# Learning Discrete and Continuous Factors of Data via Alternating Disentanglement

Yeonwoo Jeong, Hyun Oh Song

Seoul National University

ICML19

# Motivation



**Shape? square**

**Postion x? 0.3**

**Postion y? 0.7**

**Rotation? 40°**

**Size? 0.5**

▶ Our goal is to disentangle the underlying explanatory factors of data without any supervision.

# Motivation



square          square

0.3             0.3

0.7             0.7

40°             40°

0.5             0.5

# Motivation

# Motivation

# Motivation



square    square

0.3       0.3

0.7       0.7

40°        0°

0.5       0.5

3

# Motivation



square    square

0.3       0.3

0.7       0.7

40°       40°

0.5       **1**

# Motivation

▶ Most recent methods focus on learning only the continuous factors of variation.

# Motivation

► Most recent methods focus on learning only the continuous factors of variation.

► **Learning discrete representations** is known as a challenging problem. However, **learning continuous and discrete representations** is a more challenging problem.

# Outline

# Overview of our method

# Overview of our method

▶ We propose an efficient procedure for implicitly penalizing the total correlation by **controlling the information flow on each variables**.

▶ We propose a method for jointly learning discrete and continuous latent variables in an **alternating maximization framework**.

# Limitation of $\beta$-**VAE framework**

- $\beta$-VAE sets $\beta > 1$ to penalize $TC(z)$ for **disentangled representations**.

- However, it penalizes the mutual information$(= I(x, z))$ between the data and the latent variables.

## Our method

▶ We aim at penalizing $TC(z)$ by sequentially penalizing the individual summand $\mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i})$.

$$TC(z) = \sum_{i=2}^{m} \mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i}).$$

## Our method

▶ We aim at penalizing $TC(z)$ by sequentially penalizing the individual summand $\mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i})$.

$$TC(z) = \sum_{i=2}^{m} \mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i}).$$

▶ We implicitly minimizes each summand, $\mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i})$ by sequentially maximizing the left hand side $I(x; z_{1:i})$ for all $i = 2, \ldots, m$

1.

$$I(x; z_{1:i}) = I(x; z_{1:i-1}) + I(x; z_i) - \mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i}).$$
$$\uparrow$$

2.

$$I(x; z_{1:i}) = I(x; z_{1:i-1}) + I(x; z_i) - \mathbf{I}(\mathbf{z_{1:i-1}}; \mathbf{z_i}).$$
$$\uparrow \qquad\qquad \bullet \qquad\qquad \uparrow \qquad\qquad \downarrow$$

# Our method

- In practice, we maximize $I(x; z_{1:i})$ by **minimizing reconstruction term** while penalizing $z_{i+1:m}$ with high $\beta$ ($:= \beta_h$) and the others with small $\beta$ ($:= \beta_l$).

# Our method



- ▶ Every latent dimensions are heavily penalized with $\beta_h$. Each penalty on latent dimension is sequentially relieved one at a time with $\beta_l$ in a **cascading fashion**.

# Our method



▶ Every latent dimensions are heavily penalized with $\beta_h$. Each penalty on latent dimension is sequentially relieved one at a time with $\beta_l$ in a **cascading fashion**.

# Our method



- ▶ Every latent dimensions are heavily penalized with $\beta_h$. Each penalty on latent dimension is sequentially relieved one at a time with $\beta_l$ in a **cascading fashion**.

# Our method



Every latent dimensions are heavily penalized with $\beta_h$. Each penalty on latent dimension is sequentially relieved one at a time with $\beta_l$ in a **cascading fashion**.

# Graphical model



(a) $\beta$-VAE    (b) JointVAE        (c) AAE-S            (d) Ours

Figure: Graphical models view. **Solid lines** denote the **generative process** and the **dashed lines** denote the **inference process**. $x, z, d$ denotes the data, continuous latent code, and the discrete latent code respectively.

# Motviation of our method

▶ **AAE with supervised discrete variables(AAE-S)** can learn good continuous representations when the burden of simultaneously modeling the continuous and discrete factors is relieved through supervision on discrete factors unlike **jointVAE**.

# Motviation of our method

- **AAE with supervised discrete variables(AAE-S)** can learn good continuous representations when the burden of simultaneously modeling the continuous and discrete factors is relieved through supervision on discrete factors unlike **jointVAE**.

- Inspired by these findings, our idea is to **alternate** between finding the most likely discrete configuration of the variables given the continuous factors, and updating the parameters $(\phi, \theta)$ given the discrete configurations.

# Construct unary term



- The discrete latent variables are represented using one-hot encodings of each variables $d^{(i)} \in \{e_1, \ldots, e_S\}$.

# Construct unary term



▶ The discrete latent variables are represented using one-hot encodings of each variables $d^{(i)} \in \{e_1, \ldots, e_S\}$.

# Construct unary term



▶ The discrete latent variables are represented using one-hot encodings of each variables $d^{(i)} \in \{e_1, \ldots, e_S\}$.

# Construct unary term



▶ The discrete latent variables are represented using one-hot encodings of each variables $d^{(i)} \in \{e_1, \ldots, e_S\}$.

# Construct unary term



▶ The discrete latent variables are represented using one-hot encodings of each variables $d^{(i)} \in \{e_1, \ldots, e_S\}$.

# Construct unary term



- ▶ The discrete latent variables are represented using one-hot encodings of each variables $d^{(i)} \in \{e_1, \ldots, e_S\}$.

- ▶ $u_\theta^{(i)}$ denotes the vector of the likelihood $\log p_\theta(x^{(i)}|z^{(i)}, e_k)$ evaluated at each $k \in [S]$.

# Alternating minimization scheme

▶ Our goal is to maximize the variational lower bound of the following objective,

$$\mathcal{L}(\theta, \phi) = I(x; [z, d]) - \beta \mathbb{E}_{x \sim p(x)} D_{\mathsf{KL}}(q_\phi(z \mid x) \parallel p(z)) - \lambda D_{\mathsf{KL}}(q(d) \parallel p(d))$$

▶ After rearranging the terms, we arrive at the following optimization problem.

$$\underset{\theta, \phi}{\text{maximize}} \left( \underbrace{\underset{d^{(1)}, \ldots d^{(n)}}{\text{maximize}} \sum_{i=1}^{n} u_\theta^{(i)\mathsf{T}} d^{(i)} - \lambda' \sum_{i \neq j} d^{(i)\mathsf{T}} d^{(j)}}_{:= \mathcal{L}_{LB}(\theta, \phi)} \right)$$

$$- \beta \sum_{i=1}^{n} D_{KL}(q_\phi(z|x^{(i)})||p(z))$$

$$\text{subject to} \quad \|d^{(i)}\|_1 = 1, \ d^{(i)} \in \{0, 1\}^S, \ \forall i,$$

# Finding the most likely discrete configuration



- With the unary terms, we solve inner maximization problem $\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \dots, d^{(n)}]$.[1]

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations" ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations"
ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations"
ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations" ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---
[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations"
ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations" ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
  $\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations"
ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations" ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \dots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations"
ICML2018.

# Finding the most likely discrete configuration



▶ With the unary terms, we solve inner maximization problem
$\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations"
ICML2018.

# Finding the most likely discrete configuration



- ▶ With the unary terms, we solve inner maximization problem $\mathcal{L}_{LB}(\theta, \phi)$ over the discrete variables $[d^{(1)}, \ldots, d^{(n)}]$.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations" ICML2018.

# Finding the most likely discrete configuration



▶ The maximization problem can be exactly solved in polynomial time via *minimum cost flow*(mcf) without continuous relaxation.[1]

---

[1] Jeong, Y. and Song, H. O. "Efficient end-to-end learning for quantizable representations" ICML2018.

# Updating the parameters



$x^{(1)}$     $\cdots$     $x^{(i)}$     $\cdots$     $x^{(n)}$

**Min cost flow solver**

▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Updating the parameters



▶ Then, we update the parameters under this discrete configurations.

# Outline

# Notation

- We denote our full method as **CascadeVAE**.

- We evaluate with disentanglement score introduced in **FactorVAE** and unsupervised classification accuracy.

- Baselines are $\beta$-**VAE, JointVAE, FactorVAE**

# dSprites Dataset Example



- ▶ Shape (discrete) : square, ellipse, heart

- ▶ Scale: 6 values linearly spaced in $[0.5, 1]$

- ▶ Orientation: 40 values in $[0, 2\pi]$

- ▶ Position X: 32 values in $[0, 1]$

- ▶ Position Y: 32 values in $[0, 1]$

# Quantitative results on dSprites

**Disentanglement score**

| Method | m | Mean (std) | Best |
|---|---|---|---|
| $\beta$ VAE | | | |
| ($\beta = 10.0$) | 5 | 70.11 (7.54) | 84.62 |
| ($\beta = 4.0$) | 10 | 74.41 (7.68) | 88.38 |
| FactorVAE | 5 | 81.09 (2.63) | 85.12 |
| | 10 | 82.15 (0.88) | 88.25 |
| JointVAE | 6 | 74.51 (5.17) | 91.75 |
| | 4 | 73.06 (2.18) | 75.38 |
| CascadeVAE | | | |
| ($\beta_l = 1.0$) | 6 | 90.49 (5.28) | **99.50** |
| ($\beta_l = 2.0$) | 4 | **91.34 (7.36)** | 98.62 |

**Unsupervised classification accuracy**

| Method | m | Mean (std) | Best |
|---|---|---|---|
| JointVAE | 6 | 44.79 (3.88) | 53.14 |
| | 4 | 43.99 (3.94) | 54.11 |
| CascadeVAE | 6 | **78.84 (15.65)** | **99.66** |
| | 4 | 76.00 (22.16) | 98.72 |

# Outline

# Conclusion

▶ Our experiments show that information cascading and alternating maximization of discrete and continuous variables, lead to the state of the art performance in 1) **disentanglement score**, and 2) **classification accuracy**.

▶ The source code is available at https://github.com/snu-mllab/DisentanglementICML19.

**Latent dimension traversal in dSprites**

# $\beta$-**VAE**



# **FactorVAE**

## $\beta$-**VAE**

| $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ |



## **FactorVAE**

| $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ |

$\beta$-**VAE**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

**FactorVAE**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

$\beta$-**VAE**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

**FactorVAE**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

# $\beta$-VAE



$z_1$ $\quad$ $z_2$ $\quad$ $z_3$ $\quad$ $z_4$ $\quad$ $z_5$

# FactorVAE



$z_1$ $\quad$ $z_2$ $\quad$ $z_3$ $\quad$ $z_4$ $\quad$ $z_5$

$\beta$-**VAE**

$z_1$    $z_2$    $z_3$    $z_4$    $z_5$

**FactorVAE**

$z_1$    $z_2$    $z_3$    $z_4$    $z_5$

24

## $\beta$-**VAE**



## **FactorVAE**

$\beta$-**VAE**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

**FactorVAE**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

**JointVAE**

**JointVAE**

|  | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ | $z_6$ |
|---|---|---|---|---|---|---|
| $d = [1\ 0\ 0]$ | | | | | | |
| $d = [0\ 1\ 0]$ | | | | | | |
| $d = [0\ 0\ 1]$ | | | | | | |

# JointVAE

# JointVAE



|  | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ | $z_6$ |
|---|---|---|---|---|---|---|
| $d = [1\ 0\ 0]$ | | | | | | |
| $d = [0\ 1\ 0]$ | | | | | | |
| $d = [0\ 0\ 1]$ | | | | | | |

**JointVAE**

## JointVAE



|  | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ | $z_6$ |
|---|---|---|---|---|---|---|
| $d = [1\ 0\ 0]$ | | | | | | |
| $d = [0\ 1\ 0]$ | | | | | | |
| $d = [0\ 0\ 1]$ | | | | | | |

# JointVAE

**JointVAE**

|  | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ | $z_6$ |
|---|---|---|---|---|---|---|
| $d = [1\ 0\ 0]$ | | | | | | |
| $d = [0\ 1\ 0]$ | | | | | | |
| $d = [0\ 0\ 1]$ | | | | | | |

# JointVAE

**JointVAE**

# CascadeVAE

**CascadeVAE**

## CascadeVAE

# CascadeVAE

**CascadeVAE**

# CascadeVAE

**CascadeVAE**

**CascadeVAE**

## CascadeVAE

**CascadeVAE**