

Normalizing Flows

Abdul Fatir Ansari, Devamanyu Hazarika, and Remmy A. M. Zen

National University of Singapore
{abdufatir,devamanyu,remmy}@u.nus.edu

Abstract. Keywords: Normalizing flows . .

1 Introduction

2 Background

2.1 Autoencoder

An Autoencoder (AE) [4] is comprised of two neural networks: 1) The encoder network f_ϕ which transforms an input data-point $\mathbf{x} \in \mathbb{R}^d$ into a representation $\mathbf{z} \in \mathbb{R}^p$ where $p < d$ generally; 2) The decoder network g_θ which attempts to reconstruct \mathbf{x} as $\hat{\mathbf{x}} \in \mathbb{R}^d$ from the representation \mathbf{z} . Training proceeds by minimizing the reconstruction error between \mathbf{x} and $\hat{\mathbf{x}}$ and the training objective is given in Eq. (1)

$$\phi, \theta = \arg \min_{\phi, \theta} \mathcal{L}(\hat{\mathbf{x}}, \mathbf{x}) \quad (1)$$

where \mathcal{L} is a suitable reconstruction error which can be the squared error $\|\hat{\mathbf{x}} - \mathbf{x}\|^2$ for real-valued observations or binary cross-entropy $\sum_{i=1}^d -\mathbf{x}|_i \log(\hat{\mathbf{x}}|_i) - (1 - \mathbf{x}|_i) \log(1 - \hat{\mathbf{x}}|_i)$ for binary observations.

2.2 Variational Autoencoder

The Variational Autoencoder (VAE) [7] combines an inference network f_ϕ (analogous to the encoder in AE) with a generative model g_θ (analogous to the decoder in AE). VAE maximizes the *evidence lower bound* (ELBO) which can be derived as follows

$$\log p(\mathbf{x}) = \log \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \log \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \quad (2)$$

$$\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] = \mathcal{L}_{\text{ELBO}} \quad (3)$$

where the last inequality comes from Jensen's inequality.

A data-point $\mathbf{x} \in \mathbb{R}^d$ is passed through the inference network f_ϕ which outputs the parameters of the variational distribution $q(\mathbf{z}|\mathbf{x})$. A point $\mathbf{z} \in \mathbb{R}^p$ is

sampled from $q(\mathbf{z}|\mathbf{x})$ and is passed onto the generative model g_θ which attempts to reconstruct \mathbf{x} as $\hat{\mathbf{x}} \in \mathbb{R}^d$. Once training is complete, the inference network and the generative model can be used separately. New samples from the VAE can be generated by sampling \mathbf{z} from the prior $p(\mathbf{z})$ (generally set to $\mathcal{N}(\mathbf{0}, \mathbf{I})$) and plugging into the generative model.

Maximization of $\mathcal{L}_{\text{ELBO}}$ using stochastic gradient descent (SGD) requires computation of $\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$ which is done using the *reparameterization trick*. Briefly, the reparameterization trick uses an invertible standardization function $\mathcal{S}_\phi(\mathbf{z})$ such that $\mathcal{S}_\phi(\mathbf{z}) = \epsilon$ where $\epsilon \sim p(\epsilon)$ and $p(\epsilon)$ is a simple distribution. For example, $\mathcal{S}_\phi(\mathbf{z})$ for a normally distributed random variable would be $\mathcal{S}_\phi(\mathbf{z}) = \frac{\mathbf{z} - \boldsymbol{\mu}}{\boldsymbol{\sigma}} = \epsilon$ with $p(\epsilon)$ being $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Once such an invertible standardization function is available $\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [f(\mathbf{z})]$ can be written as $\nabla_\phi \mathbb{E}_{p(\epsilon)} [f(\mathcal{S}_\phi^{-1}(\epsilon))]$ which can be evaluated by using the chain rule as $\mathbb{E}_{p(\epsilon)} [\nabla_{\mathbf{z}} f(\mathcal{S}_\phi^{-1}(\epsilon)) \nabla_\phi \mathcal{S}_\phi^{-1}(\epsilon)]$. For more details please refer to [2].

3 Normalizing Flows

Before defining normalizing flows, let's consider a univariate distribution with density function $p(x)$. Define a continuous, differentiable, and increasing function f . Define $y = f(x)$ where $x \sim p(x)$. The density function of the random variable Y can then be derived analytically using the Cumulative Distribution Function (CDF) as follows.

$$F_Y(y) = P(Y \leq y) \quad (4)$$

$$= P(f(X) \leq y) \quad (5)$$

$$= P(X \leq f^{-1}(y)) = F_X(f^{-1}(y)) \quad (6)$$

We end up with the CDF of the random variable X at the point $f^{-1}(y)$. Now, $p(y) = F'_Y(y)$ by definition, where

$$F_Y(y) = F_X(f^{-1}(y)) = \int_{-\infty}^{f^{-1}(y)} p(x) dx \quad (7)$$

Differentiating Eq. (7) with respect to y (using the Fundamental Theorem of Calculus and the chain rule) we get

$$p(y) = p(f^{-1}(y)) \cdot \frac{df^{-1}}{dy} \quad (8)$$

When f is a decreasing function, we get $p(y) = -p(f^{-1}(y)) \cdot \frac{df^{-1}}{dy}$. For an invertible function in general, Eq. (8) can be written as

$$p(y) = p(f^{-1}(y)) \cdot \left| \frac{df^{-1}}{dy} \right| \quad (9)$$

Eq. (9) can be extended to the multivariate case where the derivative is replaced by the determinant of the Jacobian matrix J

$$p(\mathbf{y}) = p(f^{-1}(\mathbf{y})) \cdot \left| \det \frac{\partial f^{-1}}{\partial \mathbf{y}} \right| = p(f^{-1}(\mathbf{y})) \cdot \left| \det \frac{\partial f}{\partial f^{-1}(\mathbf{y})} \right|^{-1} \quad (10)$$

In the above equation, the second equality comes from the inverse function theorem. Successive applications of such smooth, invertible transformation on a random variable with known density is called a *normalizing flow*.

Computation of the probability density of the transformed random variable requires the computation of the determinant of the Jacobian matrix which is computationally expensive as it scales with $O(d^3)$ where d is the dimensionality of the random variable. Developing transformations with cheap determinant computation has been the primary focus of many recent works.

4 Applications

Literature on normalizing flows can be broadly classified into two parts: ones using normalizing flows for improved variational inference and ones using normalizing flows for density estimation.

4.1 Variational Inference

Variational methods perform inference by approximating the true posterior $p(z|x)$ using a simpler variational family $q_\phi(z|x)$. Recent works have focused on improving the variational posterior used in the VAE which is generally set to a multivariate normal distribution with diagonal covariance matrix $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$. It is clear that such a simplistic, unimodal choice for the posterior can be arbitrarily far away from the true posterior which can be a complex multi-modal distribution.

Recent works seek to convert samples from a simple variational posterior (such as the multivariate normal distribution) into a richer distribution by applying a series of smooth, invertible transformations or a flow. Let \mathbf{z}_0 be a sample from a simple distribution $q_0(\mathbf{z}_0)$ and \mathbf{z}_K be a sample obtained by applying a flow of length K on \mathbf{z}_0 , i.e., $\mathbf{z}_K = f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$. Using Eq. (10), the density function $q_K(\mathbf{z}_K)$ is given by

$$q_K(\mathbf{z}_K) = q_0(\mathbf{z}_0) \prod_{k=1}^K \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right|^{-1} \quad (11)$$

The variational lower bound (or evidence lower bound) in VAEs (Eq. ()) can now be modified by setting $q(\mathbf{z}|\mathbf{x}) = q_K(\mathbf{z}_K|\mathbf{x})$

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{z}_K|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}_K) - \log q(\mathbf{z}_K|\mathbf{x})] \quad (12)$$

$$= \mathbb{E}_{q(\mathbf{z}_0|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}_K) - \log q(\mathbf{z}_K|\mathbf{x})] \quad (13)$$

where $q(\mathbf{z}_0|\mathbf{x})$ is the simple initial density. Plugging in Eq. (11) into Eq. (13), we get a modified bound for flow-based VAEs

$$\mathcal{L} = \mathbb{E}_{q_0(\mathbf{z}_0|\mathbf{x})} \left[\log p(\mathbf{x}, \mathbf{z}_K) - \log q_0(\mathbf{z}_0|\mathbf{x}) + \sum_{k=1}^K \log \left| \det \frac{\partial f_k}{\partial \mathbf{z}_{k-1}} \right| \right] \quad (14)$$

Planar and Radial Flows Planar and Radial Flows [9] are one of the earliest flows proposed in the context of variational inference.

Planar flows use functions of the form

$$f(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^\top \mathbf{z} + b) \quad (15)$$

where $\mathbf{u}, \mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$, and h is an element-wise non-linearity such as \tanh . The Jacobian is then given by

$$\det \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} = (1 + h'(\mathbf{w}^\top \mathbf{z} + b)\mathbf{w}^\top \mathbf{u}) \quad (16)$$

which can be computed in $O(d)$ time.

Radial flows use functions of the form

$$f(\mathbf{z}) = \mathbf{z} + \beta h(\alpha, r)(\mathbf{z} - \mathbf{z}_0) \quad (17)$$

where $\alpha \in \mathbb{R}^+$, $\beta \in \mathbb{R}$, $h(\alpha, r) = (\alpha + r)^{-1}$ and $r = \|\mathbf{z} - \mathbf{z}_0\|$.

The Jacobian is then given by

$$\det \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} = (1 + \beta h(\alpha, r) + \beta h'(\alpha, r)r) (1 + \beta h(\alpha, r))^{d-1} \quad (18)$$

For a detailed derivation of Jacobians of Planar and Radial flows please refer Appendix. Fig. 1 shows how planar and radial flows change a standard normal distribution.

Inverse Autoregressive Flows Planar and radial flow provide a simple invertible transformation shown to be effective in a low-dimensional latent space (up to hundred dimensions). Equation 15 of Planar Flow can be seen as Multi-layer Perceptron with a single unit bottleneck hidden layer. This means that a long chain of transformations is needed to model a high-dimensional function.

Autoregressive flows [3] is a normalizing flow that scales to high-dimensional latent space by exploiting the ordering of the variables. In autoregressive flow,

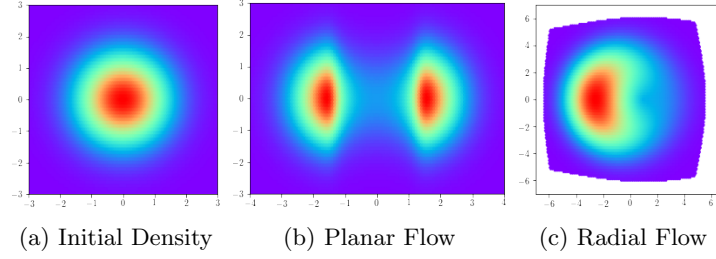


Fig. 1: Change in standard normal density on application of length 1 planar and radial flows.

given a sequence of variable $\mathbf{y} = y_{i=0}^D$, each variable is only dependent only on variables from the previous index. The distribution is then given by

$$p(\mathbf{y}) = \prod_{i=0}^D p(x_i | x_0, \dots, x_{i-1}) \quad (19)$$

Kingma et. al. [6] proposed a Gaussian version of an autoregressive flow on a noise vector $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ given as follow:

$$y_0 = \mu_0 + \sigma_0 \epsilon_0 \quad (20)$$

$$y_i = \mu_i(\mathbf{y}_{0:i-1}) + \sigma(\mathbf{y}_{0:i-1}) \epsilon_i \quad (21)$$

$$(22)$$

This flow is invertible and the noise ϵ is given by:

$$\epsilon_i = \frac{y_i - \mu_i(\mathbf{y}_{0:i-1})}{\sigma(\mathbf{y}_{0:i-1})} \quad (23)$$

Note that, *epsilon* is independent to each other so the calculation of Equation 23 can be vectorized as follow:

$$\boldsymbol{\epsilon} = \frac{\mathbf{y} - \boldsymbol{\mu}(\mathbf{y})}{\boldsymbol{\sigma}(\mathbf{y})} \quad (24)$$

This enables an efficient computation with GPU.

Due to the autoregressive structure, the transformation has a lower triangular Jacobian where diagonal is σ_i . For calculation of normalizing flows, we are interested in log-determinant of the Jacobian which is just a product of the diagonal given as:

$$\log \det \left| \frac{d\boldsymbol{\epsilon}}{d\mathbf{y}} \right| = \sum_{i=0}^D -\log \sigma_i(\mathbf{y}) \quad (25)$$

To apply Inverse Autoregressive Flows (IAF) for variational inference in VAE, we add IAF transforms after the latent variables \mathbf{z} and modify the likelihood to account for IAF transforms. Figure 2 shows the process of applying IAF to Variational Autoencoder.

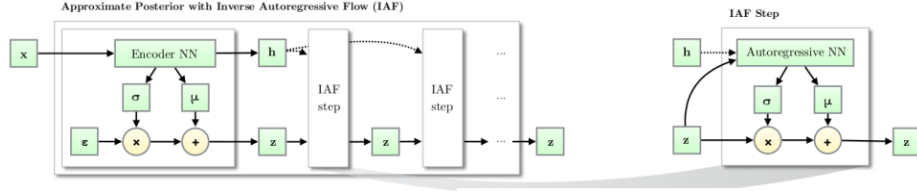


Fig. 2: The process of Inverse Autoregressive Flows in Variational Autoencoder [6].

The flow consists of a chain of T following transformations:

$$\mathbf{z}_t = \boldsymbol{\mu}_t + \boldsymbol{\sigma}_t \mathbf{z}_{t-1} \quad (26)$$

Kingma et.al. [6] proposed a more stable update based on LSTM-type update. LSTM is a type of recurrent neural network that applies autoregressive technique to the neural network. The update is given as follows:

$$\mathbf{s}_t = 1/\sigma_t \quad (27)$$

$$\mathbf{m}_t = -\boldsymbol{\mu}_t/\sigma_t \quad (28)$$

$$\sigma_t = \text{sigmoid}(\mathbf{s}_t) \quad (29)$$

$$\mathbf{z}_t = \sigma_t \mathbf{z}_{t-1} + (1 - \sigma_t) \mathbf{m}_t \quad (30)$$

where $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ learnt from an Autoregressive neural network given in [3] receiving input \mathbf{z} and \mathbf{h} from the VAE.

4.2 Density Estimation

Density Estimation techniques take a different approach from Variational Inference (VI) methods to model the complex data distribution. Unlike VI, these methods aim for exact inference, sampling and log-likelihood evaluation. The primary goal in this regime is to find a bijective function $h = f(x)$, $x \in \mathcal{X}$ to map complex data-distribution $p_X(x)$ to density $p_H(f(x))$. Given that $p_H()$ has a simpler density whose likelihood function is analytically known, the overall log-likelihood of the data can be easily calculated. The complex log-likelihood of the data can now be calculated using the change of variables as follows:

$$\log(p_X(x)) = \log(p_H(f(x))) + \log\left(\left|\det\left(\frac{\partial f(x)}{\partial x}\right)\right|\right) \quad (31)$$

Non-linear Independent Components Estimation

Real-valued Non-Volume Preserving

5 Normalizing Flows in Probabilistic Programming Languages

Normalizing flows has been implemented in mainstream probabilistic programming languages that is based on mainstream deep learning framework such as Pyro (based on Pytorch) and Tensorflow Probability (TFP, based on Tensorflow). They provide a construct called the **Bijector** for the implementation of a function in normalizing flow. We can create our own normalizing function by extending the **Bijector** construct. TFP also provides a construct called the **TransformedDistribution** which takes two input a **Bijector** and the base distribution. **TransformedDistribution** will run the flow and transform the given base distribution with the given **Bijector**. There are already numerous function (e.g. Affine, Sigmoid, RealNVP) implemented as a **Bijector** in TFP. The complete documentation of normalizing flows in TFP is given in this reference [1].

To create your own function, you need to extend the **Bijector** construct and overload these three methods. First, the `_forward` function receives samples x from the base distribution as an input and calculate the value of the function $y = f(x)$. Second, the `_inverse` function receives y from our function as the input and output x corresponds to the inverse of the function $x = f^{-1}(y)$. Third, the `_inverse_log_det_jacobian` function receives y from our function as the input and calculate the inverse log-determinant of the Jacobian $\log \det \left| \frac{\partial f^{-1}}{\partial y} \right|$. After you implement your own function, you can create a list of bijectors corresponds to the flow. TFP provides another bijector called **Chain** that will chain your list of bijectors into a flow. You can then create an instance of **TransformedDistribution** with this chain and your base distributor.

6 Recent Advances

In this section, we discuss several recent works that used normalizing flows. We only explain the high-level concept of each work. Interested reader should refer to the original paper given in the reference.

6.1 Pixel Recurrent Neural Network

Pixel Recurrent Neural Network (PixelRNN) [8] is an auto-regressive image generative model where the joint distribution over the image pixel is factorized into a product of conditional distribution. This means that the probability of pixel at position i is given as: $p(x_i | x_1, \dots, x_{i-1})$ where x_1, \dots, x_{i-1} is the previously generated pixels. This is a strong and counter-intuitive assumption for an image

generation where pixel mostly conditioned on their neighbors. However, PixelRNN is proven to work well for image completion and generation task. In image completion, the occluded pixel is generated by conditioned upon the non-occluded pixels.

Oord et. al. [8] proposed three methods model an auto-regressive image generation. First, Row LSTM 1D convolution is used to generate image row by row from top to bottom. However, Row LSTM cannot capture the whole previously generated pixels since there is a coarse-graining in the convolution method. Second, Bidirectional LSTM is used to capture all the generated pixels as the context. Every pixel is conditioned upon their neighboring pixels except the one that has not been generated. However, LSTM training is known to be expensive so the authors proposed another model based on CNN. Third, Pixel CNN used a masked convolution by setting the filter of the pixels that has not been generated as zero.

6.2 Wavenet

Wavenet [10] used the idea of Pixel CNN from [8] for raw audio generation. It used one-dimensional convolution Pixel CNN to generate raw audio. Wavenet has been proven to be the state-of-the-art model for text-to-speech model. Google Assistant in Android use wavenet to generate the audio of the assistant.

Similar to Pixel CNN, the joint probability of an input audio is modelled by a stack of causal convolutional layers. Causal convolution has similar idea to masked convolution by shifting the output of a normal convolution by a few timesteps so it does not violate the autoregressive requirement. One of the problems with causal convolution is that it requires deep layers to capture long range dependency. The authors proposed to use dilated convolution where the convolution filter is applied over an area larger than its length. This is done by skipping input values with a certain step.

6.3 Glow

Glow [5] extended and simplified the NICE and RealNVP model explained in Section 4.2 with three main extensions. First, the authors used activation normalization (act-norm) instead of batch normalization in RealNVP. Act-norm performs a channel-wise normalization and is faster than batch normalization. Second, the authors used 1×1 convolution for the channel-wise permutation. NICE and RealNVP proposed a flow containing the equivalent of a permutation that reverses the ordering of the channels. This is changed with a 1×1 convolution so that the permutation can be learned. Third, they extend RealNVP so that it can work faster by changing it into a multi-scale architecture. The authors splitted each step of the flow so it can be parallelized.

All of the additions mentioned before can be inverted and log-determinant of the Jacobian can be computed easily. Please refer to the paper for the details of the inversion and calculation of the Jacobian matrix. Glow has similar advantage as RealNVP that it can generate a high quality images. Moreover, the latent

space learned by Glow is meaningful. This means that we can control the output of the generated image by tuning the latent space.

7 Conclusion

References

1. Dillon, J.V., Langmore, I., Tran, D., Brevdo, E., Vasudevan, S., Moore, D., Patton, B., Alemi, A., Hoffman, M., Saurous, R.A.: Tensorflow distributions. arXiv preprint arXiv:1711.10604 (2017)
2. Figurnov, M., Mohamed, S., Mnih, A.: Implicit reparameterization gradients. arXiv preprint arXiv:1805.08498 (2018)
3. Germain, M., Gregor, K., Murray, I., Larochelle, H.: Made: Masked autoencoder for distribution estimation. In: International Conference on Machine Learning. pp. 881–889 (2015)
4. Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks. *science* **313**(5786), 504–507 (2006)
5. Kingma, D.P., Dhariwal, P.: Glow: Generative flow with invertible 1x1 convolutions. arXiv preprint arXiv:1807.03039 (2018)
6. Kingma, D.P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., Welling, M.: Improved variational inference with inverse autoregressive flow. In: Advances in Neural Information Processing Systems. pp. 4743–4751 (2016)
7. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013)
8. Oord, A.v.d., Kalchbrenner, N., Kavukcuoglu, K.: Pixel recurrent neural networks. arXiv preprint arXiv:1601.06759 (2016)
9. Rezende, D., Mohamed, S.: Variational inference with normalizing flows. In: International Conference on Machine Learning. pp. 1530–1538 (2015)
10. Van Den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A.W., Kavukcuoglu, K.: Wavenet: A generative model for raw audio. In: SSW. p. 125 (2016)

A Appendix (Contribution)

Abdul Fatir Ansari

- 1.

Devamanyu Hazarika

- 1.

Remmy A. M. Zen

1. Inverse autoregressive flows subsection.
2. Recent advances section.
3. Normalizing flows in PPL.

B Appendix (Abdul Fatir Ansari)

The Jacobian is then given by

$$\frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} = \mathbf{I} + \mathbf{u}h'(\mathbf{w}^\top \mathbf{z} + b)\mathbf{w}^\top$$

Now, using the matrix determinant lemma

$$\det \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} = (1 + h'(\mathbf{w}^\top \mathbf{z} + b)\mathbf{w}^\top \mathbf{I}^{-1} \mathbf{u}) \det(\mathbf{I}) \quad (32)$$

$$= (1 + h'(\mathbf{w}^\top \mathbf{z} + b)\mathbf{w}^\top \mathbf{u}) \quad (33)$$

C Appendix (Devamanyu Hazarika)**D Appendix (Remmy A. M. Zen)**