

# Amortised learning by wake sleep environment

Kevin, Ted Moskovitz, Heishiro Kanagawa, Maneesh Sahani

Gatsby Computational Neuroscience Unit, University College London

February 25, 2020

# Latent variable models

$$p_{\theta}(z, x)$$

# Latent variable models

$$p_{\theta}(z, x)$$

Caputure rich data by

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

- infer  $z$

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

- infer  $z$
- generate  $x$  (prediction, imputation)

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

- infer  $z$
- generate  $x$  (prediction, imputation)
- compression, outlier detection, and so on

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

- infer  $z$
- generate  $x$  (prediction, imputation)
- compression, outlier detection, and so on

Only if we find the **correct parameters** for data

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

- infer  $z$
- generate  $x$  (prediction, imputation)
- compression, outlier detection, and so on

Only if we find the **correct parameters** for data

- objectives: maximum likelihood (KL), score matching (Fisher), Jenson-Shannon, Wasserstein

# Latent variable models

$$p_{\theta}(z, x)$$

Capture rich data by

- model structure: hierarchical, sequential
- latent variables: continuous, discrete, positive, geometry of support

Has wide practical use

- infer  $z$
- generate  $x$  (prediction, imputation)
- compression, outlier detection, and so on

Only if we find the **correct parameters** for data

- objectives: maximum likelihood (KL), score matching (Fisher), Jenson-Shannon, Wasserstein

**This talk: maximum likelihood**

# Maximum likelihood estimation

Maximum likelihood: maximise  $\mathbb{E}_{p^*(\mathbf{x})} [\log p_{\theta}(\mathbf{x})]$

# Maximum likelihood estimation

Maximum likelihood: maximise  $\mathbb{E}_{p^*(\mathbf{x})} [\log p_{\theta}(\mathbf{x})]$

Nice properties as  $n \rightarrow \infty$ :

# Maximum likelihood estimation

Maximum likelihood: maximise  $\mathbb{E}_{p^*(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x})]$

Nice properties as  $n \rightarrow \infty$ :

- consistency:

$$p\left(\left|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\right| \geq \epsilon\right) \rightarrow 0$$

# Maximum likelihood estimation

Maximum likelihood: maximise  $\mathbb{E}_{p^*(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x})]$

Nice properties as  $n \rightarrow \infty$ :

- consistency:

$$p\left(\left|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\right| \geq \epsilon\right) \rightarrow 0$$

- asymptotic efficiency and normality (conf. bounds)

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}(0, I^{-1}) \quad I = -\mathbb{E}_p [\nabla_{\boldsymbol{\theta}}^2 \log p_{\boldsymbol{\theta}}(\mathbf{x}) |_{\boldsymbol{\theta}^*}]$$

# Maximum likelihood estimation

Maximum likelihood: maximise  $\mathbb{E}_{p^*(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x})]$

Nice properties as  $n \rightarrow \infty$ :

- consistency:

$$p\left(\left|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\right| \geq \epsilon\right) \rightarrow 0$$

- asymptotic efficiency and normality (conf. bounds)

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \xrightarrow{d} \mathcal{N}(0, I^{-1}) \quad I = -\mathbb{E}_p [\nabla_{\boldsymbol{\theta}}^2 \log p_{\boldsymbol{\theta}}(\mathbf{x}) |_{\boldsymbol{\theta}^*}]$$

Can optimise by gradient method, but requires an intractable integral

$$\Delta_{\boldsymbol{\theta}}(\mathbf{x}) := \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}) \quad \log p_{\boldsymbol{\theta}}(\mathbf{x}) = \int \log p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x}) d\mathbf{z}$$

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)] + \mathbb{H}[q] = \log p_{\theta}(x) - \text{KL}[q(z) \| p_{\theta}(z|x)] \leq \log p_{\theta}(x)$$

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)] + \mathbb{H}[q] = \log p_{\theta}(x) - \text{KL}[q(z) \| p_{\theta}(z|x)] \leq \log p_{\theta}(x)$$

Bound is tight for  $q(z) = p_{\theta}(z|x)$ .

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)] + \mathbb{H}[q] = \log p_{\theta}(x) - \text{KL}[q(z) \| p_{\theta}(z|x)] \leq \log p_{\theta}(x)$$

Bound is tight for  $q(z) = p_{\theta}(z|x)$ .

To keep maximising  $\mathcal{F}$ , when  $\theta = \theta_t$  at iteration  $t$ :

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)] + \mathbb{H}[q] = \log p_{\theta}(x) - \text{KL}[q(z) \| p_{\theta}(z|x)] \leq \log p_{\theta}(x)$$

Bound is tight for  $q(z) = p_{\theta}(z|x)$ .

To keep maximising  $\mathcal{F}$ , when  $\theta = \theta_t$  at iteration  $t$ :

- E-step: find  $q(z)$  **for each**  $x^* \sim p^*$ :

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)] + \mathbb{H}[q] = \log p_{\theta}(x) - \text{KL}[q(z) \| p_{\theta}(z|x)] \leq \log p_{\theta}(x)$$

Bound is tight for  $q(z) = p_{\theta}(z|x)$ .

To keep maximising  $\mathcal{F}$ , when  $\theta = \theta_t$  at iteration  $t$ :

- E-step: find  $q(z)$  **for each**  $x^* \sim p^*$ :

$$\arg \min_{q \in \mathcal{Q}} \text{KL}[q(z) \| p_{\theta_t}(z|x^*)]$$

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(\mathbf{z})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] + \mathbb{H}[q] = \log p_{\theta}(\mathbf{x}) - \text{KL}[q(\mathbf{z}) \| p_{\theta}(\mathbf{z}|\mathbf{x})] \leq \log p_{\theta}(\mathbf{x})$$

Bound is tight for  $q(\mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x})$ .

To keep maximising  $\mathcal{F}$ , when  $\theta = \theta_t$  at iteration  $t$ :

- E-step: find  $q(\mathbf{z})$  **for each**  $\mathbf{x}^* \sim p^*$ :

$$\arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{z}) \| p_{\theta_t}(\mathbf{z}|\mathbf{x}^*)] \quad \text{e.g.} \quad \arg \min_{\mu, \sigma} \text{KL}[\mathcal{N}(\mu, \Sigma) \| p_{\theta_t}(\mathbf{z}|\mathbf{x}^*)]$$

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(\mathbf{z})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] + \mathbb{H}[q] = \log p_{\theta}(\mathbf{x}) - \text{KL}[q(\mathbf{z}) \| p_{\theta}(\mathbf{z}|\mathbf{x})] \leq \log p_{\theta}(\mathbf{x})$$

Bound is tight for  $q(\mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x})$ .

To keep maximising  $\mathcal{F}$ , when  $\theta = \theta_t$  at iteration  $t$ :

- E-step: find  $q(\mathbf{z})$  **for each**  $\mathbf{x}^* \sim p^*$ :

$$\arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{z}) \| p_{\theta_t}(\mathbf{z}|\mathbf{x}^*)] \quad \text{e.g.} \quad \arg \min_{\mu, \sigma} \text{KL}[\mathcal{N}(\mu, \Sigma) \| p_{\theta_t}(\mathbf{z}|\mathbf{x}^*)]$$

- M-step: change  $\theta$  by

$$\nabla_{\theta} \mathcal{F}(q, \theta) = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})]$$

# Expectation-Maximisation: variational inference for learning

Instead, use the lower bound to derive the EM algorithm

$$\mathcal{F}(q, \theta) := \mathbb{E}_{q(\mathbf{z})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] + \mathbb{H}[q] = \log p_{\theta}(\mathbf{x}) - \text{KL}[q(\mathbf{z}) \| p_{\theta}(\mathbf{z}|\mathbf{x})] \leq \log p_{\theta}(\mathbf{x})$$

Bound is tight for  $q(\mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x})$ .

To keep maximising  $\mathcal{F}$ , when  $\theta = \theta_t$  at iteration  $t$ :

- E-step: find  $q(\mathbf{z})$  **for each**  $\mathbf{x}^* \sim p^*$ :

$$\arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{z}) \| p_{\theta_t}(\mathbf{z}|\mathbf{x}^*)] \quad \text{e.g.} \quad \arg \min_{\mu, \sigma} \text{KL}[\mathcal{N}(\mu, \Sigma) \| p_{\theta_t}(\mathbf{z}|\mathbf{x}^*)]$$

- M-step: change  $\theta$  by

$$\nabla_{\theta} \mathcal{F}(q, \theta) = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})] = \nabla_{\theta} \mathbb{E}_{q(\mathbf{z})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)]$$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL} [\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL} [\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL} [\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep
- Pros:

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL}[q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL}[\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

- scale to large  $x$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL}[q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL}[\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

- scale to large  $x$
  - allows flexible  $q(z)$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL} [\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

- scale to large  $x$
- allows flexible  $q(z)$
- complicated  $p_{\theta}(z, x)$  (?)

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL} [\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

- scale to large  $x$
- allows flexible  $q(z)$
- complicated  $p_{\theta}(z, x)$  (?)

- But:

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL}[q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL}[\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

- scale to large  $x$
- allows flexible  $q(z)$
- complicated  $p_{\theta}(z, x)$  (?)

- But:

- $q \stackrel{\text{KL}}{\approx} p_{\theta}(z|x) \not\Rightarrow \nabla_{\theta} \mathcal{F}(q, \theta) \stackrel{l-2}{\approx} \nabla_{\theta} \log p_{\theta}(x)$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL} [q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL} [\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

- scale to large  $x$
- allows flexible  $q(z)$
- complicated  $p_{\theta}(z, x)$  (?)

- But:

- $q \stackrel{\text{KL}}{\approx} p_{\theta}(z|x) \not\Rightarrow \nabla_{\theta} \mathcal{F}(q, \theta) \stackrel{l-2}{\approx} \nabla_{\theta} \log p_{\theta}(x)$
- reparameterisation depends on  $\mathcal{Q}$

# Amortised inference for learning

- Approximate  $\nabla_{\theta} \mathcal{F}(q, \theta) = \nabla_{\theta} \mathbb{E}_{q(z)} [\log p_{\theta}(z, x)]$  by Monte Carlo – **need samples!**
- Inferential model with parameter  $\phi$ : find  $\phi$  **for all**  $x^* \sim p^*$

$$\arg \min_{\phi} \text{KL}[q_{\phi}(z|x^*) \| p_{\theta}(z|x^*)] \quad \text{e.g.} \quad \arg \min_{\phi} \text{KL}[\mathcal{N}(\mu_{\phi}(x^*), \Sigma_{\phi}(x^*)) \| p_{\theta}(z|x^*)]$$

- Gradient w.r.t.  $\phi$ : reparameterisation trick, score trick, wake-sleep

- Pros:

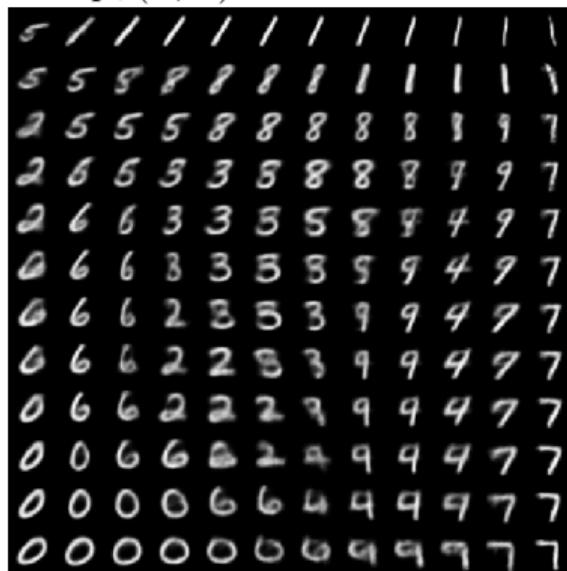
- scale to large  $x$
- allows flexible  $q(z)$
- complicated  $p_{\theta}(z, x)$  (?)

- But:

- $q \stackrel{\text{KL}}{\approx} p_{\theta}(z|x) \not\Rightarrow \nabla_{\theta} \mathcal{F}(q, \theta) \stackrel{l-2}{\approx} \nabla_{\theta} \log p_{\theta}(x)$
- reparameterisation depends on  $\mathcal{Q}$
- flexible  $p(z, x)$  can induce complex  $p_{\theta}(z|x)$

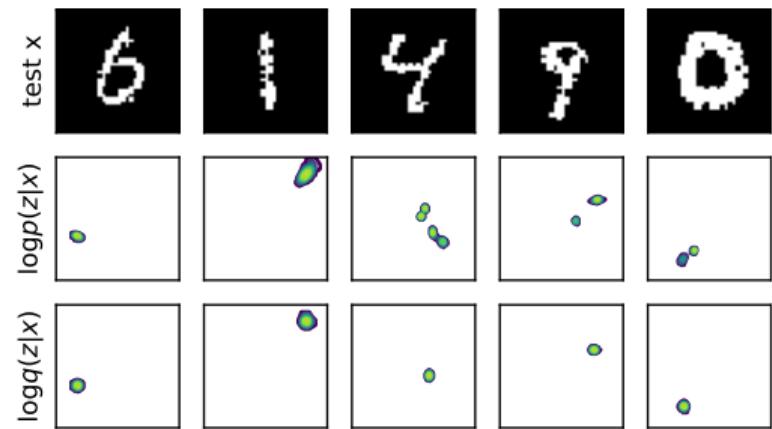
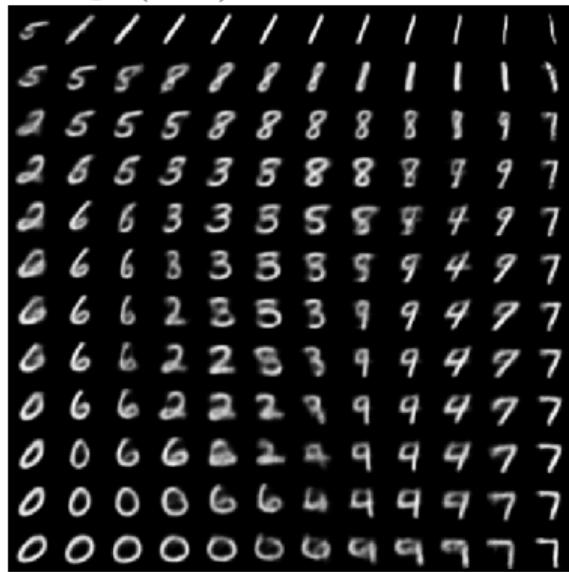
## Example: complex posterior in Gaussian VAE

$p_{\theta}(z, x)$ , where  $z \in \mathbb{R}^2$



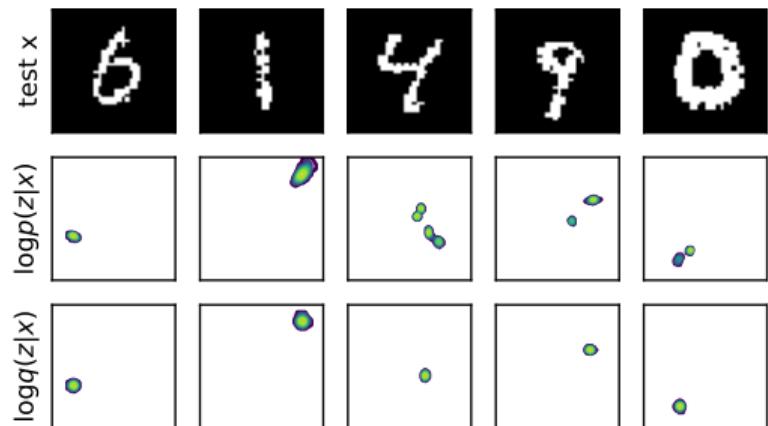
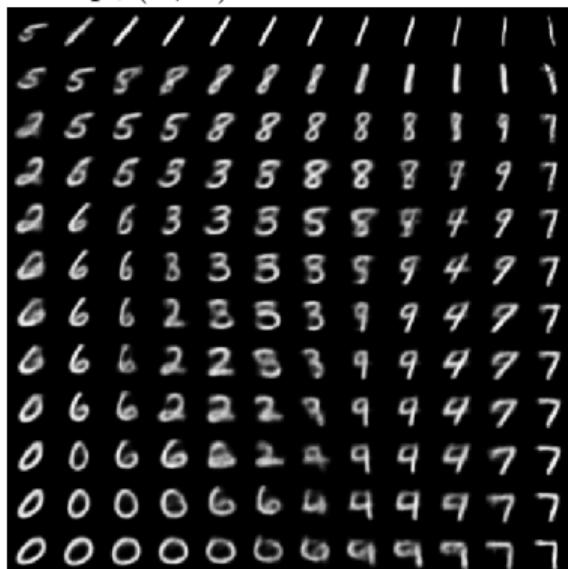
## Example: complex posterior in Gaussian VAE

$p_{\theta}(z, x)$ , where  $z \in \mathbb{R}^2$



## Example: complex posterior in Gaussian VAE

$p_{\theta}(z, x)$ , where  $z \in \mathbb{R}^2$



This talk: maximum likelihood learning without inference

# Amortised learning: focus on ML gradient

Suppose we are at iteration  $t$  with  $\theta = \theta_t$ , want

$$\Delta_{\theta_t}(\mathbf{x}) := \nabla_{\theta} \log p_{\theta}(\mathbf{x})|_{\theta_t}$$

## Amortised learning: focus on ML gradient

Suppose we are at iteration  $t$  with  $\theta = \theta_t$ , want

$$\Delta_{\theta_t}(\mathbf{x}) := \nabla_{\theta} \log p_{\theta}(\mathbf{x})|_{\theta_t}$$

using an approximate  $q$

## Amortised learning: focus on ML gradient

Suppose we are at iteration  $t$  with  $\theta = \theta_t$ , want

$$\Delta_{\theta_t}(\mathbf{x}) := \nabla_{\theta} \log p_{\theta}(\mathbf{x})|_{\theta_t}$$

using an approximate  $q$

$$\nabla_{\theta} \mathcal{F}(q, \theta)$$

## Amortised learning: focus on ML gradient

Suppose we are at iteration  $t$  with  $\theta = \theta_t$ , want

$$\Delta_{\theta_t}(\mathbf{x}) := \nabla_{\theta} \log p_{\theta}(\mathbf{x})|_{\theta_t}$$

using an approximate  $q$

$$\nabla_{\theta} \mathcal{F}(q, \theta) = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})]$$

# Amortised learning: focus on ML gradient

Suppose we are at iteration  $t$  with  $\theta = \theta_t$ , want

$$\Delta_{\theta_t}(\mathbf{x}) := \nabla_{\theta} \log p_{\theta}(\mathbf{x})|_{\theta_t}$$

using an approximate  $q$

$$\nabla_{\theta} \mathcal{F}(q, \theta) = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})]$$

When  $q(\mathbf{z}) = p_{\theta_t}(\mathbf{z}|\mathbf{x})$ , can show

$$\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathcal{F}(p_{\theta_t}(\mathbf{z}|\mathbf{x}), \theta)|_{\theta_t} = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})|_{\theta_t}]$$

# Amortised learning: focus on ML gradient

Suppose we are at iteration  $t$  with  $\theta = \theta_t$ , want

$$\Delta_{\theta_t}(\mathbf{x}) := \nabla_{\theta} \log p_{\theta}(\mathbf{x})|_{\theta_t}$$

using an approximate  $q$

$$\nabla_{\theta} \mathcal{F}(q, \theta) = \mathbb{E}_{q(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})]$$

When  $q(\mathbf{z}) = p_{\theta_t}(\mathbf{z}|\mathbf{x})$ , can show

$$\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathcal{F}(p_{\theta_t}(\mathbf{z}|\mathbf{x}), \theta)|_{\theta_t} = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})|_{\theta_t}]$$

**Amortised learning:** directly approximate  $\Delta_{\theta_t}(\mathbf{x})$  above

## Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x}) | \theta_t]$

## Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x}) | \boldsymbol{\theta}_t]$
- let  $\mathbf{y}_n = \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{z}_n, \mathbf{x}_n) | \boldsymbol{\theta}_t$

## Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})|_{\theta_t}]$
- let  $\mathbf{y}_n = \nabla_{\theta} \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)|_{\theta_t}$
- problem: given  $p(\mathbf{x}, \mathbf{y})$ , want  $f: \mathbf{x} \mapsto \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} [\mathbf{y}]$

# Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x}) | \theta_t]$
- let  $\mathbf{y}_n = \nabla_{\theta} \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n) | \theta_t$
- problem: given  $p(\mathbf{x}, \mathbf{y})$ , want  $f: \mathbf{x} \mapsto \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} [\mathbf{y}]$
- Estimate by LSR:

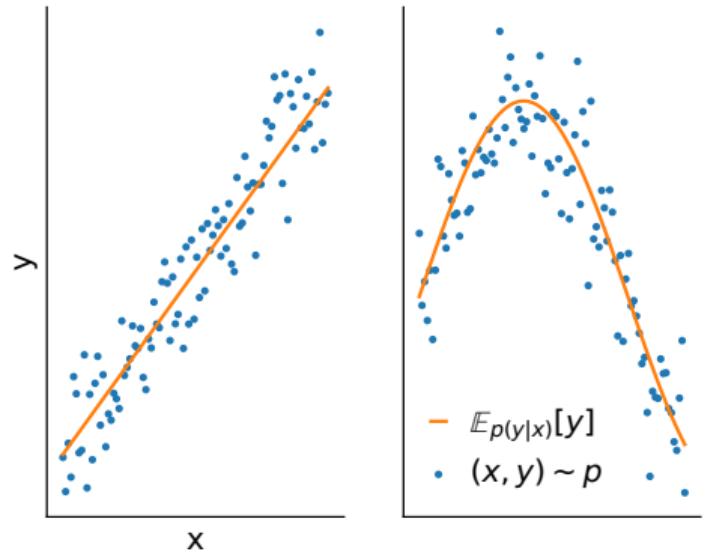
$$\min_{f \in \mathcal{H}_{\gamma}} \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{y})} [\|\mathbf{y} - f(\mathbf{x})\|_2^2]$$

# Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x}) | \theta_t]$
- let  $\mathbf{y}_n = \nabla_{\theta} \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n) | \theta_t$
- problem: given  $p(\mathbf{x}, \mathbf{y})$ , want  $f: \mathbf{x} \mapsto \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} [\mathbf{y}]$
- Estimate by LSR:

$$\min_{f \in \mathcal{H}_{\gamma}} \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{y})} [\|\mathbf{y} - f(\mathbf{x})\|_2^2]$$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p$$



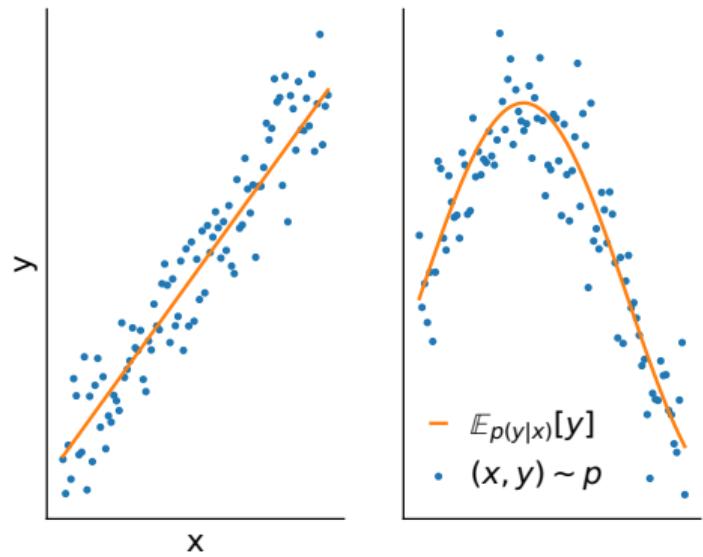
# Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x}) | \theta_t]$
- let  $\mathbf{y}_n = \nabla_{\theta} \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n) | \theta_t$
- problem: given  $p(\mathbf{x}, \mathbf{y})$ , want  $f : \mathbf{x} \mapsto \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} [\mathbf{y}]$
- Estimate by LSR:

$$\min_{f \in \mathcal{H}_{\gamma}} \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{y})} [\|\mathbf{y} - f(\mathbf{x})\|_2^2]$$

- $\hat{\Delta}_{\theta_t, \gamma}(\mathbf{x}) = \hat{f}_{\gamma}(\mathbf{x})$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p$$



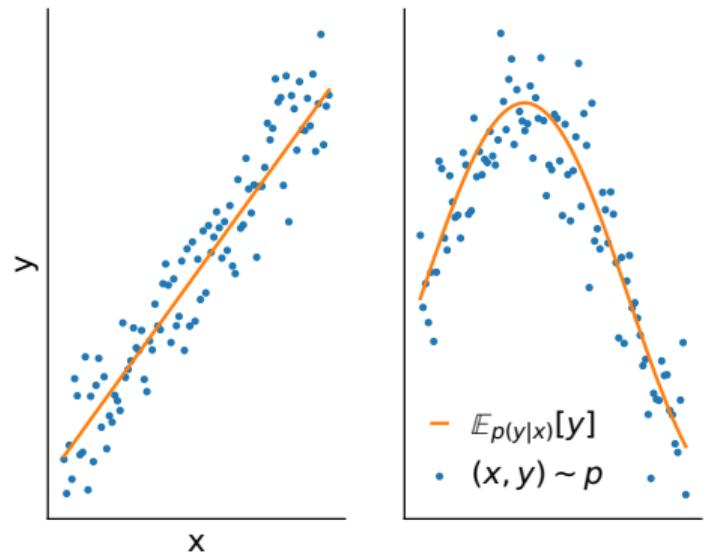
# Posterior mean from least square regression (LSR)

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x}) | \theta_t]$
- let  $\mathbf{y}_n = \nabla_{\theta} \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n) | \theta_t$
- problem: given  $p(\mathbf{x}, \mathbf{y})$ , want  $f : \mathbf{x} \mapsto \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} [\mathbf{y}]$
- Estimate by LSR:

$$\min_{f \in \mathcal{H}_{\gamma}} \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{y})} [\|\mathbf{y} - f(\mathbf{x})\|_2^2]$$

- $\hat{\Delta}_{\theta_t, \gamma}(\mathbf{x}) = \hat{f}_{\gamma}(\mathbf{x})$
- But requires  $\mathbf{y}_n = \nabla_{\theta} \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n) | \theta_t$  for all  $n$

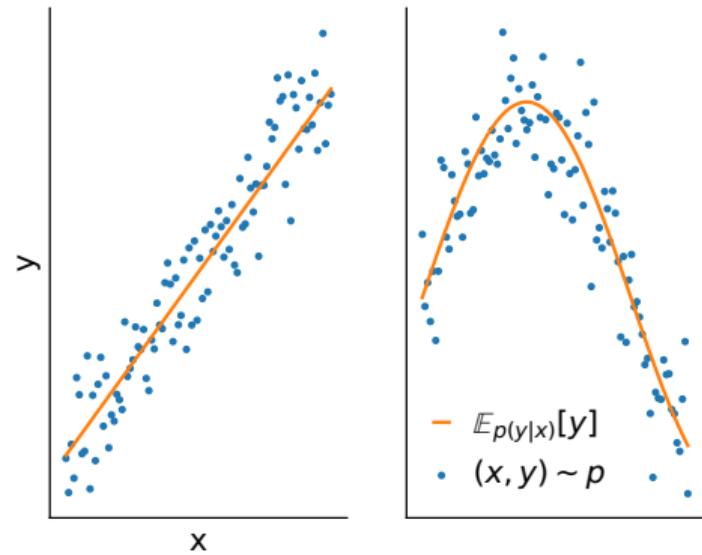
$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p$$



# Posterior mean from least square regression

- want  $\Delta_{\theta_t}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x}) | \theta_t]$

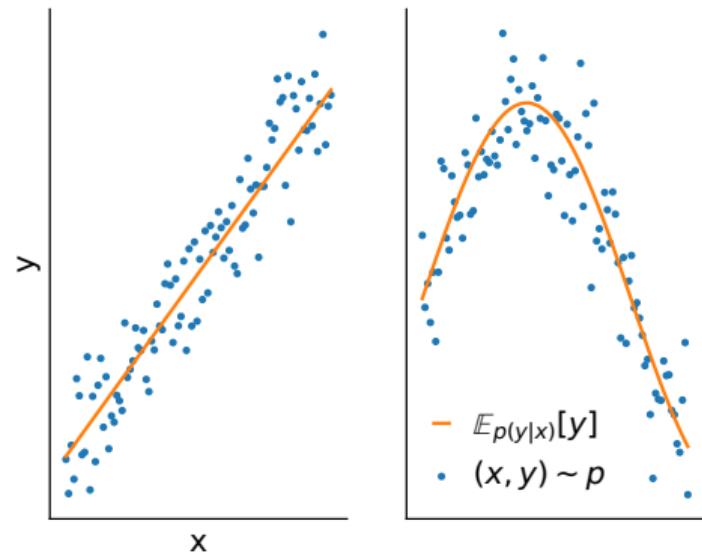
$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$



# Posterior mean from least square regression

- want  $\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] |_{\theta_t}$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$

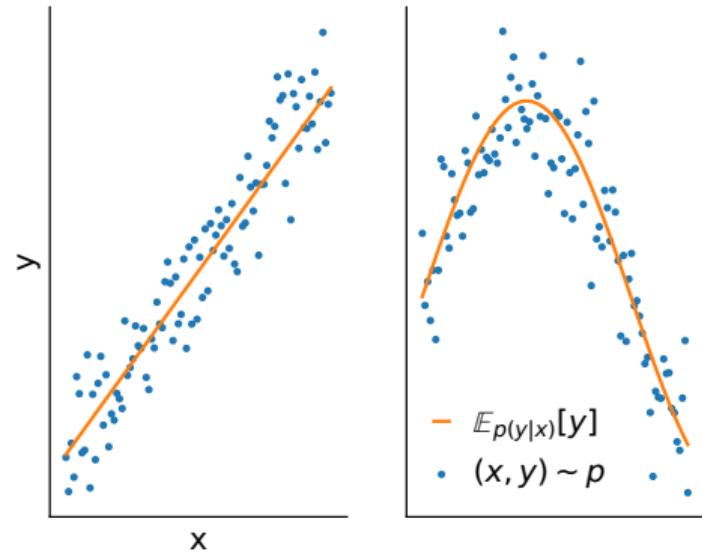


# Posterior mean from least square regression

- want  $\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] |_{\theta_t}$
- estimate  $y_{\theta,n} = \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$  and differentiate

$$\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(y, \mathbf{x})} [y_{\theta}] |_{\theta_t}$$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$

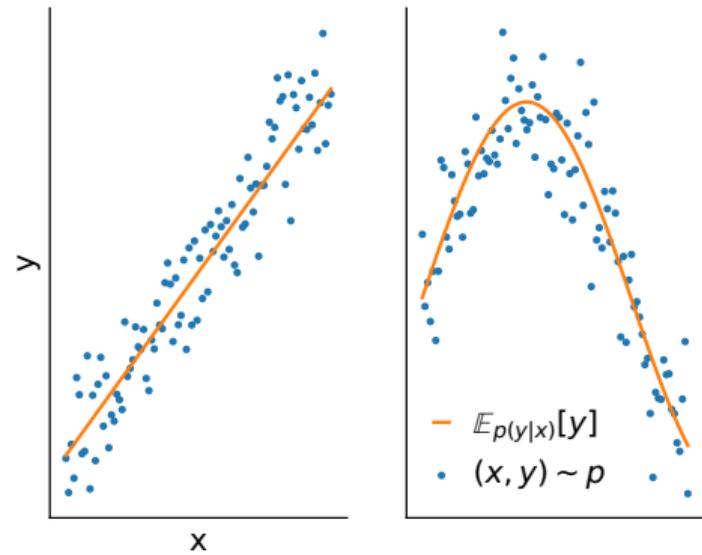


# Posterior mean from least square regression

- want  $\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] |_{\theta_t}$
- estimate  $y_{\theta,n} = \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$  and differentiate

$$\begin{aligned}\Delta_{\theta_t}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(y, \mathbf{x})} [y_{\theta}] |_{\theta_t} \\ &= \nabla_{\theta} J_{\theta}(\mathbf{x}) |_{\theta_t}\end{aligned}$$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \| \mathbf{y}_n - f(\mathbf{x}_n) \|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$



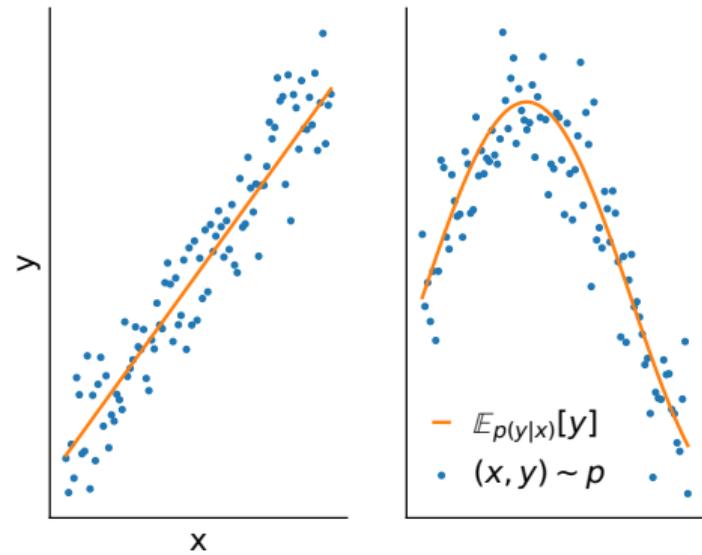
# Posterior mean from least square regression

- want  $\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})] |_{\theta_t}$
- estimate  $y_{\theta,n} = \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$  and differentiate

$$\begin{aligned}\Delta_{\theta_t}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(y, \mathbf{x})} [y_{\theta}] |_{\theta_t} \\ &= \nabla_{\theta} J_{\theta}(\mathbf{x}) |_{\theta_t}\end{aligned}$$

where  $J_{\theta}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - f(\mathbf{x}_n)\|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$



# Posterior mean from least square regression

- want  $\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(z|x)} [\log p_{\theta}(z, x)] |_{\theta_t}$
- estimate  $y_{\theta,n} = \log p_{\theta}(z_n, x_n)$  and differentiate

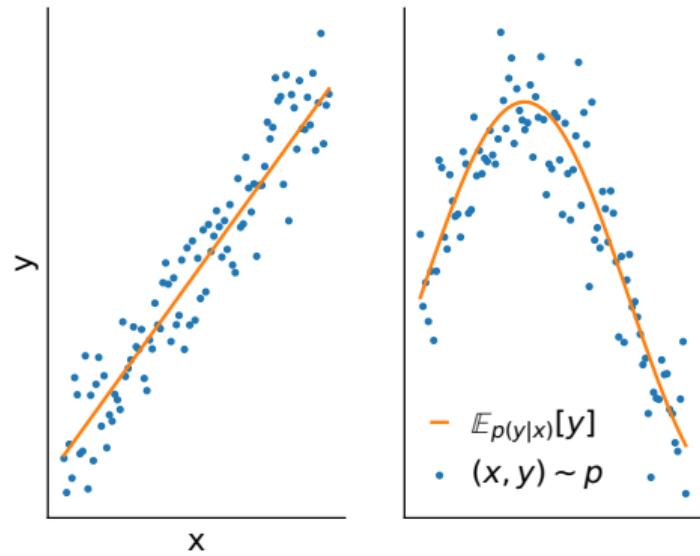
$$\begin{aligned}\Delta_{\theta_t}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(y|x)} [y_{\theta}] |_{\theta_t} \\ &= \nabla_{\theta} J_{\theta}(\mathbf{x}) |_{\theta_t}\end{aligned}$$

where  $J_{\theta}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(z|x)} [\log p_{\theta}(z, x)]$

- If estimator  $\hat{J}_{\theta,\gamma}(\mathbf{x})$ :

- depends on  $\theta$  and  $t$
- differentiable w.r.t  $\theta$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \| \mathbf{y}_n - f(\mathbf{x}_n) \|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$



# Posterior mean from least square regression

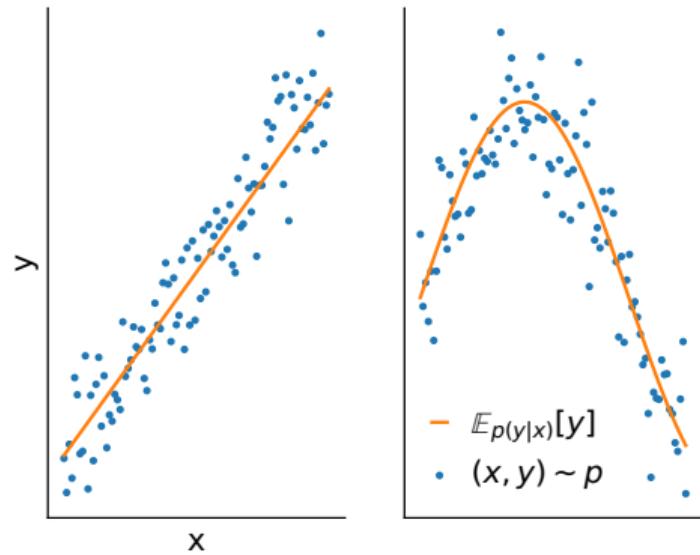
- want  $\Delta_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(z|x)} [\log p_{\theta}(z, x)] |_{\theta_t}$
- estimate  $y_{\theta,n} = \log p_{\theta}(z_n, x_n)$  and differentiate

$$\begin{aligned}\Delta_{\theta_t}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{p_{\theta_t}(y|x)} [y_{\theta}] |_{\theta_t} \\ &= \nabla_{\theta} J_{\theta}(\mathbf{x}) |_{\theta_t}\end{aligned}$$

where  $J_{\theta}(\mathbf{x}) = \mathbb{E}_{p_{\theta_t}(z|x)} [\log p_{\theta}(z, x)]$

- If estimator  $\hat{J}_{\theta,\gamma}(\mathbf{x})$ :
  - depends on  $\theta$  and  $t$
  - differentiable w.r.t  $\theta$
- Then  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta}(\mathbf{x}) |_{\theta_t}$

$$\min_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \| \mathbf{y}_n - f(\mathbf{x}_n) \|_2^2, \quad \mathbf{x}_n, \mathbf{y}_n \sim p_{\theta_t}$$



## Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x}) \right\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathcal{H}}^2$$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \|y_{\theta,n} - \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathcal{H}}^2$$

- Particles:

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x}) \right\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathcal{H}}^2$$

- Particles:

- Let  $\mathbf{z}' = S_{\gamma}(\mathbf{x}, \epsilon)$ , where  $\epsilon$  is a noise source

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x}) \right\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathcal{H}}^2$$

- Particles:

- Let  $\mathbf{z}' = S_\gamma(\mathbf{x}, \epsilon)$ , where  $\epsilon$  is a noise source
- train  $\gamma$  so that  $\mathbb{E}_{\mathbf{z}'|\mathbf{x}} [\log p_{\theta_t}(\mathbf{z}', \mathbf{x})] = \hat{J}_{\theta,\gamma}(\mathbf{x})$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x}) \right\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathcal{H}}^2$$

- Particles:

- Let  $\mathbf{z}' = S_\gamma(\mathbf{x}, \epsilon)$ , where  $\epsilon$  is a noise source
- train  $\gamma$  so that  $\mathbb{E}_{\mathbf{z}'|\mathbf{x}} [\log p_{\theta_t}(\mathbf{z}', \mathbf{x})] = \hat{J}_{\theta,\gamma}(\mathbf{x})$

$$\min_{\gamma} \mathbb{E}_{p_{\theta_t}(\mathbf{z}, \mathbf{x})} \left[ y_{\theta,n} - \mathbb{E}_{\mathbf{z}'|\mathbf{x}} [\log p_{\theta_t}(\mathbf{z}', \mathbf{x})] \right]$$

# Amortised learning

Find an estimator for  $\hat{J}_{\theta,\gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$ ,  $y_\theta = \log p_\theta(\mathbf{z}, \mathbf{x})$ , so that  $\hat{\Delta}_{\theta_t}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$

- NN regression:  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \text{NN}_{\gamma(\theta)}(\mathbf{x})$

$$\gamma_t(\theta) \leftarrow \gamma_{t-1}(\theta) + \nabla_{\theta} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \hat{J}_{\theta,\gamma}(\mathbf{x}_n) \right\|_2^2$$

- Assume a flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \sum_k \alpha_k \psi_i(\mathbf{x}) = \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x})$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \left\| y_{\theta,n} - \boldsymbol{\alpha} \cdot \boldsymbol{\psi}(\mathbf{x}) \right\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_{\mathcal{H}}^2$$

- Particles:

- Let  $\mathbf{z}' = S_\gamma(\mathbf{x}, \epsilon)$ , where  $\epsilon$  is a noise source
- train  $\gamma$  so that  $\mathbb{E}_{\mathbf{z}'|\mathbf{x}} [\log p_{\theta_t}(\mathbf{z}', \mathbf{x})] = \hat{J}_{\theta,\gamma}(\mathbf{x})$

$$\min_{\gamma} \mathbb{E}_{p_{\theta_t}(\mathbf{z}, \mathbf{x})} \left[ y_{\theta,n} - \mathbb{E}_{\mathbf{z}'|\mathbf{x}} [\log p_{\theta_t}(\mathbf{z}', \mathbf{x})] \right]$$

- But no guarantee with  $\nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x})|_{\theta_t}$  as trained on  $\hat{J}_{\theta,\gamma}(\mathbf{x})$

# Amortised learning with nonlinear regression

Flexible  $\hat{J}_{\theta, \gamma}(\mathbf{x}) = \boldsymbol{\alpha} \cdot \psi(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$

# Amortised learning with nonlinear regression

Flexible  $\hat{J}_{\theta,\gamma}(x) = \alpha \cdot \psi(x) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|x)} [y_\theta]$

$$\min_{\alpha} \frac{1}{N} \sum_{n=1}^N \|\log p_{\theta}(\mathbf{z}, \mathbf{x}) - \alpha \cdot \psi(\mathbf{x})\|_2^2 + \lambda \|\alpha\|_2^2, \quad \mathbf{x}_n, \mathbf{z}_n \sim p_{\theta_t}$$

# Amortised learning with nonlinear regression

Flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha} \cdot \psi(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \|\log p_{\theta}(\mathbf{z}, \mathbf{x}) - \boldsymbol{\alpha} \cdot \psi(\mathbf{x})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2, \quad \mathbf{x}_n, \mathbf{z}_n \sim p_{\theta_t}$$

- Given  $\mathbf{x}^* \in p^*$ , the estimator outputs

$$\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha}_{\theta,\gamma} \psi(\mathbf{x}^*)$$

# Amortised learning with nonlinear regression

Flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha} \cdot \psi(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \|\log p_{\theta}(\mathbf{z}, \mathbf{x}) - \boldsymbol{\alpha} \cdot \psi(\mathbf{x})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2, \quad \mathbf{x}_n, \mathbf{z}_n \sim p_{\theta_t}$$

- Given  $\mathbf{x}^* \in p^*$ , the estimator outputs

$$\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha}_{\theta,\gamma} \psi(\mathbf{x}^*)$$

where

$$\boldsymbol{\alpha}_{\theta,\gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1} \quad \mathbf{Y}_{\theta} = [\log p_{\theta}(\mathbf{z}_1, \mathbf{x}_1), \dots, \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)]^{\top} \quad \boldsymbol{\Psi} = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_n)]$$

# Amortised learning with nonlinear regression

Flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha} \cdot \psi(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \|\log p_{\theta}(\mathbf{z}, \mathbf{x}) - \boldsymbol{\alpha} \cdot \psi(\mathbf{x})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2, \quad \mathbf{x}_n, \mathbf{z}_n \sim p_{\theta_t}$$

- Given  $\mathbf{x}^* \in p^*$ , the estimator outputs

$$\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha}_{\theta,\gamma} \psi(\mathbf{x}^*)$$

where

$$\boldsymbol{\alpha}_{\theta,\gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1} \quad \mathbf{Y}_{\theta} = [\log p_{\theta}(\mathbf{z}_1, \mathbf{x}_1), \dots, \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)]^T \quad \boldsymbol{\Psi} = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_n)]$$

- $\hat{\Delta}_{\theta,\gamma}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x}) = \nabla_{\theta} \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1} \psi(\mathbf{x}^*)$

# Amortised learning with nonlinear regression

Flexible  $\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha} \cdot \psi(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [y_\theta]$

$$\min_{\boldsymbol{\alpha}} \frac{1}{N} \sum_{n=1}^N \|\log p_{\theta}(\mathbf{z}, \mathbf{x}) - \boldsymbol{\alpha} \cdot \psi(\mathbf{x})\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2, \quad \mathbf{x}_n, \mathbf{z}_n \sim p_{\theta_t}$$

- Given  $\mathbf{x}^* \in p^*$ , the estimator outputs

$$\hat{J}_{\theta,\gamma}(\mathbf{x}) = \boldsymbol{\alpha}_{\theta,\gamma} \psi(\mathbf{x}^*)$$

where

$$\boldsymbol{\alpha}_{\theta,\gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1} \quad \mathbf{Y}_{\theta} = [\log p_{\theta}(\mathbf{z}_1, \mathbf{x}_1), \dots, \log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)]^T \quad \boldsymbol{\Psi} = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_n)]$$

- $\hat{\Delta}_{\theta,\gamma}(\mathbf{x}) = \nabla_{\theta} \hat{J}_{\theta,\gamma}(\mathbf{x}) = \nabla_{\theta} \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1} \psi(\mathbf{x}^*)$   
same as if we regressed to  $\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})$  instead of  $\log p_{\theta}(\mathbf{z}, \mathbf{x})$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(x) \approx \mathbb{E}_{p_{\theta_t}(z|x)} [\log p_{\theta}(z, x)]$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(z|x)} [\log p_{\theta}(z, x)]$ 
  - $z_n, x_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(z_n, x_n)$
  - $\alpha_{\theta, \gamma} = \mathbf{Y}_{\theta} \Psi^{\top} (\Psi \Psi^{\top} + \lambda \mathbf{I})^{-1}$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^T (\boldsymbol{\Psi} \boldsymbol{\Psi}^T + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

- easy to sample from  $p_{\theta_t}(\mathbf{z}, \mathbf{x})$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

- easy to sample from  $p_{\theta_t}(\mathbf{z}, \mathbf{x})$
- can evaluate  $\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})$

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

- easy to sample from  $p_{\theta_t}(\mathbf{z}, \mathbf{x})$
- can evaluate  $\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})$

Non-assumptions:

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

- easy to sample from  $p_{\theta_t}(\mathbf{z}, \mathbf{x})$
- can evaluate  $\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})$

Non-assumptions:

- posterior

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

- easy to sample from  $p_{\theta_t}(\mathbf{z}, \mathbf{x})$
- can evaluate  $\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})$

Non-assumptions:

- **posterior**
- **support or domain** of latent

# Amortised learning by wake sleep (ALWS)

- Sleep phase, regression to find  $\hat{J}_{\theta, \gamma}(\mathbf{x}) \approx \mathbb{E}_{p_{\theta_t}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\mathbf{z}_n, \mathbf{x}_n \sim p_{\theta_t}$
  - evaluate  $\log p_{\theta}(\mathbf{z}_n, \mathbf{x}_n)$
  - $\boldsymbol{\alpha}_{\theta, \gamma} = \mathbf{Y}_{\theta} \boldsymbol{\Psi}^{\top} (\boldsymbol{\Psi} \boldsymbol{\Psi}^{\top} + \lambda \mathbf{I})^{-1}$
- Wake phase: update  $\theta$  according to  $\nabla \hat{J}_{\theta, \gamma}(\mathbf{x}^*) \approx \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^*)} [\log p_{\theta}(\mathbf{z}, \mathbf{x})]$ 
  - $\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*) = \nabla_{\theta} \hat{J}_{\theta, \gamma}(\mathbf{x}^*)$
  - $\theta_{t+1} = \theta_t + \text{Adam} [\hat{\Delta}_{\theta, \gamma}(\mathbf{x}^*)]$

Assumptions:

- easy to sample from  $p_{\theta_t}(\mathbf{z}, \mathbf{x})$
- can evaluate  $\nabla_{\theta} \log p_{\theta}(\mathbf{z}, \mathbf{x})$

Non-assumptions:

- **posterior**
- **support or domain** of latent
- **structure** of  $p_{\theta}(\mathbf{z}, \mathbf{x})$

# Experiment I: gradient estimation

Generative model

$$z_1, z_2 \sim \mathcal{N}(0, 1), \quad x|z \sim \mathcal{N}(\text{softplus}(\mathbf{b} \cdot \mathbf{z}) - \|\mathbf{b}\|_2^2, \sigma_x^2)$$

# Experiment I: gradient estimation

Generative model

$$z_1, z_2 \sim \mathcal{N}(0, 1), \quad x|z \sim \mathcal{N}(\text{softplus}(\mathbf{b} \cdot \mathbf{z}) - \|\mathbf{b}\|_2^2, \sigma_x^2)$$

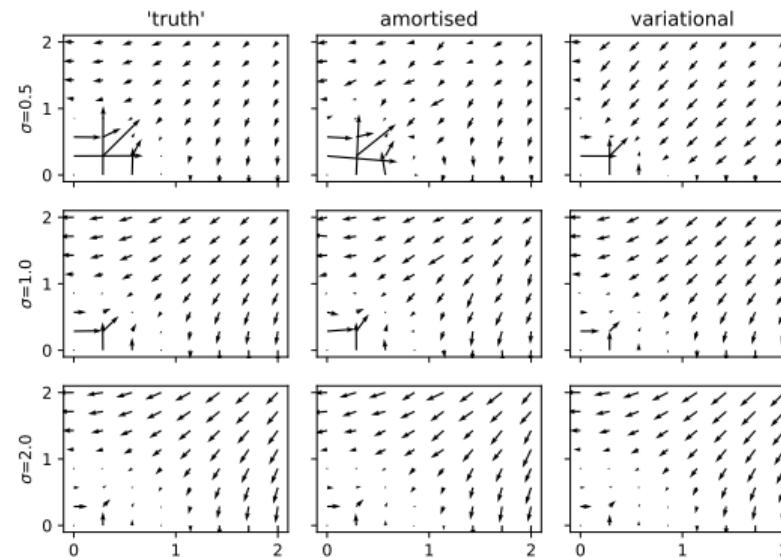
Task: estimate  $\nabla_{\mathbf{b}} \log p_{\theta}(\mathbf{x})$  for different  $\mathbf{b}$  and  $\sigma$

# Experiment I: gradient estimation

Generative model

$$z_1, z_2 \sim \mathcal{N}(0, 1), x|z \sim \mathcal{N}(\text{softplus}(\mathbf{b} \cdot \mathbf{z}) - \|\mathbf{b}\|_2^2, \sigma_x^2)$$

Task: estimate  $\nabla_{\mathbf{b}} \log p_{\theta}(\mathbf{x})$  for different  $\mathbf{b}$  and  $\sigma$



## Experiment II: non-Euclidean $\mathbf{z}$

Model:

$$\mathbf{z} = [\cos(a), \sin(a)], \quad p(a) = \mathcal{U}(a|(-\pi, \pi)), \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\text{NN}_{\mathbf{w}}(\mathbf{z}), \sigma_x^2 \mathbf{I})$$

## Experiment II: non-Euclidean $\mathbf{z}$

Model:

$$\mathbf{z} = [\cos(a), \sin(a)], \quad p(a) = \mathcal{U}(a|(-\pi, \pi)), \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\text{NN}_{\mathbf{w}}(\mathbf{z}), \sigma_x^2 \mathbf{I})$$

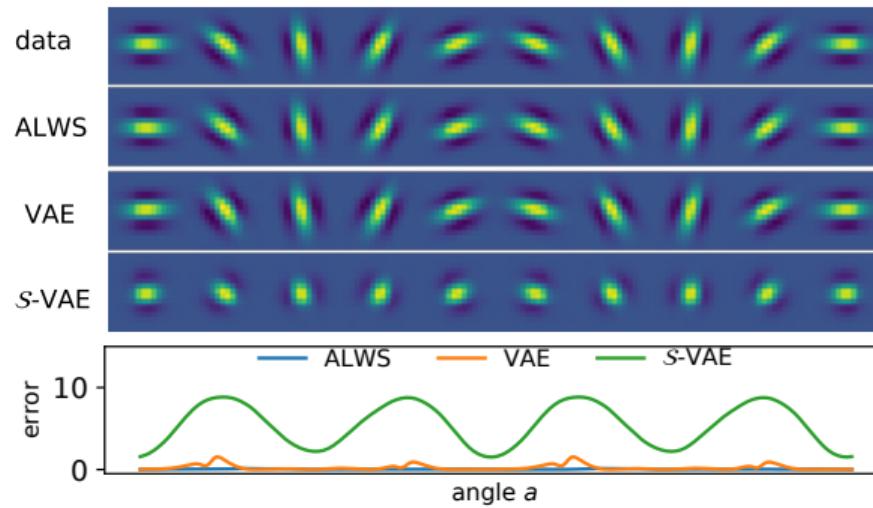
Task: generate Gabor filters of uniformly distributed orientations (no special reparameterisation)

## Experiment II: non-Euclidean $\mathbf{z}$

Model:

$$\mathbf{z} = [\cos(a), \sin(a)], \quad p(a) = \mathcal{U}(a|(-\pi, \pi)), \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\text{NN}_{\mathbf{w}}(\mathbf{z}), \sigma_x^2 \mathbf{I})$$

Task: generate Gabor filters of uniformly distributed orientations (no special reparameterisation)



## Experiment III: hierarchical models

Model:

$$p(\mathbf{z}_1) = \text{Cat}(\mathbf{z}_1 | \mathbf{m})$$

$$p(\mathbf{z}_2 | \mathbf{z}_1 = k) = \mathcal{N}(\mathbf{z}_2 | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{x} | \mathbf{z}_2) = \mathcal{N}(\mathbf{x} | \text{NN}_{\mathbf{w}}(\mathbf{z}_2), \boldsymbol{\Sigma}_x)$$

## Experiment III: hierarchical models

Model:

$$\begin{aligned} p(\boldsymbol{z}_1) &= \text{Cat}(\boldsymbol{z}_1 | \boldsymbol{m}) \\ p(\boldsymbol{z}_2 | \boldsymbol{z}_1 = k) &= \mathcal{N}(\boldsymbol{z}_2 | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ p(\boldsymbol{x} | \boldsymbol{z}_2) &= \mathcal{N}(\boldsymbol{x} | \text{NN}_{\boldsymbol{w}}(\boldsymbol{z}_2), \boldsymbol{\Sigma}_x) \end{aligned}$$

Parameters are penalised according to suitable priors

## Experiment III: hierarchical models

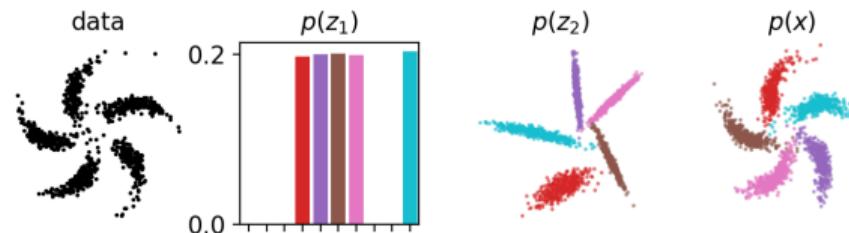
Model:

$$p(z_1) = \text{Cat}(z_1 | \mathbf{m})$$

$$p(z_2 | z_1 = k) = \mathcal{N}(z_2 | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{x} | z_2) = \mathcal{N}(\mathbf{x} | \text{NN}_{\mathbf{w}}(z_2), \boldsymbol{\Sigma}_x)$$

Parameters are penalised according to suitable priors

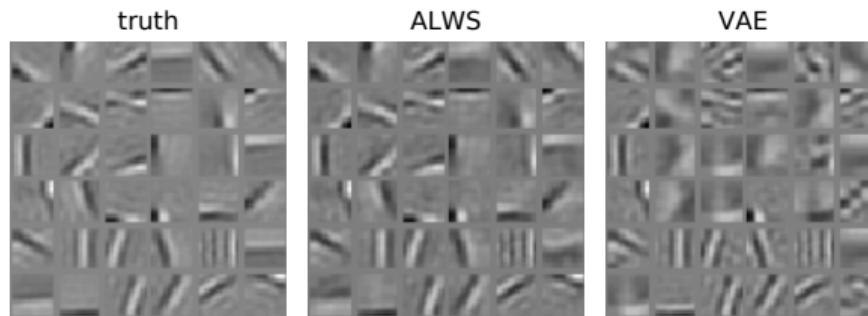


## Experiment IV: feature identification

Linear noisy ICA Model:

$$p(z_i) = \text{Lap}(z_i|0, 1)$$
$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z}, \sigma^2 I)$$

Set up true  $\mathbf{W}$ , recover from data, compare with Laplace-VAE



## Experiment V.1: noisy Hodgkin-Huxley model

$$\begin{aligned}C_m \dot{V}(t) &= -g_l[V(t) - E_l] - \bar{g}_N m^3(t) h(t) [V(t) - E_N] - \bar{g}_K n^4(t) [V(t) - E_K] + I_{in}(t) + \epsilon(t) \\&= \alpha_e(V(t)) [1 - e(t)] - \beta_e(V(t)) e(t), \quad e \in \{m, h, n\}\end{aligned}$$

## Experiment V.1: noisy Hodgkin-Huxley model

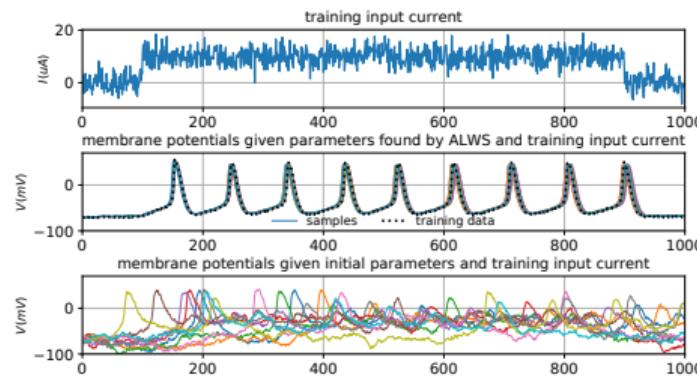
$$\begin{aligned}C_m \dot{V}(t) &= -g_l[V(t) - E_l] - \bar{g}_N m^3(t) h(t) [V(t) - E_N] - \bar{g}_K n^4(t) [V(t) - E_K] + I_{in}(t) + \epsilon(t) \\&= \alpha_e(V(t)) [1 - e(t)] - \beta_e(V(t)) e(t), \quad e \in \{m, h, n\}\end{aligned}$$

Discretised in time with  $\delta t = 0.05ms$

# Experiment V.1: noisy Hodgkin-Huxley model

$$\begin{aligned}C_m \dot{V}(t) &= -g_l[V(t) - E_l] - \bar{g}_N m^3(t)h(t)[V(t) - E_N] - \bar{g}_K n^4(t)[V(t) - E_K] + I_{in}(t) + \epsilon(t) \\&= \alpha_e(V(t))[1 - e(t)] - \beta_e(V(t))e(t), \quad e \in \{m, h, n\}\end{aligned}$$

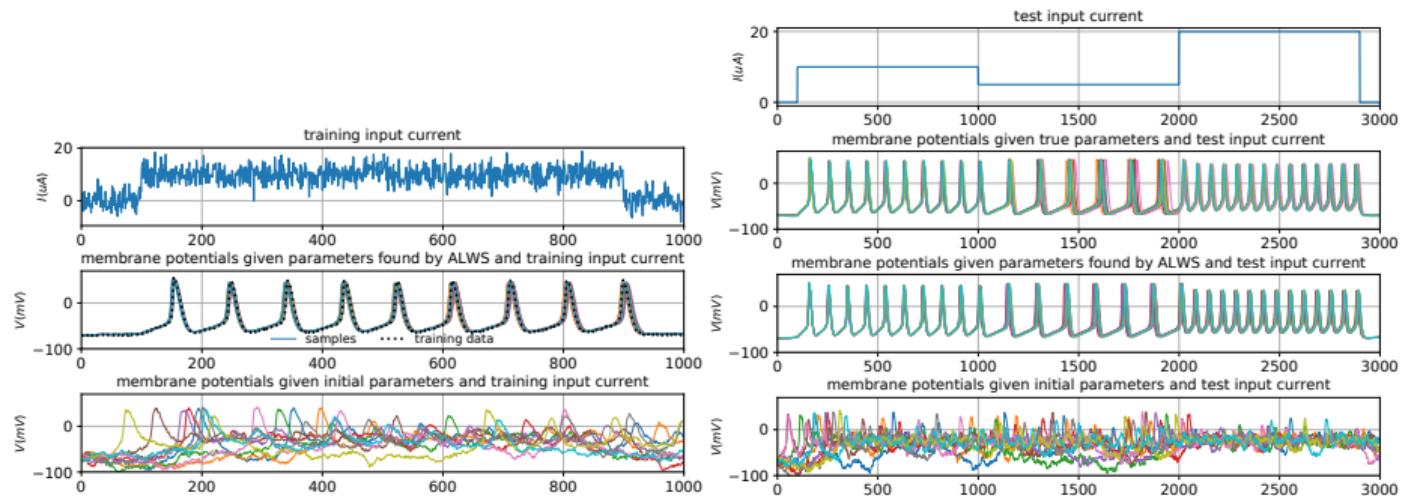
Discretised in time with  $\delta t = 0.05ms$



# Experiment V.1: noisy Hodgkin-Huxley model

$$\begin{aligned}C_m \dot{V}(t) &= -g_l[V(t) - E_l] - \bar{g}_N m^3(t) h(t)[V(t) - E_N] - \bar{g}_K n^4(t)[V(t) - E_K] + I_{in}(t) + \epsilon(t) \\&= \alpha_e(V(t))[1 - e(t)] - \beta_e(V(t))e(t), \quad e \in \{m, h, n\}\end{aligned}$$

Discretised in time with  $\delta t = 0.05ms$

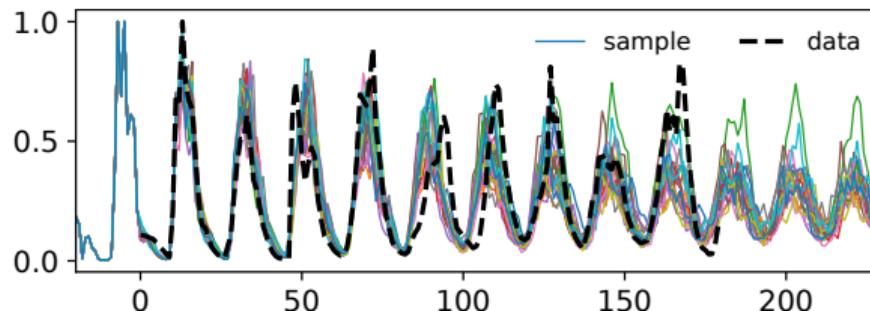


## Experiment V.2: ecology data for blowfly population size

Model

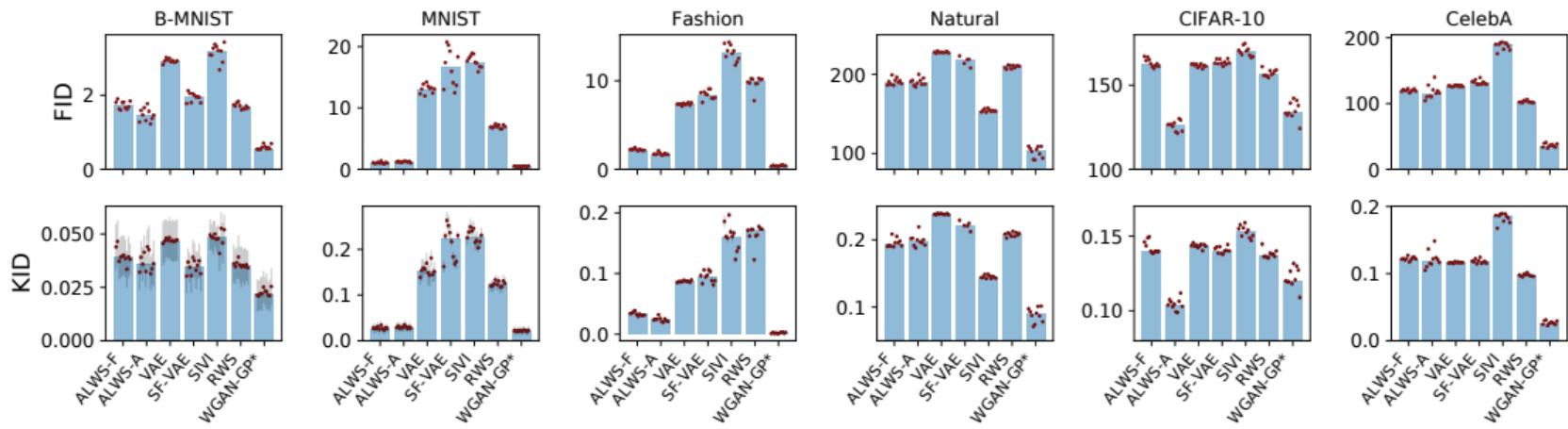
$$\tau \sim \text{Categorical}(\mathbf{m}), \tau \in \{1, \dots, 20\}, \quad e_t \sim \text{Gamma}\left(\frac{1}{\sigma_p^2}, \sigma_p^2\right), \quad \epsilon_t \sim \text{Gamma}\left(\frac{1}{\sigma_d^2}, \sigma_d^2\right),$$
$$z_t = Px_{t-\tau} \exp\left(-\frac{x_{t-\tau}}{N_0}\right) + x_t \exp(-\delta\epsilon_t), \quad p(x_t|z_t) = \text{LogNormal}(\log(z_t), \sigma_n^2)$$

Categorical is a relaxation of a discrete parameter for time lag,  $x_{-20:0}$  modelled as parameters



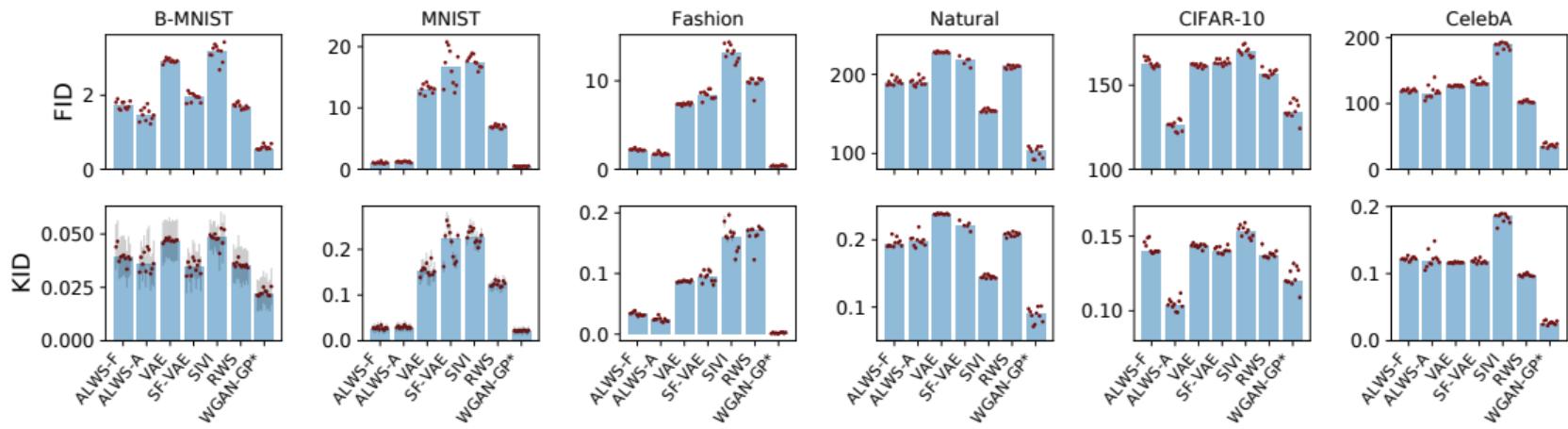
## Experiment VI: image quality

Standard benchmark for learning generative models. Quality measured as distributional discrepancies



## Experiment VI: image quality

Standard benchmark for learning generative models. Quality measured as distributional discrepancies  
Compare with: VAE, VAE+flow, VAE+implicit posterior, reweighted wake sleep

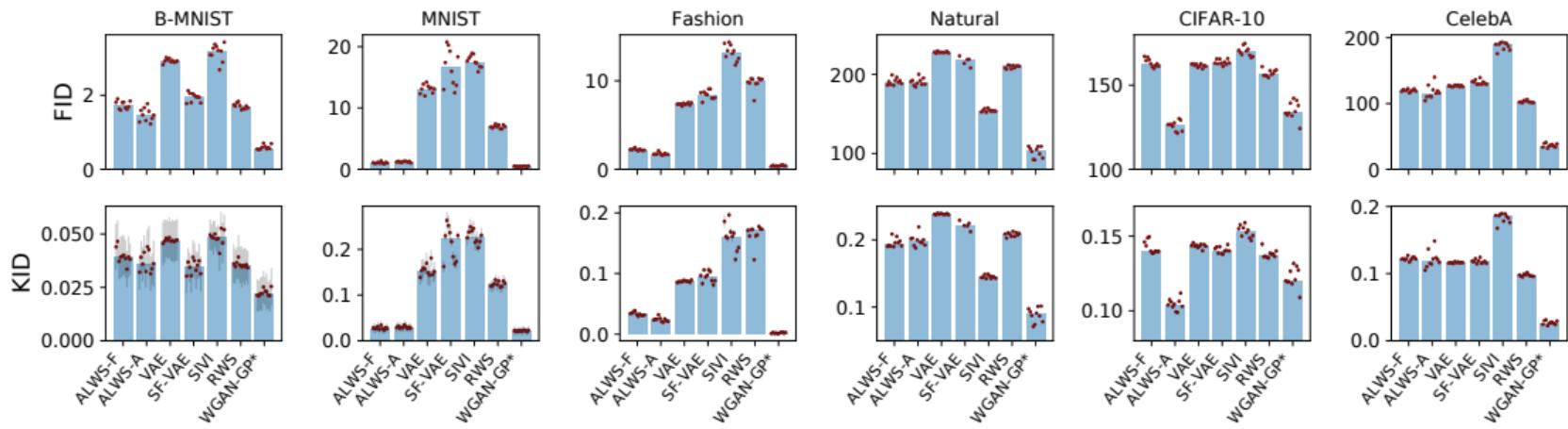


## Experiment VI: image quality

Standard benchmark for learning generative models. Quality measured as distributional discrepancies

Compare with: VAE, VAE+flow, VAE+implicit posterior, reweighted wake sleep

On datasets : B-MNIST, MNIST, Fashion, Natural, CIFAR-10, CelebA



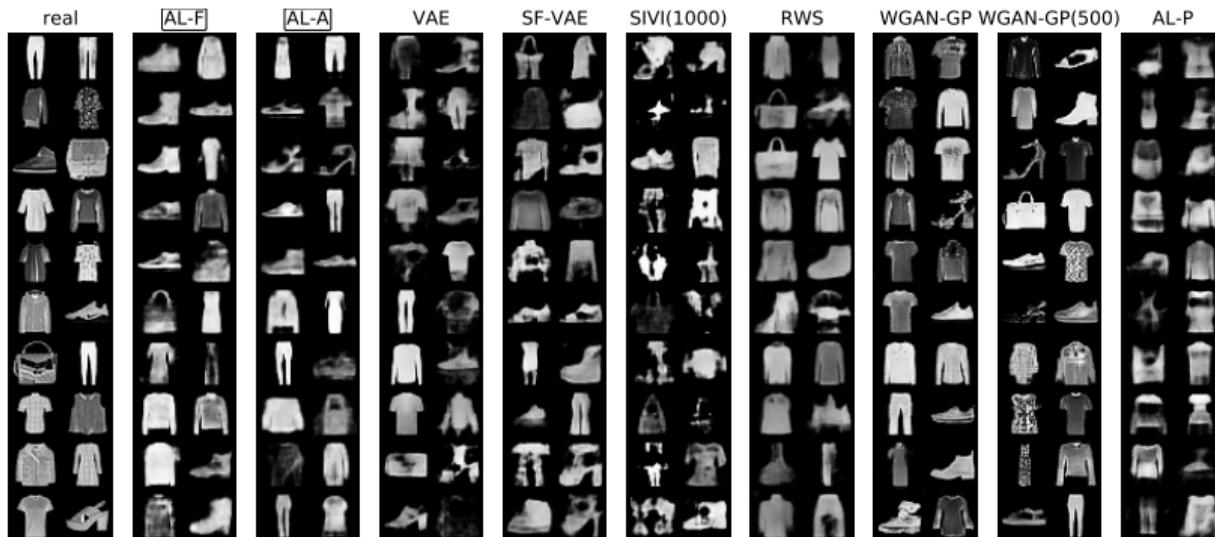
## Samples: binarised MNSIT

real	AL-F	AL-A	VAE	SF-VAE	SIVI(1000)	RWS	WGAN-GP	WGAN-GP(500)	AL-P
5 1	3 1	0 0	5 4	0 0	3 3	3 2	6 1	5 3	5 2
3 8	8 8	7 0	3 6	6 0	8 9	3 0	0 9	2 0	5 4
7 4	0 5	1 7	3 5	4 3	6 1	9 6	9 3	7 5	0 2
2 1	9 2	5 9	4 1	7 2	1 0	8 0	9 4	2 6	8 7
2 8	4 0	6 3	0 9	1 0	1 1	7 6	4 9	1 8	2 4
6 2	8 4	9 3	2 0	6 2	7 5	9 1	9 7	2 0	9 5
9 1	3 2	9 2	1 0	2 2	9 5	6 8	5 1	2 9	0 7
1 4	7 5	7 9	0 7	2 4	6 1	4 8	6 1	7 7	1 8
4 0	2 1	6 2	4 9	2 6	6 9	8 9	0 7	6 7	8 3
9 9	6 2	2 6	7 5	3 7	9 0	7 7	6 7	3 9	3 4

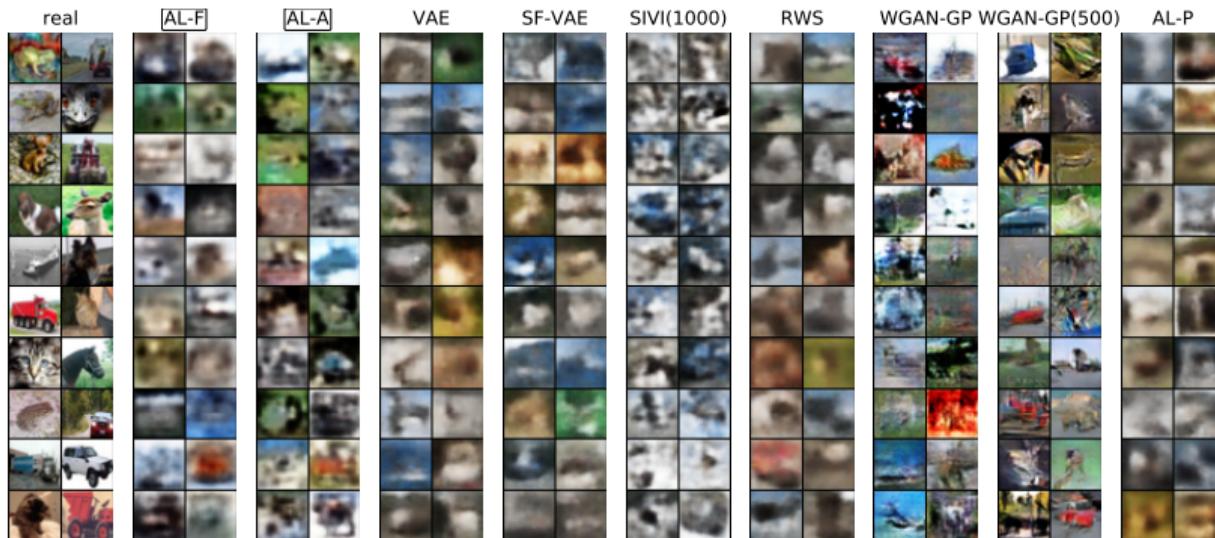
# Samples: (original) MNSIT

real	AL-F	AL-A	VAE	SF-VAE	SIVI(1000)	RWS	WGAN-GP	WGAN-GP(500)	AL-P
5 5	5 3	1 3	2 3	3 2	2 4	8 1	6 7	5 3	5 3
2 0	9 9	3 4	4 5	1 3	1 2	6 7	4 1	4 3	2 5
0 4	2 3	2 5	3 4	5 3	2 0	7 3	0 8	2 5	2 7
7 7	8 2	3 7	8 1	7 8	8 2	2 3	2 6	2 0	3 8
0 3	3 1	7 7	7 3	5 2	2 5	7 1	7 0	3 0	2 0
4 2	1 0	9 8	4 5	3 4	8 0	9 2	7 8	1 1	4 8
9 5	3 3	9 6	1 2	3 4	8 8	0 9	8 4	8 8	7 7
1 4	7 9	1 6	7 3	1 4	5 3	3 7	3 6	4 1	3 7
9 4	7 7	6 9	5 3	2 5	2 4	9 2	8 5	9 3	3 6
3 1	5 7	0 4	1 4	4 8	2 3	9 3	6 1	0 1	9 7

# Samples: Fashion



# Samples: CIFAR-10 (bad)



# Samples: Celeb-A (bad)

