
RNADE: The real-valued neural autoregressive density-estimator

Supplementary material

Benigno Uria and **Iain Murray**
 School of Informatics
 University of Edinburgh
 {b.uria, i.murray}@ed.ac.uk

Hugo Larochelle
 Département d'informatique
 Université de Sherbrooke
 hugo.larochelle@usherbrooke.ca

In this document we provide pseudo-code for the calculation of densities and learning gradients. No new material is presented. A Python implementation of RNADE is available from <http://www.benignouria.com/en/research/RNADE>.

1 Density estimation

In Algorithm 1 we detail the pseudocode for calculating the density of a datapoint under an RNADE with mixture of Gaussian conditionals. The model has parameters: $\rho \in \mathbb{R}^D$, $\mathbf{W} \in \mathbb{R}^{H \times D-1}$, $\mathbf{c} \in \mathbb{R}^H$, $\mathbf{b}^\alpha \in \mathbb{R}^{D \times K}$, $\mathbf{V}^\alpha \in \mathbb{R}^{D \times H \times K}$, $\mathbf{b}^\mu \in \mathbb{R}^{D \times K}$, $\mathbf{V}^\mu \in \mathbb{R}^{D \times H \times K}$, $\mathbf{b}^\sigma \in \mathbb{R}^{D \times K}$, $\mathbf{V}^\sigma \in \mathbb{R}^{D \times H \times K}$

Algorithm 1 Computation of $p(\mathbf{x})$

```

 $\mathbf{a} \leftarrow \mathbf{c}$ 
 $p(\mathbf{x}) \leftarrow 1$ 
for  $d$  from 1 to  $D$  do
   $\psi_d \leftarrow \rho_d \mathbf{a}$  ▷ Rescaling factors
   $\mathbf{h}_d \leftarrow \psi_d \mathbf{1}_{\psi_d > 0}$  ▷ Rectified linear units
   $\mathbf{z}_d^\alpha \leftarrow \mathbf{V}_d^{\alpha \top} \mathbf{h}_d + \mathbf{b}_d^\alpha$ 
   $\mathbf{z}_d^\mu \leftarrow \mathbf{V}_d^{\mu \top} \mathbf{h}_d + \mathbf{b}_d^\mu$ 
   $\mathbf{z}_d^\sigma \leftarrow \mathbf{V}_d^{\sigma \top} \mathbf{h}_d + \mathbf{b}_d^\sigma$ 
   $\alpha_d \leftarrow \text{softmax}(\mathbf{z}_d^\alpha)$  ▷ Enforce constraints
   $\mu_d \leftarrow \mathbf{z}_d^\mu$ 
   $\sigma_d \leftarrow \exp(\mathbf{z}_d^\sigma)$ 
   $p(\mathbf{x}) \leftarrow p(\mathbf{x}) p_{MoG}(x_d; \alpha_d, \mu_d, \sigma_d)$  ▷  $p_{MoG}$  is the density of a mixture of Gaussians
   $\mathbf{a} \leftarrow \mathbf{a} + x_d \mathbf{W}_{\cdot, d}$  ▷ Activations are calculated recursively,  $x_d$  is a scalar
end for
return  $p(\mathbf{x})$ 

```

2 Learning gradients

Training of an RNADE model can be done using a gradient ascent algorithm on the loglikelihood of the model given the training data. Gradients can be calculated using automatic differentiation libraries (e.g. Theano [1]). However we found our manual implementation to work faster in practice, possibly due to our recomputation of the \mathbf{a} terms in the second *for* loop in Algorithm 2, which is more cache-friendly than storing them during the first loop.

Here we show the derivation of the gradients of each parameter of a NADE model with MoG conditionals. Following [2], we define $\phi_i(x_d | \mathbf{x}_{<d})$ as the density of x_d under the i -th component of

the conditional:

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

and $\pi_i(x_d | \mathbf{x}_{<d})$ as the “responsability” of the i -th component for x_d :

$$\pi_i(x_d | \mathbf{x}_{<d}) = \frac{\alpha_{d,i} \phi_i(x_d | \mathbf{x}_{<d})}{\sum_{j=1}^K \alpha_{d,j} \phi_j(x_d | \mathbf{x}_{<d})}. \quad (2)$$

It is easy to find just by taking their derivatives that:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} = \pi_i(x_d | \mathbf{x}_{<d}) - \alpha_{d,i} \quad (3)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} = \pi_i(x_d | \mathbf{x}_{<d}) \frac{x_d - \mu_{d,i}}{\sigma_{d,i}^2} \quad (4)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} = \pi_i(x_d | \mathbf{x}_{<d}) \left\{ \frac{(x_d - \mu_{d,i})^2}{\sigma_{d,i}^2} - 1 \right\} \quad (5)$$

Using the chain rule we can calculate the derivative of the parameters of the output layer parameters:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{V}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \mathbf{h} \quad (6)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{b}_d^\alpha} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \quad (7)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{V}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \mathbf{h} \quad (8)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{b}_d^\mu} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \quad (9)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{V}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{V}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \mathbf{h} \quad (10)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{b}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{b}_d^\sigma} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \quad (11)$$

By “backpropagating” the we can calculate the partial derivatives with respect to the output of the hidden units:

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \frac{\partial \mathbf{z}_{d,i}^\alpha}{\partial \mathbf{h}_d} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \frac{\partial \mathbf{z}_{d,i}^\mu}{\partial \mathbf{h}_d} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \frac{\partial \mathbf{z}_{d,i}^\sigma}{\partial \mathbf{h}_d} \quad (12)$$

$$= \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\alpha} \mathbf{V}_d^\alpha + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\mu} \mathbf{V}_d^\mu + \frac{\partial p(\mathbf{x})}{\partial \mathbf{z}_{d,i}^\sigma} \mathbf{V}_d^\sigma \quad (13)$$

and calculate the partial derivatives with respect to all other parameters in RNADE:

$$\frac{\partial p(\mathbf{x})}{\partial \psi_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} \mathbf{1}_{\psi_d > 0} \quad (14)$$

$$\frac{\partial p(\mathbf{x})}{\partial \rho_d} = \sum_j \frac{\partial p(\mathbf{x})}{\partial \psi_{d,j}} \mathbf{a}_{d,j} \quad (15)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{a}_d} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{a}_{d+1}} + \frac{\partial p(\mathbf{x})}{\partial \mathbf{h}_d} \rho_d \mathbf{1}_{\psi_d > 0} \quad (16)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{W}_{\cdot,d}} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{a}_d} x_d \quad (17)$$

$$\frac{\partial p(\mathbf{x})}{\partial \mathbf{c}} = \frac{\partial p(\mathbf{x})}{\partial \mathbf{a}_1} \quad (18)$$

Note that gradients are calculated recursively, due to (16), starting at $d = D$ and progressing down to $d = 1$.

Algorithm 2 Computation of the learning gradients for a datapoint x

```

 $a \leftarrow c$ 
for  $d$  from 1 to  $D$  do                                     ▷ Compute the activation of the last dimension
     $a \leftarrow a + x_d W_{\cdot,d}$ 
end for
for  $d$  from  $D$  to 1 do                                       ▷ Backpropagate errors
     $\psi \leftarrow \rho_d a$                                        ▷ Rescaling factors
     $h \leftarrow \psi \mathbf{1}_{\psi > 0}$                                 ▷ Rectified linear units
     $z^\alpha \leftarrow V_d^{\alpha \top} h + b_d^\alpha$ 
     $z^\mu \leftarrow V_d^{\mu \top} h + b_d^\mu$ 
     $z^\sigma \leftarrow V_d^{\sigma \top} h + b_d^\sigma$ 
     $\alpha \leftarrow \text{softmax}(z^\alpha)$                                ▷ Enforce constraints
     $\mu \leftarrow z^\mu$ 
     $\sigma \leftarrow \exp(z^\sigma)$ 
     $\phi \leftarrow \frac{1}{2} \frac{(\mu - x_d)^2}{\sigma^2} - \log \sigma - \frac{1}{2} \log(2\pi)$    ▷ Calculate gradients
     $\pi \leftarrow \frac{\alpha \phi}{\sum_{j=1}^K \alpha_j \phi_j}$ 
     $\partial z^\alpha \leftarrow \pi - \alpha$ 
     $\partial V_d^\alpha \leftarrow \partial z^\alpha h$ 
     $\partial b_d^\alpha \leftarrow \partial z^\alpha$ 
     $\partial z^\mu \leftarrow \pi(x_d - \mu)/\sigma^2$ 
     $\partial z^\mu \leftarrow \partial z^\mu * \sigma$                                ▷ Move tighter components slower, allows higher learning rates
     $\partial V_d^\mu \leftarrow \partial z^\mu h$ 
     $\partial b_d^\mu \leftarrow \partial z^\mu$ 
     $\partial z^\sigma \leftarrow \pi\{(x_d - \mu)^2/\sigma^2 - 1\}$ 
     $\partial V_d^\sigma \leftarrow \partial z^\sigma h$ 
     $\partial b_d^\sigma \leftarrow \partial z^\sigma$ 
     $\partial h \leftarrow \partial z^\alpha V_d^\alpha + \partial z^\mu V_d^\mu + \partial z^\sigma V_d^\sigma$ 
     $\partial \psi \leftarrow \partial h \mathbf{1}_{\psi > 0}$                                ▷ Second factor: indicator function with condition  $\psi > 0$ 
     $\partial \rho_d \leftarrow \sum_j \partial \psi_j a_j$ 
     $\partial a \leftarrow \partial a + \partial \psi \rho$ 
     $\partial W_{\cdot,d} \leftarrow \partial a x_d$ 
    if  $d = 1$  then
         $\partial c \leftarrow \partial a$ 
    else
         $a \leftarrow a - x_d W_{\cdot,d}$ 
    end if
end for
return  $\partial \rho, \partial W, \partial c, \partial b^\alpha, \partial V^\alpha, \partial b^\mu, \partial V^\mu, \partial b^\sigma, \partial V^\sigma$ 

```

References

- [1] James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*, June 2010. Oral Presentation.
- [2] C. M. Bishop. Mixture density networks. Technical Report NCRG 4288, Neural Computing Research Group, Aston University, Birmingham, 1994.