

[Paper review 28]

Dennsity Estimation using Real NVP

(Laurent Dinh, et al., 2017)

[Contents]

1. Abstract
2. Introduction

1. Abstract

Advantages of "Flow based generative models "

- 1) tractability of the exact log-likelihood
- 2) tractability of exact latent-variable inference
- 3) parallelizability of both training and syntehsis

Glow

- simple type of generative flow, using "invertible 1 x 1 convolution"
- significant improvement in log-likelihood on standard benchmarks

2. Introduction

2 major problems in ML

- 1) data efficiency (ability to learn from few data points)
- 2) generalization (robustness to changes of the task)

Promise of generative models : overcome these 2 problems by

- learning realistic world models
- learning meaningful features of the input

Generative Modeling have advanced with likelihood-based methods

Likelihood-based methods : three categories

- 1) Autoregressive models
- 2) VAEs

- 3) Flow-based generative models (ex. NICE, RealNVP)

3) Flow-based generative model's merit

- exact latent-variable inference and log-likelihood evaluation
- efficient inference and efficient synthesis
- useful latent space for downstream tasks
- significant potential for memory savings

3. Background : Flow-based Generative Models

x : high-dimensional random vector

$x \sim p^*(x)$: unknown true distribution

log-likelihood objective : minimizing....

- (discrete x)

$$\mathcal{L}(\mathcal{D}) = \frac{1}{N} \sum_{i=1}^N -\log p_{\theta}(\mathbf{x}^{(i)})$$

- (continuous x)

$$\mathcal{L}(\mathcal{D}) \simeq \frac{1}{N} \sum_{i=1}^N -\log p_{\theta}(\tilde{\mathbf{x}}^{(i)}) + c$$

- $\tilde{\mathbf{x}}^{(i)} = \mathbf{x}^{(i)} + u$ with $u \sim \mathcal{U}(0, a)$, and $c = -M \cdot \log a$
- a : determined by the discretization level of the data
- M : dimensionality of x

Generative process of most flow-based generative models

$$\mathbf{z} \sim p_{\theta}(\mathbf{z})$$

$$\mathbf{x} = g_{\theta}(\mathbf{z})$$

- z : latent variable
- $p_{\theta}(z)$: tractable density
- $g_{\theta}(\cdot)$: invertible (=bijective)

Change of variables + triangular matrix

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &= \log p_{\theta}(\mathbf{z}) + \log |\det(d\mathbf{z}/d\mathbf{x})| \\ &= \log p_{\theta}(\mathbf{z}) + \sum_{i=1}^K \log |\det(d\mathbf{h}_i/d\mathbf{h}_{i-1})| \\ &= \log p_{\theta}(\mathbf{z}) + \sum_{i=1}^K \text{sum}(\log |\text{diag}(d\mathbf{h}_i/d\mathbf{h}_{i-1})|) \end{aligned}$$

4. Proposed Generative Flow

we propose a new flow

- built on NICE and RealNVP
- consists of a series of steps of flows
- combined with multi-scale architecture

4.1 Actnorm : scale & bias layer with data dependent initialization

actnorm layer :

- performs an affine transformation of the activations, using a "scale and bias" parameters per channel
- data dependent initialization

4.2 Invertible 1 x 1 convolution

permutation that reverses the ordering of the channels

(1x1 convolution with equal number of input & output channels = generalization of permutation operation)

log-determinant of an invertible 1x1 convolution of $h \times w \times c$ tensor \mathbf{h} with $c \times c$ weight matrix \mathbf{W}

$$\log \left| \det \left(\frac{d \text{conv } 2D(\mathbf{h}; \mathbf{W})}{d\mathbf{h}} \right) \right| = h \cdot w \cdot \log |\det(\mathbf{W})|$$

Computation cost

- before : $\text{conv } 2D(\mathbf{h}; \mathbf{W}) \rightarrow \mathcal{O}(h \cdot w \cdot c^2)$
- after : $\det(\mathbf{W}) \rightarrow \mathcal{O}(c^3)$,

proof) LU Decomposition

$$\mathbf{W} = \mathbf{P}\mathbf{L}(\mathbf{U} + \text{diag}(\mathbf{s}))$$

$$\log |\det(\mathbf{W})| = \text{sum}(\log |\mathbf{s}|)$$

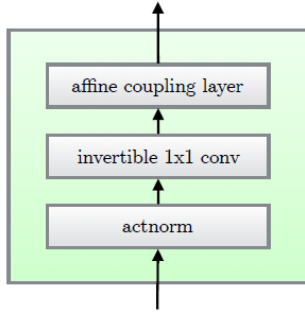
where P : permutation matrix & W random rotation matrix

4.3 Affine Coupling Layers

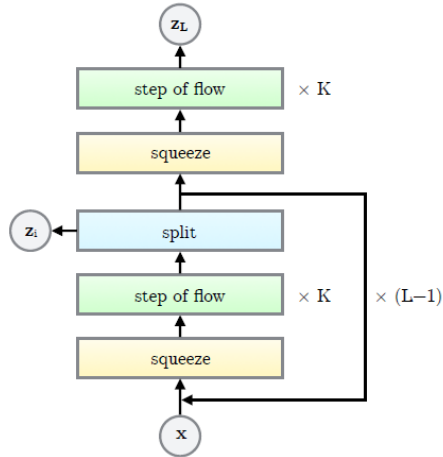
ex) $s = 1$: additive coupling layer

- Zero initialization
(initialize the last convolution of each NN() with zeros)
- Split and Concatenation
(splits \mathbf{h} the input tensor into 2 halves)

Description	Function	Reverse Function	Log-determinant
Actnorm. See Section 3.1	$\forall i, j : y_{i,j} = s \odot x_{i,j} + \mathbf{b}$	$\forall i, j : x_{i,j} = (y_{i,j} - \mathbf{b})/s$	$h \cdot w \cdot \text{sum}(\log s)$
Invertible 1×1 convolution. $\mathbf{W} : [c \times c]$. See Section 3.2	$\forall i, j : y_{i,j} = \mathbf{W}x_{i,j}$	$\forall i, j : x_{i,j} = \mathbf{W}^{-1}y_{i,j}$	$h \cdot w \cdot \log \det(\mathbf{W}) $ or $h \cdot w \cdot \text{sum}(\log s)$ (see eq. (10))
Affine coupling layer. See Section 3.3 and (Dinh et al., 2014)	$x_a, x_b = \text{split}(x)$ $(\log s, t) = \text{NN}(x_b)$ $s = \exp(\log s)$ $y_a = s \odot x_a + t$ $y_b = x_b$ $y = \text{concat}(y_a, y_b)$	$y_a, y_b = \text{split}(y)$ $(\log s, t) = \text{NN}(y_b)$ $s = \exp(\log s)$ $x_a = (y_a - t)/s$ $x_b = y_b$ $x = \text{concat}(x_a, x_b)$	$\text{sum}(\log(s))$



(a) One step of our flow.



(b) Multi-scale architecture (Dinh et al., 2016).