

# Introduction

## TL;DR

- Text is represented as
  - a sequence of integers: an integer index into the list of tokens in the vocabulary
- We want to represent other types of data (e.g., images) as
  - a sequence of integers: an integer index into a "code book": finite list of vectors
- By making the "shape" (sequence) and "type" (integer) compatible between text and other data types
  - We facilitate mixing text and other data types
  - Will enable an implementation of Text to Image as a simple extension of the Language Modeling objective

We now present an Autoencoder with a twist

- the latent representation produced for an input
- is limited to be one member of a *finite list* of vectors
- enabling us to describe the latent by the *integer index* in the list

Why is an integer encoding of an input interesting ?

- It is analogous to the way we treat words (tokens) in Natural Language Processing
  - an index into a finite Vocabulary of words
- This opens the possibility of dealing with sequences that are a *mixture* of text and other data types (e.g., images)

Rather than pre-specifying the finite list, we will *learn* the list by training a Neural Network.

In a subsequent module, we will use a similar technique for the task of Text to Image

- given the description of an image in words
- create an image matching the description

But there is a significant problem with a Neural Network that learns discrete values

- the network may need to make a "hard" (as compared to "soft") choice
  - a true `if` statement ("hard") versus a "soft" conditional (sigmoid)
  - a Python `dict` ("hard" lookup) versus a "soft" lookup (Context Sensitive Memory)
- "hard" means derivatives are not continuous
- Gradient Descent won't work

We will introduce a new Deep Learning operator (*Stop Gradient*) to deal with "hard" operators.

## References

- [paper: vanilla VQ-VAE \(https://arxiv.org/pdf/1711.00937.pdf\)](https://arxiv.org/pdf/1711.00937.pdf)
- [paper: VQ-VAE-2 \(https://arxiv.org/pdf/1906.00446.pdf\)](https://arxiv.org/pdf/1906.00446.pdf)

# From PCA to VQ-VAE

The common element in the design of any Autoencoder method is

- to create a latent representation  $\mathbf{z}$  of input  $\mathbf{x}$
- such that  $\mathbf{z}$  can be (approximately) inverted to reconstruct  $\mathbf{x}$ .

Principal Components Analysis is a type of Autoencoder that produces a latent representation  $\mathbf{z}$  of  $\mathbf{x}$

- $\mathbf{x}$  is a vector of length  $n$ :  $\mathbf{x} \in \mathbb{R}^n$
- $\mathbf{z}$  is a vector of length  $n' \leq n$ :  $\mathbf{z} \in \mathbb{R}^{n'}$

Usually  $n' \ll n$ : achieving *dimensionality reduction*

This is accomplished by decomposing  $\mathbf{x}$  into a weighted product of  $n$  *Principal Components*

- $\mathbf{V} \in \mathbb{R}^{n \times n}$

$$\mathbf{x} = \mathbf{z}' \mathbf{V}^T$$

- where  $\mathbf{z}' \in \mathbb{R}^n$
- rows of  $\mathbf{V}^T$  are the components

So  $\mathbf{x}$  can be decomposed into the weighted sum (with  $\mathbf{z}'$  specifying the weights)

- of  $n$  component vectors
- each of length  $n$



That is:

- $\mathbf{V}^T$  is a basis space that can define all vectors in  $n$ -dimensional space
- with loadings  $\mathbf{z}'$  on the basis vectors.
- to illustrate
  - think of the familiar basis space  $I \in \mathbb{R}^{n \times n} = \text{diagonal}(n)$
  - $\mathbf{x} = \mathbf{z}'$  are the "coordinates" of a point within the  $n$ -dimensional space spanned by  $I$

Since  $\mathbf{z}' \in \mathbb{R}^n$ : there is **no** dimensionality reduction just yet.

One can also view  $\mathbf{V}^T$  as a kind of *code book*

- any  $\mathbf{x}$  can be represented (as a linear combination) of the *codes* (components) in  $\mathbf{V}^T$

$$\mathbf{x} = \mathbf{z}' \mathbf{V}^T$$

$\mathbf{z}'$  is like a translation of  $\mathbf{x}$ , using  $\mathbf{V}$  as the vocabulary.

- weights in the codebook
- rather than weights in the standard basis space  $I \in \mathbb{R}^{n \times n} = \text{diagonal}(n)$

$$\mathbf{x} = \mathbf{x} I$$

Dimensionality reduction is achieved by defining  $\mathbf{z}$  as a length  $n'$  prefix of  $\mathbf{z}'$

- $\mathbf{z} = \mathbf{z}'_{1:n'}$
- $\mathbf{z} \in \mathbb{R}^{n'}$

Similarly, we needed only  $n'$  components from  $\mathbf{V}$

- $\mathbb{V}^T = \mathbf{V}_{1:n'}^T$
- $\mathbb{V}^T \in \mathbb{R}^{n' \times n}$

We can construct an *approximation*  $\hat{\mathbf{x}}$  of  $\mathbf{x}$  using *reduced dimension*  $\mathbf{z}'$  and  $\mathbb{V}$

$$\hat{\mathbf{x}} = \mathbf{z}\mathbb{V}^T$$

The Autoencoder (and variants such as VAE) produces  $\mathbf{z}^{(i)}$ , the latent representation of  $\mathbf{x}^{(i)}$

- directly
- independent of any other training example  $\mathbf{x}^{(i')}$  for  $i \neq i'$

One of our goals in using AE's is in generating synthetic data

- the dimensionality reduction achieved thus far was a necessity, not a goal
- "bottle-neck" in the AE enforces the dimensionality reduction

Our goal in introducing the Vector Quantized Autoencoder is **not** in synthesizing data

- it is to show how non-language data types
- can be represented as *single* codes (vectors) in a code-book

This makes the non-language data types

- similar to language tokens
- which are elements from a finite "code book": the Vocabulary
- and we can create non-language "tokens"

If we are able to do this

- both language and non-language data types
- become sequences of tokens
- and we can use an LLM
  - which processes token sequences
  - to handle *mixed* types of data (multi-modal)

# Vector Quantized Autoencoder

A *Vector Quantized VAE* is a VAE with similarities to PCA. It creates  $\mathbf{z}$

- which is an **integer**
- that is the index of a row
- in a codebook with  $K$  rows

That is: the input is represented by one of  $K$  possible vectors.

The result is a **discrete** (rather than continuous) representation of a non-language token.

The representation is a

- *reference* (via an index)
- to a continuous vector of length  $n$
- where the continuous vector is *one* of  $K$  possible vectors of length  $n$  within a fixed codebook



In essence

- the index
- is a OHE of a token
- with a non-language "vocabulary" of  $K$  possible tokens

Thus, a sequence of  $T$  continuously valued vectors

- can be represented as a sequence of  $T$  integers
- over a "vocabulary" defined by the code book

This is analogous to text

- sequence of words
- represented as a sequence of integer indices in a vocabulary of tokens

Once we put complex objects

- like images
- timeseries
- speech

into a representation similar to text

- we can have *mixed type* sequences
  - e.g., words, images

In a subsequent module we will take advantage of mixed type sequences

- to produce an image
- from a text *description* of the image
- using the "predict the next" element of a sequence technique of Large Language Models

---

DALL-E: Text to Image

---

Text input: "An illustration of a baby daikon radish in a tutu walking a dog"

---

Image output:

---



# Details

A VQ-VAE can produce a sequence of integers that encodes many different types of data

- including data with "shape": *non-feature* dimensions
  - image
  - audio

Each "location" in the space of non-feature dimensions

- is a vector consisting only of features

So an image, for example,

- is a collection of feature vectors
- whose locations are arranged into a 2D grid
  - a location is a pair of row number/column number
- the feature vector at each location
  - has 3 features: Red, Green, Blue

We illustrate the VQ-VAE using Image examples.

Here is diagram of a VQ-VAE

- that creates a latent representation of a 3-dimensional image ( $w \times h \times 3$ )
- as a 2-dimensional matrix of integers
  - each location
  - has an integer index
    - which references a vector of length  $n$  in the codebook (with  $K$  codes in the codebook)

There is a bit of notation: referring to the diagram should facilitate understanding the notation.

VQ-VAE





## The picture: in words

- the 3D *input* image (dog)  $\mathbf{x}$  of shape  $(h \times w \times 3)$  is processed by a CNN
- into a 3D alternate representation  $\mathbf{z}_e(\mathbf{x})$ 
  - called the *Encoder Output*
- each location in the alternate representation
  - is converted into an integer index  $k$  of a vector  $\mathbf{e}_k$  in the codebook
- resulting in a 2D alternate representation  $q(\mathbf{z}|\mathbf{x})$ 
  - with the same non-feature dimensions as  $\mathbf{z}_e(\mathbf{x})$
- replacing each integer with its length  $n$  vector
  - is the 3D equivalent
  - called the *quantized* equivalent of  $\mathbf{x}$ , denoted  $\mathbf{z}_q(\mathbf{x})$

The details follow

In general, we assume the input has  $\#S$  non-feature (*spatial* in the case of images) dimensions

- where each location in the spatial dimension is a vector of features with length  $n$
- input shape  $(n_1 \times n_2 \dots \times n_{\#S} \times n)$

We will explain this diagram in steps.

First, we summarize the notation in a single spot for easy subsequent reference.

# Notation summary

| term                         | shape                                   | meaning  |   |
|------------------------------|---|--|---|
| $S$                          | $(n_1 \times n_2 \dots \times n_{\#S})$ | Spatial dimensions of $\#S$ -dimensional input   |   |
| $\mathbf{x}$                 | $\mathbb{R}^{S \times n}$               | Input  |   |
| $D$                          |   | length of latent vectors (Encoder output, Quantized Encoder output, Codebook entry)              |   |
| $\mathcal{E}$                |   | Encoder function   |   |
| $\mathbf{z}_e(\mathbf{x})$   | $\mathbb{R}^{S \times D}$               | Encoder output over each location of spatial dimension   |   |
|                              |   | $\mathbf{z}_e(\mathbf{x}) = \mathcal{E}(\mathbf{x})$   |   |
| $\mathbf{z}_e(\mathbf{x})$   | $\mathbb{R}^D$                          | Encoder output at a <b>single</b> representative spatial location                                |   |
|                              |   | $\mathbf{z}_e(\mathbf{x}) = \mathcal{E}(\mathbf{x})$   |   |
| $K$                          |   | number of codes  |   |
| $\mathbf{E}$                 | $\mathbb{R}^{K \times D}$               | Codebook/Embedding   |   |
|                              |   | $K$ codes, each of length $D$  |   |
| $e \in \mathbf{E}$           | $\mathbb{R}^D$                          | code/embedding   |   |
| $\mathbf{z}$                 | $\{1, \dots, K\}^{S \times D}$          | latent representation over all spatial dimensions  |   |
| $\mathbf{z}$                 | $\{1, \dots, K\}$                       | Latent representation at a <b>single</b> representative spatial location                         |   |
|                              |   | one integer per spatial location   |   |
| $\lfloor \mathbf{z} \rfloor$ |   | integer $\in [1 \dots K]$  | Index $k$ of $e_k \in \mathbf{E}$ that is closest to $\mathbf{z}_e(\mathbf{x})$ |
|                              |   | $k = \underset{j \in [1, K]}{\operatorname{argmin}}  \mathbf{z}_e(\mathbf{x}) - \mathbf{e}_j _2$ |   |
|                              |   | actually: encoded as a OHE vector of length $K$  |   |
| $\mathbf{z}_q(\mathbf{x})$   | $\mathbb{R}^D$                          | Quantized $\mathbf{z}_e(\mathbf{x})$   |   |
|                              |   | $\mathbf{z}_q(\mathbf{x}) = e_k$ where $k = \lfloor \mathbf{z} \rfloor$                          | $\lfloor \mathbf{x} \rfloor$  |
|                              |   | i.e, the element of codebook that is closest to $\mathbf{z}_e(\mathbf{x})$                       |   |
|                              |   | $\mathbf{z}_q(\mathbf{x}) \approx \mathbf{z}_e(\mathbf{x})$                                      |   |
| $\tilde{\mathbf{x}}$         | $n$                                     | Output: reconstructed $\mathbf{x}$   |   |
|                              |   | $\lfloor \mathbf{z}_q(\mathbf{x}) \rfloor$   | $\lfloor \mathbf{z}_q(\mathbf{x}) \rfloor$                                      |

| term                                    | shape                                      | meaning |
|---|--|---------|
| $\mathcal{D}$                           | $\mathbb{R}^{n'} \rightarrow \mathbb{R}^n$ | Decoder |
| input: element of codebook $\mathbf{E}$ |  |         |

# Quantization

Let  $S$  denote the spatial dimensions, e.g.  $S = (n_1 \times n_2)$  for 2D

So input  $\mathbf{x} \in \mathbb{R}^{S \times n}$

- $n$  features over  $S$  spatial locations

The input  $\mathbf{x}$  is transformed in a sequence of steps

- Encoder output (continuous value)
- Latent representation (discrete value)
  - Quantized (continuous value)

In the first step, the *Encoder* maps input  $\mathbf{x}$

- to Encoder output  $\mathbf{z}_e(\mathbf{x})$
- an alternate representation of  $D$  features over  $S'$  spatial locations

(For simplicity, we will assume  $S' = S$ )

## Notational simplification

In the sequel, we will apply the same transformation **to each element** of the spatial dimension

Rather than explicitly iterating over each location we write

$$\mathbf{z}_e(\mathbf{x}) \in \mathbb{R}^D$$

to denote a representative element of  $\mathbf{z}_e(\mathbf{x})$  at a single location  $s = (i_1, \dots, i_{\#S})$

$$\mathbf{z}_e(\mathbf{x}) = \mathbf{z}_e(\mathbf{x})_s$$

We will continue the transformation at the single representative location

- and implicitly iterate over all locations  $s \in S$



The continuous (length  $D$ ) Encoder output vector  $\mathbf{z}_e(\mathbf{x})$

- is mapped to a *latent representation*  $q(\mathbf{z}|\mathbf{x})$
- which is a **discrete** value (integer)

$$k = q(\mathbf{z}|\mathbf{x}) \in \{1, \dots, K\}$$

where  $k$  is the *index* of a row  $\mathbf{e}_k$  in codebook  $\mathbf{E}$

$$\mathbf{e}_k = \mathbf{E}_k \in \mathbb{R}^D$$

The codebook is also called an *Embedding* table.

$k$  is chosen such that  $\mathbf{e}_k$  is the row in  $\mathbf{E}$  closest to  $\mathbf{z}_e(\mathbf{x})$

$$\begin{aligned} k &= q(\mathbf{z}|\mathbf{x}) \\ &= \operatorname{argmin}_{j \in \{1, \dots, K\}} \|\mathbf{z}_e(\mathbf{x}) - \mathbf{e}_j\|_2 \end{aligned}$$

We denote the codebook vector

- closest to representative encoder output  $\mathbf{z}_e(\mathbf{x})$
  - as  $\mathbf{z}_q(\mathbf{x})$
- $$\mathbf{z}_q(\mathbf{x}) = \mathbf{e}_k \text{ where } k = q(\mathbf{z}|\mathbf{x})$$

The Decoder tries to invert the codebook entry  $\mathbf{e}_k = \mathbf{z}_q(\mathbf{x})$  so that

$$\begin{aligned}\tilde{\mathbf{x}} &= \mathcal{D}(\mathbf{z}_q(\mathbf{x})) \\ &\approx \mathbf{x}\end{aligned}$$

# Discussion

## Why do we need the CNN Encoder ?

The input  $\mathbf{x}$  is first transformed into an *alternate representation*

- the **number** and length of the non-feature dimensions (i.e.,  $h \times w$ ) are preserved
  - convenient but not necessary
- but the number of features is transformed from  $n$  raw features to  $D \geq n$  synthetic features
  - typical behavior for, e.g., an image classifier

The part of the VQ-VAE after the initial CNN

- reduces the size of the **feature dimension** from  $D$  to 1
- this is the primary source of dimensionality reduction
  - the raw  $n$  of image input is usually only  $n = 3$  channels

Although our illustration preserved the non-feature dimensions

- it may be useful for the CNN to *down-sample* spatial dimension  $S$  to a smaller  $S'$
- resulting in shorter sequences when we eliminate the non-feature dimensions by "flattening"

For example

- 3 layers of stride 2 CNN layers
- will reduce a 2D image of spatial dimension  $(n_1 \times n_2)$
- to spatial dimension  $(\frac{n_1}{8} \times \frac{n_2}{8})$

This replaces each  $(8 \times 8 \times n)$  *patch* of raw input

- into a single vector of length  $D$
- that summarizes the  $(8 \times 8)$  the patch

One possible role (not strictly necessary) for the CNN Encoder

- is to replace a large spatial dimensions
- by smaller "summaries" of local neighborhoods (patches)



# Why quantize ?

Quantization

- converts the continuous  $\mathbf{z}_e(\mathbf{x})$
- into discrete  $q(\mathbf{z}|\mathbf{x})$
- representing the approximation  $\mathbf{z}_q(\mathbf{x}) \approx \mathbf{z}_e(\mathbf{x})$

The Decoder inverts the approximation.

Why bother when the Quantization/De-Quantization is Lossy ?

One motivation comes from observing what happens if we *quantize and flatten*

- the  $\#S'$ -dimensional spatial locations
- to a sequence of integers

Quantizing replaces each patch with a single integer index.

- the integer is the index of an *image token* within a list of  $K$  possible tokens

By flattening the quantized higher dimensional matrix of patches, we convert the input

- into a sequence of image tokens
- over a "vocabulary" defined by the codebook **E**.

This yields an image representation

- similar to the representation of text

Thus, we open the possibility of processing sequences of mixed text and image tokens.

## Quantized image embeddings mixed with Text: preview of DALL-E

The Large Language Model operates on a sequence of text tokens

- where the text tokens are fragments of words
- when run autoregressively
  - concatenating each output to the initial input sequence
  - the LLM shows an ability to produce a "sensible" continuation of an initial "thought"

Suppose we train a LLM on input sequences

- that start with a sequence of *text* tokens describing an image
- followed by a separator [SEP] token
- followed by a sequence of quantized image tokens

<text token> <text token> ... <text token> [SEP] <image token> <image token> ...

What continuation will our trained LLM produce given prompt

<text token> <text token> ... <text token> [SEP]

Hopefully:

- a sequence of *image tokens*
- that can be reconstructed
- into an image matching the description given by the text tokens !

That is the key idea behind a Text to Image model called DALL-E that we will discuss in a later module.

There remains an important technical detail

- the embedding space of text and image are distinct
- they need to be merged into a common embedding space

We will visit these issues in the module on CLIP.

# Loss function

The Loss function for the VQ-VAE entails several parts

- Reconstruction loss
  - enforcing constraint that reconstructed image is similar to input
$$\tilde{\mathbf{x}} \approx \mathbf{x}$$
- Vector Quantization (VQ) Loss:
  - enforcing similarity of quantized encoder output and actual encoder output
$$\mathbf{z}_q(\mathbf{x}) \approx \mathbf{z}_e(\mathbf{x})$$
- Commitment Loss
  - a constraint that prevents the Quantization of  $\mathbf{z}_e(\mathbf{x})$  from alternating rapidly between code book entries



## Stop Gradient operator

Some of the Loss terms will involve an operator that we have not yet seen:

- The sg operator is the *Stop Gradient* operator.

The need for this operator stems from the Gradient Descent update process

- the partial derivative assumes "everything else" other than the denominator (variable being updated) remains constant
- if "everything else" *also* changes: the gradient update step may not reduce Loss

But

- changing the Encoder parameters affects Encoder output  $\mathbf{z}_e(\mathbf{x})$ 
  - which may affect the Embeddings

Similarly

- changing the Embeddings "code book" will affect the Encoder output (the key used for lookup in the code book)

In addition

- Quantization (selecting a discrete code from the code book) **is not** differentiable
- "hard" rather than "soft" choice

We need to be able to pass gradients backward through the non-differentiable operator.

The Stop Gradient operator will facilitate all these cases.

On the Forward Pass, it acts as an Identity operator

$$\text{sg}(\mathbf{x}) = \mathbf{x}$$

But on the Backward Pass of Backpropagation: *it stops the gradient* from flowing backwards

$$\frac{\partial \text{sg}(\mathbf{x})}{\partial \mathbf{y}} = 0 \text{ for all } \mathbf{y}$$

# Reconstruction Loss

The Reconstruction Loss term is our familiar: Maximize Likelihood

- written to minimize the negative of the log likelihood, as usual

$$p(\mathbf{x}|\mathbf{z}_q(\mathbf{x}))$$

It will serve to update parameters for the Encoder and Decoder.

As we will see (section: "Quantization is not differentiable")

- during the backward pass
- the Loss gradient from the (quantized) Decoder input  $\mathbf{z}_q(\mathbf{x})$
- flows directly to the (continuous) Encoder output  $\mathbf{z}_e(\mathbf{x})$

Effectively, for the purpose of gradient/weight update due to Reconstruction Loss:

$$\mathbf{z}_q(\mathbf{x}) = \mathbf{z}_e(\mathbf{x})$$

So the Reconstruction Loss will **not** cause the code book embeddings to be updated

- update to the Encoder parameters thus satisfies "all else being constant" in that the Code book does not change

The Code Book updates will be the job of the two other Loss terms.

# Vector Quantization Loss

The Vector Quantization Loss and Commitment Loss are similar.

- differ only in the placement of the Stop Gradient

Vector Quantization Loss:

$$\| \text{sg}(\mathbf{z}_e(\mathbf{x})) - \mathbf{z}_q(\mathbf{x}) \|$$

where  $\text{sg}$  is the *Stop Gradient* Operator (details to follow).

The purpose of the Vector Quantization Loss is to update the Embedding (codebook)  $\mathbf{E}$

- by moving "code"  $\mathbf{z}_q(\mathbf{x}) = e_k$  closer to Encoder output  $\mathbf{z}_e(\mathbf{x})$
- assuming Encoder output is held constant
  - the stop-gradient on the Encoder output  $\mathbf{z}_e(\mathbf{x})$  prevents the Encoder output from being changed by the VQ Loss



# Commitment Loss

Commitment Loss:

$$\|\mathbf{z}_e(\mathbf{x}) - \text{sg}(\mathbf{z}_q(\mathbf{x}))\|$$

It is similar to the Vector Quantization loss except for the placement of the Stop Gradient operator.

The Stop Gradient in the Commitment Loss prevents a change in the Embeddings from affecting the Encoder weights (and thus,  $z_e(\mathbf{x})$ ).

The reason for the different placement of the Stop Gradient operators in the VQ and Commitment Loss terms

- there is a circular dependency
  - encoder output  $\mathbf{z}_e(\mathbf{x})$  is affected by the Code Book quantization  $\mathbf{z}_q(\mathbf{x})$
  - the Code Book quantization  $\mathbf{z}_q(\mathbf{x})$  is affected by the encoder output  $\mathbf{z}_e(\mathbf{x})$

By holding one constant while updating the other

- we ensure that the embeddings  $\mathbb{E}$  converge.

# Total Loss

Loss function

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \mathcal{D}(\mathbf{e})) = & \|\mathbf{x} - \mathcal{D}(\mathbf{e})\|_2^2 && \text{Reconstruction Loss} \\ & + \|\text{sg}[\mathcal{E}(\mathbf{x})] - \mathbf{e}\|_2^2 && \text{VQ loss, codebook loss: train codebook} \\ & + \beta \|\text{sg}[\mathbf{e}] - \mathcal{E}(\mathbf{x})\|_2^2 && \text{Commitment Loss: force } \mathcal{E}(\mathbf{x}) \text{ to be close to } \mathbf{e} \\ & \text{where } \mathbf{e} = \mathbf{z}_q(\mathbf{x})\end{aligned}$$

Need the stop gradient operator  $\text{sg}$  to control the mutual dependence

- of the Encoder output  $\mathcal{E}(\mathbf{x})$  and the chosen code  $\mathbf{e}$
-

# Quantization is not differentiable

There is a subtle but important problem.

The Quantization operation

$$\begin{aligned} k &= q(\mathbf{z}|\mathbf{x}) \\ &= \operatorname{argmin}_{j \in \{1, \dots, K\}} \|\mathbf{z}_e(\mathbf{x}) - \mathbf{e}_j\|_2 \end{aligned}$$

is not differentiable because of  $\operatorname{argmin}$

`argmin` is a problematic operation because

- it contains a "hard choice" so is not differentiable
  - output may change dis-continuously from index  $k$  to index  $k' \neq k$
  - for small changes in the input
  - not continuous at the point of change
- it may also be non-deterministic
  - when minimum value occurs at more than one index
  - when  $\mathbf{e}_k = \mathbf{e}_{k'}$  for  $k \neq k'$

There is a work-around

- implement a VectorQuantizer layer
- using a [Straight Through Estimator \(Straight Through Estimator.ipynb\)](#).
  - see that module for details of the technique

We see this in the [Colab \(https://keras.io/examples/generative/vq\\_vae/\)](https://keras.io/examples/generative/vq_vae/) implementation of Vector Quantization (the VectorQuantizer layer)

```
class VectorQuantizer(layers.Layer):  
    ...  
    def call(self, x):  
    ...  
        # Straight-through estimator.  
        quantized = x + tf.stop_gradient(quantized - x)
```

Code similar to the [VectorQuantizer of the paper's authors \(https://github.com/deepmind/sonnet/blob/v1/sonnet/python/modules/nets/vqvae.py\)](https://github.com/deepmind/sonnet/blob/v1/sonnet/python/modules/nets/vqvae.py)

# Code

Here ([https://colab.research.google.com/github/keras-team/keras-io/blob/master/examples/generative/ipynb/vq\\_vae.ipynb#scrollTo=LWYJf1MYvzap](https://colab.research.google.com/github/keras-team/keras-io/blob/master/examples/generative/ipynb/vq_vae.ipynb#scrollTo=LWYJf1MYvzap)) is a Colab notebook.



# Creating synthetic examples

Consider an example that is an Image.

It is a structured arrangement of feature vectors

- a 2D grid
- each element of the grid is a vector of features

In general: our examples may have an arbitrary number of dimensions.

By convention, we will refer

- to the last dimensions as the "feature" dimension
- all the preceding dimensions as *non-feature* dimensions
  - for an Image: *spatial* dimensions

Denoting  $\#S$  as the number of non-feature dimensions and  $n_{(0)}$  as the number of features, examples have shape

$$(n_1 \times \dots \times n_{\#S} \times n_{(0)})$$

where  $\#S$  denotes the number of non-feature dimensions (  $\#S = 2$  for an Image)

The Encoder creates an output  $\mathbf{z}_e(\mathbf{x})$  of shape

$$(n_1 \times \dots \times n_{\#S} \times n_e)$$

which is quantized into  $\mathbf{z}_q(\mathbf{x})$  of shape

$$(n_1 \times \dots \times n_{\#S} \times 1)$$

That is: the feature dimension of  $\mathbf{z}_q(\mathbf{x})$  is a single integer index into a learned codebook.

In order for us to generate synthetic examples

- we must create a Tensor  $\mathbf{z}_q(\mathbf{x})$  of  $(n_1 \times \dots \times n_{\#S} \times 1)$  integers (*latents*)
- feed this to the Decoder
- get a synthetic example as output

# Learning the distribution of latents

But we can't create  $\mathbf{z}_q(\mathbf{x})$  *completely at random*.

- The elements at adjacent locations in the non-feature dimensions may not be **independent**

For example: consider an Image

- the pixels in an image are related to one another

We must learn a distribution of  $\mathbf{z}_q(\mathbf{x})$  that respects the dependencies.

One solution

- flatten  $\mathbf{z}_q(\mathbf{x})$
- from shape  $(n_1 \times \dots \times n_{\#S} \times 1)$
- into a sequence  
 $\mathbf{z}_{(1)}, \mathbf{z}_{(2)}, \dots, \mathbf{z}_{(n_1 * n_2 \dots * n_{\#S})}$   
of integer indices.

We can then learn the distribution of the latents

- through auto-regressive modeling of the sequence sequence  $\mathbf{z}$

$$p(\mathbf{z}_{(k+1)} | \mathbf{z}_{(1)}, \dots, \mathbf{z}_{(k)})$$

This is just like the Language Model objection for NLP.

## Aside

*Learning* a distribution is less restrictive than *assuming* a distribution

For the case of "structured" examples: we have no choice.

- adjacent elements are not independent

But recall that, for the VAE, we *assumed* that the latents came from a Normal distribution.

Use Auto-regressive modeling is a nice "trick" to learn distributions rather than having to assume a "convenient" functional form.



In [2]: `print("Done")`

Done

