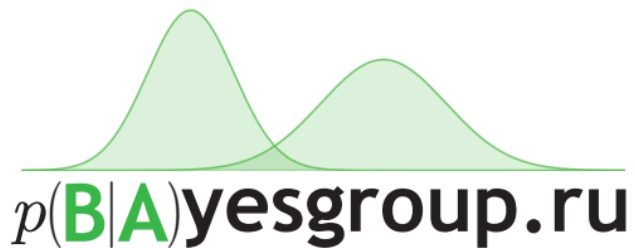# Universal Conditional Machine
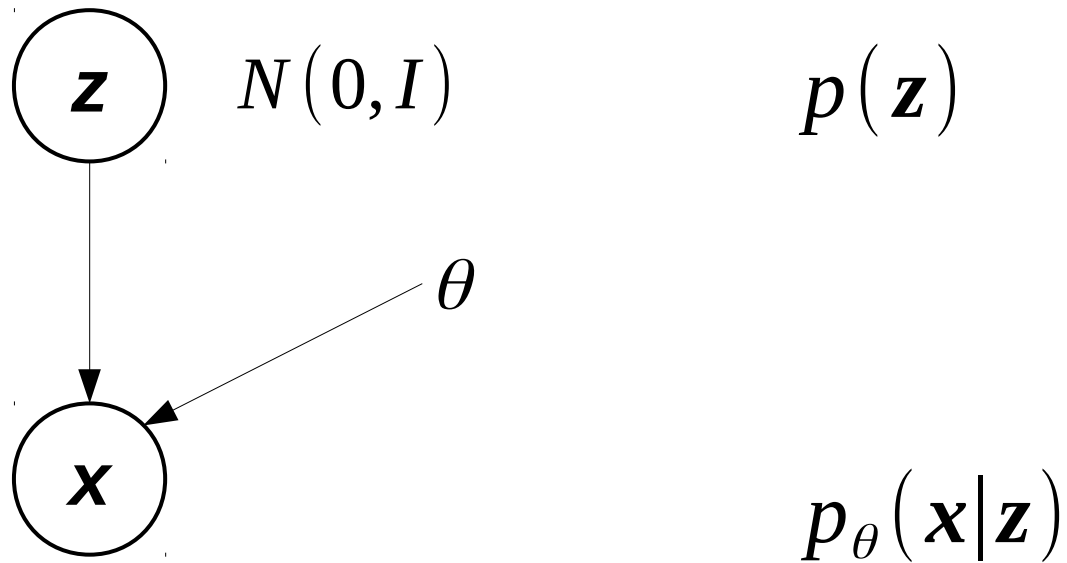
Oleg Ivanov
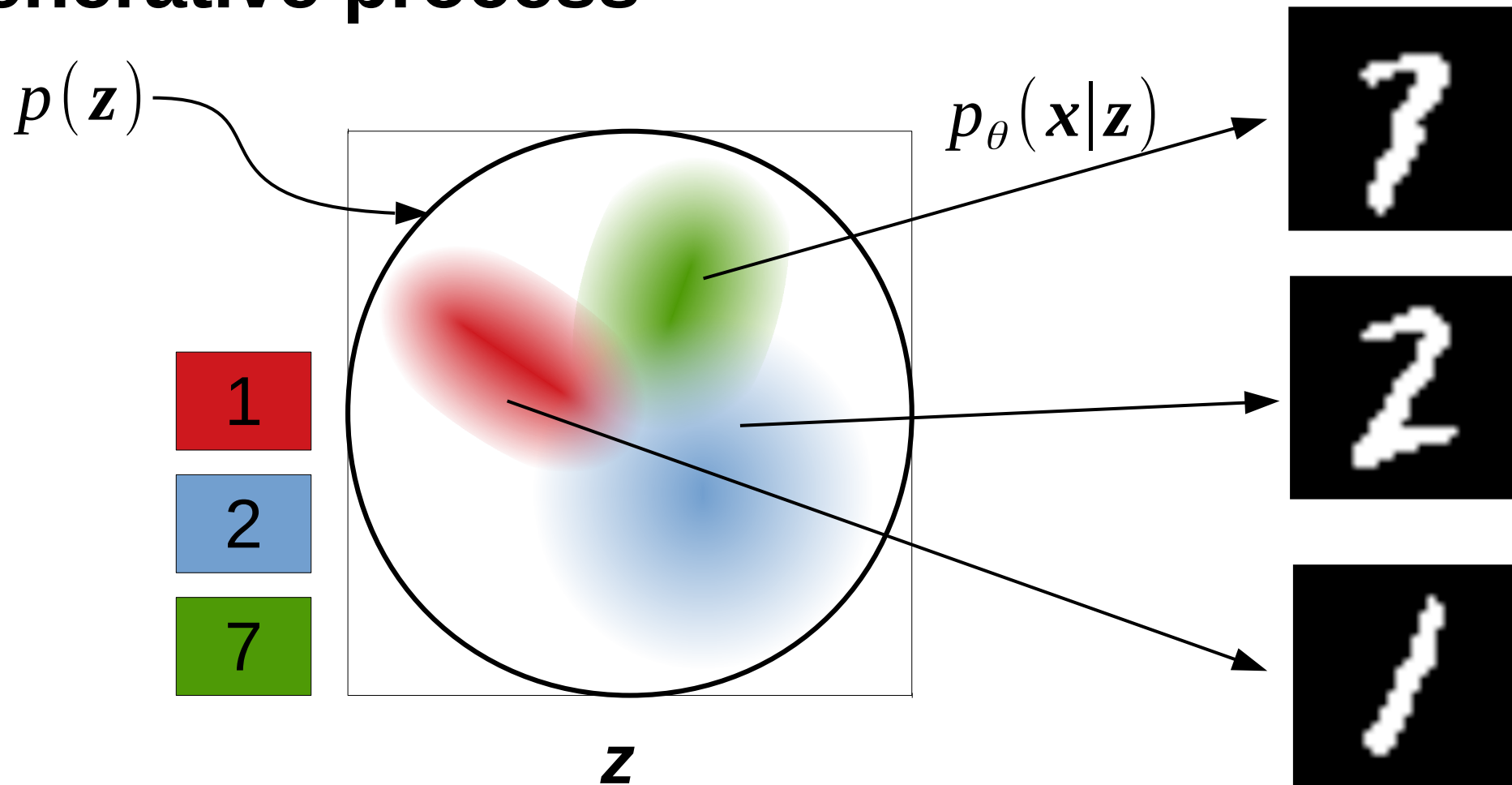
# Variational Autoencoder

$$p_\theta(x)$$

# Generative process

# Generative process

$p(\mathbf{z})$

$p_\theta(\mathbf{x}|\mathbf{z})$



| | |
|---|---|
| 🟥 | 1 |
| 🟦 | 2 |
| 🟩 | 7 |

$\mathbf{z}$

# Generative process

Transform prior into distribution over objects:



$p_\theta(\boldsymbol{x}|\boldsymbol{z})$

$p(\boldsymbol{z})$

PDF

$$p_\theta(\boldsymbol{x}) = \mathop{\mathbb{E}}_{\boldsymbol{z} \sim p(\boldsymbol{z})} p_\theta(\boldsymbol{x}|\boldsymbol{z})$$

# Model Log-Likelihood

$$p_\theta(x) = \mathop{\mathbf{E}}_{z \sim p(z)} p_\theta(x|z)$$

$$L(\theta) = \mathop{\mathbf{E}}_{x \sim p_d(x)} \log p_\theta(x)$$

# Variational Lower Bound

$$L(\theta) \geqslant L(\phi, \theta) =$$

$$= \mathop{\mathbf{E}}_{x \sim p_d(x)} \left[ \underbrace{\mathop{\mathbf{E}}_{z \sim q_\phi(z|x)} \log p_\theta(x|z)}_{\text{Reconstruction term}} - \underbrace{KL(q_\phi(z|x) \| p(z))}_{\text{Regularization term}} \right]$$

Reconstruction term  Regularization term

$$L(\phi, \theta) \rightarrow \max_{\phi, \theta}$$

# Variational Autoencoder

# Latent space visualizations



d = 2

d = 10

# Variational Autoencoder

- Probabilistic generative model
  - Learns distribution $p_d(\boldsymbol{x})$
  - Works with both categorical and real features
  - No assumptions on the data distribution
- Semantic latent space
- Fast inference and generation

# Conditional Variational Autoencoder

$$p_{\psi,\theta}(\boldsymbol{y}|\boldsymbol{x})$$

# Generative process

# Generative process

# Generative process



$$p_\psi(\mathbf{z}|\mathbf{x})$$

$$p_\theta(\mathbf{y}|\mathbf{z},\mathbf{x})$$

| | |
|---|---|
| 🟥 | 1 |
| 🟦 | 2 |
| 🟩 | 7 |

$\mathbf{z}$

# Model Log-Likelihood

$$p_\theta(\boldsymbol{y}) = \mathop{\mathrm{E}}_{\boldsymbol{z} \sim \boldsymbol{p}(\boldsymbol{z})} p_\theta(\boldsymbol{x}|\boldsymbol{z})$$

$$L(\theta) = \mathop{\mathrm{E}}_{\boldsymbol{y} \sim \boldsymbol{p}_d(\boldsymbol{y})} \log p_\theta(\boldsymbol{y})$$

# Model Log-Likelihood

$$p_{\psi,\theta}(\boldsymbol{y}|\boldsymbol{x}) = \operatorname*{E}_{\boldsymbol{z} \sim \boldsymbol{p_\psi}(\boldsymbol{z}|\boldsymbol{x})} p_\theta(\boldsymbol{y}|\mathbf{z}, \boldsymbol{x})$$

$$L(\psi,\theta) = \operatorname*{E}_{\boldsymbol{x}, \boldsymbol{y} \sim \boldsymbol{p_d}(\boldsymbol{x}, \boldsymbol{y})} \log p_\theta(\boldsymbol{y}|\boldsymbol{x})$$

# Variational Lower Bound

$$L(\theta) \geqslant L(\phi, \theta) =$$

Reconstruction loss

$$= \mathop{\mathbf{E}}_{\boldsymbol{y} \sim \boldsymbol{p}_d(\boldsymbol{y})} \left[ \mathop{\mathbf{E}}_{\boldsymbol{z} \sim \boldsymbol{q}_\phi(\boldsymbol{z}|\boldsymbol{x})} \log p_\theta(\boldsymbol{y}|\boldsymbol{z}) - \right.$$

$$\left. - KL(q_\phi(\boldsymbol{z}|\boldsymbol{y}) \| p(\boldsymbol{z})) \right]$$

$$L(\phi, \theta) \to \max_{\phi, \theta}$$

Regularization term

# Variational Lower Bound

$$L(\psi,\theta) \geqslant L(\phi,\psi,\theta) =$$

Reconsruction loss

$$= \mathop{\mathbf{E}}_{\boldsymbol{x},\boldsymbol{y} \sim \boldsymbol{p_d(x,y)}} \left[ \mathop{\mathbf{E}}_{\boldsymbol{z} \sim \boldsymbol{q_\phi(z|x,y)}} \log p_\theta(\boldsymbol{y}|\boldsymbol{z},\boldsymbol{x}) - \right.$$

$$\left. - KL\big(q_\phi(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{y})\|p_\psi(\boldsymbol{z}|\boldsymbol{x})\big) \right]$$

Regularization term

$$L(\phi,\psi,\theta) \to \max_{\phi,\psi,\theta}$$

# Conditional Variational Autoencoder



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

$\boldsymbol{x}$

$\boldsymbol{y}$

Proposal network

$\boldsymbol{\mu}$

$\boldsymbol{\sigma}$

$\boldsymbol{\mu}$

$\boldsymbol{\sigma}$

Prior network

$\boldsymbol{x}$

$p_\psi(\mathbf{z}|\boldsymbol{x})$

$\mathbf{z}$

Generative network

$\boldsymbol{\mu}$

$\boldsymbol{\sigma}$

$p_\theta(\boldsymbol{y}|\mathbf{z}, \boldsymbol{x})$

# Gaussian Stochastic Neural Network

Motivation:
1. Train/test procedure inconsistency
2. Gaps in latent space
3. Better Monte-Carlo log-likelihood estimations

$$q_\phi(\mathbf{z}|\boldsymbol{x},\boldsymbol{y}) = p_\psi(\mathbf{z}|\boldsymbol{x_{1-b}},\boldsymbol{b})$$

$$L_{GSNN}(\phi,\psi,\theta) = \operatorname*{E}_{\boldsymbol{x},\boldsymbol{y}\sim \boldsymbol{p_d}(x,y)} \operatorname*{E}_{\boldsymbol{z}\sim \boldsymbol{p_\phi}(z|x)} \log p_\theta(\boldsymbol{y}|\mathbf{z},\boldsymbol{x})$$

$$L_{hybrid} = \alpha\, L_{CVAE} + (1-\alpha)\, L_{GSNN}$$

# Gaussian Stochastic Neural Network



Gaussian Stochastic Layer

$x$ Prior network $p_\psi(\mathbf{z}|\boldsymbol{x})$ $\boldsymbol{\mu}$ $\boldsymbol{\sigma}$ $\boldsymbol{z}$ Generative network $p_\theta(\boldsymbol{y}|\mathbf{z},\boldsymbol{x})$ $\boldsymbol{\mu}$ $\boldsymbol{\sigma}$

# Some motivation pictures



Input

Samples

# Conditional Variational Autoencoder

- Learns conditional distribution $p_d\left(\boldsymbol{y}\middle|\boldsymbol{x}\right)$
  - Essential when $p_d\left(\boldsymbol{y}\middle|\boldsymbol{x}\right)$ has several local maximums
- Obtains by conditioning Variational Autoencoder
  - Inherits the majority of its properties
- Has prior network to model prior latent distribution
- Modifications: GSNN, hybrid model

# Universal Conditional Machine

$$p_{\psi,\theta}\left(\boldsymbol{x_I}|\boldsymbol{x_{U \setminus I}}\right)$$

# Problem statement

Test set

Features

Objects



Missing value

Observed value

# Missing features mask

Features



Objects

$\boldsymbol{b} \in \{0,1\}^D$

(0, 0, 0, 0, 0)
(0, 0, 0, **1**, 0)
(0, **1**, 0, 0, 0)
(0, 0, 0, 0, 0)
(0, 0, **1**, **1**, 0)
(0, **1**, **1**, 0, 0)

Indexation example

$\boldsymbol{b} = (0,0,1,1,0)$
$\boldsymbol{x_b} = (x_3, x_4)$
$\boldsymbol{x_{1-b}} = (x_1, x_2, x_5)$

# **Problem statement**

## Features

## Objects



$$p\left(x_4 \mid x_1, x_2, x_3, x_5\right)$$

$$p\left(x_2 \mid x_1, x_3, x_4, x_5\right)$$

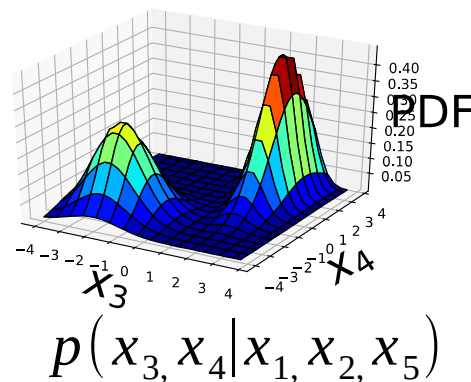$$p_{\psi,\theta}\left(\boldsymbol{x_b} \mid \boldsymbol{x_{1-b}}\right)$$

$$p\left(x_2, x_3 \mid x_1, x_4, x_5\right)$$

$$p\left(x_3, x_4 \mid x_1, x_2, x_5\right)$$

# Why distributions?

Input

Average

Most probable



Single imputation causes information loss

# Why distributions?

Input

Average
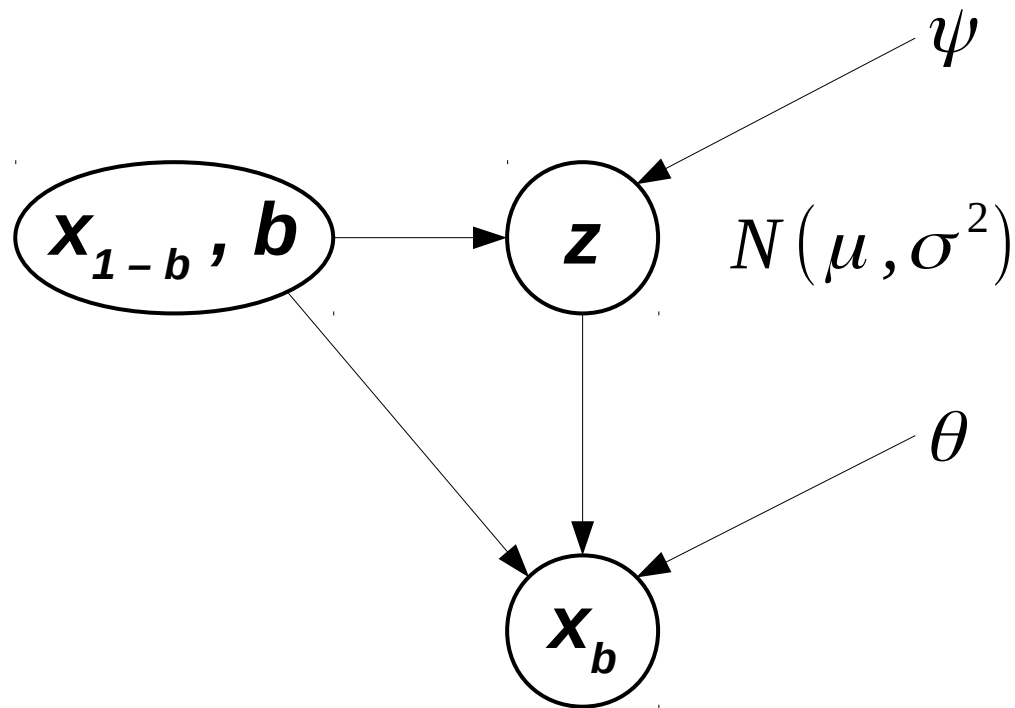
Most probable



Distribution samples

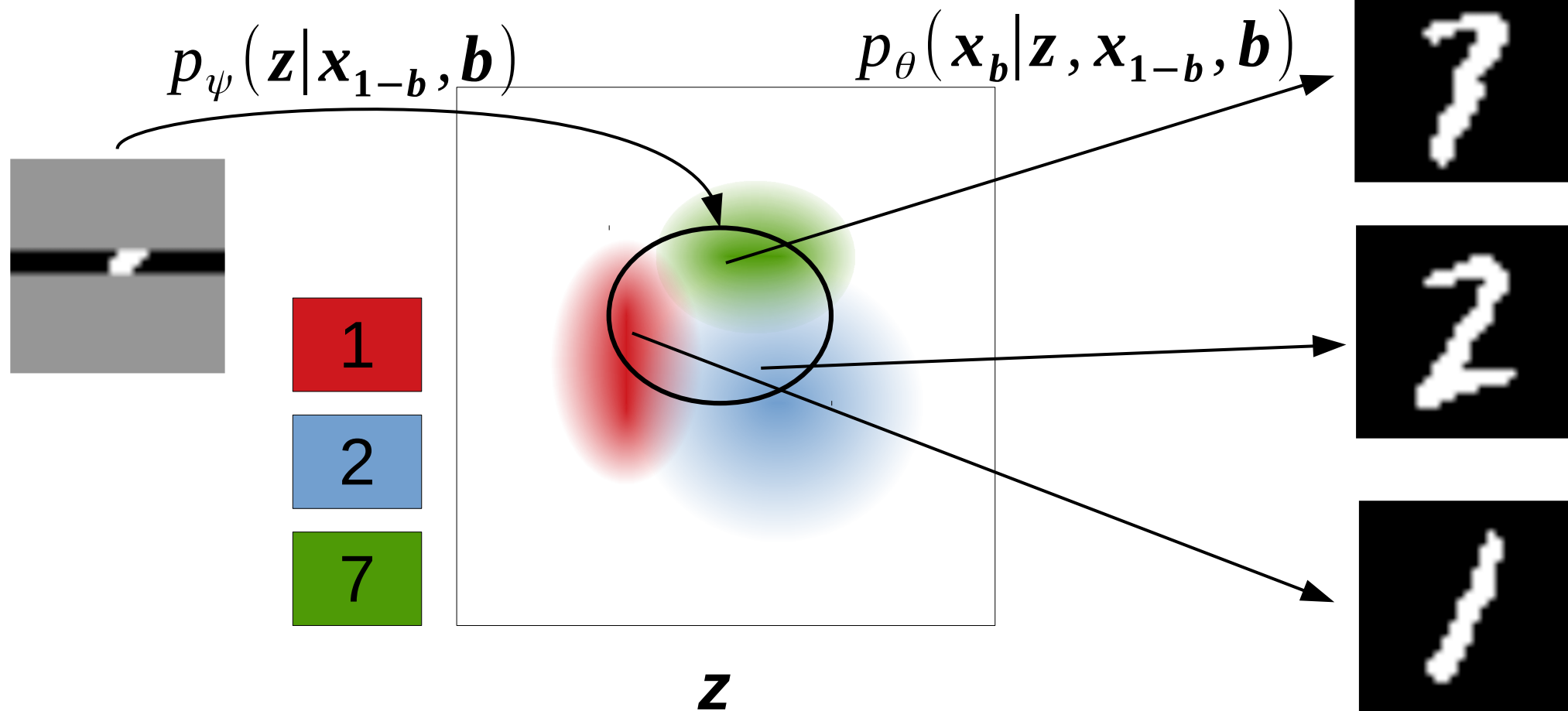# Generative process



$$p_\psi(\mathbf{z}|\mathbf{x}_{1-b}, \mathbf{b})$$

$$p_\theta(\mathbf{x}_b|\mathbf{z}, \mathbf{x}_{1-b}, \mathbf{b})$$

# Generative process



$$p_\psi\left(\boldsymbol{z}|\boldsymbol{x_{1-b}}, \boldsymbol{b}\right)$$

$$p_\theta\left(\boldsymbol{x_b}|\boldsymbol{z}, \boldsymbol{x_{1-b}}, \boldsymbol{b}\right)$$

| | |
|---|---|
| 1 | |
| 2 | |
| 7 | |

$\boldsymbol{z}$

# Mask distribution and model likelihood

- $p_{\psi,\theta}(\boldsymbol{x_b}|\boldsymbol{x_{1-b}},\boldsymbol{b}) = \underset{\boldsymbol{z} \sim \boldsymbol{p_\psi}(\boldsymbol{z}|\boldsymbol{x_{1-b}},\boldsymbol{b})}{\mathrm{E}} p_\theta(\boldsymbol{x_b}|\boldsymbol{z},\boldsymbol{x_{1-b}},\boldsymbol{b})$

- User-defined mask distribution $p_b(\boldsymbol{b})$

- "Train set":

$$(\boldsymbol{x_i},\boldsymbol{b_i})_{i=1}^{N} : \boldsymbol{x} \sim p_d(\boldsymbol{x}), \boldsymbol{b} \sim p_b(\boldsymbol{b})$$

- Model log-likelihood

$$L(\psi,\theta) = \underset{\substack{\boldsymbol{x} \sim \boldsymbol{p_d}(\boldsymbol{x}) \\ \boldsymbol{b} \sim \boldsymbol{p_b}(\boldsymbol{b})}}{\mathrm{E}} \log p_{\psi,\theta}(\boldsymbol{x_b}|\boldsymbol{x_{1-b}},\boldsymbol{b})$$

# Variational Lower Bound

$$L(\psi,\theta) \geqslant L(\phi,\psi,\theta) =$$

$$= \mathop{\mathbf{E}}_{\substack{x \sim p_d(x) \\ b \sim p_b(b)}} \left[ \mathop{\mathbf{E}}_{z \sim q_\phi(z|x,b)} \log p_\theta(x_b|z,x_{1-b},b) - \right.$$

$$\left. - KL(q_\phi(z|x,b) \| p_\psi(z|x_{1-b},b)) \right]$$

# Universal Conditional Machine

# Universal Conditional Machine



Proposal network

ResNet Blocks

Fully Connected

$q_\phi(z|x,b)$

KL divergence

Reconstruction loss

$p_\psi(z|x\circ(1-b),b)$

Fully Connected

$z$

ResNet Blocks

$p_\theta(x\circ b|z,x\circ(1-b),b)$

Skip connections

Prior network

Generative network

# Gaussian Stochastic Neural Network



Gaussian Stochastic Layer

$p_\psi(\boldsymbol{z}|\boldsymbol{x})$

$x$

Prior network

$\boldsymbol{\mu}$

$\boldsymbol{\sigma}$

$\boldsymbol{z}$

Generative network

$\boldsymbol{\mu}$

$\boldsymbol{\sigma}$

$p_\theta(\boldsymbol{y}|\boldsymbol{z},\boldsymbol{x})$

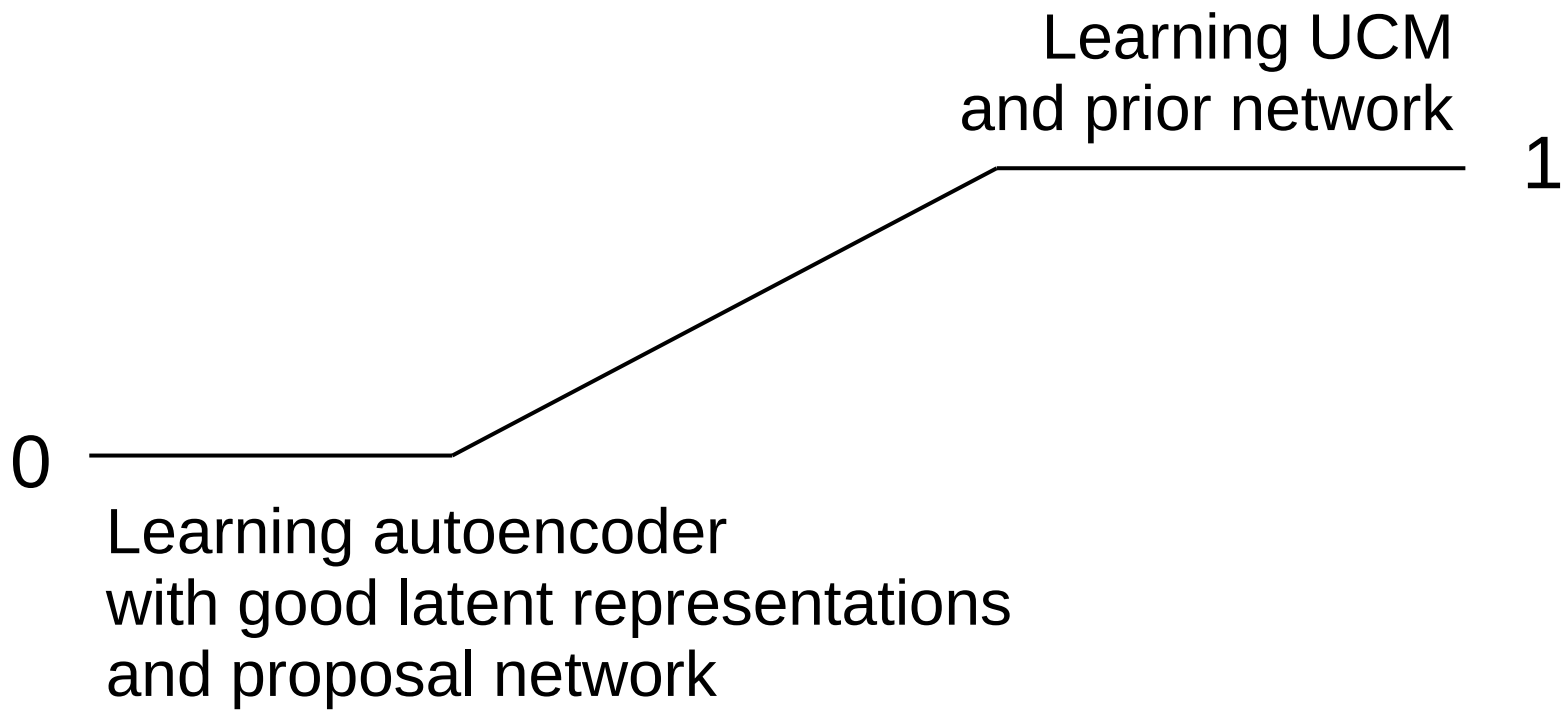$$L_{hybrid} = \alpha\, L_{UCM} + (1-\alpha)\, L_{GSNN}$$

# Missing features in train set

- Missing feature $x_i = \omega$

- Conditioned mask distribution $p_b(\boldsymbol{b}|\boldsymbol{x})$

- $x_i = \omega \;\Rightarrow\; p_b(b_i|\boldsymbol{x}) = 1$

- Missing features marginalization: $x_i = \omega \;\Rightarrow$

$$\log p_\theta(x_i|\boldsymbol{z}, \boldsymbol{x_{1-b}}, \boldsymbol{b}) \;=\; \log \int p_\theta(\hat{x}_i|\boldsymbol{z}, \boldsymbol{x_{1-b}}, \boldsymbol{b})\, d\hat{x}_i \;=\; \log 1 \;=\; 0$$

# KL coefficient

Found necessary only for syntetic data

Learning UCM
and prior network

1

0

Learning autoencoder
with good latent representations
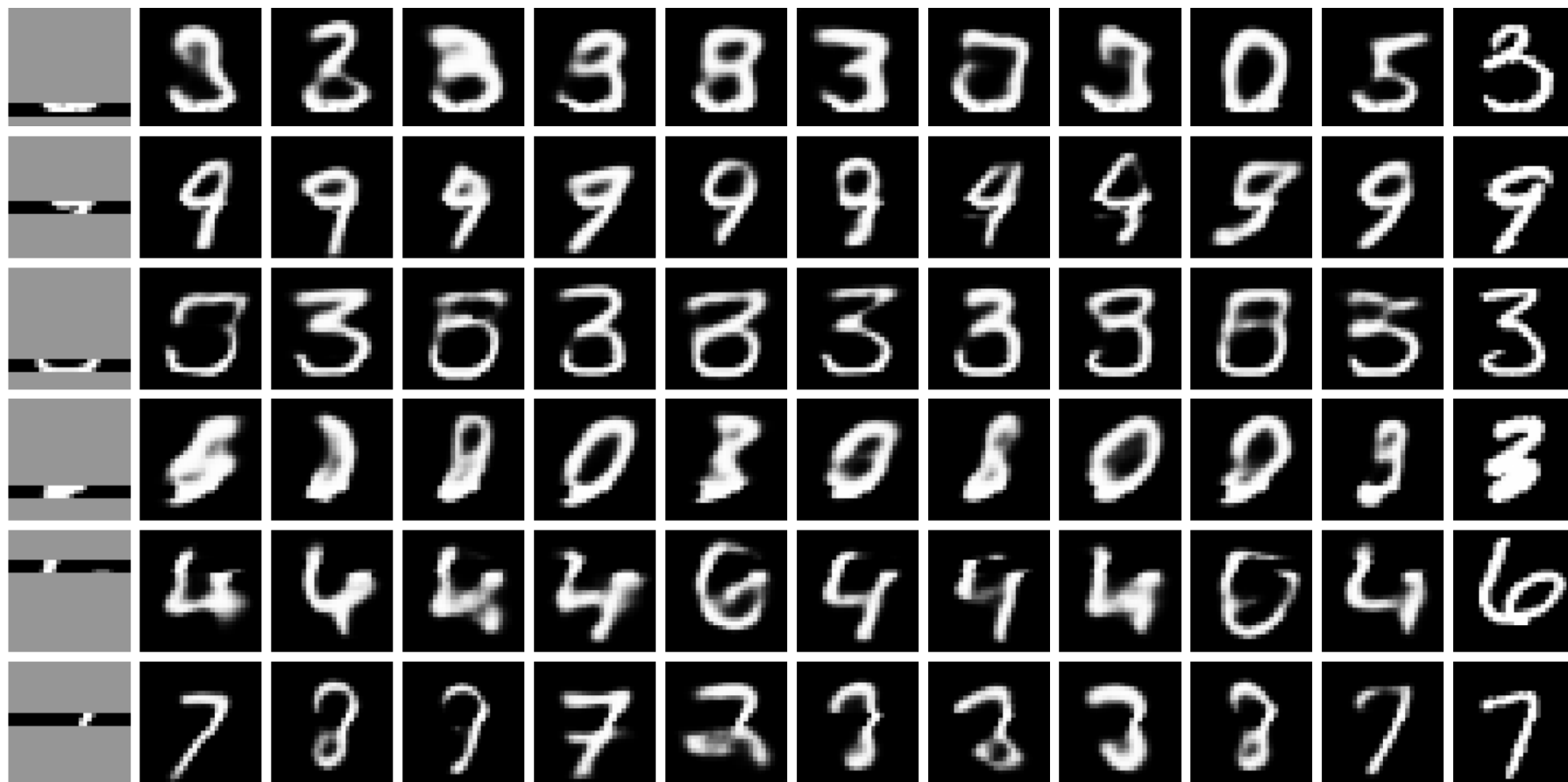and proposal network

# Experiment 1

Missing feature multiple imputation for supervised learning.
50% features of dataset are missed.
10 imputations for UCM, GSNN, MG.

| DATASET | AVERAGE | XGBOOST | UMC | GSNN | MG |
|---------|---------|---------|-----|------|-----|
| BOSTON | $0.505 \pm 0.061$ | $0.502 \pm 0.056$ | $\mathbf{0.577} \pm 0.069$ | $0.563 \pm 0.069$ | $0.564 \pm 0.055$ |
| CONCRETE | $0.452 \pm 0.042$ | $0.458 \pm 0.040$ | $\mathbf{0.494} \pm 0.032$ | $0.453 \pm 0.030$ | $0.484 \pm 0.030$ |
| CASP | $0.840 \pm 0.002$ | $0.842 \pm 0.002$ | $\mathbf{0.856} \pm 0.002$ | $\mathbf{0.856} \pm 0.002$ | $0.850 \pm 0.003$ |
| WINE | $0.230 \pm 0.012$ | $0.236 \pm 0.008$ | $0.232 \pm 0.016$ | $\mathbf{0.243} \pm 0.014$ | $0.238 \pm 0.011$ |
| YEAST | $0.423 \pm 0.025$ | $0.426 \pm 0.026$ | $\mathbf{0.436} \pm 0.019$ | $0.419 \pm 0.025$ | $0.430 \pm 0.019$ |

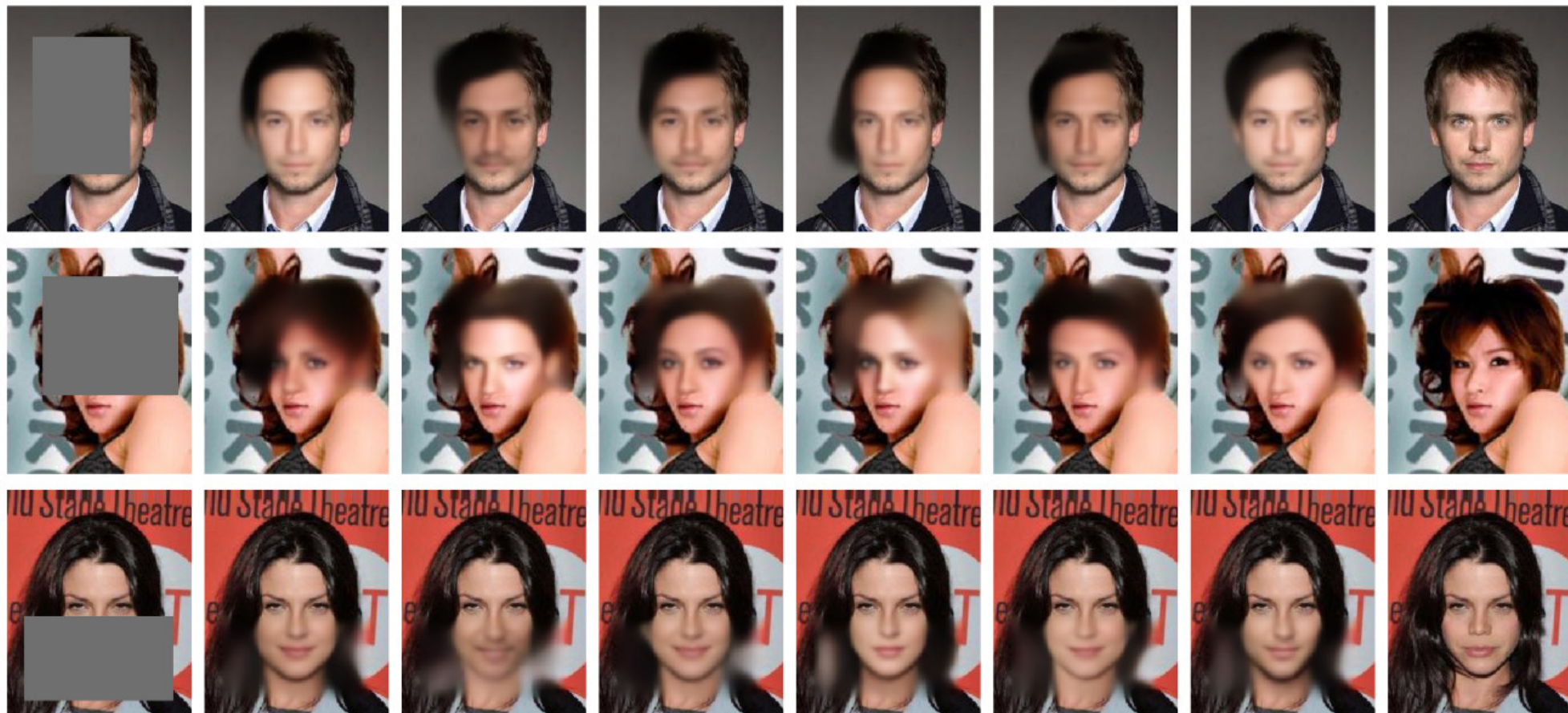MG – Multivariate Gaussian

# Experiment 2: image inpainting, MNIST
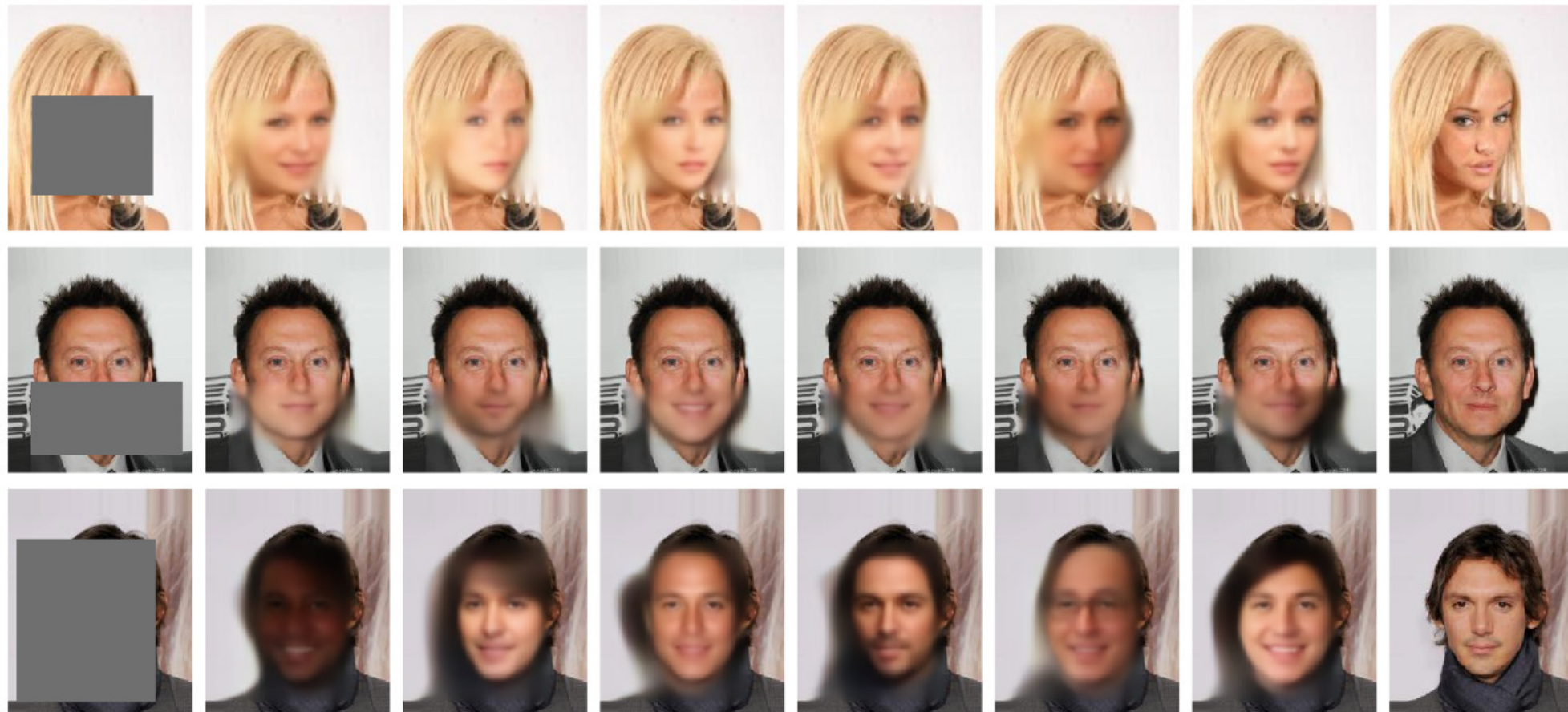
# Experiment 2: image inpainting, Omniglot
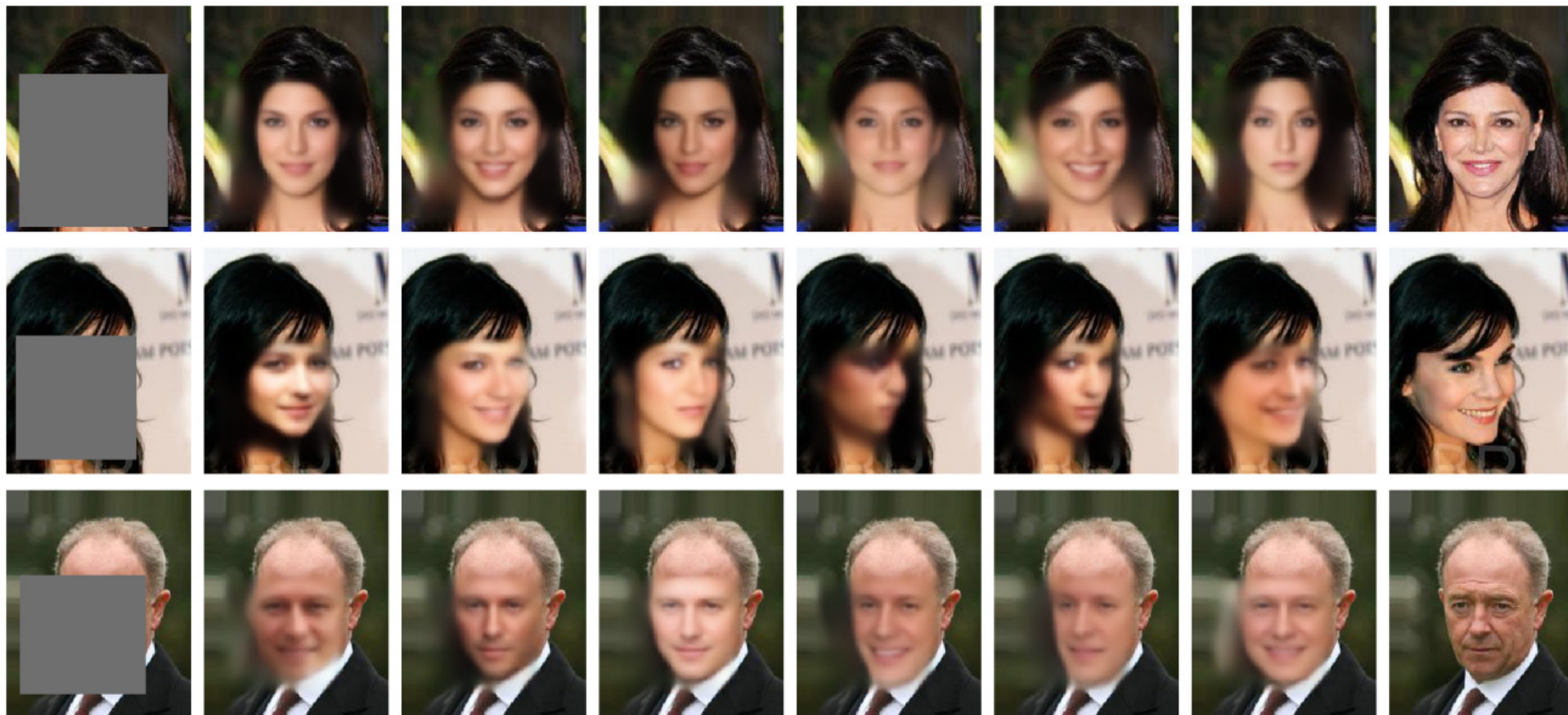
Experiment 2: image inpainting, CelebA

# Experiment 2: image inpainting, CelebA

Experiment 2: image inpainting, CelebA

# Universal Conditional Machine

- Learns all conditional distributions $p_d\left(x_b \middle| x_{1-b}\right)$
  - The importance of the conditioning is given by $p_b(b)$
- Further extension of VAE and CVAE
  - Same conditioning technique as in CVAE
  - Lots of slight modifications
- It works!
  - As preprocessing (multiple imputation) for datasets with missing data
  - For image inpainting
  - TBD: image colourization

# Saga of hybrid model

$$\alpha\, L_{CVAE} + (1 - \alpha)\, L_{GSNN}$$

# GSNN – good or evil?

$$L_{hybrid} = \alpha\, L_{UCM} + (1 - \alpha)\, L_{GSNN}$$

Motivation:
1. Train/test procedure inconsistency
2. Gaps in latent space
3. Better Monte-Carlo log-likelihood estimations

# Log-likelihood estimations

$$\log p_{\psi,\theta}(\boldsymbol{x_b}|\boldsymbol{x_{1-b}},\boldsymbol{b}) = \log \operatorname*{E}_{\boldsymbol{z} \sim \boldsymbol{p_\psi(z|x_{1-b},b)}} p_\theta(\boldsymbol{x_b}|\boldsymbol{z},\boldsymbol{x_{1-b}},\boldsymbol{b})$$

Monte-Carlo

$$\approx \log \frac{1}{S} \sum_{i=1}^{S} p_\theta(\boldsymbol{x_b}|\boldsymbol{z_i},\boldsymbol{x_{1-b}},\boldsymbol{b})$$

$$\boldsymbol{z_i} \sim p_\psi(\boldsymbol{z}|\boldsymbol{x_{1-b}},\boldsymbol{b})$$

Importance Sampling

$$\approx \log \frac{1}{S} \sum_{i=1}^{S} \frac{p_\theta(\boldsymbol{x_b}|\boldsymbol{z_i},\boldsymbol{x_{1-b}},\boldsymbol{b})\, p_\psi(\boldsymbol{z_i}|\boldsymbol{x_{1-b}},\boldsymbol{b})}{q_\phi(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{b})}$$

$$\boldsymbol{z_i} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x},\boldsymbol{b})$$

# Log-likelihood estimations

| Method | MNIST | Omniglot | CelebA |
|---|---|---|---|
| UMC IS-$10^2$ | $\mathbf{83} \pm 2$ | $\mathbf{275} \pm 17$ | $\mathbf{34035} \pm 1609$ |
| UMC MC-$10^4$ | $98 \pm 4$ | $1452 \pm 109$ | $41513 \pm 2163$ |
| UMC MC-$10^2$ | $135 \pm 6$ | $2203 \pm 150$ | $53904 \pm 3121$ |
| GSNN MC-$10^4$ | $139 \pm 3$ | $1199 \pm 62$ | $53427 \pm 2208$ |
| GSNN MC-$10^2$ | $139 \pm 3$ | $1200 \pm 62$ | $53486 \pm 2210$ |
| Naive Bayes | $205$ | $2490$ | $269480$ |

# Gaps & overlapping in latent space: XOR

$x_1$

$x_2$

True distribution
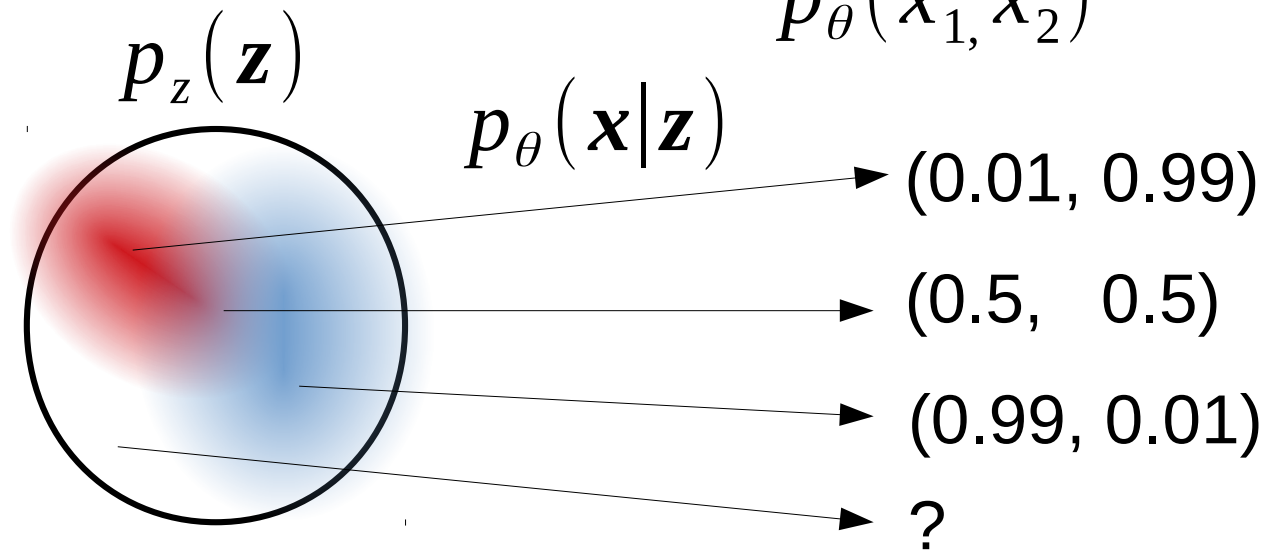
$$x_1, x_2 \in \{0, 1\}$$
$$p_d(x_1, x_2) = 0.5(x_1 \oplus x_2)$$
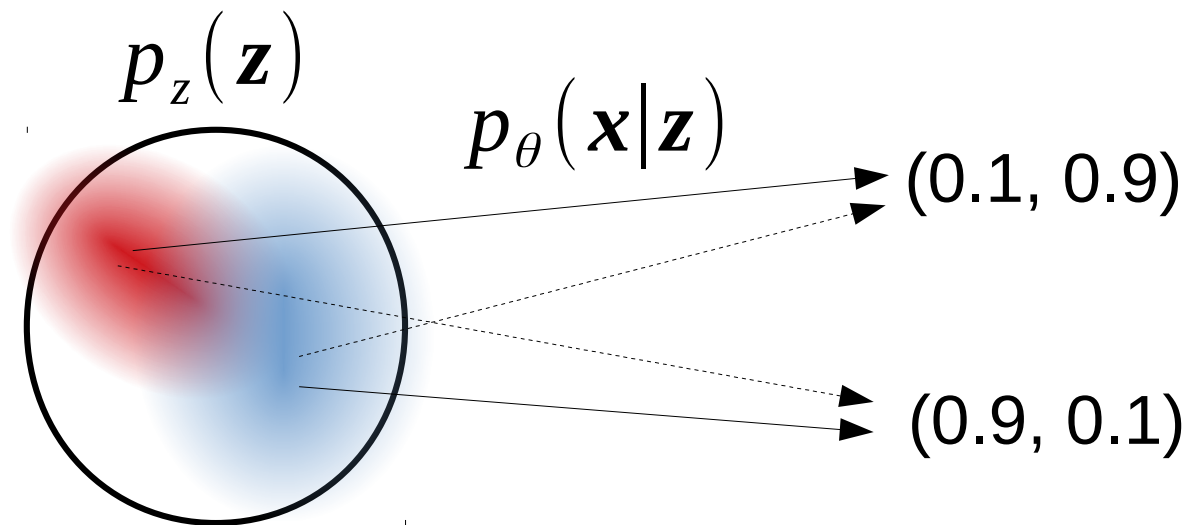
$x_1$

$x_2$

$$p_\theta(x_1, x_2)$$

$$q_\phi(z|(0,1))$$

$$q_\phi(z|(1,0))$$

$$p_z(z)$$

$$p_\theta(x|z)$$

(0.01, 0.99)

(0.5,   0.5)

(0.99, 0.01)

?

# Why don't use Monte-Carlo?



$p_z(\mathbf{z})$

$p_\theta(\boldsymbol{x}|\mathbf{z})$

$q_\phi(\mathbf{z}|(0,1))$

$q_\phi(\mathbf{z}|(1,0))$

(0.1, 0.9)

(0.9, 0.1)

# Why don't use Monte-Carlo?

$q_\phi(\mathbf{z}|(0,1))$

$q_\phi(\mathbf{z}|(1,0))$

$p_z(\mathbf{z})$

$p_\theta(\mathbf{x}|\mathbf{z})$

(0.5, 0.5)

(0.5, 0.5)

# GSNN – good or evil?

True distribution

# Saga of hybrid model

- Gaps and overlapping are problems for Gaussian latent space

    - Might be with normalizing flows, etc

- GSNN Monte-Carlo estimations with a few samples are better

    - Monte-Carlo estimator is not precise

        - Needs too many samples to find the region in latent space suitable for the given object

- GSNN can't learn multimodal distribution

    - GSNN might work better with big S at the training stage

- True log-likelihood is better for UCM without any GSNN

# Universal Marginalizer

Improved version

$$p_\theta\big(\boldsymbol{x_i}\big|\boldsymbol{x_{U\setminus I}}\big)$$

# Mask distribution and model likelihood

- User-defined mask distribution $p_b(\boldsymbol{b})$

- "Train set":

$$\left(\boldsymbol{x}_i, \boldsymbol{b}_i\right)_{i=1}^{N} : \boldsymbol{x} \sim p_d(\boldsymbol{x}), \boldsymbol{b} \sim p_b(\boldsymbol{b})$$

- Model log-likelihood

$$L(\theta) = \mathop{\mathbf{E}}_{\substack{\boldsymbol{x} \sim \boldsymbol{p_d}(\boldsymbol{x}) \\ \boldsymbol{b} \sim \boldsymbol{p_b}(\boldsymbol{b})}} \sum_{i=1}^{D} b_i \log p_\theta\left(x_i | \boldsymbol{x_{1-b}}, \boldsymbol{b}\right)$$

# Universal Marginalizer



$$p_\theta\left(x_i | \boldsymbol{x_{1-b}}, \boldsymbol{b}\right)$$

# Joint distribution: chain rule

- Choose $i \in \boldsymbol{b}$

- Sample $x_i \sim p_\theta \left( x_i \middle| \boldsymbol{x_{1-b}}, \boldsymbol{b} \right)$

- Update $\boldsymbol{b} \Leftarrow \boldsymbol{b} - \boldsymbol{e_i}$

- Repeat while $\boldsymbol{b} \neq \boldsymbol{0}$

- Log-likelihood: don't sample $x_i$, but compute the product of conditional probabilities instead

Here ends the original paper

# Joint distribution

- Choose $i \in \boldsymbol{b}$

  – Sequential left to right:

$$\boldsymbol{e}_{1..i} = \big(\overbrace{0, 0, ..., 0}^{i}, 1, 1, ..., 1\big)$$

$$p_\theta(\boldsymbol{x_b} | \boldsymbol{x_{1-b}}, \boldsymbol{b}) = \prod_{i \in \boldsymbol{b}} p_\theta(x_i | \boldsymbol{x_{1-b \wedge e_{1..i}}}, \boldsymbol{b} \wedge \boldsymbol{e}_{1..i})$$

# Joint distribution

- Choose $i \in \boldsymbol{b}$
  - Sequential left to right
  - At random uni-probable

# Joint distribution

- Choose $i \in \boldsymbol{b}$

  - Sequential left to right

  - At random uni-probable

- The distribution over $\boldsymbol{b}$ at test stage is not $p_b(\boldsymbol{b})$

- This inconsistency ruins everything

- Need generative process for induced $\hat{p}_b(\boldsymbol{b})$

# Consistent mask generative process

Generative process for $\hat{p}_b(\boldsymbol{b})$

Choose $i \in \boldsymbol{b}$

**Sequential left to right**

- $\boldsymbol{b} \sim p_b(\boldsymbol{b})$
- $u \sim U_D\left[0, 1, ..., \sum \boldsymbol{b}\right]$
- $j$: u$^{\text{th}}$ 1 in $\boldsymbol{b}$
- $\hat{\boldsymbol{b}} = \boldsymbol{b} \wedge \boldsymbol{e}_{1..j}$

**At random uni-probable**

- $\boldsymbol{b} \sim p_b(\boldsymbol{b})$
- $u \sim U[0, 1]$
- $\boldsymbol{b_0} \sim Bernoulli(u)$
- $\hat{\boldsymbol{b}} = \boldsymbol{b_0} \circ \boldsymbol{b}$

# MNIST inpaintings

# Universal Marginalizer

- Needs fast generative process for induced $\hat{p}_b(\boldsymbol{b})$ for given $p_b(\boldsymbol{b})$
  - Allows to keep the training speed for one epoch
- Needs O(D) time to generate sample or estimate likelihood
- Takes into account local dependencies
- The relation to UCM is similar to the relation between VAE and PixelCNN