



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

The VAE was our first example of a generative model that is capable of sampling from $P(x)$. A VAE can also estimate $P(x)$ by going from the encoder to z , and then using the known distribution $P(z)$. However, VAEs typically are not great at either tasks. Their samples are often described as "blurry" because they have mean-seeking behavior. That is, the individual samples do not correspond to a very likely example (not mode-seeking). Instead the distribution of samples is good (mean-seeking). **Generative adversarial networks** (GANs) are an alternative to VAEs which have high probability samples. However, both methods suffer often have training problems and also give poor estimates of $P(x)$ because of both lack of normalization and assumptions of normal distributions. An alternative is a **normalizing flow*** that has better stability training and better estimates of $P(x)$. Normalizing flows are also used as components in other networks, like it can act as $P(z)$ of a latent space for a VAE instead of standard normal distributions.

NOTE:

You can construct interesting bijective functions between R^2 and R , but in the case of finite precision on a computer this is not possible so we can work with the slightly imprecise heuristic that functions need to have the same number of inputs and outputs to be bijective.

A **normalizing flow** is similar to a VAE in that we try to build up $P(x)$ by starting from a simple known distribution $P(z)$. We use functions, like the decoder from a VAE, to go from x to z . However, we make sure that the functions we choose keep the probability mass normalized ($\sum P(x) = 1$) and can be used forward (to sample from x) and backward (to compute $P(x)$). We call these functions **bijectors** because they are bijective (surjective and injective). Recall surjective (onto) means every output has a corresponding input and injective (into) means each output has exactly one corresponding input.

An example of a bijector is an element-wise cosine $y_i = \cos x_i$ (assuming x_i is between 0 and π). A non-bijective function would be $y_i = \cos x_i$ on the interval from 0 to 2π , because it outputs all values from $[0,1]$ twice and hence is not injective. Any function which changes the number of elements is automatically not bijective (see margin note). A consequence of using only bijectors in constructing our normalizing flow is that the size of the latent space must be equal to the size of the feature space. Remember the VAE used a smaller latent space than the feature space.

This lab sheet builds on VAE and assumes the same background of probability theory. This lab sheet is an introduction to the key ideas, but has many things beyond this lab sheet left to explore. Some knowledge of vector calculus (Jacobians) is assumed as well. After completing it, you should be able to:



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

- Understand the trade-offs between a VAE, GAN, and normalizing flow.
- Identify a bijector and construct a bijector chain
- Construct a normalizing flow using common bijectors types and train it
- Sample from a normalizing flow and compute sample probabilities

You can find a recent review of normalizing [flows here](<https://arxiv.org/pdf/1908.09257.pdf>). Although generating images and sound is the most popular application of normalizing flows, some of their biggest scientific impact has been on more efficient sampling from posteriors or likelihoods and other complex probability distributions. You find details on how to do normalizing flows on categorical (discrete) data in Hooeboom et al.

Flow Equation

Recall for the VAE decoder, we had an explicit formula for $p(x|z)$. This allowed us to compute $p(x) = \int dz p(x|z)p(z)$ which is the quantity of interest. The VAE decoder is a conditional probability density function. In the normalizing flow, we do not use probability density functions. We use bijective functions. So we cannot just compute an integral to change variables. We can use the change of variable formula. Consider our normalizing flow to be defined by our bijector $x = f(z)$, its inverse $z = g(x)$, and the starting probability distribution $P_z(z)$. Then the formula for probability of x is

$$P(x) = P_z(g(x)) |\det[\mathbf{J}_g]|$$

where the term on the right is the absolute value of the determinant of the Jacobian of g . [Jacobians](https://en.wikipedia.org/wiki/Jacobian_matrix_and_determinant) are matrices that describe how infinitesimal changes in each domain dimension change each range dimension. This term corrects for the volume change of the distribution. For example, if $f(z) = 2z$, then $g(x) = x/2$, and the Jacobian determinant is $1/2$. The intuition is that we are stretching out z by 2, so we need to account for the increase in volume to keep the probability normalized. You can read more about the change of variable formula for [probability distributions here](<https://cranmer.github.io/stats-ds-book/distributions/change-of-variables.html>)



Bijectors

A bijector is a function that is [injective](https://en.wikipedia.org/wiki/Injective_function) (1 to 1) and [surjective](https://en.wikipedia.org/wiki/Surjective_function) (onto). An equivalent way to view a bijective function is if it has an inverse. For example, a sum reduction has no inverse and is thus not bijective. $\sum[1,0] = 1$ and $\sum[-1,2] = 1$. Multiplying by a matrix which has an inverse is bijective. $y = x^2$ is not bijective, since $y = 4$ has two solutions.

Remember that we must compute the determinant of the bijector Jacobian. If the Jacobian is dense (all output elements depend on all input elements), computing this quantity will be $O(|x|_0^3)$ where $|x|_0$ is the number of dimensions of x because a determinant scales by $O(n^3)$. This would make computing normalizing flows impractical in high-dimensions. However, in practice we restrict ourselves to bijectors that have easy to calculate Jacobians. For example, if the bijector is $x_i = \cos z_i$ then the Jacobian will be diagonal. Such a diagonal Jacobian means that each dimension is independent of the other though.

One way to get faster determinants without just making each dimension independent is to get a triangular Jacobian. Then x_0 only depends on z_0 , x_1 depends on z_0, z_1 , and x_2 depends on z_0, z_1, z_2 , etc. *This enables fitting high-dimensional correlations for some of the dimensions (like x_n).* The matrix determinant of a triangular matrix is computed in linear time with respect to the number of dimensions -- because it is just the product of the matrix diagonal.

Bijector Chains

Just like in deep neural networks, multiple bijectors are chained together to increase how complex of the final fit distribution $\hat{P}(x)$ can be. The change of variable equation can be repeatedly applied:

$$P(x) = P_z[g_1(g_0(x))] |\det[\mathbf{J}_{g_1}]| |\det[\mathbf{J}_{g_0}]|$$

where we would compute x with $f_0(f_1(z))$. One critical point is that you should also include a **permute bijector** that swaps the order of dimensions. Since the bijectors typically have triangular Jacobians, certain output dimensions will depend on many input dimensions and others will only depend on a single one. By applying a permutation, you allow each dimension to influence each other.



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

Training

At this point, you may be wondering how you could possibly train a normalizing flow. The trainable parameters appear in the bijectors. They have adjustable parameters. The loss equation is quite simple: the negative log-likelihood (negative to make it minimization). Explicitly:

$$\ell = -\log P_z [g_1(g_0(x))] - \sum_i \log |\det[\mathbf{J}_{g_i}]|$$

where x is the training point and when you take the gradient of the loss, it is with respect to the parameters of the bijectors.

Common Bijectors

The choice of bijector functions is a fast changing area. We will thus only mention a few. You can of course use any bijective function or matrix, but these become inefficient at high-dimension due to the Jacobian calculation. One class of efficient bijectors are autoregressive bijectors. These have triangular Jacobians because each output dimension can only depend on the dimensions with a lower index. There are two variants: masked autoregressive flows (MAF) and inverse autoregressive flows (IAF). MAFs are efficient at training and computing probabilities, but are slow for sampling from $P(x)$. IAFs are slow at training and computing probabilities but efficient for sampling. Wavenets combine the advantages of both. We'll mention one other common bijector which is not autoregressive: real non-volume preserving (RealNVPs). RealNVPs are less expressive than IAFs/MAFs, meaning they have trouble replicating complex distributions, but are efficient at all three tasks: training, sampling, and computing probabilities. Another interesting variant is the Glow bijector, which is able to expand the rank of the normalizing flow, for example going from a matrix to an RGB image. What are the equations for all these bijectors? Most are variants of standard neural network layers but with special rules about which outputs depend on which inputs.

NOTE

Remember to add permute bijectors between autoregressive bijectors to ensure the dependence between dimensions is well-mixed.



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

```
import tensorflow as tf
import tensorflow_probability as tfp
import matplotlib.pyplot as plt
import sklearn.datasets as datasets
import numpy as np
```

```
np.random.seed(0)
tf.random.set_seed(0)
```

```
tfd = tfp.distributions
tfb = tfp.bijectors
```

Moon Example

We'll start with a basic 2D example to learn the two moons distribution with a normalizing flow. Two moons is a common example dataset that is hard to cluster and model as a probability distribution.

When doing normalizing flows you have two options to implement them. You can do all the Jacobians, inverses, and likelihood calculations analytically and implement them in a normal ML framework like Jax, PyTorch, or TensorFlow. This is actually most common. The second option is to utilize a probability library that knows how to use bijectors and distributions. The packages for that are PYMC, TensorFlow Probability (which has a non-tensorflow JAX version confusingly), and Pyro (Pytorch). We'll use TensorFlow Probability for this work.

Generating Data

In the code below, we set-up imports and sample points which will be used for training. Remember, this code has nothing to do with normalizing flows -- it's just to generate data.

```
moon_n = 10000
ndim = 2
data, _ = datasets.FUNCTION(moon_n, noise=0.05) # TODO: Replace
FUNCTION with appropriate function call to make moon
```



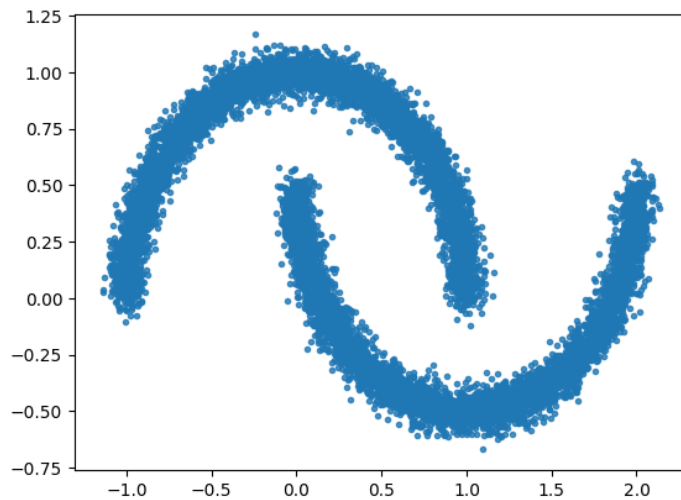
BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
plt.plot(data[:, 0], data[:, 1], ".", alpha=0.8)
```



Z Distribution

Our Z distribution should always be as simple as possible. We'll create a 2D Gaussian with unit variance, no covariance, and centered at 0. We'll be using the tensorflow probability package for this example. The key new concept is that we organize our tensors that were *sampled* from a probability distribution in a specific way. We, by convention, make the first axes be the 'sample' shape, the second axes be the 'batch' shape, and the final axes be the 'event' shape. The sample shape is the number of times we sampled from our distribution. It is a *shape* and not a single dimension because occasionally you'll want to organize your samples into some shape. The batch shape is a result of possibly multiple distributions batched together. For example, you might have 2 Gaussians, instead of a single 2D Gaussian. Finally, the event shape is the shape of a single sample from the distribution. This is overly complicated for our example, but you should be able to read information about the distribution now by understanding this nomenclature. You can find a tutorial on these [shapes here](https://www.tensorflow.org/probability/examples/Understanding_TensorFlow_Distributions_Shapes) and more tutorials on [tensorflow probability here](https://www.tensorflow.org/probability/examples/A_Tour_of_TensorFlow_Probability).

```
zdist = tfd.MultivariateNormalDiag(loc=[0.0] * ndim)
zdist
```

With our new understanding of shapes, you can see that this distribution has no 'batch_shape' because there is only one set of parameters and the 'event_shape' is '[2]' because it's a 2D Gaussian. Let's now sample from this distribution and view it

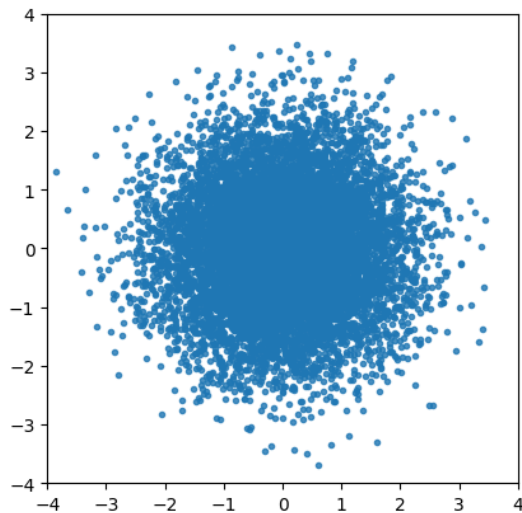


BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

```
zsamples = FUNCTION(moon_n) # TODO: Make call to z-distribution  
sampling function (Refer docs)  
plt.plot(zsamples[:, 0], zsamples[:, 1], ".", alpha=0.8)  
plt.xlim(-4, 4)  
plt.ylim(-4, 4)  
plt.gca().set_aspect("equal")
```



As expected, our starting distribution looks nothing like our target distribution. Let's demonstrate a bijector now. We'll implement the following bijector:

$$x = \vec{z} \times (1, 0.5)^T + (0.5, 0.25)$$

This is bijective because the operations are element-wise and invertible. Rather than just write this out using operations like `@` or `*`, we'll use the built-in bijectors from TensorFlow probability. The reason we do this is that they have their inverses and Jacobian determinants already defined. We first create a bijector that scales by `[1, 0.5]` and one then one that shifts by `[0.5, 0.25]`.

```
shift_bij = ... # TODO  
scale_bij = ... # TODO  
# Hint: Use tfd, tfb  
# make composite via function convolution  
b = shift_bij(scale_bij)
```

To now apply the change of variable formula, we create a **transformed distribution**. What is important about this choice is that we can compute likelihoods, probabilities, and sample from it.



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

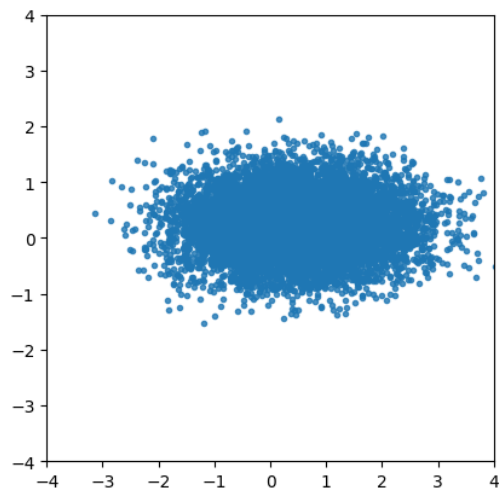
GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
td = tfd.TransformedDistribution(zdist, bijector=b)
td
```

```
zsamples = td.sample(moon_n)
plt.plot(zsamples[:, 0], zsamples[:, 1], ".", alpha=0.8)
plt.xlim(-4, 4)
plt.ylim(-4, 4)
plt.gca().set_aspect("equal")
```



We show above the sampling from this new distribution. We can also plot it's probability, which is impossible for a VAE-like model!



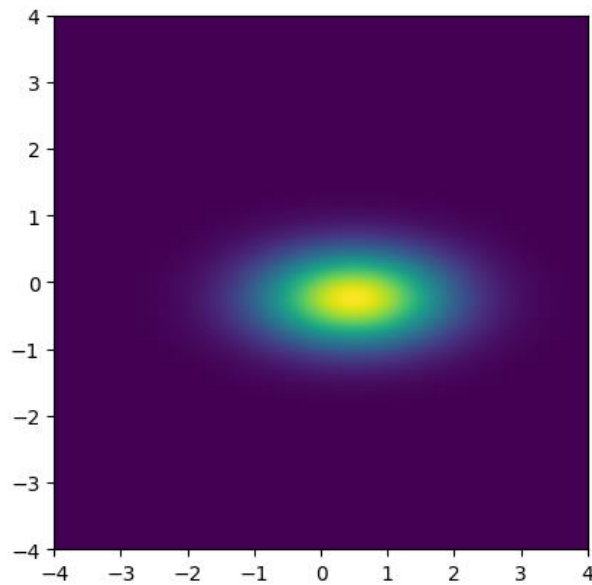
BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
# make points for grid
zpoints = np.linspace(-4, 4, 150)
(
    z1,
    z2,
) = np.meshgrid(zpoints, zpoints)
zgrid = np.concatenate((z1.reshape(-1, 1), z2.reshape(-1, 1)), axis=1)
# compute P(x)
p = np.exp(td.log_prob(zgrid))
fig = plt.figure()
# plot and set axes limits
plt.imshow(p.reshape(z1.shape), aspect="equal", extent=[-4, 4, -4, 4])
plt.show()
```



The Normalizing Flow

Now we will build bijectors that are expressive enough to capture the moon distribution. I will use 3 sets of a MAF and permutation for 6 total bijectors. MAF's have dense neural network layers in them, so I will also set the usual parameters for a neural network: dimension of hidden layer and activation.



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
num_layers = 3
my_bijects = []
# loop over desired bijectors and put into list
for i in range(num_layers):
    # Syntax to make a MAF
    anet = ...
    # TODO: Autoregressive Network with hidden units [128, 128] and
    relu activation

    ab = tfb.MaskedAutoregressiveFlow(anet)
    # Add bijector to list
    my_bijects.append(ab)
    # Now permuate (!important!)
    permute = ... # TODO: tfb Permute
    my_bijects.append(permute)
# put all bijectors into one "chain bijector"
# that looks like one
big_bijector = ... # TODO: tfb Chain
# make transformed dist
td = tfd.TransformedDistribution(zdist, bijector=big_bijector)
```

At this point, we have not actually trained but we can still view our distribution.



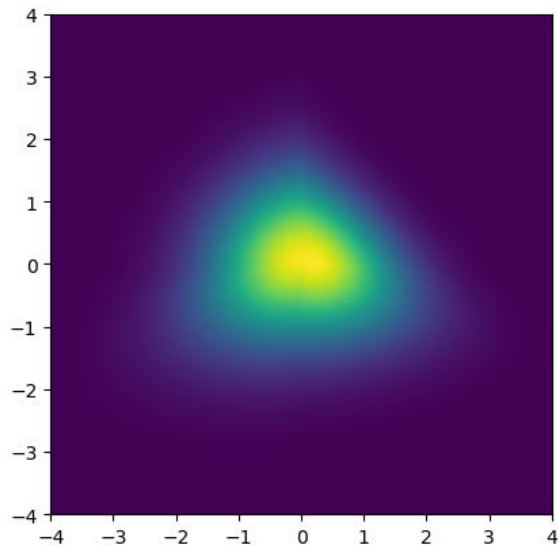
BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
zpoints = np.linspace(-4, 4, 150)
(
    z1,
    z2,
) = np.meshgrid(zpoints, zpoints)
zgrid = np.concatenate((z1.reshape(-1, 1), z2.reshape(-1, 1)), axis=1)
p = np.exp(td.log_prob(zgrid))
fig = plt.figure()
plt.imshow(p.reshape(z1.shape), aspect="equal", extent=[-4, 4, -4, 4])
plt.show()
```



You can already see that the distribution looks more complex than a Gaussian.

Training

To train, we'll use TensorFlow Keras, which just handles computing derivatives and the optimizer.



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

```
# declare the feature dimension
x = tf.keras.Input(...) # TODO: Appropriate shape and dtype
# create a "placeholder" function
# that will be model output
log_prob = td.log_prob(x)
# use input (feature) and output (log prob)
# to make model
model = ... # TODO: tf keras Model

# define a loss
def neg_loglik(yhat, log_prob):
    # losses always take in label, prediction
    # in keras. We do not have labels,
    # but we still need to accept the arg
    # to comply with Keras format
    return -log_prob

# now we prepare model for training
model.compile(optimizer=...), loss=...) # TODO: Adam optimizer,
appropriate loss function call
```

One detail is that we have to create fake labels (zeros) because Keras expects there to always be training labels. Thus our loss we defined above (negative log-likelihood) takes in the labels but does nothing with them.

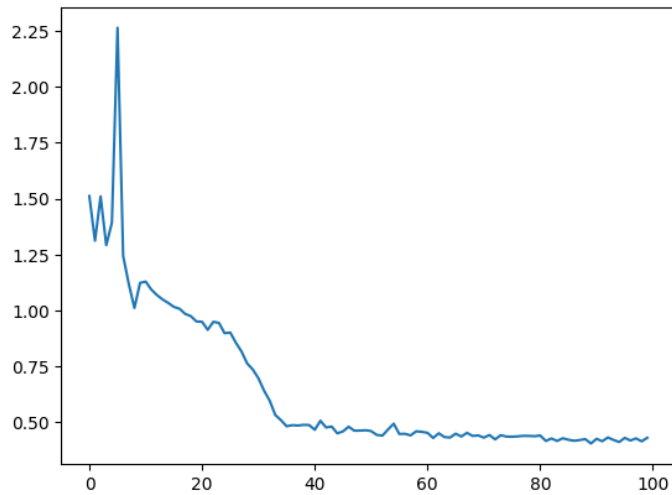


BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB [NORMALIZING FLOWS]

TIME: 02 Hours

```
result = ...  
# TODO: Train with x=data, y from moon_n, for 80 epochs  
plt.plot(result.history["loss"])  
plt.show()
```



Training looks reasonable. Let's now see our distribution.



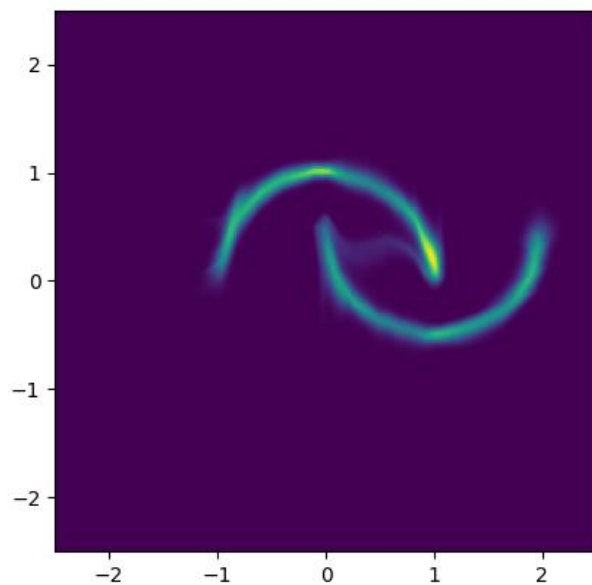
BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
zpoints = np.linspace(-2.5, 2.5, 200)
(
    z1,
    z2,
) = np.meshgrid(zpoints, zpoints)
zgrid = np.concatenate((z1.reshape(-1, 1), z2.reshape(-1, 1)), axis=1)
p = np.exp(td.log_prob(zgrid))
fig = plt.figure()
plt.imshow(
    p.reshape(z1.shape), aspect="equal", origin="lower", extent=[-2.5,
2.5, -2.5, 2.5]
)
plt.show()
```



Wow! We now can compute the probability of any point in this distribution. You can see there are some oddities that could be fixed with further training. One issue that cannot be overcome is the connection between the two curves -- it is not possible to get fully disconnected densities. This is because of our requirement that the bijectors are invertible and volume preserving -- you can only squeeze volume so far but cannot completely disconnect. Some work has been done on addressing this issue by adding sampling to the flow and this gives more expressive normalizing flows.

Finally, we'll sample from our model just to show that indeed it is generative.



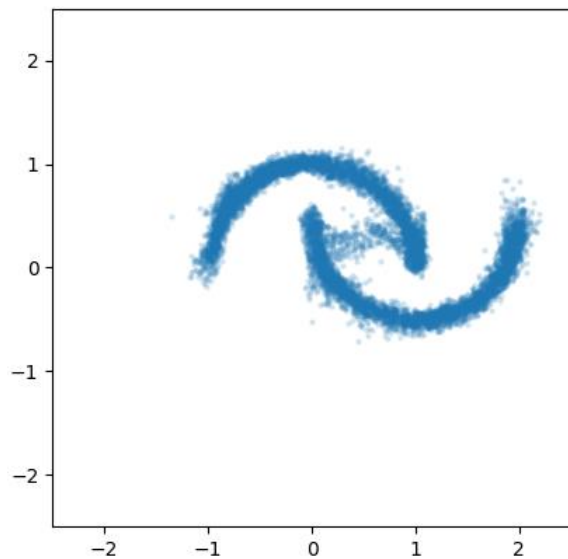
BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

```
zsamples = td.sample(moon_n)
plt.plot(zsamples[:, 0], zsamples[:, 1], ".", alpha=0.2,
         markeredgewidth=0.0)
plt.xlim(-2.5, 2.5)
plt.ylim(-2.5, 2.5)
plt.gca().set_aspect("equal")
```



Labsheet Summary

- A normalizing flow builds up a probability distribution of x by starting from a known distribution on z . Bijective functions are used to go from z to x .
- Bijections are functions that keep the probability mass normalized and are used to go forward and backward (because they have well-defined inverses).
- To find the probability distribution of x we use the change of variable formula, which requires a function inverse and Jacobian.
- The bijection function has trainable parameters, which can be trained using a negative log-likelihood function.
- Multiple bijections can be chained together, but typically must include a permute bijection to swap the order of dimensions.



BIRLA INSTITUTE OF TECHNOLOGY AND SCIENCE

GENERATIVE AI LAB

[NORMALIZING FLOWS]

TIME: 02 Hours

Exercise:

1. Complete all the TODOs given in the labsheet
2. Further reading – Attempt Normalizing Flows for Image Models (https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial_notebooks/tutorial11/NF_image_modeling.html), not required to be submitted, but give it a read