
Supplementary Material for the RNN-RNADE model

Anonymous Author(s)

Affiliation

Address

email

Abstract

This document contains some of the mathematical derivations for the RNN-RNADE model.

1 Introduction

The RNADE fprop:

$$p(x) = \prod_{d=1}^D p(x_d | \mathbf{x}_{< \mathbf{d}}) \text{ with } p(x_d | \mathbf{x}_{< \mathbf{d}}) = p_{\mathcal{M}(x_d | \theta_d)}$$

$$\mathbf{a}_{d+1} = \mathbf{a}_d + x_d \mathbf{W}_{\cdot, d}$$

$$\mathbf{h}_d = \text{sigm}(\rho_d \mathbf{a}_d)$$

$$\boldsymbol{\alpha}_d = \text{softmax}(\mathbf{V}_d^\alpha \mathbf{h}_d + \mathbf{b}_d^\alpha)$$

$$\mu_d = \mathbf{V}_d^\mu \mathbf{h}_d + \mathbf{b}_d^\mu$$

$$\sigma_d = \exp(\mathbf{V}_d^\sigma \mathbf{h}_d + \mathbf{b}_d^\sigma)$$

Algorithm 1 RNADE fprop

procedure FPROP

$a \leftarrow c$

$p \leftarrow 1$

for d from 1 to D **do**

$\psi_d = \rho_d \mathbf{a}$

$\mathbf{h}_d = \text{sigm}(\psi_d)$

$\mathbf{z}_d^\alpha = \mathbf{V}_d^\alpha \mathbf{h}_d + \mathbf{b}_d^\alpha$

$\mathbf{z}_d^\mu = \mathbf{V}_d^\mu \mathbf{h}_d + \mathbf{b}_d^\mu$

$\mathbf{z}_d^\sigma = \mathbf{V}_d^\sigma \mathbf{h}_d + \mathbf{b}_d^\sigma$

$\boldsymbol{\alpha}_d = \text{softmax}(\mathbf{z}_d^\alpha)$

$\mu_d = \mathbf{z}_d^\mu$

$\sigma_d = \exp(\mathbf{z}_d^\sigma)$

$p(\mathbf{x}) = p(\mathbf{x}) p_{\mathcal{M}}(x_d; \boldsymbol{\alpha}_d; \mu_d; \sigma_d)$

$\mathbf{a} = \mathbf{a} + x_d \mathbf{W}_{\cdot, d}$

Gradients for the RNADE:

$$\phi_i(x_d|\mathbf{x}_{<d}) = \frac{1}{\sqrt{2\pi}\sigma_{d,i}} \exp \left\{ -\frac{(x_d - \mu_{d,i})^2}{2\sigma_{d,i}^2} \right\}$$

Posterior/Responsibility:

$$\pi_i(x_d|\mathbf{x}_{<d}) = \frac{\alpha_{d,i}\phi_i(x_d|\mathbf{x}_{<d})}{\sum_{j=1}^K \alpha_{d,j}\phi_j(x_d|\mathbf{x}_{<d})}$$

Gradients with respect to the gaussian parameters:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{V}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \mathbf{h}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{b}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{V}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \mathbf{h}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{b}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{V}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \mathbf{h}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{b}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma}$$

Gradient with respect to the hidden activations:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{h}_d} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{h}_d} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{h}_d}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \mathbf{V}_d^\alpha + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \mathbf{V}_d^\mu + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \mathbf{V}_d^\sigma$$

2 RNN-RNADE fprop

$$p(x_1^T) = \prod_1^T p(x^t | \mathcal{A}^t) \text{ where } \mathcal{A}^t \equiv \{x^\tau | \tau < t\}$$

Maximum-likelihood Training:

$$\begin{aligned} \log p(x_1^T) &= \sum_1^T \log p(x^t | \mathcal{A}^t) \\ \log p(x_1^T) &= \sum_{t=1}^T \sum_{d=1}^D \log p(x_d^t | \mathbf{x}_{<d}^t) \end{aligned}$$

W' is the weight matrix from the recurrent layer to the mixing coefficients.

W'' is the weight matrix from the recurrent layer to the means.

W''' is the weight matrix from the recurrent layer to the sigmas.

The rest of the architecture is similar to the RNN-RBM paper.

$$\begin{aligned} \hat{h}^{(t)} &= \sigma(W_2 \mathbf{x}^t + W_3 \hat{h}^{(t-1)} + b_{\hat{h}}) \\ r_\alpha &= W^\alpha \hat{h}^{(t-1)} \\ r_\mu &= W^\mu \hat{h}^{(t-1)} \\ r_\sigma &= W^\sigma \hat{h}^{(t-1)} \end{aligned}$$

Algorithm 2 RNN-RNADE fprop

procedure FPROP

$a \leftarrow c$

$p \leftarrow 1$

for t from 1 to T **do**

for d from 1 to D **do**

$\psi_d = \rho_d \mathbf{a}$

$\mathbf{h}_d = \text{sigm}(\psi_d)$

$\mathbf{z}_d^\alpha = \mathbf{V}_d^\alpha \mathbf{T} \mathbf{h}_d + \mathbf{b}_d^\alpha + \mathbf{r}_d^\alpha = \mathbf{V}_d^\alpha \mathbf{T} \mathbf{h}_d + \mathbf{b}_d^\alpha + \mathbf{r}_d^\alpha$

$\mathbf{z}_d^\mu = \mathbf{V}_d^\mu \mathbf{T} \mathbf{h}_d + \mathbf{b}_d^\mu + \mathbf{r}_d^\mu$

$\mathbf{z}_d^\sigma = \mathbf{V}_d^\sigma \mathbf{T} \mathbf{h}_d + \mathbf{b}_d^\sigma + \mathbf{r}_d^\sigma$

$\alpha_d = \text{softmax}(\mathbf{z}_d^\alpha)$

$\mu_d = \mathbf{z}_d^\mu$

$\sigma_d = \exp(\mathbf{z}_d^\sigma)$

$p(\mathbf{x}) = p(\mathbf{x}) p_{\mathcal{M}}(x_d; \alpha_d; \mu_d; \sigma_d)$

$\mathbf{a} = \mathbf{a} + x_d^t \mathbf{W}_{\cdot, d}$

Gradients for the RNN-RNADE:

$$\phi_i(x_d | \mathbf{x}_{<d}) = \frac{1}{\sqrt{2\pi}\sigma_{d,i}} \exp \left\{ -\frac{(x_d - \mu_{d,i})^2}{2\sigma_{d,i}^2} \right\}$$

Posterior/Responsibility:

$$\pi_i(x_d | \mathbf{x}_{<d}) = \frac{\alpha_{d,i} \phi_i(x_d | \mathbf{x}_{<d})}{\sum_{j=1}^K \alpha_{d,j} \phi_j(x_d | \mathbf{x}_{<d})}$$

Gradients with respect to the gaussian parameters:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{V}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \mathbf{h}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{b}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{V}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \mathbf{h}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{b}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{V}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \mathbf{h}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{b}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma}$$

Gradient with respect to the hidden activations:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{h}_d} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{h}_d} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{h}_d}$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \mathbf{V}_d^\alpha + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \mathbf{V}_d^\mu + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \mathbf{V}_d^\sigma$$