

Exercise 3 from GNN

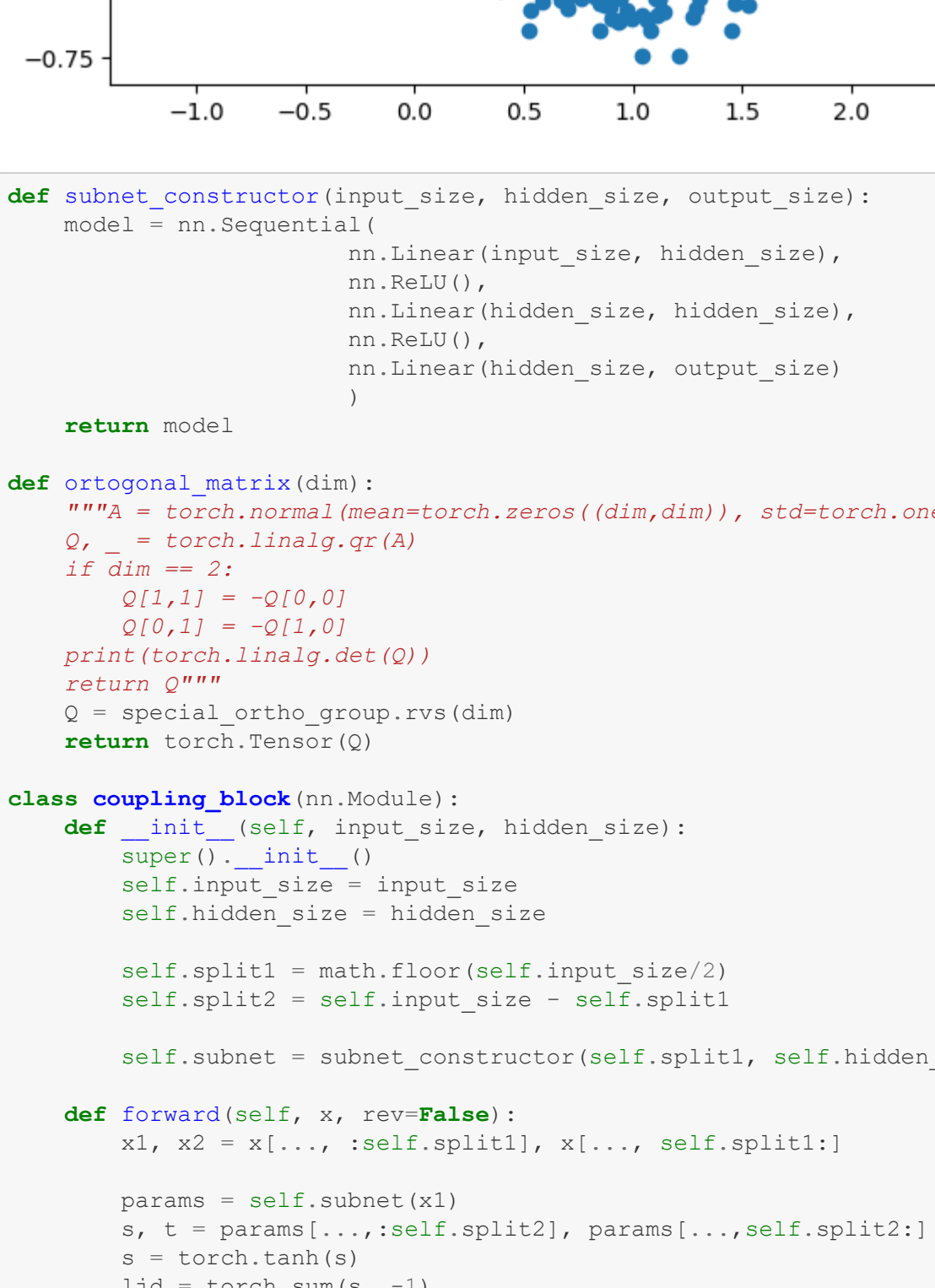
```
In [1]: import torch
import torch.nn as nn
import torch.nn.functional as F
from sklearn.datasets import make_moons, load_digits
from scipy.stats import special_ortho_group
import matplotlib.pyplot as plt
import random
```

1 Two Moons and Gaussian Mixture Model with INN

In this task we build a RealNVP invertible architecture to try and learn the two moons distribution. We inspect what influence model parameters and learning parameters have on the samples the networks produce. We check if the latent vectors which the network assigns to a Gaussian data is gaussian distributed. We also evaluate our synthetic data via MMD calculation. Once this is done, we also train our model on a Gaussian mixture distribution and compare the performance to the two moons case.

```
In [2]: test_x, _ = make_moons(n_samples=1000, shuffle=True, noise=0.1, random_state=42)
plt.scatter(test_x[:,0], test_x[:,1])

Out[2]: <matplotlib.collections.PathCollection at 0x2118cae5980>
```



```
In [3]: def subnet_constructor(input_size, hidden_size, output_size):
    m = nn.Sequential(
        nn.Linear(input_size, hidden_size),
        nn.ReLU(),
        nn.Linear(hidden_size, hidden_size),
        nn.ReLU(),
        nn.Linear(hidden_size, output_size)
    )
    return model

def orthogonal_matrix(dim):
    """ = torch.normal(mean=torch.zeros(dim,dim), std=torch.ones(dim,dim))
    Q, _ = torch.linalg.qr(A)
    if dim == 2:
        Q[0,1] = -Q[1,0]
        Q[1,1] = Q[0,0]
    print(torch.linalg.det(Q))
    return Q"""
    Q = special_ortho_group.randn(dim)
    return torch.Tensor(Q)

class coupling_block(nn.Module):
    def __init__(self, input_size, hidden_size):
        super().__init__()
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.split1 = math.floor(self.input_size/2)
        self.split2 = self.input_size - self.split1
        self.subnet = subnet_constructor(self.split1, self.hidden_size, 2*self.split2)

    def forward(self, x, rev=False):
        x1, x2 = x[:, :self.split1], x[:, self.split1:]
        params = self.subnet(x1)
        s, t = params[...,:self.split2], params[...,:self.split2:]
        s = torch.tanh(s)
        ljd = torch.randn(s, -1)
        if not rev:
            s = torch.exp(s)
            x2 = s*x2 + t
            return torch.cat([x1,x2], -1), ljd
        if rev:
            s = torch.exp(-s)
            x2 = s * (x2-t)
            return torch.cat([x1,x2], -1)

class realNVP(nn.Module):
    def __init__(self, input_size, hidden_size, n_blocks):
        super().__init__()
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.n_blocks = n_blocks
        self.coupling_blocks = nn.ModuleList([coupling_block(input_size, hidden_size) for _ in range(n_blocks)])
        self.orthogonal_matrices = [orthogonal_matrix(input_size) for _ in range(n_blocks-1)]

    def forward(self, x, rev=False):
        if rev: return self.inverse(x)
        return self.forward(x)

    def forward(self, x):
        ljd = torch.zeros(x.shape[0])
        for l in range(self.n_blocks-1):
            x, partial_ljd = self.coupling_blocks[l](x)
            ljd += partial_ljd
            x = torch.matmul(x, self.orthogonal_matrices[l])
            x, partial_ljd = self.coupling_blocks[-l](x)
            ljd += partial_ljd
        return x, ljd

    def inverse(self, self, x):
        for l in range(self.n_blocks-1, 0, -1):
            x = self.coupling_blocks[l](x, rev=True)
            x = torch.matmul(x, self.orthogonal_matrices[-l-1].T)
            x = self.coupling_blocks[0](x, rev=True)
        return x

def sample(self, num_samples):
    z = torch.normal(mean=torch.zeros(num_samples, self.input_size), std=torch.ones(num_samples, self.input_size))
    return self.inverse(z)
```

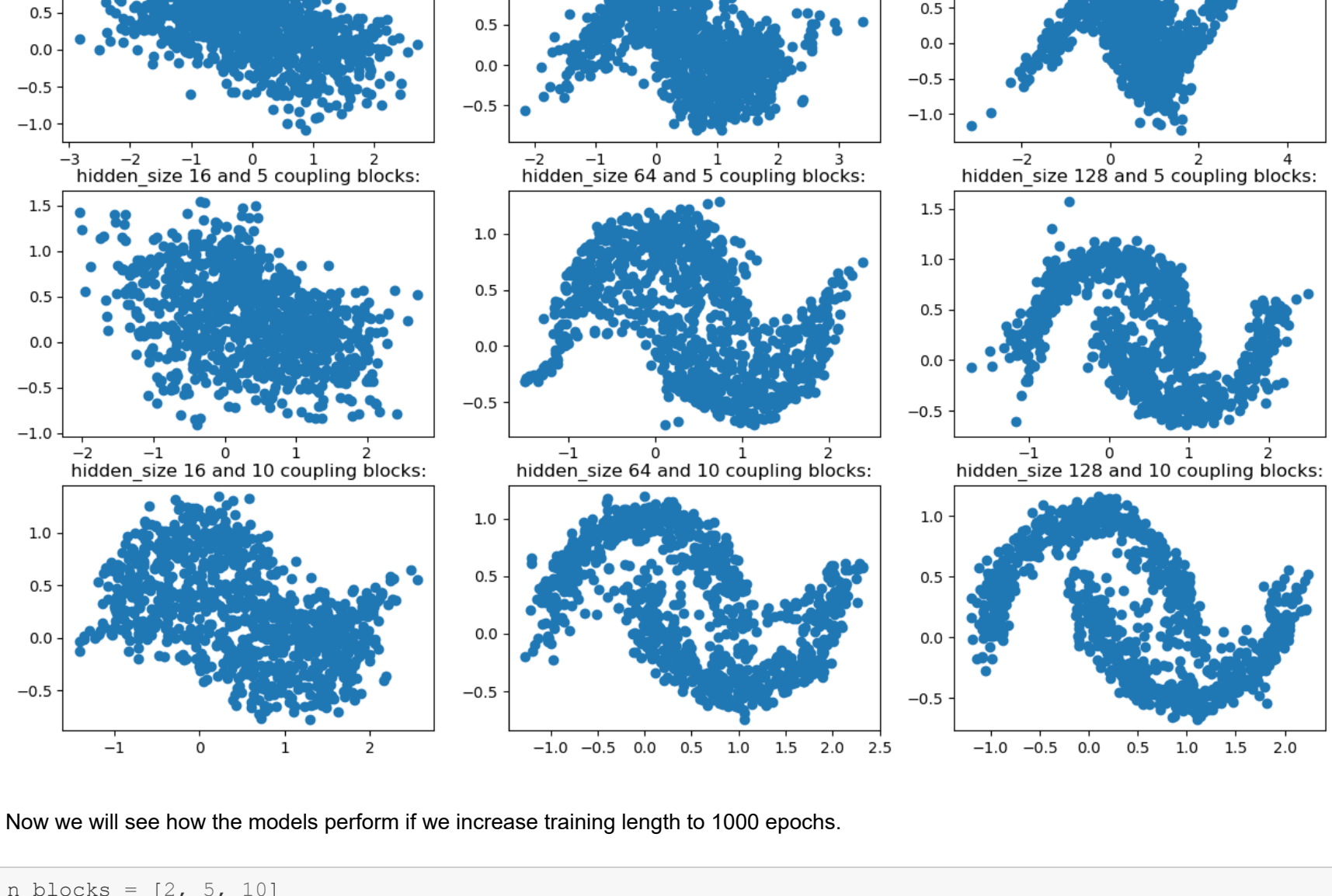
```
In [4]: def train_inn(model, batchsize=1000, epochs=1000, lr=0.001):
    optimizer = torch.optim.Adam(params=model.parameters(), lr=lr)
    for epoch in range(epochs):
        optimizer.zero_grad()
        x_data, y_data = make_moons(n_samples=batchsize, shuffle=True, noise=0.1, random_state=42)
        x_data, y_data = torch.Tensor(x_data), torch.Tensor(y_data)
        z, ljd = model(x_data)
        loss = torch.sum(0.5*torch.sum(z**2, -1)-ljd) / batchsize
        loss.backward()
        optimizer.step()
        if (epoch+1) % (epochs//3) == 0:
            print(f"Epoch {(epoch+1)/(epochs)}, Loss: {loss.item():.4f}")
```

```
In [5]: def plot_samples(model, title="Generated samples from INN", axes=None):
    samples = model.sample(1000)
    samples = samples.detach().numpy()
    if axes is None:
        fig, axs = plt.subplots(1,1)
        axs.scatter(samples[:,0], samples[:,1])
        axs.set_title(title)
```

First we will compare the effect of changing the model parameters (width of the coupling nets and the number of coupling blocks), while training for a fixed number of 100 epochs and a fixed learning rate.

```
In [6]: n_blocks = [2, 5, 10]
hidden_sizes = [16, 64, 128]
fig, axs = plt.subplots(3,3,figsize=(15,10))
for n in range(3):
    for h in range(3):
        model = realNVP(2, hidden_size[h], n_blocks[n])
        print(f"Training of INN with hidden_size {hidden_sizes[h]} and {n_blocks[n]} coupling blocks started.")
        train_inn(model, epochs=100, lr=0.001)
        plot_samples(model, title=f"hidden_size {hidden_sizes[h]} and {n_blocks[n]} coupling blocks:",
            axes=axs[n, h])
```

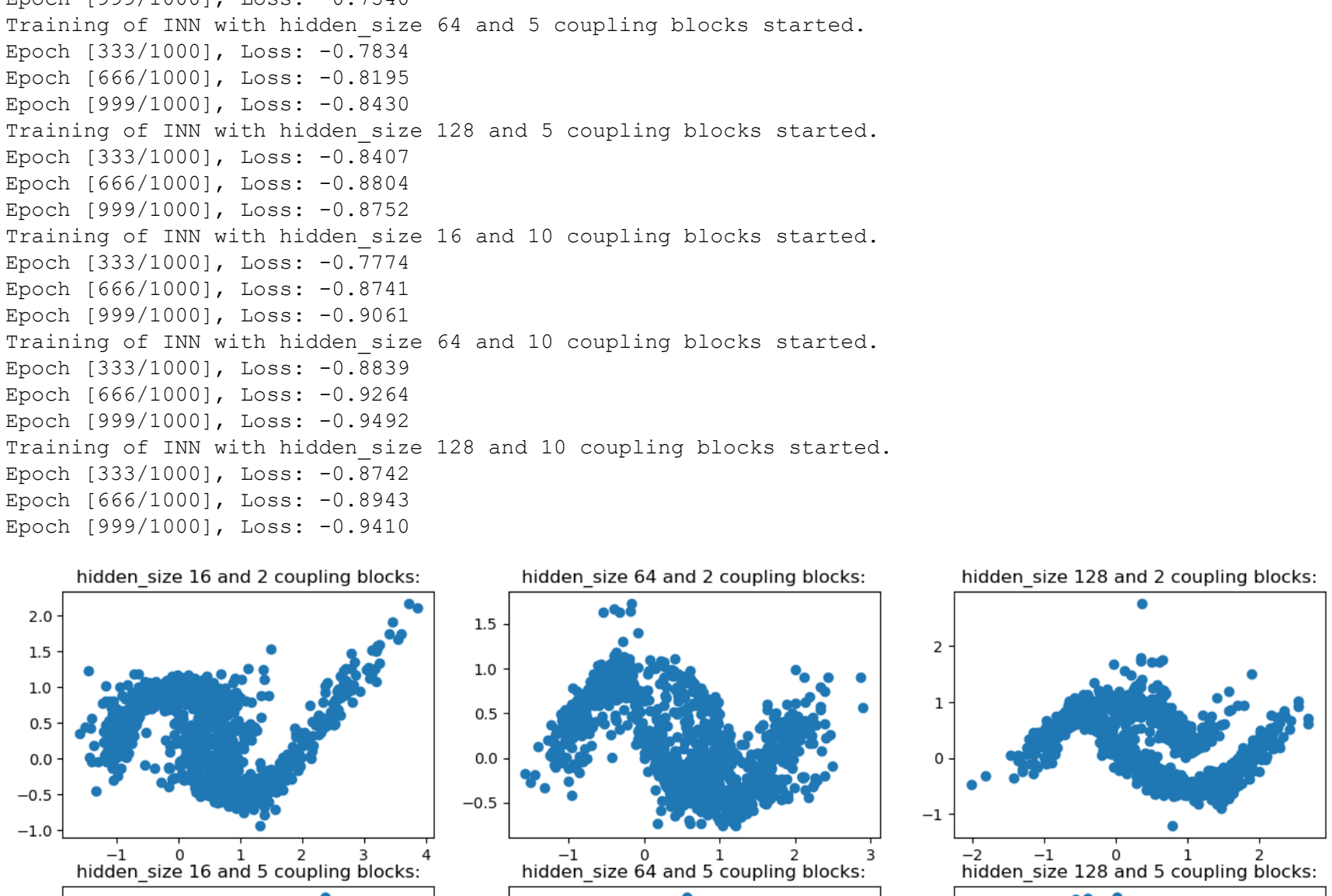
Training of INN with hidden_size 16 and 2 coupling blocks started.
Epoch (33/100), Loss: 0.2504
Epoch (66/100), Loss: 0.1504
Epoch (99/100), Loss: 0.0978
Epoch (33/100), Loss: 0.1423
Epoch (66/100), Loss: -0.2243
Epoch (99/100), Loss: -0.2700
Training of INN with hidden_size 64 and 2 coupling blocks started.
Epoch (33/100), Loss: 0.0403
Epoch (66/100), Loss: -0.1402
Epoch (99/100), Loss: -0.2243
Epoch (33/100), Loss: -0.1602
Epoch (66/100), Loss: -0.2389
Epoch (99/100), Loss: -0.2700
Training of INN with hidden_size 128 and 2 coupling blocks started.
Epoch (33/100), Loss: -0.1602
Epoch (66/100), Loss: 0.3824
Epoch (99/100), Loss: 0.0974
Epoch (33/100), Loss: -0.0433
Epoch (66/100), Loss: -0.1423
Epoch (99/100), Loss: -0.2243
Epoch (33/100), Loss: -0.4579
Epoch (99/100), Loss: -0.5681
Training of INN with hidden_size 128 and 5 coupling blocks started.
Epoch (33/100), Loss: -0.4114
Epoch (66/100), Loss: -0.6675
Epoch (99/100), Loss: -0.7544
Training of INN with hidden_size 16 and 10 coupling blocks started.
Epoch (33/100), Loss: -0.0796
Epoch (66/100), Loss: -0.1491
Epoch (99/100), Loss: -0.3319
Training of INN with hidden_size 64 and 10 coupling blocks started.
Epoch (33/100), Loss: -0.5040
Epoch (66/100), Loss: -0.7008
Epoch (99/100), Loss: -0.7645
Training of INN with hidden_size 128 and 10 coupling blocks started.
Epoch (33/100), Loss: -0.6532
Epoch (66/100), Loss: -0.8337
Epoch (99/100), Loss: -0.8337



Now we will see how the models perform if we increase training length to 1000 epochs.

```
In [7]: n_blocks = [2, 5, 10]
hidden_sizes = [16, 64, 128]
fig, axs = plt.subplots(3,3,figsize=(15,10))
for n in range(3):
    for h in range(3):
        model = realNVP(2, hidden_size[h], n_blocks[n])
        print(f"Training of INN with hidden_size {hidden_sizes[h]} and {n_blocks[n]} coupling blocks started.")
        train_inn(model, epochs=1000, lr=0.001)
        plot_samples(model, title=f"hidden_size {hidden_sizes[h]} and {n_blocks[n]} coupling blocks:",
            axes=axs[n, h])
```

Training of INN with hidden_size 16 and 2 coupling blocks started.
Epoch (333/1000), Loss: -0.2168
Epoch (666/1000), Loss: -0.3991
Epoch (999/1000), Loss: -0.4637
Training of INN with hidden_size 64 and 2 coupling blocks started.
Epoch (333/1000), Loss: -0.3990
Epoch (666/1000), Loss: -0.4374
Epoch (999/1000), Loss: -0.4461
Training of INN with hidden_size 128 and 2 coupling blocks started.
Epoch (333/1000), Loss: -0.5535
Epoch (666/1000), Loss: -0.5643
Epoch (999/1000), Loss: -0.5451
Training of INN with hidden_size 16 and 5 coupling blocks started.
Epoch (333/1000), Loss: -0.5159
Epoch (666/1000), Loss: -0.7038
Epoch (999/1000), Loss: -0.7540
Training of INN with hidden_size 64 and 5 coupling blocks started.
Epoch (333/1000), Loss: -0.7834
Epoch (666/1000), Loss: -0.8195
Epoch (999/1000), Loss: -0.8430
Training of INN with hidden_size 128 and 5 coupling blocks started.
Epoch (333/1000), Loss: -0.8407
Epoch (666/1000), Loss: -0.8804
Epoch (999/1000), Loss: -0.8752
Training of INN with hidden_size 16 and 10 coupling blocks started.
Epoch (333/1000), Loss: -0.7774
Epoch (666/1000), Loss: -0.8741
Epoch (999/1000), Loss: -0.9061
Training of INN with hidden_size 64 and 10 coupling blocks started.
Epoch (333/1000), Loss: -0.7084
Epoch (666/1000), Loss: -0.8839
Epoch (999/1000), Loss: -0.9264
Training of INN with hidden_size 128 and 10 coupling blocks started.
Epoch (333/1000), Loss: -0.9492
Epoch (666/1000), Loss: -0.8742
Epoch (999/1000), Loss: -0.8943
Epoch (333/1000), Loss: -0.9340
Epoch (666/1000), Loss: -0.9214
Epoch (999/1000), Loss: -0.9214

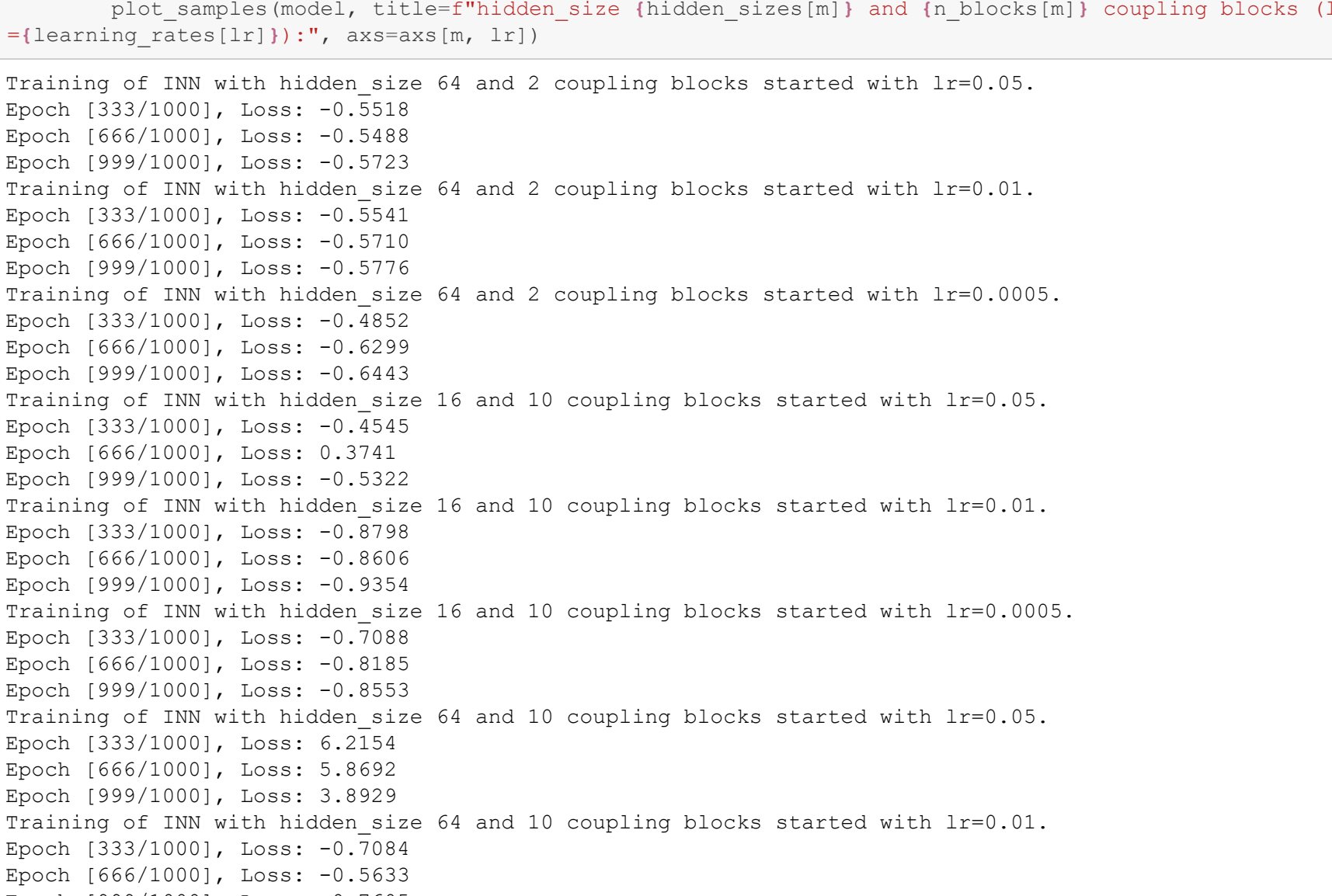


Interestingly, increasing the width of the coupling nets beyond 16 does not yield further improvements, but rather decreases the quality due to over-fitting. 10 seems to be a good choice for the number of coupling blocks.

We can also inspect changing the learning rate for training of different models.

```
In [8]: learning_rates = [0.05, 0.001, 0.0005]
hidden_sizes = [64, 16, 64]
n_blocks = [2, 10, 10]
fig, axs = plt.subplots(3,3,figsize=(15,10))
for lr in range(3):
    for h in range(3):
        model = realNVP(2, hidden_size[h], n_blocks[h])
        print(f"Training of INN with hidden_size {hidden_sizes[h]} and {n_blocks[h]} coupling blocks started with lr={learning_rates[lr]}")
        train_inn(model, epochs=1000, lr=learning_rates[lr])
        plot_samples(model, title=f"hidden_size {hidden_sizes[h]} and {n_blocks[h]} coupling blocks (lr={learning_rates[lr]}", axes=axs[lr, h])
```

Training of INN with hidden_size 64 and 2 coupling blocks started with lr=0.05.
Epoch (333/1000), Loss: -0.5518
Epoch (666/1000), Loss: -0.5991
Epoch (999/1000), Loss: -0.5723
Training of INN with hidden_size 64 and 2 coupling blocks started with lr=0.01.
Epoch (333/1000), Loss: -0.5541
Epoch (666/1000), Loss: -0.5710
Epoch (999/1000), Loss: -0.5772
Training of INN with hidden_size 64 and 2 coupling blocks started with lr=0.0005.
Epoch (333/1000), Loss: -0.6299
Epoch (666/1000), Loss: -0.6443
Epoch (999/1000), Loss: -0.6455
Training of INN with hidden_size 16 and 10 coupling blocks started with lr=0.05.
Epoch (333/1000), Loss: -0.3741
Epoch (666/1000), Loss: -0.5322
Epoch (999/1000), Loss: -0.6141
Training of INN with hidden_size 16 and 10 coupling blocks started with lr=0.01.
Epoch (333/1000), Loss: -0.8798
Epoch (666/1000), Loss: -0.8606
Epoch (999/1000), Loss: -0.9354
Training of INN with hidden_size 16 and 10 coupling blocks started with lr=0.0005.
Epoch (333/1000), Loss: -0.7088
Epoch (666/1000), Loss: -0.8185
Epoch (999/1000), Loss: -0.8533
Training of INN with hidden_size 64 and 10 coupling blocks started with lr=0.05.
Epoch (333/1000), Loss: 6.2154
Epoch (666/1000), Loss: 5.8692
Epoch (999/1000), Loss: -0.5633
Training of INN with hidden_size 64 and 10 coupling blocks started with lr=0.01.
Epoch (333/1000), Loss: -0.7084
Epoch (666/1000), Loss: -0.7605
Epoch (999/1000), Loss: -0.7460
Training of INN with hidden_size 64 and 10 coupling blocks started with lr=0.0005.
Epoch (333/1000), Loss: -0.8943
Epoch (666/1000), Loss: -0.8982
Epoch (999/1000), Loss: -0.9214



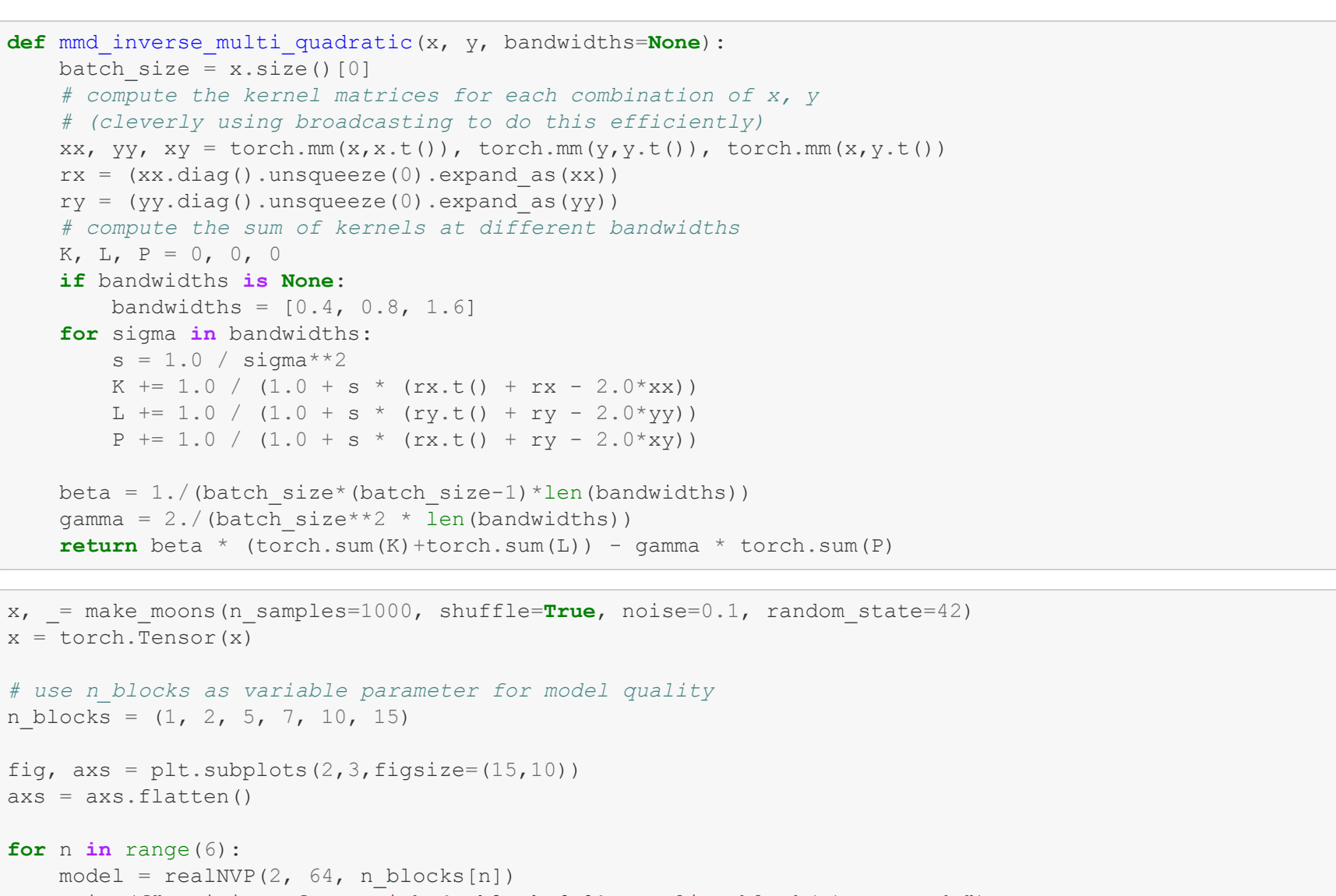
A learning rate which is too high can cause training instabilities as we see here. The default lr=0.001 seems to be a good choice.

Comparison of MMD scores with the function from the last sample solution:

```
In [9]: def mmd_inverse_multi_quadratic(x, y, bandwidth=None):
    batch_size = x.size()[0]
    # compute the kernel matrices for each combination of x, y
    # (cleverly using broadcasting to do this efficiently)
    xy, yy, xx = torch.mm(x,x.t()), torch.mm(y,y.t()), torch.mm(x,y.t())
    rx = (xx.diag().unsqueeze(0).expand_as(xx))
    ry = (yy.diag().unsqueeze(0).expand_as(yy))
    # compute the sum of kernels at different bandwidths
    Kx, Ly, P = 0, 0, 0
    for bandwidth in bandwidths:
        s = 1.0 / (bandwidth**2)
        P += 1.0 / (1.0 + s * (rx.t() + rx - 2.0*xx))
        L += 1.0 / (1.0 + s * (ry.t() + ry - 2.0*yy))
        P += 1.0 / (1.0 + s * (rx.t() + ry - 2.0*xy))
    beta = 1. / (batch_size*(batch_size-1)*len(bandwidths))
    gamma = 2. / (batch_size**2 * len(bandwidths))
    return beta * (torch.sum(Kx)*torch.sum(L) - gamma * torch.sum(P))
```

```
In [10]: x, _ = make_moons(n_samples=1000, shuffle=True, noise=0.1, random_state=42)
x = torch.Tensor(x)
# use n_blocks as variable parameter for model quality
n_blocks = (1, 2, 5, 7, 10, 15)
fig, axs = plt.subplots(2,3,figsize=(15,10))
axs = axs.flatten()
for n in range(6):
    model = realNVP(2, 64, n_blocks[n])
    print(f"Training of INN with {n_blocks[n]} coupling blocks started.")
    train_inn(model)
    samples = model.sample(1000)
    mmd = mmd_inverse_multi_quadratic(samples, x)
    plot_samples(model, title=f"{n_blocks[n]}-block INN, MMD = {mmd:.6f}", axes=axs[n])
```

Training of INN with 1 coupling block(s) started.
Epoch (333/1000), Loss: -0.0052
Epoch (666/1000), Loss: -0.0167
Epoch (999/1000), Loss: -0.0186
Training of INN with 2 coupling block(s) started.
Epoch (333/1000), Loss: -0.6141
Epoch (666/1000), Loss: -0.6483
Epoch (999/1000), Loss: -0.6556
Training of INN with 5 coupling block(s) started.
Epoch (333/1000), Loss: -0.8569
Epoch (666/1000), Loss: -0.8859
Epoch (999/1000), Loss: -0.8986
Training of INN with 7 coupling block(s) started.
Epoch (333/1000), Loss: -0.8649
Epoch (666/1000), Loss: -0.9079
Epoch (999/1000), Loss: -0.9407
Training of INN with 10 coupling block(s) started.
Epoch (333/1000), Loss: -0.8840
Epoch (666/1000), Loss: -0.9127
Epoch (999/1000), Loss: -0.9369
Training of INN with 15 coupling block(s) started.
Epoch (333/1000), Loss: -0.8846
Epoch (666/1000), Loss: -0.9174
Epoch (999/1000), Loss: -0.9660

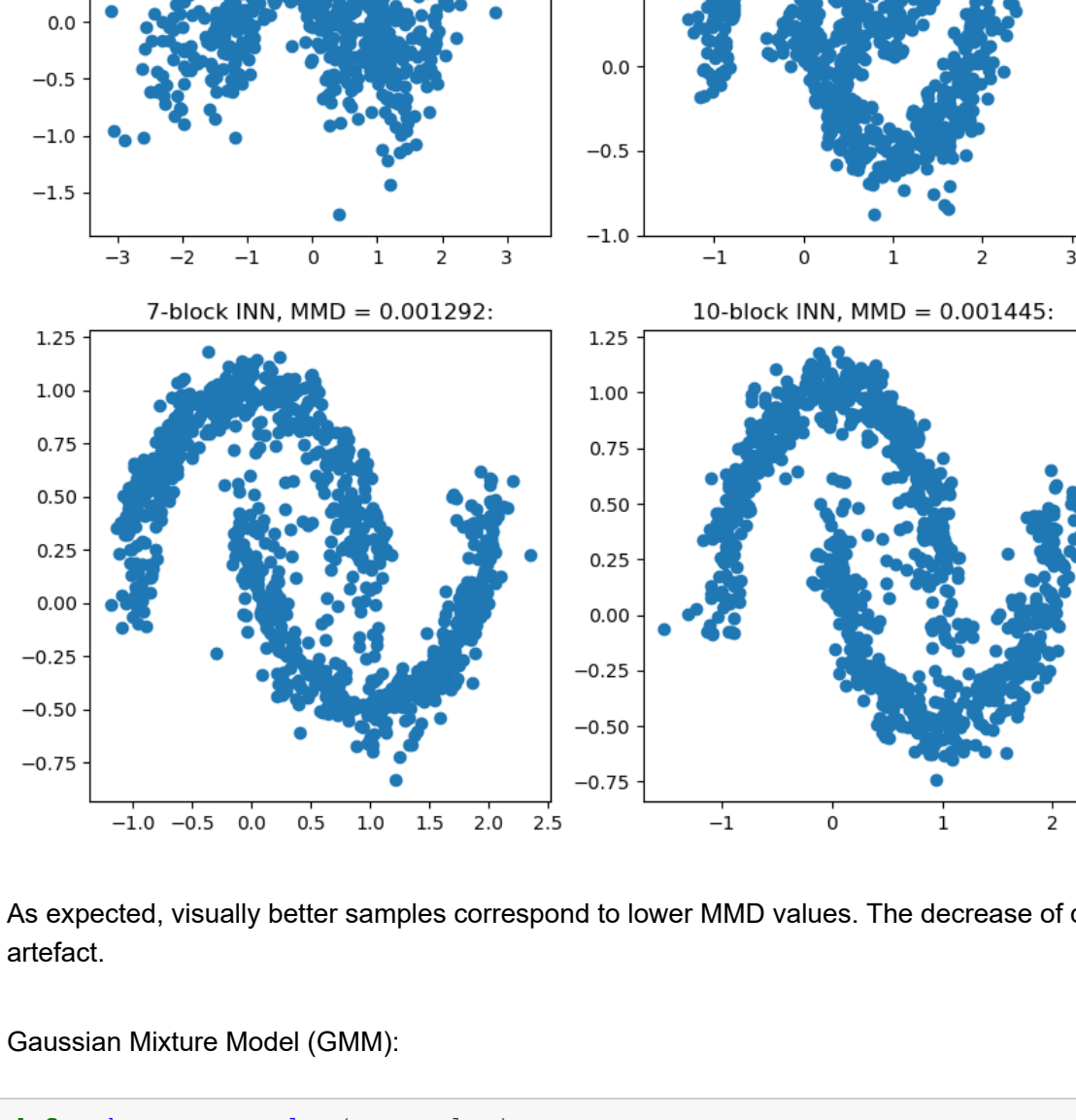


As expected, visually better samples correspond to lower MMD values. The decrease of quality for 10 coupling blocks seems to be a random artefact.

Gaussian Mixture Model (GMM):

```
In [11]: def make_gmm_samples(n_samples):
    a = 0.5*torch.exp(3)
    centers = torch.Tensor([[-1, 0],
        [-0.5, -a],
        [0.5, -a],
        [1, 0],
        [0.5, a],
        [-0.5, a]])
    stddev = torch.Tensor([0.1, 0.1])
    samples = []
    indices = []
    for i in range(n_samples):
        i = random.randint(0,5)
        indices.append(i)
        samples.append(torch.normal(centers[i], stddev))
    return torch.stack(samples), torch.Tensor(indices)
```

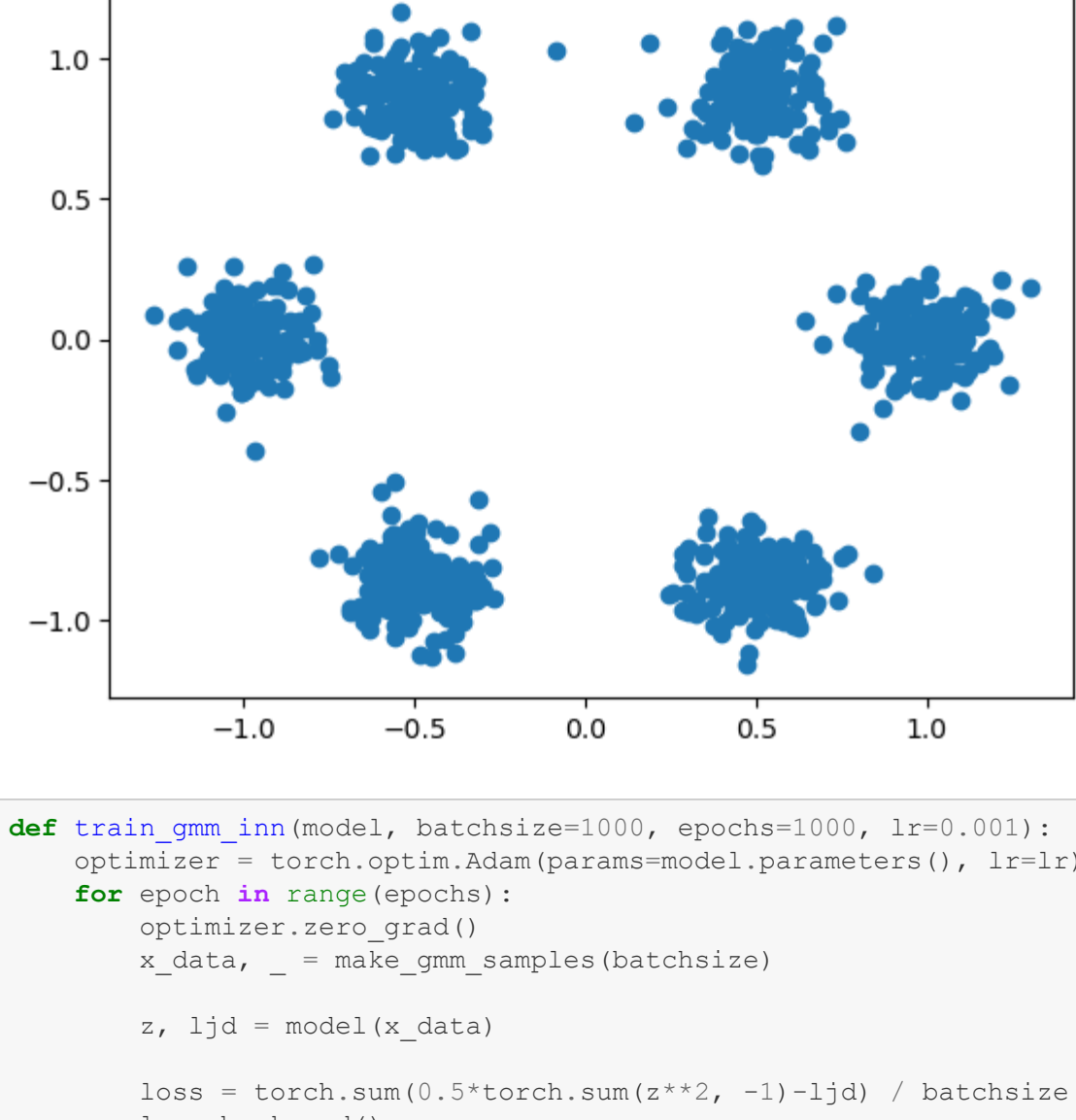
```
In [12]: samples, _ = make_gmm_samples(1000)
plt.scatter(samples[:,0], samples[:,1])
fig, axs = plt.subplots(2,3,figsize=(15,10))
axs = axs.flatten()
```



```
In [13]: def train_gmm_inn(model, batchsize=1000, epochs=1000, lr=0.001):
    optimizer = torch.optim.Adam(params=model.parameters(), lr=lr)
    for epoch in range(epochs):
        optimizer.zero_grad()
        x_data, _ = make_gmm_samples(batchsize)
        z, ljd = model(x_data)
        loss = torch.sum(0.5*torch.sum(z**2, -1)-ljd) / batchsize
        loss.backward()
        optimizer.step()
        if (epoch+1) % (epochs//3) == 0:
            print(f"Epoch {(epoch+1)/(epochs)}, Loss: {loss.item():.4f}")
```

```
In [14]: gmm_model = realNVP(2, 64, 10)
train_gmm_inn(gmm_model)
samples = gmm_model.sample(1000)
plt.scatter(samples.detach().numpy()[0:10], samples.detach().numpy()[1:1])
true_data, _ = make_gmm_samples(1000)
print(f"MMD of generated data with training data: {mmd_inverse_multi_quadratic(samples, true_data):.6f}")
```

Epoch (333/1000), Loss: -1.4501
Epoch (666/1000), Loss: -1.6455
Epoch (999/1000), Loss: -1.5911
MMD of generated data with training data: 0.003657



The INN is able to learn the distribution with no further changes. But let's see if we can improve the performance with a larger network and more training.

```
In [15]: n_epochs = [1000, 3000]
true_data, _ = make_gmm_samples(1000)
fig, axs = plt.subplots(1,2,figsize=(10, 4))
axs = axs.flatten()
larger_gmm_model = realNVP(2, 128, 15)
train_gmm_inn(larger_gmm_model, epochs=n_epochs[0])
samples = larger_gmm_model.sample(1000)
axs[0].scatter(samples[:,0], samples[:,1])
axs[1].set_title(f"Large INN, {n_epochs[0]} epochs, MMD = {mmd_inverse_multi_quadratic(samples, true_data):.6f}")
```

Epoch (333/1000), Loss: -1.4595
Epoch (666/1000), Loss: -1.6455
Epoch (999/1000), Loss: -1.5911
Epoch (1666/3000), Loss: -1.7460
Epoch (3332/3000), Loss: -1.6744
Epoch (4998/3000), Loss: -1.7179

Large INN, 1000 epochs, MMD = 0.004831: Large INN, 5000 epochs, MMD = 0.002602:



The larger architecture seems to be sufficient for an increase in visual quality, but further increasing the number of training epochs seems to have an effect on the amount of outliers. The MMD is not improved though.

2 Two Moons and GMM with a Conditional INN

In this exercise we are going to add functionality to our RealNVP model, to turn it into a conditional INN. We will then compare how the cINN performs when learning conditional distributions $p(x|y)$ in comparison to learning the marginalized distribution $p(x) = \sum_y p(x|y)$. We will inspect the two moon case, the full GMM case and the GMM case where we cluster together 2 and 4 of the modes to make it a binary classification task.

We reuse most of the code for our classes coupling_block and realNVP.

```
In [16]: class conditional_coupling_block(nn.Module):
    def __init__(self, input_size, hidden_size, condition_size):
        super().__init__()
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.condition_size = condition_size
        self.split1 = math.floor(self.input_size/2)
        self.split2 = self.input_size - self.split1
        self.subnet = subnet_constructor(self.split1 + self.condition_size, self.hidden_size, 2*self.split2)

    def forward(self, x, cond, rev=False):
        x1, x2 = x[:, :self.split1], x[:, self.split1:]
        params = self.subnet(torch.cat([x1, cond], -1))
        s, t = params[...,:self.split2], params[...,:self.split2:]
        s = torch.tanh(s)
        ljd = torch.randn(s, -1)
        if not rev:
            s = torch.exp(s)
            x2 = s*x2 + t
            return torch.cat([x1,x2], -1), ljd
        if rev:
            s = torch.exp(-s)
            x2 = s * (x2-t)
            return torch.cat([x1,x2], -1)

class conditional_realNVP(nn.Module):
    def __init__(self, input_size, hidden_size, n_blocks, condition_size):
        super().__init__()
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.n_blocks = n_blocks
        self.condition_size = condition_size
        self.coupling_blocks = nn.ModuleList([conditional_coupling_block(input_size, hidden_size, self.condition_size) for _ in range(n_blocks-1)])
        self.orthogonal_matrices = [orthogonal_matrix(input_size) for _ in range(n_blocks-1)]

    def forward(self, x, cond, rev=False):
        if rev: return self.inverse(x, cond)
        return self.forward(x, cond)

    def forward(self, x, cond):
        cond = nn.functional.one_hot(cond.to(torch.int64), num_classes=self.condition_size)
        ljd = torch.zeros(x.shape[0])
        for l in range(self.n_blocks-1):
            x, partial_ljd = self.coupling_blocks[l](x, cond)
            ljd += partial_ljd
            x = torch.matmul(x, self.orthogonal_matrices[l])
            x, partial_ljd = self.coupling_blocks[-l](x, cond)
            ljd += partial_ljd
        return x, ljd

    def inverse(self, self, x, cond):
        cond = nn.functional.one_hot(cond.to(torch.int64), num_classes=self.condition_size)
        for l in range(self.n_blocks-1, 0, -1):
            x = self.coupling_blocks[l](x, cond, rev=True)
            x = torch.matmul(x, self.orthogonal_matrices[-l-1].T)
            x = self.coupling_blocks[0](x, cond, rev=True)
        return x

def sample(self, num_samples, cond=None):
    samples = []
    if cond is None:
        z = torch.normal(mean=torch.zeros(num_samples, self.input_size), std=torch.ones(num_samples, self.input_size))
        samples.append(self.inverse(z, cond="torch.ones(num_samples)"))
    else:
        z = torch.normal(mean=torch.zeros(num_samples, self.input_size), std=torch.ones(num_samples, self.input_size))
        return torch.cat([self.inverse(z, cond=cond*torch.ones(num_samples))])
```

```
In [17]: def train_cinn_moons(model, batchsize=1000, epochs=1000, lr=0.001):
    optimizer = torch.optim.Adam(params=model.parameters(), lr=lr)
    for epoch in range(epochs):
        optimizer.zero_grad()
        x_data, y_data = make_moons(n_samples=batchsize, shuffle=True, noise=0.1, random_state=42)
        x_data, y_data = torch.Tensor(x_data), torch.Tensor(y_data)
        z, ljd = model(x_data, y_data)
        loss = torch.sum(0.5*torch.sum(z**2, -1)-ljd) / batchsize
        loss.backward()
        optimizer.step()
        if (epoch+1) % (epochs//3) == 0:
            print(f"Epoch {(epoch+1)/(epochs)}, Loss: {loss.item():.4f}")
```