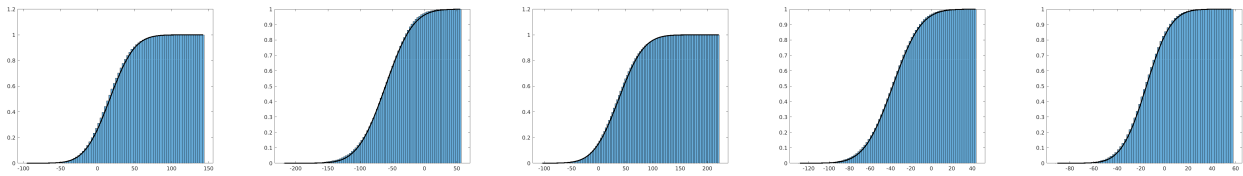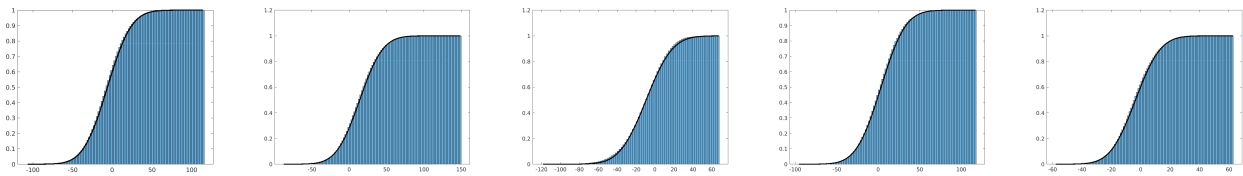Figure 2: $D = 512$



Figure 3: $D = 128$



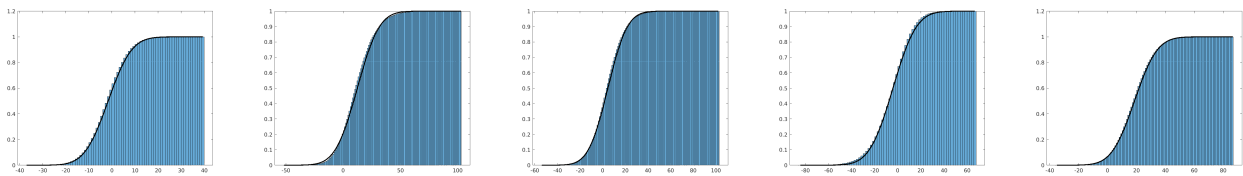Figure 4: $D = 32$
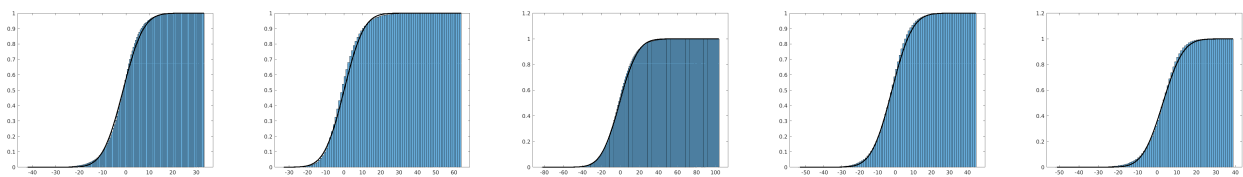


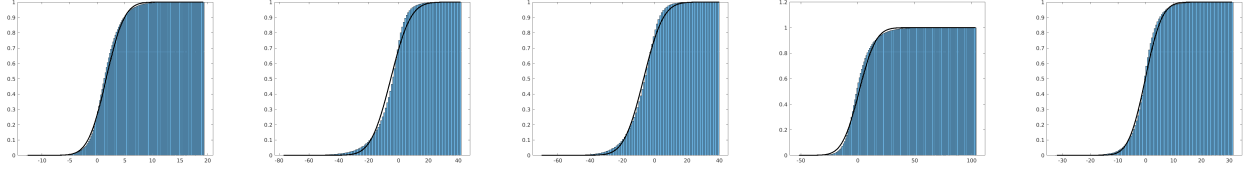Figure 5: $D = 16$



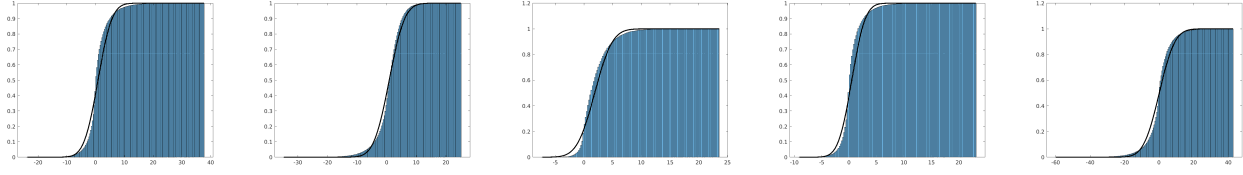Figure 6: $D = 8$



Figure 7: $D = 4$

Figure 8: $D=2$



Figure 9: $D=1$

## 3. Variational Approximation

In this section we derive the Expected Lower Bound. We start with the marginal probability of a triplet label $I(X,Y,Z)$.

$$logP(I(X,Y,Z))=log\int P(I(X,Y,Z),(X,Y,Z))d(X,Y,Z) \tag{27}$$

$$=log\int\int\int P(I(X,Y,Z),X,Y,Z)dXdYdZ \tag{28}$$

$$=log\int\int\int P(I(X,Y,Z)|X,Y,Z)P(X,Y,Z)dXdYdZ \tag{29}$$

$$=log\int\int\int P(I(X,Y,Z)|X,Y,Z)P(X)P(Y)P(Z)dXdYdZ \tag{30}$$

$$=log\int\int\int \frac{P(I(X,Y,Z)|X,Y,Z)P(X)P(Y)P(Z)}{q(X)q(Y)q(Z)}q(X)q(Y)q(Z)dXdYdZ \tag{31}$$

$$\geq\int\int\int log\frac{P(I(X,Y,Z)|X,Y,Z)P(X)P(Y)P(Z)}{q(X)q(Y)q(Z)}q(X)q(Y)q(Z)dXdYdZ \tag{32}$$

$$=\int\int\int logP(I(X,Y,Z)|X,Y,Z)q(X)q(Y)q(Z) \tag{33}$$

$$+log\frac{P(X)P(Y)P(Z)}{q(X)q(Y)q(Z)}q(X)q(Y)q(Z)dXdYdZ \tag{34}$$

$$=\int\int\int logP(I(X,Y,Z)|X,Y,Z)q(X)q(Y)q(Z)dXdYdZ \tag{35}$$

$$+\int\int\int log\frac{P(X)P(Y)P(Z)}{q(X)q(Y)q(Z)}q(X)q(Y)q(Z)dXdYdZ \tag{36}$$

$$\tag{37}$$

We have two terms. The first part becomes

$$\int\int\int logP(I(X,Y,Z)|X,Y,Z)q(X)q(Y)q(Z)dXdYdZ \tag{38}$$

$$=\mathbb{E}_{q(X)q(Y)q(Z)}[logP(I(X,Y,Z)|X,Y,Z)] \tag{39}$$

The second part becomes

4

$$\int\int\int log\frac{P(X)P(Y)P(Z)}{q(X)q(Y)q(Z)}q(X)q(Y)q(Z)dXdYdZ \tag{40}$$

$$=\int\int\int[logP(X)-logq(X)+logP(Y)-logq(Y)+logP(Z)-logq(Z)]q(X)q(Y)q(Z)dXdYdZ \tag{41}$$

$$=\int\int\int[logP(X)-logq(X)]q(X)q(Y)q(Z)dXdYdZ \tag{42}$$

$$+\int\int\int[logP(Y)-logq(Y)]q(X)q(Y)q(Z)dXdYdZ \tag{43}$$

$$+\int\int\int[logP(Z)-logq(Z)]q(X)q(Y)q(Z)dXdYdZ \tag{44}$$

$$=\int[logP(X)-logq(X)]q(X)dX+\int[logP(Y)-logq(Y)]q(Y)dY+\int[logP(Z)-logq(Z)]q(Z)dZ \tag{45}$$

$$=-KL(q(X)||P(X))-KL(q(Y)||P(Y))-KL(q(Z)||P(Z)) \tag{46}$$

This means that if we choose $P(X)=P(Y)=P(Z)=N(0,1)$, we can remove the normalization layer and optimize the ELBO directly.

## 4. Implementation

We present our Pytorch implementation of the Bayesian triplet loss with both the von Mises Fischer and Gaussian embeddings. We will make the source code available upon acceptance.

```
1  import torch
2  import torch.nn as nn
3  import functional as LF
4
5  class BayesianTripletLoss(nn.Module):
6
7      def __init__(self, margin, varPrior, kl_scale_factor = 1e-6, distribution = 'gauss'):
8          super(BayesianTripletLoss, self).__init__()
9
10         self.margin = margin
11         self.varPrior = varPrior
12         self.kl_scale_factor = kl_scale_factor
13         self.distribution = distribution
14
15     def forward(self, x, label):
16
17         # divide x into anchor, positive, negative based on labels
18         D, N = x.shape
19         nq = torch.sum(label.data == -1).item()  # number of tuples
20         S = x.size(1) // nq  # number of images per tuple including query: 1+1+n
21         A = x[:, label.data == -1].permute(1, 0).repeat(1, S - 2).view((S - 2) * nq, D).permute(1, 0)
22         P = x[:, label.data == 1].permute(1, 0).repeat(1, S - 2).view((S - 2) * nq, D).permute(1, 0)
23         N = x[:, label.data == 0]
24
25         varA = A[-1:, :]
26         varP = P[-1:, :]
27         varN = N[-1:, :]
28
29         muA = A[:-1, :]
30         muP = P[:-1, :]
31         muN = N[:-1, :]
32
33         # calculate nll
34         nll = LF.negative_loglikelihood(muA, muP, muN, varA, varP, varN, margin=self.margin)
35
36         # KL(anchor|| prior) + KL(positive|| prior) + KL(negative|| prior)
37         if self.distribution == 'gauss':
38
```

5

```
39            muPrior = torch.zeros_like(muA, requires_grad = False)
40            varPrior = torch.ones_like(varA, requires_grad = False) * self.varPrior
41
42            kl = (LF.kl_div_gauss(muA, varA, muPrior, varPrior) + \
43                    LF.kl_div_gauss(muP, varP, muPrior, varPrior) + \
44                    LF.kl_div_gauss(muN, varN, muPrior, varPrior))
45
46        elif self.distribution == 'vMF':
47            kl = (LF.kl_div_vMF(muA, varA) + \
48                    LF.kl_div_vMF(muP, varP) + \
49                    LF.kl_div_vMF(muN, varN))
50
51        return nll + self.kl_scale_factor * kl
52
53    def __repr__(self):
54        return self.__class__._Name__ + '(' + 'margin=' + '{:.4f}'.format(self.margin) + ')'
```

Listing 1: loss.py

```
1  import torch
2  from torch.distributions.normal import Normal
3
4  def negative_loglikelihood(muA, muP, muN, varA, varP, varN, margin = 0.0):
5
6      muA2 = muA**2
7      muP2 = muP**2
8      muN2 = muN**2
9      varP2 = varP**2
10     varN2 = varN**2
11
12     mu = torch.sum(muP2 + varP - muN2 - varN - 2*muA*(muP - muN), dim=0)
13     T1 = varP2 + 2*muP2 * varP + 2*(varA + muA2)*(varP + muP2) - 2*muA2 * muP2 - 4*muA*muP*varP
14     T2 = varN2 + 2*muN2 * varN + 2*(varA + muA2)*(varN + muN2) - 2*muA2 * muN2 - 4*muA*muN*varN
15     T3 = 4*muP*muN*varA
16     sigma2 = torch.sum(2*T1 + 2*T2 - 2*T3, dim=0)
17     sigma = sigma2**0.5
18
19     probs = Normal(loc = mu, scale = sigma + 1e-8).cdf(margin)
20     nll = -torch.log(probs + 1e-8)
21
22     return nll.mean()
23
24 def kl_div_gauss(mu_q, var_q, mu_p, var_p):
25
26     N, D = mu_q.shape
27
28     # kl diverence for isotropic gaussian
29     kl = 0.5 * ((var_q / var_p) * D + \
30             1.0 / (var_p) * torch.sum(mu_p**2 + mu_q**2 - 2*mu_p*mu_q, axis=1) - D + \
31             D*(torch.log(var_p) - torch.log(var_q)))
32
33     return kl.mean()
34
35 def kl_div_vMF(mu_q, var_q):
36
37     N, D = mu_q.shape
38
39     # we are estimating the variance and not kappa in the network.
40     # They are propertional
41     kappa_q = 1.0 / var_q
42     kl = kappa_q - D * torch.log(2.0)
43
44     return kl.mean()
```

Listing 2: functional.py