

Generative Models - DRE7053

Lecture 3

NORA Summer School 2024

Rogelio A Mancisidor
Assistant Professor
Department of Data Science and Analytics
BI Norwegian Business School

June 10-14, 2024

Outline

1 Generative Clustering

- Theoretical Background
- ELBO

2 Multimodal Learning with VAEs

- Multimodal Learning
- Learning Setting
- Scalable Multimodal VAEs
- Estimating Joint Posterior Distributions

3 Cross-Modal Generation

- Generative modeling in multimodal VAEs
- Examples of Cross-modal Generation

Outline

1 Generative Clustering

- Theoretical Background
- ELBO

2 Multimodal Learning with VAEs

- Multimodal Learning
- Learning Setting
- Scalable Multimodal VAEs
- Estimating Joint Posterior Distributions

3 Cross-Modal Generation

- Generative modeling in multimodal VAEs
- Examples of Cross-modal Generation

Unsupervised Clustering

- The main idea in [12] is to induce a flexible latent space in an unsupervised setting
- This can be achieved in different ways, for example

$$p(z) = \sum_c p(c)p(z|c),$$

where $p(c) \sim \text{Cat}(\pi_0)$, and $p(z|c)$ is a conditional prior distribution.

- The marginal distribution $p(z)$ is a mixture model!

Generative and Inference Models

- The generative model is factorized as $p(\mathbf{x}, \mathbf{z}, c) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z}|c)p(c)$, with distributions

$$p(c) \sim \text{Cat}(\pi_0) \quad (1)$$

$$p(\mathbf{z}|c) \sim \mathcal{N}(\boldsymbol{\mu}_c, \sigma_c^2) \quad (2)$$

$$p(\mathbf{x}|\mathbf{z}) \sim f_{\theta}(\theta) \quad (3)$$

- The inference model is factorized as $q(\mathbf{z}, c) = q(\mathbf{z}|c)q(c|\mathbf{x})$, where

$$q(\mathbf{z}|\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2) \quad (4)$$

$$q(c|\mathbf{x}) \sim \text{Cat}(\boldsymbol{\pi}) \quad (5)$$

Deriving the ELBO

Qualitative Results

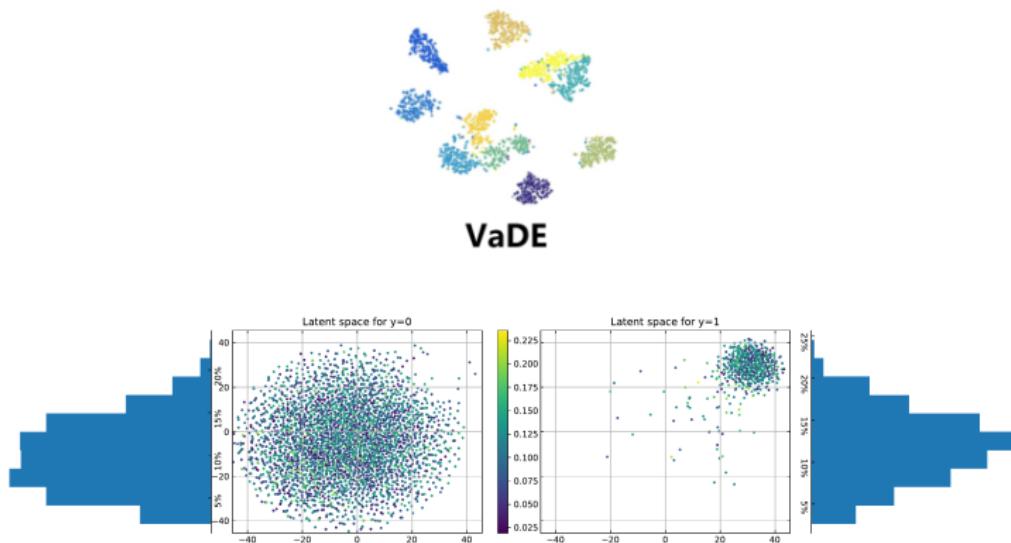


Figure: Deep generative models for reject inference in credit scoring

Things to keep in mind

- Why do we model the prior distribution as a GMM and not the posterior?
- Note the effect of the KL divergence term in the ELBO
- Unsupervised generative clustering is difficult to train!
- VaDE initialize the values for μ_c and σ_c^2 using the expectation-maximization algorithm for GMMs, otherwise the ELBO cannot be optimized.

Outline

1 Generative Clustering

- Theoretical Background
- ELBO

2 Multimodal Learning with VAEs

- Multimodal Learning
- Learning Setting
- Scalable Multimodal VAEs
- Estimating Joint Posterior Distributions

3 Cross-Modal Generation

- Generative modeling in multimodal VAEs
- Examples of Cross-modal Generation

What is Multimodal Learning?

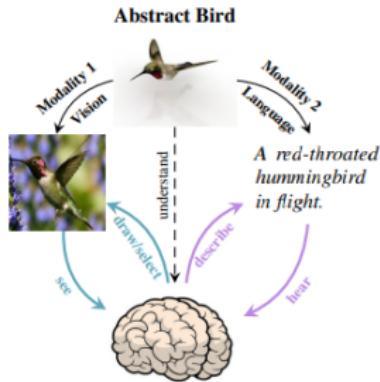
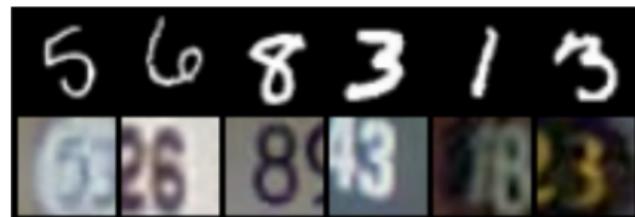


Figure: Borrowed from Shi et al. 2019

Multimodal Learning

- Multimodal learning (MML) uses different measurement modalities of the same object to learn *shared representations*.



$$f : (\textcolor{red}{x}_1, \textcolor{blue}{x}_2) \rightarrow z$$

- Two major approaches:
 - ① Joint representation: projects all modalities into a common space.
 - ② Coordinated representation: assumes that each representation exists in their own space.

Canonical Correlation Analysis

- First coordinated representation approach introduced by Hotelling (1936):
- Canonical Correlation Analysis seeks the vectors \mathbf{a} and \mathbf{b} in

$$\mathbf{U} = \mathbf{a}^T \mathbf{x}_1$$

$$\mathbf{V} = \mathbf{b}^T \mathbf{x}_2$$

such that $\text{corr}(\mathbf{U}, \mathbf{V})$ is maximized. That is,

$$\arg \max_{\mathbf{a}, \mathbf{b}} \text{corr}(\mathbf{U}, \mathbf{V})$$

- Note: \mathbf{V} and \mathbf{U} are shared (*linear*) representations of \mathbf{x}_1 and \mathbf{x}_2 .

Learning Setting in Multimodal VAEs

- Fusion setting: data from all modalities are available at all phases.
- **Cross modality learning:** data from multiple modalities is available only during feature learning; during the testing phase, a(some) modality(ies) is(are) missing.
- The goal is to approximate the joint posterior

$$p(\mathbf{z} | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$$

and to maximize the marginal likelihood of all modalities

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$$

Theoretical Framework

- We assume M modalities $\mathbb{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M)$ that are conditionally independent given the latent variable \mathbf{z}
- Therefore, the generative model becomes

$$\begin{aligned} p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M, \mathbf{z}) &= p(\mathbf{z})p(\mathbf{x}_1|\mathbf{z})p(\mathbf{x}_2|\mathbf{z}) \cdots p(\mathbf{x}_M|\mathbf{z}) \\ p(\mathbb{X}, \mathbf{z}) &\equiv p(\mathbf{z})p(\mathbb{X}|\mathbf{z}) \end{aligned} \tag{6}$$

- Given M modalities, there are 2^M subsets in the powerset $\mathcal{P}(\mathbb{X})$, where only $K = 2^M - 1$ includes any of the modalities
- We need to consider all K subsets for cross-modality generation in any direction (the prior takes care of the empty set)

The Multimodal ELBO

- Assume that we observe 2 modalities, i.e $\mathbb{X} = (\mathbf{x}_1, \mathbf{x}_2)$
- We need to minimize

$$\sum_k KL[q(\mathbf{z}|\mathbb{X}_k) || p(\mathbf{z}|\mathbb{X})] \quad (7)$$

for all $\mathbb{X}_k \in \mathcal{P}(\mathbb{X})$ (excluding the empty set)

Joint Posterior Distributions

- Note the joint posterior distribution $q(\mathbf{z}|\mathbf{x}_1, \mathbf{x}_2)$
- We need to estimate $K - M$ joint posterior distributions, which are called *consensus distributions*
- If we observe, say, 5 modalities, we need to estimate $31 - 5 = 26$ consensus distributions
- Scalable multimodal learning deals with efficient methods that
 - scale linearly on M
 - estimate consensus distributions based on unimodal distributions, which are called *expert distributions*
- Mixture of Experts (MoE) and Product of Experts (PoE) are the common choices to approximate the joint posterior distributions

Product of Experts

- Product of Experts (PoE) simply multiply expert distributions

$$q(z|x_1, x_2) = \frac{1}{Z} q(z|x_1)q(z|x_2) \quad (8)$$

- If experts are Gaussian (as they always are in VAEs):

$$\mu = \left(\sum_i^M \mu_i \tau_i \right) \left(\sum_i^M \tau_i \right)^{-1}$$

$$\sigma^2 = \left(\sum_i^M \tau_i \right)^{-1}$$

where $\tau = 1/\sigma^2$ is the precision

Mixture of Experts

- MoE simply estimates consensus distributions as mixture models

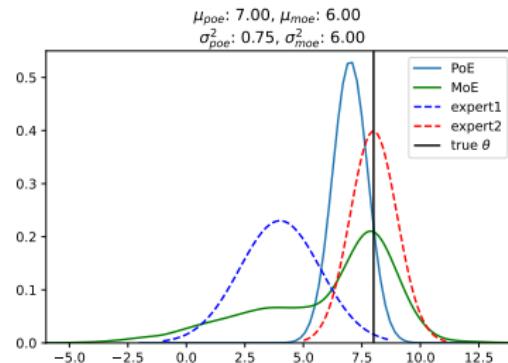
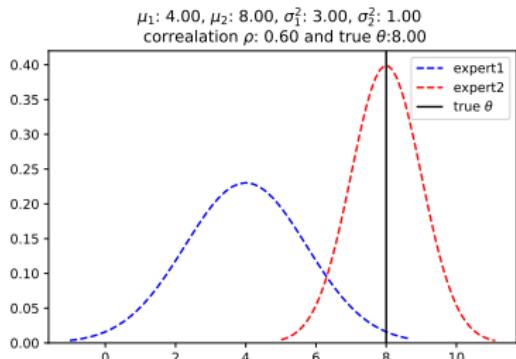
$$q(\mathbf{z}|\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{2}q(\mathbf{z}|\mathbf{x}_1) + \frac{1}{2}q(\mathbf{z}|\mathbf{x}_2) \quad (9)$$

- If experts are Gaussian (as they always are in VAEs):

$$\boldsymbol{\mu} = \frac{1}{2}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2$$

$$\sigma^2 = \frac{1}{2}\sigma_1^2 + \frac{1}{2}\sigma_2^2 + \frac{1}{2}\frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^2$$

MoE vs PoE



- PoE has always smaller variance than any of the experts
- MoE is never "sharper" than any of the experts

Consensus of Dependent Experts

- Both PoE and MoE assume independence between expert distributions
- Consensus of Dependent Experts (CoDE) accounts for the dependency between experts!
- For $\rho = 0$, CoDE recovers the PoE parameters!
- Therefore, CoDE is a more general method

- The main idea is to think about estimates from expert distributions as observations from a likelihood function
- We assume a (non-informative) prior distribution and derive the posterior distribution in principled Bayesian manner
- The posterior is the consensus distribution that we are interested in!

Definition

All M distributions $q(z|\bar{\bar{\mathbb{X}}}_k = 1)$, where $\bar{\bar{\mathbb{X}}}_k$ denotes the cardinality of \mathbb{X}_k , are considered expert distributions, estimating the remaining unknown distributions $q(z|\bar{\bar{\mathbb{X}}}_k > 1)$, which are called consensus distributions.

Definition

Let $\boldsymbol{\theta}^k = (\theta_1^k, \theta_2^k, \dots, \theta_D^k)^T$ denote the latent variable $\mathbf{z} \in \mathbb{R}^D$ for the k -th consensus distribution $q(\mathbf{z}|\mathbb{X}_k)$. Each j -th expert distribution provides a point estimate μ_j^d on the d -th dimension of θ_d^k , and the uncertainty about the estimation is expressed in the parameter σ_j^d of each expert distribution.

Definition

The error of estimation of the j -th expert in the d -th dimension is $e_j^d = \mu_j^d - \theta_d^k$ and the error of estimation of all experts in the d -th dimension is $\mathbf{e}_d = (e_1^d, e_2^d, \dots, e_{M'}^d)^T$, where M' is the number of experts who evaluate the subset \mathbb{X}_k . The overall error of estimation on θ^k is $\mathbf{e}^k = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_D)^T$ with distribution $\mathbf{e}^k \sim \mathcal{N}(\mathbf{0}, \Sigma^k)$ where Σ^k

$$\Sigma^k = \begin{bmatrix} \Sigma_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma_D \end{bmatrix} \quad \text{where } \Sigma_d = \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} & \cdots & \sigma_{1,M'} \\ \sigma_{2,1} & \sigma_2^2 & \cdots & \sigma_{2,M'} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{M',1} & \sigma_{M',2} & \cdots & \sigma_{M'}^2 \end{bmatrix}$$

and $\mathbf{0}$ is a $[M' \times M']$ matrix.

Lemma 2

Assuming a flat prior distribution on the parameter θ^k of the k-th consensus distribution and known covariance matrix Σ , the posterior distribution is

$$h(\theta^k | \mu^k) \sim \mathcal{N}(\mathcal{A}_k^{-1} \mathcal{B}_k, \mathcal{A}_k^{-1}), \quad (10)$$

where $\mathcal{A}_k = \mathbf{u}^t \Sigma_k^{-1} \mathbf{u}$, $\mathcal{B}_k = \mathbf{u}^t \Sigma_k^{-1} \mu^k$, Σ_k^{-1} is the inverse matrix of Σ^k , and \mathbf{u} is a $[M' \cdot D \times D]$ design matrix with size M' vectors of 1s along the diagonal and 0s elsewhere, and $\mu^k = (\mu^1, \mu^2, \dots, \mu^D)^T$ be the vector containing estimates of all expert distributions about all dimensions of θ^k , where $\mu^d = (\mu_1^d, \mu_2^d, \dots, \mu_{M'}^d)^T$.

CoDE - A Simple Example

Outline

1 Generative Clustering

- Theoretical Background
- ELBO

2 Multimodal Learning with VAEs

- Multimodal Learning
- Learning Setting
- Scalable Multimodal VAEs
- Estimating Joint Posterior Distributions

3 Cross-Modal Generation

- Generative modeling in multimodal VAEs
- Examples of Cross-modal Generation

Generative modeling with the posterior distribution

Algorithm 1 Generative Modeling with Multimodal VAEs

$\theta, \phi \leftarrow$ Optimized trainable parameters

repeat for $i = 1, \dots, N$

$\mathbb{X}^i | \text{available} \leftarrow$ Random multimodal sample from available modalities

for $\mathbb{X}_k^i \in \mathcal{P}(\mathbb{X}^i)$ **do**

$\epsilon^i \leftarrow$ Random samples from $\mathcal{N}(\mathbf{0}, \mathbf{1})$

$\mu^k, \sigma^k = q_k(z | \mathbb{X}_k^i)$

$z_k^i = \mu^k + \sigma^k \epsilon^i \leftarrow$ latent from posterior

for $m \in M$ **do**

$x_{m|k} \sim p_\theta(x_m | z_k^i) \leftarrow$ Generate from the m -th likelihood

end for

end for

until

return $x_{m|k} \dots x_{M|K}$

Generative modeling with the prior distribution

Algorithm 1 Generative Modeling with Multimodal VAEs

```
 $\theta, \phi \leftarrow$  Optimized trainable parameters  
repeat    for  $i = 1, \dots, N$   
           $z^i \sim p(z) \leftarrow$  latent from prior  
          for  $m \in M$  do  
               $x_m \sim p_\theta(x_m | z^i) \leftarrow$  Generate from the  $m$ -th likelihood  
          end for  
until  
return  $x_1 \dots x_M$ 
```

- Point-estimate: μ_θ . The mode or median is also possible.
- Random samples: $x^i = \mu_\theta^i + \sigma_\theta^i \epsilon^i$

Image to Text

Tags describing images are generated with the multi-modal learning deep Boltzmann machine (DBM) ([Srivastava & Salakhutdinov, 2012](#)) and with CMMMD. DBM fails to generate coherent tags in the first 3 images. CMMMD is, however, able to generate meaningful tags. In the last image, both models generate coherent tags.

				
Generated tags DBM	water, glass, wine, drink, beer, bubbles, splash, drops, drop	portrait, women, soldier, postcard soldiers, army	nikon, d200, tamron, d300, f28, sb600, d60 nikkor, d50, d90	foliage, autumn, trees, leaves, fall, forest, woods, path
Generated tags MVAE	-	statue	car, performance	a700
Generated tags MMVAE IWAE	canon, night, 2007	nikon, green, lion	flower	trees, autumn
Generated tags <u>CMMMD</u>	sign, fisheye	animal, lion, outdoors, zoo, k10d, challengeyouwinner, boston, wildlife	apple, food	nature, light autumn, leaves wood, path, forest

Image to Text and Text to Image

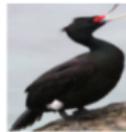
language → vision



this small, white bellied bird
has a brown head and a
red tipped beak.



this bird has a grey
back, a white and black
spotted belly and breast
and a yellow eyebrow.



a black bird with a curved
neck, and a long silver beak
with blue eyes.



a small sized bird that has
a yellow facial marking
with a pointed bill.

vision → language



this bird has wings that are
black and has a white bill.



this bird is brown with white
and a very short beak.



this bird is black black
with a and a red
beak.



this bird is yellow and
yellow, white and
and a and a beak.

Text to Text (Translation)

Type	Sentence
$x_{en} \sim p_{data}$ $x_{vi} \sim p(x_{vi} z(x_{en}))$ GOOGLE(x_{vi})	this was one of the highest points in my life. Đó là một gian tôi vời của cuộc đời tôi. It was a great time of my life.
$x_{en} \sim p_{data}$ $x_{vi} \sim p(x_{vi} z(x_{en}))$ GOOGLE(x_{vi})	the project's also made a big difference in the lives of the people . tôi án này được ra một Điều lớn lao cuộc sống của chúng người sống chữa hưỡng . this project is a great thing for the lives of people who live and thrive .
$x_{vi} \sim p_{data}$ $x_{en} \sim p(x_{en} z(x_{vi}))$ GOOGLE(x_{vi})	trước tiên , tại sao chúng lại có ăn tượng xấu như vậy ? first of all, you do not a good job ? First, why are they so bad?
$x_{vi} \sim p_{data}$ $x_{en} \sim p(x_{en} z(x_{vi}))$ GOOGLE(x_{vi})	Ông ngoại của tôi là một người thật đáng <unk> phục vào thời ấy . grandfather is the best experience of me family . My grandfather was a worthy person at the time .

CelebA Faces generated from the prior



Figure: Left: common latent variable. Right: modality-specific and common latent variables

CelebA Faces generated from the posterior



Figure: Left: common latent variable. Right: modality-specific and common latent variables

Generating Descriptions



archedeyebrows
bigrps
bagsundereyes
biglps
blackhair
brownhair
bushyeyebrows
highcheekbones
male
mouthslightlyopen
mustacstrctbusthyeyebrows
straighthair
male
young

Soclocksh
ey
gbusheyebroes
male
mouthslightlyopen
ovalface
pointynose
smiling

bignps
ey
gbusheyebroes
attracttgo
tm
emile
roweyes
sidebungs

Soclockshadow
attracttgo
recedinghairline
smiling
straighthair
wearinglipstick
young

archedeyebrows
biglps
bignose
bowhair
heavymakeup
highcheekbones
moutslightlyopen
narroweyes
recedinghairline
smiling
straighthair
wearinglipstick
young

Soclockshaiw
blackhair
beury
highcheekbones
moutslightlyopen
rebeard
smiling
sideburns
young

ttactive
biglps
biglps
bignose
bowhair
heavymakeup
highcheekbones
moutslightlyopen
rebeard
smiling
wearinglipstick
young

Soclockshaiw
blackhair
beury
highcheekbones
moutslightlyopen
rebeard
smiling
sideburns
young

ttactive
biglps
bignose
bowhair
heavymakeup
highcheekbones
moutslightlyopen
rebeard
smiling
wearinglipstick
young

archedeyerowbrows
attractive
brownhair
heavymakeup
narroweyes
rebdnoed

blondhair
noebeard
ointy
straighthair
young

Soclockshadow
bagsundereyes
bignoseblondhair
male
noebeard
young



Soclockshadow
blackhair
male
pointynose
smiling
young

bignose
highcheekbones
male
noebeard
ovalface
pointynose
wearingnecktie

archedeyebrows
attractive
blondhair
heavymakeup
highcheekbones
moutslightlyopen
rebeard
smiling
wavyhair
wearinglipstick
young

archedeyebrows
attractive
blondhair
heavymakeup
highcheekbones
moutslightlyopen
rebeard
smiling
wavyhair
wearinglipstick
young

blackhair
biglps
bignose
bowhair
heavymakeup
highcheekbones
moutslightlyopen
narroweyes
rebeard
smiling
wavyhair
wearingearrings
wearinglipstick
young

blackhair
biglps
bignose
bowhair
heavymakeup
highcheekbones
moutslightlyopen
narroweyes
rebeard
smiling
wavyhair
wearingearrings
wearinglipstick
young

atbgive
blondhair
heavymakeup
highcheekbones
moutslightlyopen
narroweyes
rebeard
smiling
wavyhair
wearingearrings
wearinglipstick
young

blackhair
bushyeyebrows
heavymakeup
highcheekbones
moutslightlyopen
narroweyes
rebeard
ovalface
paleskin
wearinglipstick
young

archedeyebrows
aees
ewrownhair
bushyeyebrows
highcheekbones
male
moutslightlyopen
rebeard
ovalface
recedinghairline
ointynose
young

blackhair
bushyeyebrowy
moustache
ovalface
recedinghairline
ointynose
young

Figure: Top: text generated from image. Right: text generated from posterior

Varying Style and Common Latent Variables



Figure: Left: fixed common latent variables. Right: fixed modality-specific latent variables.