

---

## Supplementary Material: Max-margin learning with the Bayes factor

---

Rahul G. Krishnan  
MIT

Arjun Khandelwal  
MIT

Rajesh Ranganath  
NYU

David Sontag  
MIT

## 1 MODEL ARCHITECTURES

We detail the neural architectures for each of the experimental setups. We use Keras (Chollet *et al.*, 2015) to implement the models described.

### 1.1 Pinwheel Dataset

**Encoder:**  $p(z|x)$ :

- $x \rightarrow \text{Dense}(20, \text{'relu'})$
- $h_1 \rightarrow \text{Dense}(20, \text{'relu'}) \rightarrow h_2$
- $h_2 \rightarrow \text{Dense}(1) \rightarrow \mu$
- $h_2 \rightarrow \text{Dense}(1) \rightarrow \log \Sigma$

**Decoder:**  $p(x|z)$ :

- $z \rightarrow \text{Dense}(20, \text{'relu'})$
- $h_1 \rightarrow \text{Dense}(20, \text{'relu'})$
- $h_2 \rightarrow \text{Dense}(2) \rightarrow \mu_{\text{obs}}$

**Reasoning Model:**  $p(z|Q)$ :

- $\{x_1, \dots, x_Q\} \rightarrow p(z|x)$  (**Elementwise**)
- $\{[\mu_1, \log \Sigma_1], \dots, [\mu_Q, \log \Sigma_Q]\} \rightarrow \text{PermutationEquivariant}(20, \text{'elu'})$
- $\{h_1^1, \dots, h_Q^1\} \rightarrow \text{PermutationEquivariant}(20, \text{'elu'})$
- $\{h_1^2, \dots, h_Q^2\} \rightarrow \text{PermutationInvariant}(1) \rightarrow \mu$
- $\{h_1^2, \dots, h_Q^2\} \rightarrow \text{PermutationInvariant}(1) \rightarrow \log \Sigma$

### 1.2 MiniImagenet Dataset

**Embedding Network**  $f(x) \rightarrow x'$ :

- $x \rightarrow \text{ResNet18}$  (He *et al.*, 2016) Conv Layers (see below)  $\rightarrow h_1$
- $h_1 \rightarrow \text{AveragePooling} \rightarrow x'$

**Encoder:**  $p(z|x')$ :

- $x' \rightarrow \text{Dense}(512, \text{'relu'}) \rightarrow h_1$
- $h_1 \rightarrow \text{Dense}(128, \text{'linear'}) \rightarrow \mu$
- $h_1 \rightarrow \text{Dense}(128, \text{'linear'}) \rightarrow \sigma$

**Decoder:**  $p(x'|z)$ :

- $z \rightarrow \text{Dense}(512, \text{'relu'}) \rightarrow h_1$
- $h_1 \rightarrow \text{Dense}(256, \text{'linear'}) \rightarrow \mu_{\text{obs}}$

**Reasoning Model:**  $p(z|Q)$ :

- $\{x_1, \dots, x_Q\} \rightarrow p(z|x)$  (**Elementwise**)
- $\{[\mu_1, \log \Sigma_1], \dots, [\mu_Q, \log \Sigma_Q]\} \rightarrow \text{PermutationEquivariant}(2048, \text{'linear'})$
- $\{h_1^2, \dots, h_Q^2\} \rightarrow \text{PermutationInvariant}(128) \rightarrow \mu$
- $\{h_1^2, \dots, h_Q^2\} \rightarrow \log \Sigma$

**Training Details:**

We take  $|Q_s| = 1$ ,  $|Q_{ns}| = 5$ , learning rate =  $5e - 5$ .

### 1.3 MNIST Dataset

**Encoder:**  $p(z|x)$ :

- $x \rightarrow \text{Flatten}() \rightarrow h_1$
- $h_1 \rightarrow \text{Dense}(500, \text{'relu'}) \rightarrow h_2$
- $h_2 \rightarrow \text{Dense}(500, \text{'relu'}) \rightarrow h_3$
- $h_3 \rightarrow \text{Dense}(2) \rightarrow \mu$
- $h_3 \rightarrow \text{Dense}(2) \rightarrow \sigma$

**Decoder:**  $p(x|z)$ :

- $z \rightarrow \text{Dense}(500, \text{'relu'}) \rightarrow h_1$
- $h_1 \rightarrow \text{Dense}(784, \text{'sigmoid'}) \rightarrow h_2$
- $h_2 \rightarrow \text{Reshape}((28,28)) \rightarrow \mu$

**Reasoning Model:**  $p(z|Q)$ :

- $\{x_1, \dots, x_Q\} \rightarrow p(z|x)$  (**Elementwise**)
- $\{[\mu_1, \log \Sigma_1], \dots, [\mu_Q, \log \Sigma_Q]\}$   
 $\rightarrow \text{PermutationEquivariant}(20, \text{'relu'})$
- $\{h_1^2, \dots, h_Q^2\} \rightarrow \text{PermutationInvariant}(2) \rightarrow \mu$
- $\{h_1^2, \dots, h_Q^2\} \rightarrow \text{PermutationInvariant}(2) \rightarrow \log \Sigma$

**Training Details:** We take  $|Q_s| = 5$ ,  $|Q_{ns}| = 5$ , learning rate =  $1e - 4$ .

## References

- Chollet, François, et al. . 2015. *Keras*. <https://github.com/keras-team/keras>.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, & Sun, Jian. 2016. Deep residual learning for image recognition. In: *CVPR*.