

▼ VQ-VAE by Aäron van den Oord et al. in PyTorch

Introduction

Variational Auto Encoders (VAEs) can be thought of as what all but the last layer of a neural network is doing, namely feature extraction or separating out the data. Thus given some data we can think of using a neural network for representation generation.

Recall that the goal of a generative model is to estimate the probability distribution of high dimensional data such as images, videos, audio or even text by learning the underlying structure in the data as well as the dependencies between the different elements of the data. This is very useful since we can then use this representation to generate new data with similar properties. This way we can also learn useful features from the data in an unsupervised fashion.

The VQ-VAE uses a discrete latent representation mostly because many important real-world objects are discrete. For example in images we might have categories like "Cat", "Car", etc. and it might not make sense to interpolate between these categories. Discrete representations are also easier to model since each category has a single value whereas if we had a continuous latent space then we will need to normalize this density function and learn the dependencies between the different variables which could be very complex.

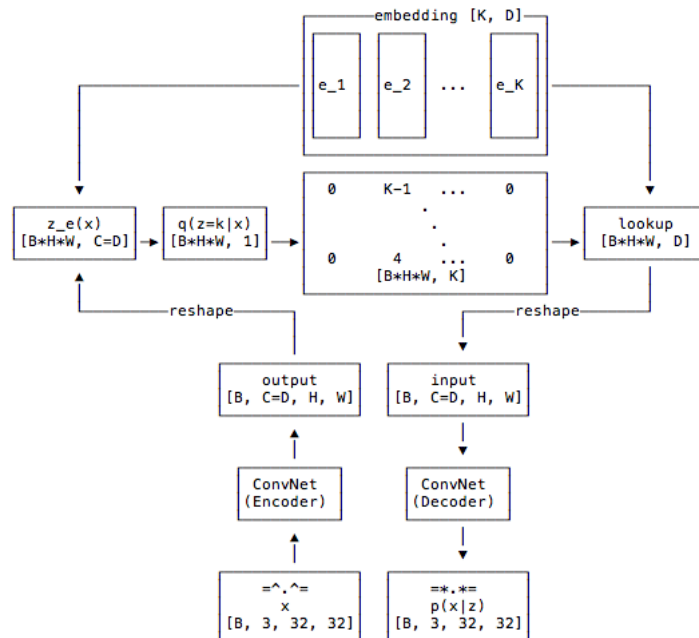
Code

I have followed the code from the TensorFlow implementation by the author which you can find here [vqvae.py](#) and [vqvae_example.ipynb](#).

Another PyTorch implementation is found at [pytorch-vqvae](#).

Basic Idea

The overall architecture is summarized in the diagram below:



We start by defining a latent embedding space of dimension $[K, D]$ where K are the number of embeddings and D is the dimensionality of each latent embedding vector, i.e. $e_i \in \mathbb{R}^D$. The model is comprised of an encoder and a decoder. The encoder will map the input to a sequence of discrete latent variables, whereas the decoder will try to reconstruct the input from these latent sequences.

More precisely, the model will take in batches of RGB images, say x , each of size 32×32 for our example, and pass it through a ConvNet encoder producing some output $E(x)$, where we make sure the channels are the same as the dimensionality of the latent embedding vectors. To calculate the discrete latent variable we find the nearest embedding vector and output its index.

The input to the decoder is the embedding vector corresponding to the index which is passed through the decoder to produce the reconstructed image.

Since the nearest neighbour lookup has no real gradient in the backward pass we simply pass the gradients from the decoder to the encoder unaltered. The intuition is that since the output representation of the encoder and the input to the decoder share the same D channel dimensional space, the gradients contain useful information for how the encoder has to change its output to lower the reconstruction loss.

▼ Loss

The total loss is actually composed of three components

1. **reconstruction loss**: which optimizes the decoder and encoder
2. **codebook loss**: due to the fact that gradients bypass the embedding, we use a dictionary learning algorithm which uses an l_2 error to move the embedding vectors e_i towards the encoder output
3. **commitment loss**: since the volume of the embedding space is dimensionless, it can grow arbitrarily if the embeddings e_i do not train as fast as the encoder parameters, and thus we add a commitment loss to make sure that the encoder commits to an embedding

TODOS:

1. Do PCA on the z's to see if it's low-dimensional
2. Maybe z's are in a union of low-dimensional subspaces.
3. Maybe z entries can be represented by sparse dictionary items.
4. What is the set of the purple vectors to allow efficient reconstruction?
5. The codebook can be quadratic on D instead of exponential.
6. Try the regular dictionary learning first.

```
1 !nvidia-smi
```

```
Tue Oct 10 23:23:07 2023
+-----+
| NVIDIA-SMI 525.105.17    Driver Version: 525.105.17    CUDA Version: 12.0    |
+-----+-----+-----+-----+-----+-----+
| GPU   Name               Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
|                                           MIG M. |
+-----+-----+-----+-----+-----+-----+
|   0   NVIDIA A100-SXM...    Off      | 00000000:00:04.0 Off |                    0 |
| N/A   33C    P0      47W / 400W |  0MiB / 40960MiB |           0%    Default |
+-----+-----+-----+-----+-----+-----+
|
+-----+-----+-----+-----+-----+-----+
| Processes: |
| GPU   GI    CI          PID    Type    Process name                  GPU Memory |
|   ID   ID                 |                   |            Usage              |
+-----+-----+-----+-----+-----+-----+
|  No running processes found |
+-----+-----+-----+-----+-----+-----+
|
```

```
1 !pip install umap-learn
```

```
2 !pip install six
```

```
Collecting umap-learn
  Downloading umap-learn-0.5.4.tar.gz (90 kB)
    90.8/90.8 kB 2.8 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from umap-learn) (1.23.5)
Requirement already satisfied: scipy>=1.3.1 in /usr/local/lib/python3.10/dist-packages (from umap-learn) (1.11.3)
Requirement already satisfied: scikit-learn>=0.22 in /usr/local/lib/python3.10/dist-packages (from umap-learn) (1.2.2)
Requirement already satisfied: numba>=0.51.2 in /usr/local/lib/python3.10/dist-packages (from umap-learn) (0.56.4)
Collecting pynndescent>=0.5 (from umap-learn)
  Downloading pynndescent-0.5.10.tar.gz (1.1 MB)
    1.1/1.1 MB 43.7 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from umap-learn) (4.66.1)
Requirement already satisfied: tbb>=2019.0 in /usr/local/lib/python3.10/dist-packages (from umap-learn) (2021.10.0)
Requirement already satisfied: llvmlite<0.40,>=0.39.0dev0 in /usr/local/lib/python3.10/dist-packages (from numba>=0.51.2->umap-learn) (0.40.0)
Requirement already satisfied: setuptools in /usr/local/lib/python3.10/dist-packages (from numba>=0.51.2->umap-learn) (67.7.2)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.10/dist-packages (from pynndescent>=0.5->umap-learn) (1.3.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn>=0.22->umap-learn) (3.2.0)
Building wheels for collected packages: umap-learn, pynndescent
  Building wheel for umap-learn (setup.py) ... done
  Created wheel for umap-learn: filename=umap_learn-0.5.4-py3-none-any.whl size=86770 sha256=fdf0f25b29fb1ea95ced08fc01551b86b9128893fd66
  Stored in directory: /root/.cache/pip/wheels/fb/66/29/199acf5784d0f7b8add6d466175ab45506c96e386ed5dd0633
  Building wheel for pynndescent (setup.py) ... done
  Created wheel for pynndescent: filename=pynndescent-0.5.10-py3-none-any.whl size=55615 sha256=0dd4af283afa169e3157f4ff66aebf589bf8a3b5
  Stored in directory: /root/.cache/pip/wheels/4a/38/5d/f60a40a66a9512b7e5e83517ebc2d1b42d857be97d135f1096
Successfully built umap-learn pynndescent
Installing collected packages: pynndescent, umap-learn
Successfully installed pynndescent-0.5.10 umap-learn-0.5.4
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (1.16.0)
```

```

1 from __future__ import print_function
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5 from scipy.signal import savgol_filter
6
7
8 from six.moves import xrange
9
10 import umap
11
12 import torch
13 import torch.nn as nn
14 import torch.nn.functional as F
15 from torch.utils.data import DataLoader
16 import torch.optim as optim
17
18 import torchvision.datasets as datasets
19 import torchvision.transforms as transforms
20 from torchvision.utils import make_grid
21
22 import time

1 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
2 print(device)

    cuda

```

▼ Load Data

```

1 training_data = datasets.CIFAR10(root="data", train=True, download=True,
2                                 transform=transforms.Compose([
3                                     transforms.ToTensor(),
4                                     transforms.Normalize((0.0,0.0,0.0), (1.0,1.0,1.0))
5                                 ]))
6
7 validation_data = datasets.CIFAR10(root="data", train=False, download=True,
8                                    transform=transforms.Compose([
9                                        transforms.ToTensor(),
10                                       transforms.Normalize((0.0,0.0,0.0), (1.0,1.0,1.0))
11                                   ]))

Downloading https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz to data/cifar-10-python.tar.gz
100%|██████████| 170498071/170498071 [00:13<00:00, 12992914.81it/s]
Extracting data/cifar-10-python.tar.gz to data
Files already downloaded and verified

1 # training_data = datasets.SVHN(root="data", split='train', download=True,
2 #                               transform=transforms.Compose([
3 #                                   transforms.ToTensor(),
4 #                                   transforms.Normalize((0.0,0.0,0.0), (1.0,1.0,1.0))
5 #                               ]))
6
7 # validation_data = datasets.SVHN(root="data", split='test', download=True,
8 #                                 transform=transforms.Compose([
9 #                                     transforms.ToTensor(),
10 #                                     transforms.Normalize((0.0,0.0,0.0), (1.0,1.0,1.0))
11 #                                 ]))

1 # size of the training data
2 training_data_size = len(training_data)
3 print(training_data_size)
4
5 # size of the validation data
6 validation_data_size = len(validation_data)
7 print(validation_data_size)

50000
10000

```

```

1 # variance of the training data
2 data_variance_training = np.var(training_data.data / 255.0)
3 print(data_variance_training)
4
5 # variance of the validation data
6 data_variance_validation = np.var(validation_data.data / 255.0)
7 print(data_variance_validation)

0.06328692405746414
0.06311123249672204

```

▼ Vector Quantizer Layer

This layer takes a tensor to be quantized. The channel dimension will be used as the space in which to quantize. All other dimensions will be flattened and will be seen as different examples to quantize.

The output tensor will have the same shape as the input.

As an example for a BCHW tensor of shape `[16, 64, 32, 32]`, we will first convert it to an BHCW tensor of shape `[16, 32, 32, 64]` and then reshape it into `[16384, 64]` and all 16384 vectors of size 64 will be **quantized independently**. In other words, the channels are used as the space in which to quantize. All other dimensions will be flattened and be seen as different examples to quantize, 16384 in this case.

```

1 class VectorQuantizer(nn.Module):
2     def __init__(self, num_embeddings, embedding_dim, commitment_cost):
3         super(VectorQuantizer, self).__init__()
4
5         self._embedding_dim = embedding_dim
6         self._num_embeddings = num_embeddings
7
8         self._embedding = nn.Embedding(self._num_embeddings, self._embedding_dim)
9         self._embedding.weight.data.uniform_(-1/self._num_embeddings, 1/self._num_embeddings)
10        self._commitment_cost = commitment_cost
11
12    def forward(self, inputs):
13        # convert inputs from BCHW -> BHCW
14        inputs = inputs.permute(0, 2, 3, 1).contiguous()
15        input_shape = inputs.shape
16
17        # Flatten input
18        flat_input = inputs.view(-1, self._embedding_dim)
19
20        # Calculate distances
21        distances = (torch.sum(flat_input**2, dim=1, keepdim=True)
22                     + torch.sum(self._embedding.weight**2, dim=1)
23                     - 2 * torch.matmul(flat_input, self._embedding.weight.t()))
24
25        # Encoding
26        encoding_indices = torch.argmax(distances, dim=1).unsqueeze(1)
27        encodings = torch.zeros(encoding_indices.shape[0], self._num_embeddings, device=inputs.device)
28        encodings.scatter_(1, encoding_indices, 1)
29
30        # Quantize and unflatten
31        quantized = torch.matmul(encodings, self._embedding.weight).view(input_shape)
32
33        # Loss
34        e_latent_loss = F.mse_loss(quantized.detach(), inputs)
35        q_latent_loss = F.mse_loss(quantized, inputs.detach())
36        loss = q_latent_loss + self._commitment_cost * e_latent_loss
37
38        quantized = inputs + (quantized - inputs).detach()
39        avg_probs = torch.mean(encodings, dim=0)
40        perplexity = torch.exp(-torch.sum(avg_probs * torch.log(avg_probs + 1e-10)))
41
42        # convert quantized from BHCW -> BCHW
43        return loss, quantized.permute(0, 3, 1, 2).contiguous(), perplexity, encodings

```

We will also implement a slightly modified version which will use exponential moving averages to update the embedding vectors (see appendix in the original paper) instead of an auxiliary loss. This has the **advantage** that the embedding updates are independent of the choice of optimizer for the encoder, decoder and other parts of the architecture. For most experiments the EMA version *trains faster* than the non-EMA version.

```

1 class VectorQuantizerEMA(nn.Module):
2     def __init__(self, num_embeddings, embedding_dim, commitment_cost, decay, epsilon=1e-5):
3         super(VectorQuantizerEMA, self).__init__()
4
5         self._embedding_dim = embedding_dim
6         self._num_embeddings = num_embeddings
7
8         self._embedding = nn.Embedding(self._num_embeddings, self._embedding_dim)
9         self._embedding.weight.data.normal_()
10        self._commitment_cost = commitment_cost
11
12        '''
13        Here the register_buffer() method is typically used to register a buffer
14        that should not to be considered a model parameter. For example,
15        BatchNorm's running_mean is not a parameter, but is part of the module's
16        state. Buffers, by default, are persistent and will be saved alongside
17        parameters. This behavior can be changed by setting persistent to False.
18        The only difference between a persistent buffer and a non-persistent
19        buffer is that the latter will not be a part of this module's state_dict.
20        '''
21        self.register_buffer('_ema_cluster_size', torch.zeros(num_embeddings)) # N
22        self._ema_w = nn.Parameter(torch.Tensor(num_embeddings, self._embedding_dim))
23        self._ema_w.data.normal_()
24
25        self._decay = decay # gamma term in the paper
26        self._epsilon = epsilon # numerical stability constant (laplace smoothing parameter)
27
28    def forward(self, inputs):
29        # convert inputs from BCHW -> BHWC
30        inputs = inputs.permute(0, 2, 3, 1).contiguous()
31        input_shape = inputs.shape
32
33        # Flatten input
34        flat_input = inputs.view(-1, self._embedding_dim)
35
36        # Calculate distances
37        distances = (torch.sum(flat_input**2, dim=1, keepdim=True)
38                    + torch.sum(self._embedding.weight**2, dim=1)
39                    - 2 * torch.matmul(flat_input, self._embedding.weight.t()))
40
41        # Encoding
42        encoding_indices = torch.argmax(distances, dim=1).unsqueeze(1)
43        encodings = torch.zeros(encoding_indices.shape[0], self._num_embeddings, device=inputs.device)
44        encodings.scatter_(1, encoding_indices, 1)
45
46        # Quantize and unflatten
47        quantized = torch.matmul(encodings, self._embedding.weight).view(input_shape)
48
49        # Use EMA to update the embedding vectors
50        if self.training:
51            self._ema_cluster_size = self._ema_cluster_size * self._decay + \
52                (1 - self._decay) * torch.sum(encodings, 0)
53
54            # Laplace smoothing of the cluster size
55            n = torch.sum(self._ema_cluster_size.data)
56            self._ema_cluster_size = (
57                (self._ema_cluster_size + self._epsilon)
58                / (n + self._num_embeddings * self._epsilon) * n) # N_i
59
60            dw = torch.matmul(encodings.t(), flat_input) # z_ij
61            self._ema_w = nn.Parameter(self._ema_w * self._decay + (1 - self._decay) * dw) # m_i
62
63            self._embedding.weight = nn.Parameter(self._ema_w / self._ema_cluster_size.unsqueeze(1)) # e_i
64
65        # Loss
66        e_latent_loss = F.mse_loss(quantized.detach(), inputs)
67        loss = self._commitment_cost * e_latent_loss
68
69        # Straight Through Estimator
70        quantized = inputs + (quantized - inputs).detach()
71        avg_probs = torch.mean(encodings, dim=0)
72        perplexity = torch.exp(-torch.sum(avg_probs * torch.log(avg_probs + 1e-10)))
73
74        # convert quantized from BHWC -> BCHW
75        return loss, quantized.permute(0, 3, 1, 2).contiguous(), perplexity, encodings

```

▼ The Dictionary Learning Algorithm

Dictionary learning can be thought of as the **vector space encoding** problem. Specifically, consider the input data matrix consisting of N data points, i.e.,

$$\mathbf{X} = [\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_N] \in \mathbb{R}^{D \times N}, \quad \mathbf{x}_n \in \mathbb{R}^D \text{ for } n = 1, 2, \dots, N$$

where D is the data dimension. We say that a **dictionary** is a set of K vectors (atoms) that lie in the same vector space as the original data vectors,

$$\mathbf{A} = [\mathbf{a}_1 | \mathbf{a}_2 | \dots | \mathbf{a}_K] \in \mathbb{R}^{D \times K}, \quad \mathbf{a}_k \in \mathbb{R}^D \text{ for } k = 1, 2, \dots, K$$

and K is the number of dictionary atoms. We then define the reconstruction error,

$$\|\mathbf{X} - \mathbf{A}\mathbf{R}\|_F^2,$$

where the learned representation $\mathbf{R} \in \mathbb{R}^{K \times N}$ can be thought of as the soft assignments (linear combination) of each point with the dictionary atoms,

$$\mathbf{R} = [\mathbf{r}_1 | \mathbf{r}_2 | \dots | \mathbf{r}_N] \in \mathbb{R}^{K \times N}, \quad \text{where } \mathbf{r}_n \in \mathbb{R}^K \text{ for } n = 1, 2, \dots, N.$$

Another thing to note that we would like the learned representations to be sparse, thus we use an ℓ_1 regularizer on \mathbf{R} to enforce sparsity, with sparsity level control parameter λ , yielding the objective function for the dictionary learning problem,

$$\|\mathbf{X} - \mathbf{A}\mathbf{R}\|_F^2 + \lambda \|\mathbf{R}\|_1.$$

```

1 class DictionaryLearning(nn.Module):
2     """A simple dictionary learning algorithm.
3
4     The algorithm is L1 regularized and optimized with SGD.
5
6     Attributes:
7         dim: data dimension D.
8         num_atoms: number of dictionary atoms K.
9         dictionary: dictionary matrix A with dimension K x D.
10        representation: representaion matrix R with dimension N x K,
11        N as the number of data samples
12        sparsity: sparsity control parameter, i.e., the lambda.
13    """
14    def __init__(self, dim, num_atoms, sparsity, sparsity_level) -> None:
15        super(DictionaryLearning, self).__init__()
16        self.dim = dim
17        self.num_atoms = num_atoms
18        self.dictionary = nn.Parameter(
19            torch.rand((num_atoms, dim), dtype=torch.float, requires_grad=True))
20        self.representation = self.representation_builder()
21        self.sparsity = sparsity
22        self.sparsity_level = sparsity_level
23
24    def representation_builder(self):
25        layers = nn.ModuleList()
26        layers.append(nn.Linear(self.dim, self.num_atoms)) # output dim for one sample: K
27        layers.append(nn.Softmax()) # convert the output to [0, 1] range
28
29        return nn.Sequential(*layers)
30
31    def loss(self, x, representation):
32        """Calculate the loss.
33
34        Args:
35            x: input data, for image data, the dimension is: N x C x H x W
36            representation: representation matrix R with dimension N x K,
37            N as the number of data samples
38
39        Returns:
40            combined loss, reconstruction, representation
41        """
42        # print(f'DEBUG:shape of the input: {x.shape}')
43        # print(f'DEBUG:shape of the representation: {representation.shape}')
44        # print(f'DEBUG:shape of the dictionary: {self.dictionary.shape}')
45        batch_size, num_channels, height, width = x.shape
46        reconstruction = torch.matmul(representation, self.dictionary)
47        reconstruction = reconstruction.view(batch_size, self.dim, height, width).contiguous()
48
49        reconstruction = reconstruction + (reconstruction - x).detach() # straight-through gradient
50        # recon_loss = nn.MSELoss()(x, reconstruction)
51        recon_loss = nn.MSELoss()(x, reconstruction.detach()) * 0.25 + nn.MSELoss()(x.detach(), reconstruction)
52        regularization = torch.sum(torch.abs(representation))
53

```

```

54     return recon_loss + self.sparsity * regularization, reconstruction, representation
55
56 def forward(self, x):
57     """Forward pass.
58
59     Args:
60         x: input data, for image data, the dimension is:
61             batch_size x num_channels x height x width
62     """
63     batch_size, num_channels, height, width = x.shape
64     x = x.permute(0, 2, 3, 1).contiguous()
65     # TODO: Might need to flatten the representation instead of x
66     x_flattened = x.view(-1, self.dim) # data dimension: N x D
67     # print(f'DEBUG:shape of the flattened input: {x_flattened.shape}')
68     # print(f'DEBUG:shape of the dictionary: {self.dictionary.shape}')
69     representation = self.representation(x_flattened)
70     # print(f'DEBUG:shape of the representation: {representation.shape}')
71     sparse_operator = torch.zeros(representation.shape, requires_grad=True).cuda()
72
73     # compute the l2 distances (N x K) between x_flattened (N x D) and the dictionary (K x D)
74     distances = -(torch.sum(x_flattened**2, dim=1, keepdim=True) + torch.sum(self.dictionary**2, dim=1) - 2 * torch.matmul(x_flattened, self.dictionary))
75
76     _, indices = torch.topk(distances, self.sparsity_level, dim=1)
77     # _, indices = torch.topk(distances, 1, dim=1)
78     # print(f'shape of the indices: {indices.shape}')
79
80     sparse_operator.scatter_(1, indices, 1)
81     representation = representation * sparse_operator
82
83     reconstruction = torch.matmul(representation, self.dictionary)
84     reconstruction = reconstruction.view(batch_size, self.dim, height, width).contiguous()
85
86     x = x.permute(0, 3, 1, 2).contiguous() # permute back
87     x = x.view(batch_size, self.dim, height, width)
88     # print(f'DEBUG:shape of the input: {x.shape}')
89     # recon_loss = nn.MSELoss()(x, reconstruction)
90     recon_loss = nn.MSELoss()(x, reconstruction.detach()) * 0.25 + nn.MSELoss()(x.detach(), reconstruction)
91     regularization = torch.sum(torch.abs(representation))
92
93     reconstruction = reconstruction + (reconstruction - x).detach() # straight-through gradient
94
95     # return representation
96
97     return recon_loss + self.sparsity * regularization, reconstruction, representation

```

▼ Encoder & Decoder Architecture

The encoder and decoder architecture is based on a ResNet and is implemented below:

```

1 class Residual(nn.Module):
2     def __init__(self, in_channels, num_hiddens, num_residual_hiddens):
3         super(Residual, self).__init__()
4         self._block = nn.Sequential(
5             nn.ReLU(True),
6             nn.Conv2d(in_channels=in_channels,
7                       out_channels=num_residual_hiddens,
8                       kernel_size=3, stride=1, padding=1, bias=False),
9             nn.ReLU(True),
10            nn.Conv2d(in_channels=num_residual_hiddens,
11                     out_channels=num_hiddens,
12                     kernel_size=1, stride=1, bias=False)
13        )
14
15    def forward(self, x):
16        return x + self._block(x)
17
18
19 class ResidualStack(nn.Module):
20     def __init__(self, in_channels, num_hiddens, num_residual_layers, num_residual_hiddens):
21         super(ResidualStack, self).__init__()
22         self._num_residual_layers = num_residual_layers
23         self._layers = nn.ModuleList([Residual(in_channels, num_hiddens, num_residual_hiddens)
24                                         for _ in range(self._num_residual_layers)])
25
26    def forward(self, x):

```

```

27         for i in range(self._num_residual_layers):
28             x = self._layers[i](x)
29         return F.relu(x)

1 class Encoder(nn.Module):
2     def __init__(self, in_channels, num_hiddens, num_residual_layers, num_residual_hiddens):
3         super(Encoder, self).__init__()
4
5         self._conv_1 = nn.Conv2d(in_channels=in_channels,
6                                   out_channels=num_hiddens//2,
7                                   kernel_size=4,
8                                   stride=2, padding=1) # output dim: (256, 64, 16, 16)
9         self._conv_2 = nn.Conv2d(in_channels=num_hiddens//2,
10                                   out_channels=num_hiddens,
11                                   kernel_size=4,
12                                   stride=2, padding=1) # output dim: (256, 128, 8, 8)
13         self._conv_3 = nn.Conv2d(in_channels=num_hiddens,
14                                   out_channels=num_hiddens,
15                                   kernel_size=3,
16                                   stride=1, padding=1) # output dim: (256, 128, 8, 8)
17         self._residual_stack = ResidualStack(in_channels=num_hiddens,
18                                               num_hiddens=num_hiddens,
19                                               num_residual_layers=num_residual_layers,
20                                               num_residual_hiddens=num_residual_hiddens)
21
22     def forward(self, inputs):
23         x = self._conv_1(inputs)
24         x = F.relu(x)
25
26         x = self._conv_2(x)
27         x = F.relu(x)
28
29         x = self._conv_3(x)
30         return self._residual_stack(x)

1 class Decoder(nn.Module):
2     def __init__(self, in_channels, num_hiddens, num_residual_layers, num_residual_hiddens):
3         super(Decoder, self).__init__()
4
5         self._conv_1 = nn.Conv2d(in_channels=in_channels,
6                                   out_channels=num_hiddens,
7                                   kernel_size=3,
8                                   stride=1, padding=1)
9
10        self._residual_stack = ResidualStack(in_channels=num_hiddens,
11                                              num_hiddens=num_hiddens,
12                                              num_residual_layers=num_residual_layers,
13                                              num_residual_hiddens=num_residual_hiddens)
14
15        self._conv_trans_1 = nn.ConvTranspose2d(in_channels=num_hiddens,
16                                                  out_channels=num_hiddens//2,
17                                                  kernel_size=4,
18                                                  stride=2, padding=1)
19
20        self._conv_trans_2 = nn.ConvTranspose2d(in_channels=num_hiddens//2,
21                                                  out_channels=3,
22                                                  kernel_size=4,
23                                                  stride=2, padding=1)
24
25    def forward(self, inputs):
26        x = self._conv_1(inputs)
27        x = self._residual_stack(x)
28
29        x = self._conv_trans_1(x)
30        x = F.relu(x)
31
32        return self._conv_trans_2(x)

```

▼ Train

We use the hyperparameters from the author's code:

```

1 batch_size = 256
2

```



```

3 num_hiddens = 128
4 num_residual_hiddens = 32
5 num_residual_layers = 2
6
7 embedding_dim = 256 # D
8 num_embeddings = 512 # number of e's, which might be much smaller for dictionary learning
9 sparsity = 0
10 sparsity_level = 100 # number of atoms selected
11
12 commitment_cost = 0.25
13
14 decay = 0 # set to 0 to disable EMA dictionary updates
15
16 learning_rate = 1e-4

1 training_loader = DataLoader(training_data,
2                               batch_size=batch_size,
3                               shuffle=True,
4                               pin_memory=True)

1 validation_loader = DataLoader(validation_data,
2                                 batch_size=32,
3                                 shuffle=False,
4                                 pin_memory=True)

1 class DictionaryLearningVAE(nn.Module):
2     def __init__(self, num_hiddens, num_residual_layers, num_residual_hiddens,
3                   num_embeddings, embedding_dim, sparsity=5, sparsity_level=0.01):
4         super(DictionaryLearningVAE, self).__init__()
5
6         # encoder input dimension: 3 (RGB image channel dimension),
7         # output dimension: num_hiddens
8         self._encoder = Encoder(3, num_hiddens,
9                                  num_residual_layers,
10                                 num_residual_hiddens)
11
12         # 2D convolution operation with kernel_size = 1 to change the size in the channel
13         # dimension (number of convolutional kernels) from 128 (for RGB images) to embedding_dim
14         # input channel dimension: num_hiddens (Encoder output dimension)
15         # output channel dimension: embedding_dim (dimension of the dictionary atoms)
16         self._pre_vq_conv = nn.Conv2d(in_channels=num_hiddens,
17                                         out_channels=embedding_dim,
18                                         kernel_size=1,
19                                         stride=1)
20
21         self._dl = DictionaryLearning(embedding_dim, num_embeddings, sparsity, sparsity_level)
22
23         self._decoder = Decoder(embedding_dim,
24                                   num_hiddens,
25                                   num_residual_layers,
26                                   num_residual_hiddens)
27     def forward(self, x):
28         z = self._encoder(x)
29         z = self._pre_vq_conv(z)
30         # representation = self._dl(z)
31         # dlloss, z_recon, representation = self._dl.loss(z, representation)
32
33         dlloss, z_recon, representation = self._dl(z)
34
35         x_recon = self._decoder(z_recon)
36
37         return dlloss, x_recon, representation

1 model = DictionaryLearningVAE(num_hiddens, num_residual_layers, num_residual_hiddens,
2                               num_embeddings, embedding_dim, sparsity=sparsity, sparsity_level=sparsity_level).to(device)

1 model.eval()

DictionaryLearningVAE(
  (_encoder): Encoder(
    (_conv_1): Conv2d(3, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
    (_conv_2): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
    (_conv_3): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (_residual_stack): ResidualStack(

```

```

        (_layers): ModuleList(
          (0-1): 2 x Residual(
            (_block): Sequential(
              (0): ReLU(inplace=True)
              (1): Conv2d(128, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
              (2): ReLU(inplace=True)
              (3): Conv2d(32, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
            )
          )
        )
      )
    )
  )
  (_pre_vq_conv): Conv2d(128, 256, kernel_size=(1, 1), stride=(1, 1))
  (_dl): DictionaryLearning(
    (representation): Sequential(
      (0): Linear(in_features=256, out_features=512, bias=True)
      (1): Softmax(dim=None)
    )
  )
  (_decoder): Decoder(
    (_conv_1): Conv2d(256, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (_residual_stack): ResidualStack(
      (_layers): ModuleList(
        (0-1): 2 x Residual(
          (_block): Sequential(
            (0): ReLU(inplace=True)
            (1): Conv2d(128, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
            (2): ReLU(inplace=True)
            (3): Conv2d(32, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
          )
        )
      )
    )
  )
  (_conv_trans_1): ConvTranspose2d(128, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
  (_conv_trans_2): ConvTranspose2d(64, 3, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
)

```

▼ Examining the model parameters.

```

1 optimizer = optim.Adam(model.parameters(), lr=learning_rate, amsgrad=False)

1 import time
2
3 num_training_updates = 200000
4
5 model.train() # set the model in training mode
6
7 train_res_recon_error = []
8
9 start_time = time.time()
10
11 z = None
12
13 for i in xrange(num_training_updates):
14     (data, _) = next(iter(training_loader))
15     data = data.to(device)
16     optimizer.zero_grad()
17
18     dlloss, data_recon, representation = model(data)
19     recon_error = F.mse_loss(data_recon, data) / data_variance_training
20     loss = recon_error + dlloss
21     loss.backward()
22
23     optimizer.step()
24
25     train_res_recon_error.append(recon_error.item())
26
27     if (i+1) % 100 == 0:
28         print('%d iterations' % (i+1), end="\n")
29         print('recon_error: %.3f' % np.mean(train_res_recon_error[-100:]))
30         print(f'DEBUG: sparsity level of the representation matrix: {torch.sum(representation != 0) / torch.numel(representation)}')
31         print()
32 print(f"Time elapsed: {time.time() - start_time}")
33 print(f"Experiment with sparsity control parameter = {sparsity} completed :)")

```

```

DEBUG:sparsity level of the representation matrix: 0.1953125

198300 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

198400 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

198500 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

198600 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

198700 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

198800 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

198900 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

199000 iterations      recon_error: 0.034
DEBUG:sparsity level of the representation matrix: 0.1953125

199100 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

199200 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

199300 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

199400 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

199500 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

199600 iterations      recon_error: 0.034
DEBUG:sparsity level of the representation matrix: 0.1953125

199700 iterations      recon_error: 0.034
DEBUG:sparsity level of the representation matrix: 0.1953125

199800 iterations      recon_error: 0.034
DEBUG:sparsity level of the representation matrix: 0.1953125

199900 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

200000 iterations      recon_error: 0.035
DEBUG:sparsity level of the representation matrix: 0.1953125

Time elapsed: 21723.437965869904
Experiment with sparsity control parameter = 0 completed :)

```

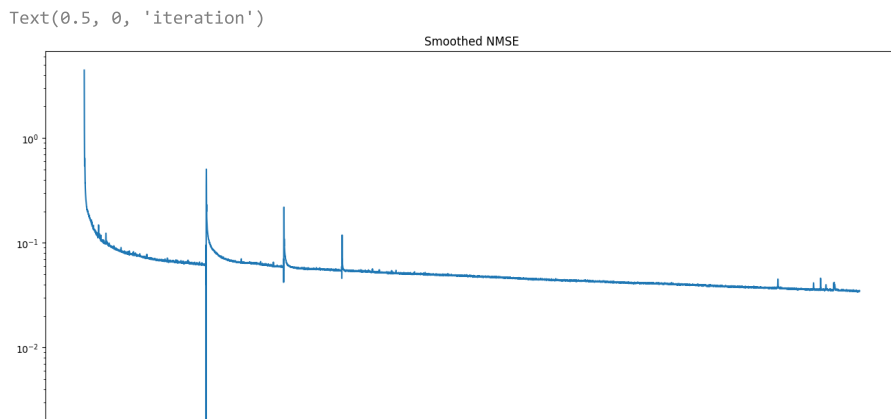
▼ Plot Loss

```

1  train_res_recon_error_smooth = savgol_filter(train_res_recon_error, 201, 7)

1  f = plt.figure(figsize=(16,8))
2  ax = f.add_subplot(1,1,1)
3  ax.plot(train_res_recon_error_smooth)
4  ax.set_yscale('log')
5  ax.set_title('Smoothed NMSE')
6  ax.set_xlabel('iteration')

```



View Reconstructions

```

1 model.eval()
2
3 (validation_originals, _) = next(iter(validation_loader))
4 validation_originals = validation_originals.to(device)

1 vq_output_eval = model._pre_vq_conv(model._encoder(validation_originals))
2 H, W = vq_output_eval.shape[2], vq_output_eval.shape[3]
3 # validation_representation = model._dl(vq_output_eval) # dictionary learning input dim: (batch_size, N, D), N = H * W (number of pixels)
4 # loss, validation_dict_learned, representation = model._dl.loss(vq_output_eval, validation_representation)
5 loss, validation_dict_learned, representation = model._dl(vq_output_eval)
6 print(validation_dict_learned.shape)

    torch.Size([32, 256, 8, 8])

1 plt.rcParams["figure.figsize"] = (10,8)

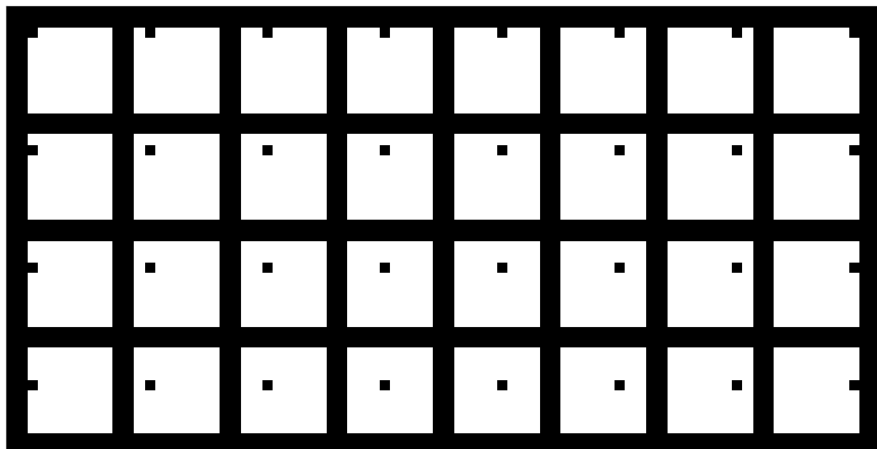
1 def show(img):
2     npimg = img.numpy()
3     fig = plt.imshow(np.transpose(npimg, (1,2,0)), interpolation='nearest')
4     fig.axes.get_xaxis().set_visible(False)
5     fig.axes.get_yaxis().set_visible(False)

1 # Latent space images
2 latents = validation_dict_learned.view(32 * 64, embedding_dim)
3 U, S, V = torch.pca_lowrank(latents)
4 projections = torch.matmul(latents, V[:, :3]).reshape(32, 3, 8, 8) # project to 3 channel space
5 latent_imgs = projections.cpu().data
6 for i in range(len(latent_imgs)):
7     for j in range(3):
8         tmp = latent_imgs[i, j, :, :]
9         latent_imgs[i, j, :, :] = (tmp - torch.min(tmp)) / (torch.max(tmp) - torch.min(tmp))
10 show(make_grid(latent_imgs), )

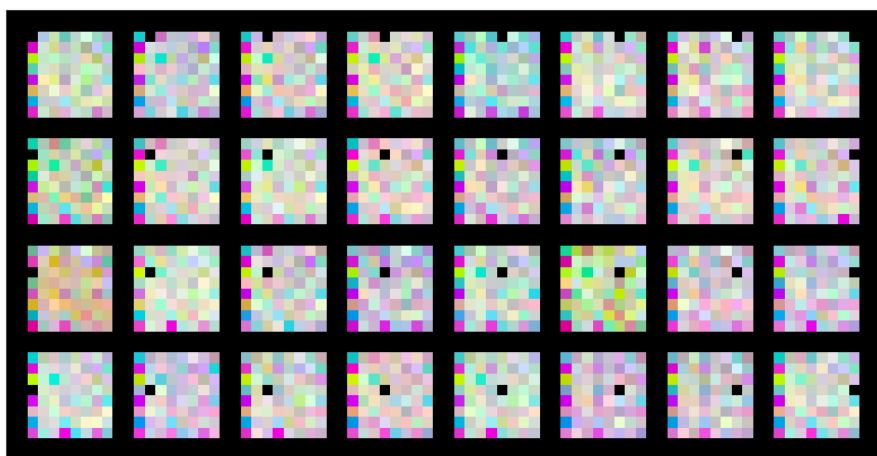
```



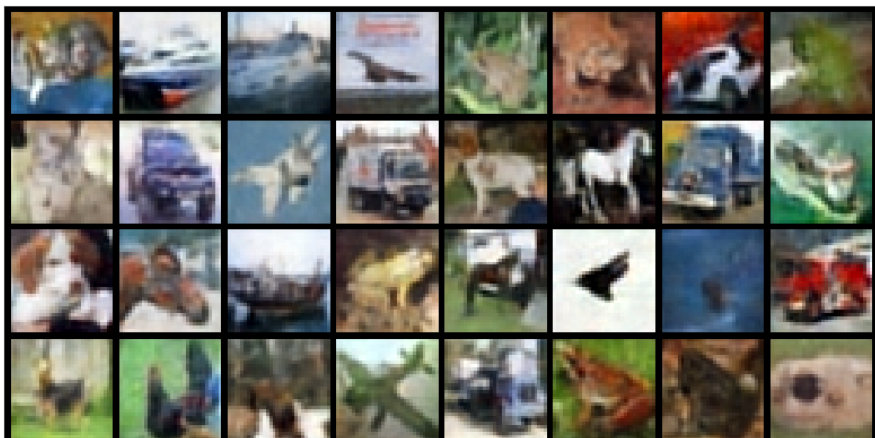
```
1 # create masks
2 masks = torch.ones((32, 3, 8, 8))
3 for i, mask in enumerate(masks):
4     mask[:, i // 8, i % 8] = 0
5 show(make_grid(masks))
```



```
1 masked_latent_imgs = torch.mul(masks, latent_imgs)
2 tmp = masked_latent_imgs
3
4 show(make_grid(masked_latent_imgs))
```



```
1 validation_reconstructions = (model._decoder(validation_dict_learned.reshape(32, embedding_dim, H, W))) # decoder input dim: (batch_size,
2
3 reconstructed_imgs = validation_reconstructions.cpu().data
4
5 for i in range(len(reconstructed_imgs)):
6     for j in range(3):
7         tmp = torch.clamp(reconstructed_imgs[i, j, :, :], 0, 1)
8         reconstructed_imgs[i, j, :, :] = tmp
9
10 show(make_grid(reconstructed_imgs), )
```

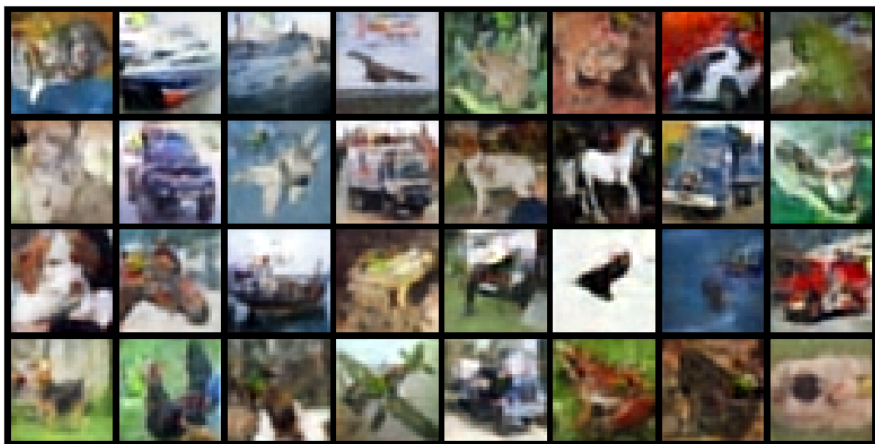


```

1 masks = torch.ones((32, embedding_dim, 8, 8))
2 for i, mask in enumerate(masks):
3     mask[:, i // 8, i % 8] = 0
4 masked_latents = torch.mul(masks.cuda(), latents.reshape(32, embedding_dim, 8, 8))
5 masked_validation_reconstructions = (model._decoder(masked_latents)) # decoder input dim: (batch_size, D, H, W), D = embedding_dim (data

1 for i in range(len(masked_validation_reconstructions)):
2     for j in range(3):
3         tmp = torch.clamp(masked_validation_reconstructions[i, j, :, :], 0, 1)
4         masked_validation_reconstructions[i, j, :, :] = tmp
5
6 show(make_grid(masked_validation_reconstructions.cpu()))

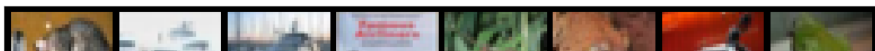
```



```

1 show(make_grid(validation_originals.cpu()))

```



We can see the results are decent but a little blurry.



View Embedding



```
1 proj = umap.UMAP(n_neighbors=3,  
2                 min_dist=0.1,  
3                 metric='cosine').fit_transform(model._dl.dictionary.cpu().detach().numpy())
```



```
1 plt.scatter(proj[:,0], proj[:,1], alpha=0.3)
```

<matplotlib.collections.PathCollection at 0x7f3da539d2d0>

