



UNIVERSITY OF CAPE TOWN

STA5003W

MULTIVARIATE STATISTICS

---

# Multivariate Report

---

*Authors:*  
Gawronsky, Marcus  
Singo, Una

*Student Number:*  
GWRMAR002  
SNGUNA003

October 7, 2019

### Abstract

Variational Autoencoders (VAEs) are a class of generative artificial neural network commonly used for dimensionality reduction in machine learning applications. In this paper, we implement a stable Variational Autoencoder and benchmark its performance against traditional methods in multivariate statistics using common open data sets. We build on our implementation with two experiments; the first is an application to noncolumnar data in the medical field, and the second is an original generative sampling problem.

## Contents

<b>Literature Review</b>	<b>2</b>
<b>Setup</b>	<b>3</b>
<b>Data</b>	<b>3</b>
Iris Flowers . . . . .	3
Wine . . . . .	3
Breast Cancer . . . . .	4
Heartbeats Audio Recordings . . . . .	4
<b>Analysis</b>	<b>4</b>
Benchmarking . . . . .	4
Dimensionality reduction for clustering . . . . .	4
Manifold Learning . . . . .	5
Application to Medical Heartbeats Audio Data . . . . .	6
Generative Modeling . . . . .	8
<b>Appendixes</b>	<b>8</b>
Benchmarking . . . . .	8
Manifold Learning . . . . .	9
Heartbeats Audio . . . . .	10
Generative Sampling . . . . .	11
Code . . . . .	12
Exploratory Analysis . . . . .	12
Dataset Benchmark Code . . . . .	19
Manifold Learning Comparison . . . . .	23
Latent Space Sampling . . . . .	26
Heartbeats Dataset . . . . .	29
Library Code . . . . .	36
Environment . . . . .	55

## Literature Review

Variational Autoencoders (VAEs) are a class of generative artificial neural network commonly used for dimensionality reduction in machine learning applications. The technique extends on the application of traditional autoencoders and deep belief networks by introducing techniques from Variational Bayesian Methods.

An autoencoder consists of two neural networks, an encoder network and decoder network. The encoder network takes an input and produces a lower-dimensional representation of the data. The decoder network takes the lower-dimensional representation of the data and learns a reconstruction of the data back into its original vector-space. These neural networks can have a varying number of hidden-layers, activation functions and neurons to learn non-linear representations of the data (Bengio et al., 2007, 2013).

Variational Bayesian Methods are a class of techniques used in approximating intractable integrals found in Bayesian Modeling. Using an approximate distribution,  $q^*(x)$ , a loss function is used to minimize the difference between the approximate distribution and true posterior. This loss function is often chosen to be the Evidence Lower Bound (ELBO) or Kullback–Leibler divergence and is weighted in the optimization procedure to ensure the validity of the applied approximate distribution.

The Variational Autoencoder aims to extend on encoder-decoder model architectures by placing distributional assumptions on the low-dimensional latent space. By making this assumption, these models can be generative, allowing one to sample data across a complex data manifold. Figure 1 aims to show the basic architecture of a VAE. The input data  $\mathbf{X}$  is assumed to be an *i.i.d.* high dimensional dataset and the latent variables are denoted as  $\mathbf{z}$  are assumed to follow an unobservable distribution. The encoder network is defined as the conditional probability distribution  $q_\theta(\mathbf{z}|\mathbf{x})$  and outputs estimates to  $q_\theta(\mathbf{z}|\mathbf{x})$  which is assumed to follow a parametric probability distribution. Using the estimated parameters from the encoder neural network, sample values of  $\mathbf{z}$  are computed and used as an input to the decoder neural network  $p_\theta(\mathbf{x}|\mathbf{z})$  which outputs estimates to the conditional distribution  $p_\theta(\mathbf{x}|\mathbf{z})$  which can be used to sample values of  $\mathbf{x}$ . A key aspect to the model is the distributional assumption on the latent space. While several methods can be used to ensure the model distributional assumption, the original authors include the Kullback–Leibler divergence into the loss function of the model to ensure a distribution over latent variable which is approximately standard normal (Kingma and Welling, 2013; Zhao et al., 2017; Tolstikhin et al., 2017).

Using these techniques, Variational Autoencoders provide a flexible method which can be applied to solve a wide range of problems across domains. In the original

paper, Kingma and Welling (2013) use the popular MNIST dataset to both embed and generate new images of handwritten digits. A significant number of subsequent research on VAEs has focused on the learning latent distributions of images and generating new images. Generative Adversarial Networks are a popular example of such research outputs where the objective is to generate fake images that appear to be authentic to human observers (Goodfellow et al., 2014).

Apart from image generation VAEs have been shown to solve anomaly detection problems. An and Cho (2015) propose a method using the reconstruction probability from a variational autoencoder to identify outliers. Their experimental results show that this method outperforms traditional auto-encoder and principal components based methods.

## Setup

Code for this assignment was written in Python version 3.7.1, using a random seed of 1234, a full list of the Python dependencies are listed in the dependencies in a YML file. This can be used, along with the code <sup>1</sup>, to reproduce the analysis provided.

## Data

Five open data sets and two simulated are used in this assignment. Of the open data sets, three are column-oriented and two audio based. This section will provide a short description of the data, its source, and some information about its characteristics.

### Iris Flowers

The Iris flower dataset (Fisher, 1936) measures 4 distinguishing features of three related species of Iris flowers, Setosa, Versicolour and Virginica. It contains 150 observations with 5 variables, a classifier, the sepal length, sepal width, petal length and petal width. It is commonly used as an introduction to solving classification, clustering and dimensional reduction problems.

### Wine

The Wine dataset (Aeberhard et al., 1994) contains results from a chemical study of wine grown in the same region in Italy but grown by three different farmers.

---

<sup>1</sup>The code is accessible on github: <https://github.com/marcusinthesky/super-spirals.git>

The data contains 178 observations, 13 predictive numerical features that measure chemicals properties.

### **Breast Cancer**

The Breast Cancer (Street et al., 1993) data contains features computed from an image of a fine needle aspirate (FNA). Fine-needle aspiration (FNA) is a diagnostic procedure used to investigate lumps. Using a fine needle cells from the lump are extracted and studied under a microscope, from the image 30 features can be measured that characterise the mass and distinguish it from being Benign or Malignant.

### **Heartbeats Audio Recordings**

The Heartbeats dataset (Bentley et al.) will form part of our non column-oriented datasets.

Heartbeats is an audio based data set that contains recordings of heartbeats from a stethoscope. There are four classes of audio, artifact, extrahls, murmer and normal. The data set was originally created to identify S1 (dub) and S2 (dub) sounds and classify beats into one of the four classes.

## **Analysis**

In the literature, we have seen a variety of VAE implementations to solve multivariate problems. In this section, we present benchmark the applications of VAEs to problems in dimensionality reduction for clustering and manifold learning. We also demonstrate the application of VAEs to real-world datasets from the medical field and test an original generative modelling example.

### **Benchmarking**

#### **Dimensionality reduction for clustering**

Clustering high dimensional data can be difficult as regions in the space become increasingly sparse and data tends to equidistance. To solve this 'curse of dimensionality', dimensionality reduction is often applied in order to control for the correlation between dimensions and to aid in identifying separable clusters. We explored the feasibility of using a VAE to map high dimensional data to a lower-dimensional space for use in clustering.

We trained VAEs, on the Iris, Cancer and Wine dataset and visualised their latent space. As a benchmark, we also implemented Principal Component Analysis (PCA),

Independent Component Analysis (ICA) and Kernel Principal Component Analysis (KPCA). Figures 2, 3 and 4 show the input data mapped to the latent space, point is coloured according to class.

By observation, the PCA for all data sets produces latent spaces with significantly more variation within classes compared to the other models. The KPCA seems to be the most consistent across datasets in the way that it maps inputs to the latent space. The ICA seems to produce tightly condensed clusters that will be difficult to separate.

To get a quantitative sense of the feasibility of clustering data in the latent space, we needed to create a ground truth clustering measure. We used the original class labels computed silhouette scores. These scores summarise the level of cohesion within clusters and separation between clusters. Scores range from -1 to 1, and high values mean that objects of a cluster are cohesive and well separated to other clusters.

Figure 5 shows the silhouette scores for each method trained on Iris, Cancer and Wine. The KPCA methods consistently produce scores that are above 0.4. The KPCA(RBF) method produces the highest score with the Wine dataset. The VAEs performance varies significantly between datasets. The VAE(Tanh) performs the best when trained on the Iris and Cancer dataset but is the weakest performer on the wine dataset. The VAE(ReLU) only performs well with the Cancer dataset and performs poorly with the Iris and Wine dataset. It appears that the VAEs performance is quite sensitive to the type of activation function used in training. Given the lack of consistent performance, the feasibility of using a VAE for efficient dimensionality reduction as a pre-processing step for clustering is questionable. While methods such as KPCA produced consistent results, there is a trade-off. The KPCA method is a basis expansion method, and so time complexity is likely to be a problem.

## Manifold Learning

## Application to Medical Heartbeats Audio Data

For the application section of this report, a public medical dataset of patient heartbeats was chosen. This dataset, provided by Bentley et al., was gathered from both general public via the iStethoscope Pro iPhone app, and from clinical trials in hospitals using the digital stethoscope DigiScope. Many of these 832 are noisy and contain artefacts in the data. The dataset has labels for the recordings, marking them as recording either normal heartbeats, recordings with heart murmurs or extrahls or recordings with noise artefacts which make classification by physicians challenging. Even though the distributions of subsets of the data may be different the aims of this task was to use the data collected from the clinical trial using the digital stethoscope DigiScope to train a Variational Autoencoder to encode data from the iStethoscope Pro iPhone app collected from the general public and identify whether, using this data we could accurately identify recordings with artefacts, as well as heart conditions.

### Data processing

The original data structure was stored in wav files. A wav file is a raw uncompressed audio file that stores audio data, sample rate, and bit rate. The sample rate was 44,100Hz per second, and so in the time domain, there are too many observations to efficiently perform machine learning. Data was read from the wav files and normalized in order to ensure that the audio tracks were roughly the same volume. A simple filter was used in order to clip loud pops from the recordings and split in individual heartbeats by searching for extrema in the signal from within a specified neighbourhood of one-thousand time-steps and then subsampling to ensure the recordings were of equal length. By splitting the data into individual heartbeats, we massively increased the sample size and also managed to better control for elements of the waveform which determined heart-rate as opposed to artefacts, heart murmurs or normal heartbeats. We used a Fast Fourier Transform (FFT) method to transform the high dimensional audio data from an amplitude-time domain to an amplitude-frequency domain of 40 dimensions. The final dataset yielded 11483 samples of heartbeats for in training our model.

While we did experiment with various hyper-parameters our final Variational Autoencoder featured 594 parameters on the encoder with two hidden layers of 15 and 10 nodes each using the tanh activation function, the model was trained over 1000 epochs using the ADAM optimizer with an l2 regularization term of 0.0001 and a one-to-one weighting between our reconstruction loss and the weighting of our KL Divergence term.

### Results

Figure 8 shows the latent space of the trained VAE across three random initializa-

tions. While the data is characterized by imperfect sample recordings, based on the figure 8 and results of tables 1 and 2, we can see that the model tends to separate our the data artefacts at the tails of the distribution, with normal heartbeats closest to the mean. This suggests that for data-cleaning and anomaly detection, the Variational Autoencoder may present value to further downstream analysis.

<b>Component 1</b>		
	mean	std
label		
artifact	-0.565199	0.736880
extrahls	-0.634950	0.719695
murmur	-0.981001	0.404676
normal	-0.795812	0.568068

Table 1: Summary statistics on component 1 for each heartbeat label

<b>Component 2</b>		
	mean	std
label		
artifact	0.172482	1.096362
extrahls	-0.313284	1.023796
murmur	-0.077684	0.546320
normal	-0.198015	0.840402

Table 2: Summary statistics on component 2 for each heartbeat label



## Generative Modeling

Since their popularization, a widely publicized application of Variational Autoencoders has been in generating new images. Using the prior distribution of the latent space, researchers will commonly sample from this distribution and use the decoder model to generate sampled from the high dimensional data manifold. While this method is novel, this method does make several assumptions around the size of the data, the capacity of the network and the underlying distribution of the data.

In order to experiment and test the notion of sampling from complicated distributions, we took to experiment with an extremely shallow neural network architecture- with only 6 and 3 hidden nodes with Rectified Linear Unit (ReLU) activation functions- and a simple distribution- the truncated Bivariate Gaussian. In this example, we aimed to compare the ability of the Variational Autoencoder- given enough samples- to embed this truncated bivariate data onto a one-dimensional line in a manner which maintained the original distributional properties of the data.

Figure 9 shows the newly sampled data superimposed with the original Gaussian dataset. The random samples generated from the VAE seem to have low variance and low bias- sampling points from around the mean of the truncated distribution. This turns out to be a complicated yet straightforward problem to solve for methods in dimensionality reduction as similar properties arise when applying methods such as Principle Component Analysis, as shown in figure 10, and Kernel Principal Component Analysis, as shown in figure 11. When trying to use linear dimensionality reduction, it is impossible to find a reconstruction which accurately reconstructs the data and when using a kernel method methods either tend to high variance or the mean of the distribution as in the case of the kernel method.

While the model architecture and choice of loss function weighting make a huge difference in the stability and capacity of the model to generate samples accurately, it is clear that these methods rely on the correlations in the data in order to accurately embed the data in a lower-dimensional space. This is not easy for all datasets, as shown when given non-symmetric data from the truncated Bivariate Gaussian.

## Appendixes

### Benchmarking

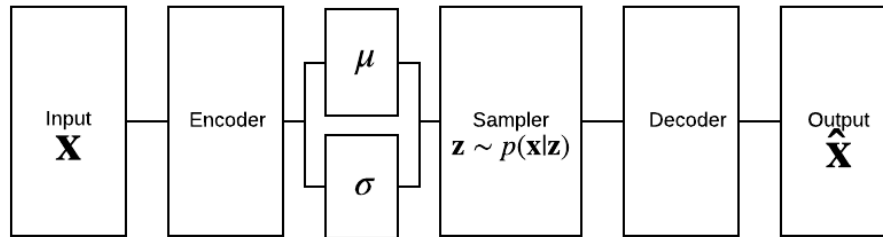


Figure 1: Basic architecture of a Variational Autoencoder

Iris Dataset

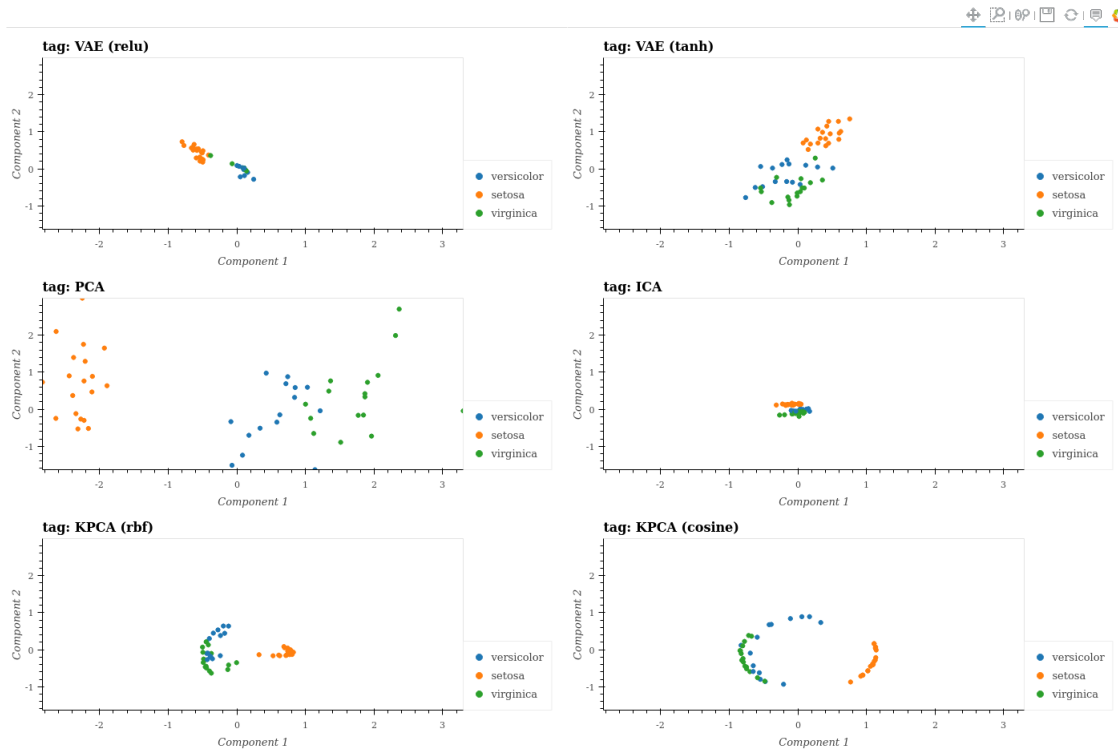


Figure 2: Iris dataset Latent Space representation

## Manifold Learning

# Multivariate Report

Cancer Dataset

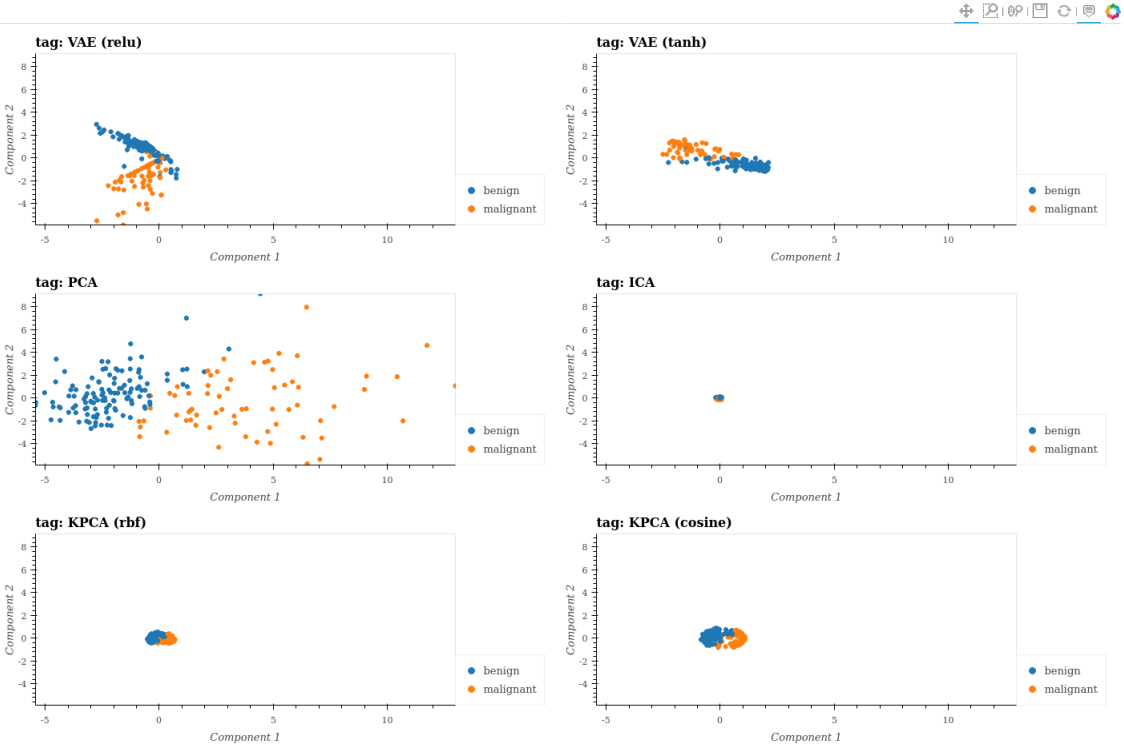


Figure 3: Cancer dataset Latent Space representation

## Heartbeats Audio

## Multivariate Report

---

Wine Dataset

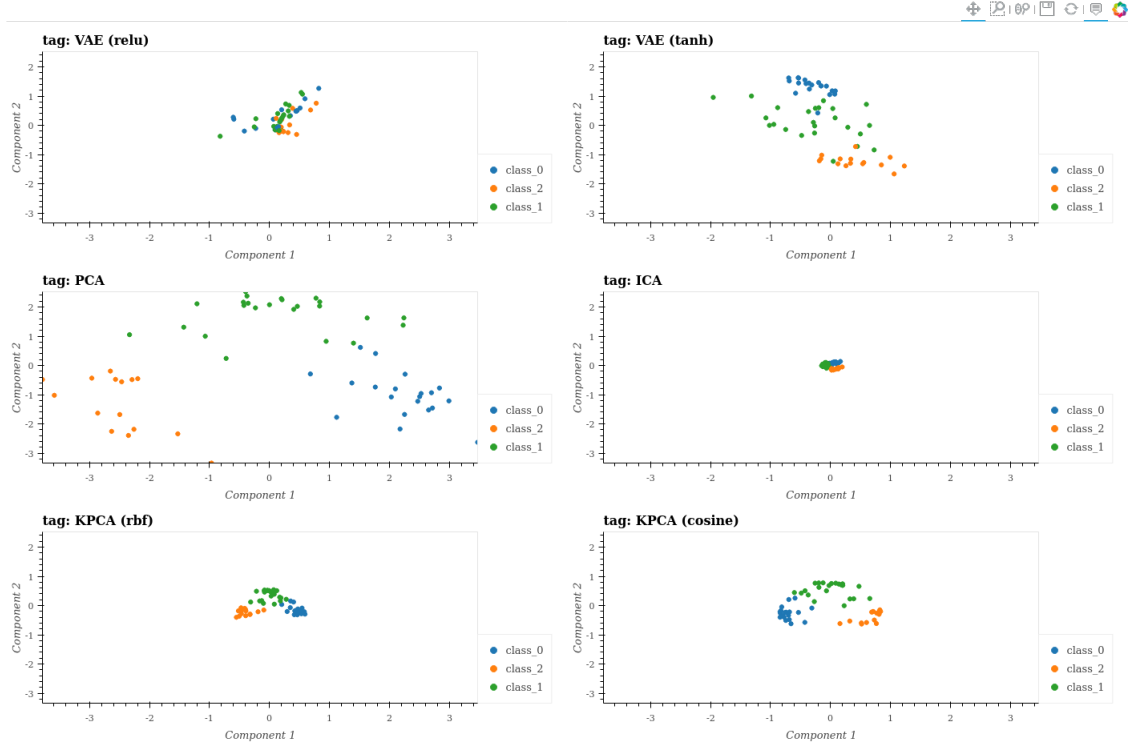


Figure 4: Wine dataset Latent Space representation

## Generative Sampling

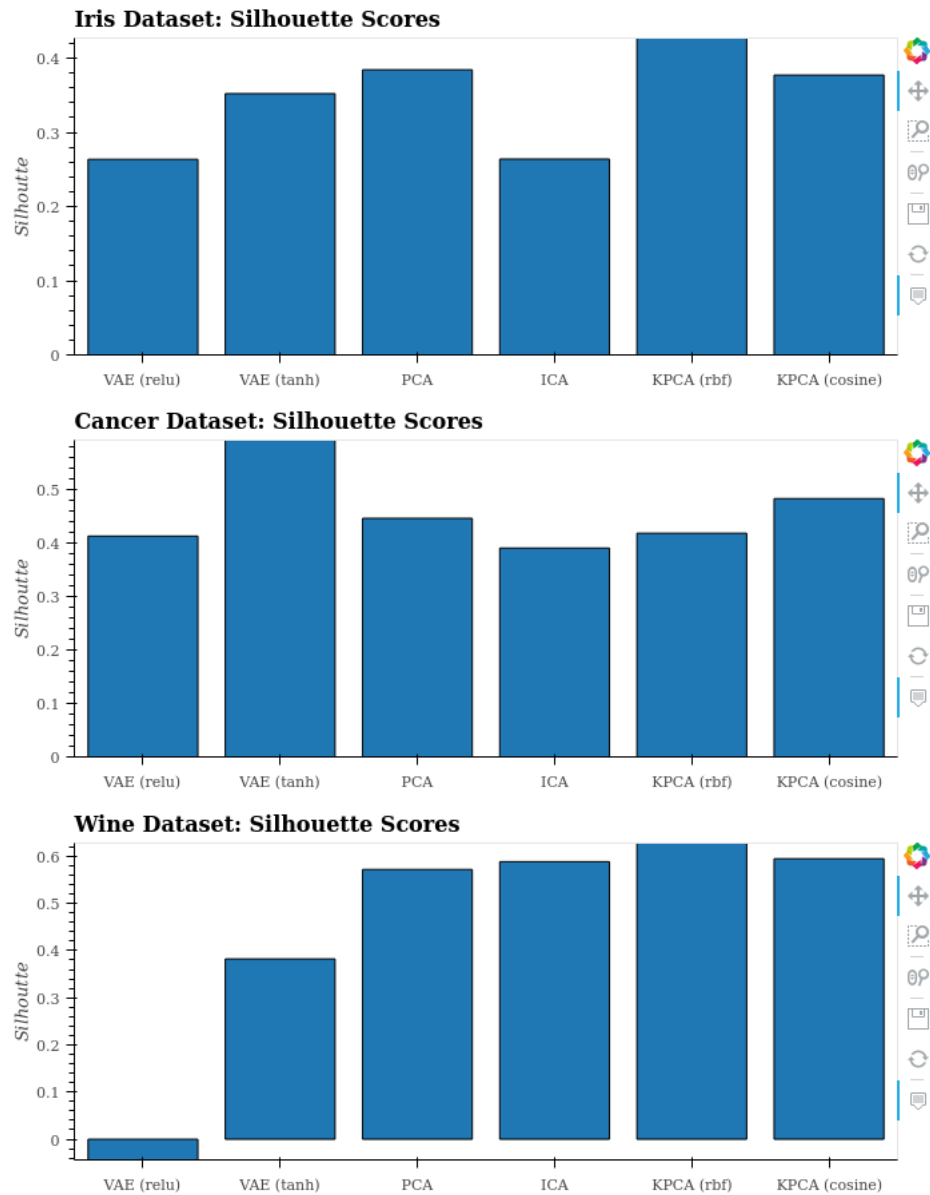


Figure 5: Silhouette Scores

## Code

### Exploratory Analysis

```
1 # ---  
2 # jupyter:  
3 #   jupyter:  
4 #     text_representation:
```

## Multivariate Report

S-Curve Manifold

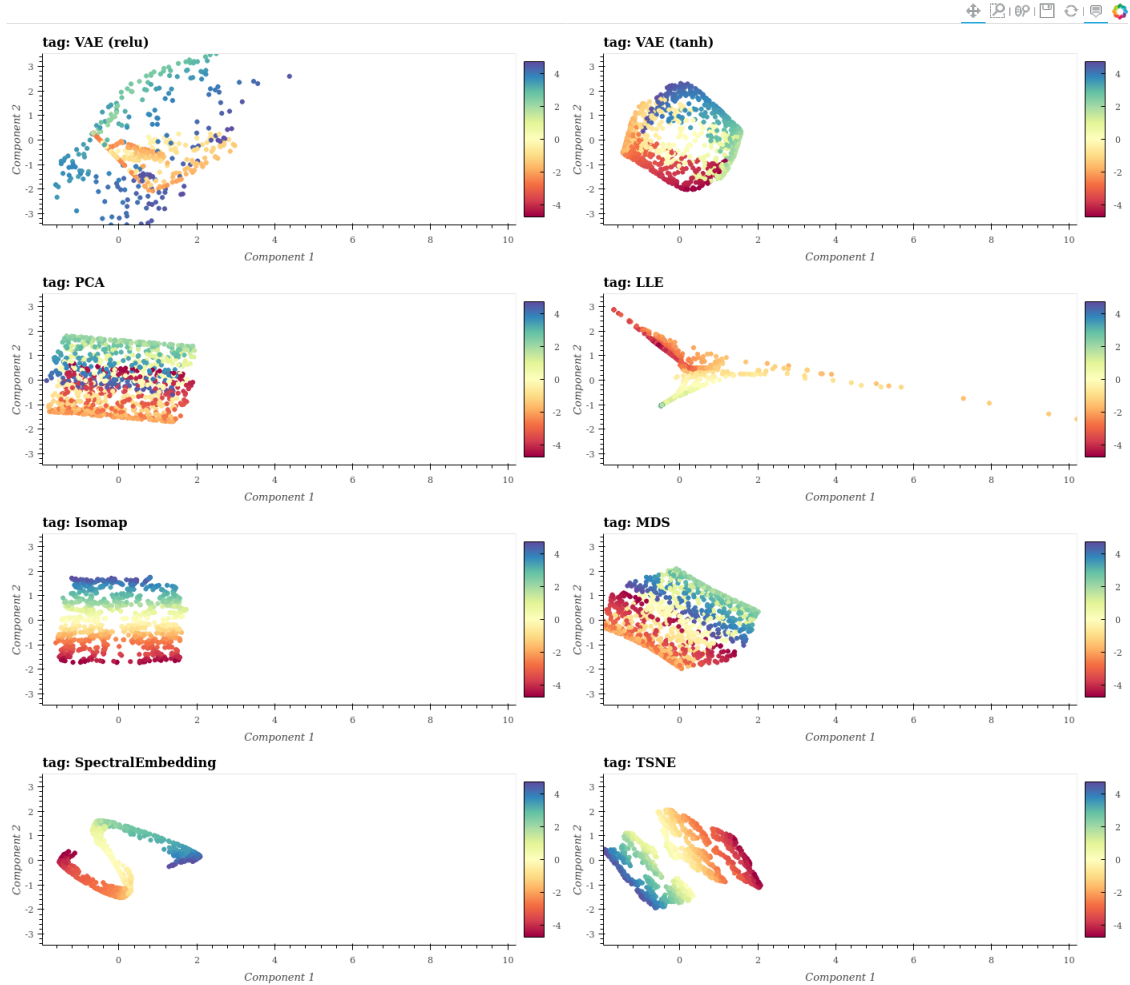


Figure 6: S-Curve Manifold Learning

```
5 # extension: .py
6 # format_name: percent
7 # format_version: '1.2'
8 # jupyter_text_version: 1.2.4
9 # kernelspec:
10 #   display_name: Python 3
11 #   language: python
12 #   name: python3
13 # ---
14
15 # %% {"slideshow": {"slide_type": "skip"}, "language": "html"}
16 # <style>
17 # div.input {
```

## Multivariate Report

```
18 #     display:none;
19 # }
20 # </style>
21
22 # %% {"slideshow": {"slide_type": "skip"}, "language": "html"}
23 # <style>
24 #   div.input {
25 #     display:contents;
26 #   }
27 # </style>
28
29 # %% {"slideshow": {"slide_type": "skip"}}
30 import sys
```

Swill-roll Manifold

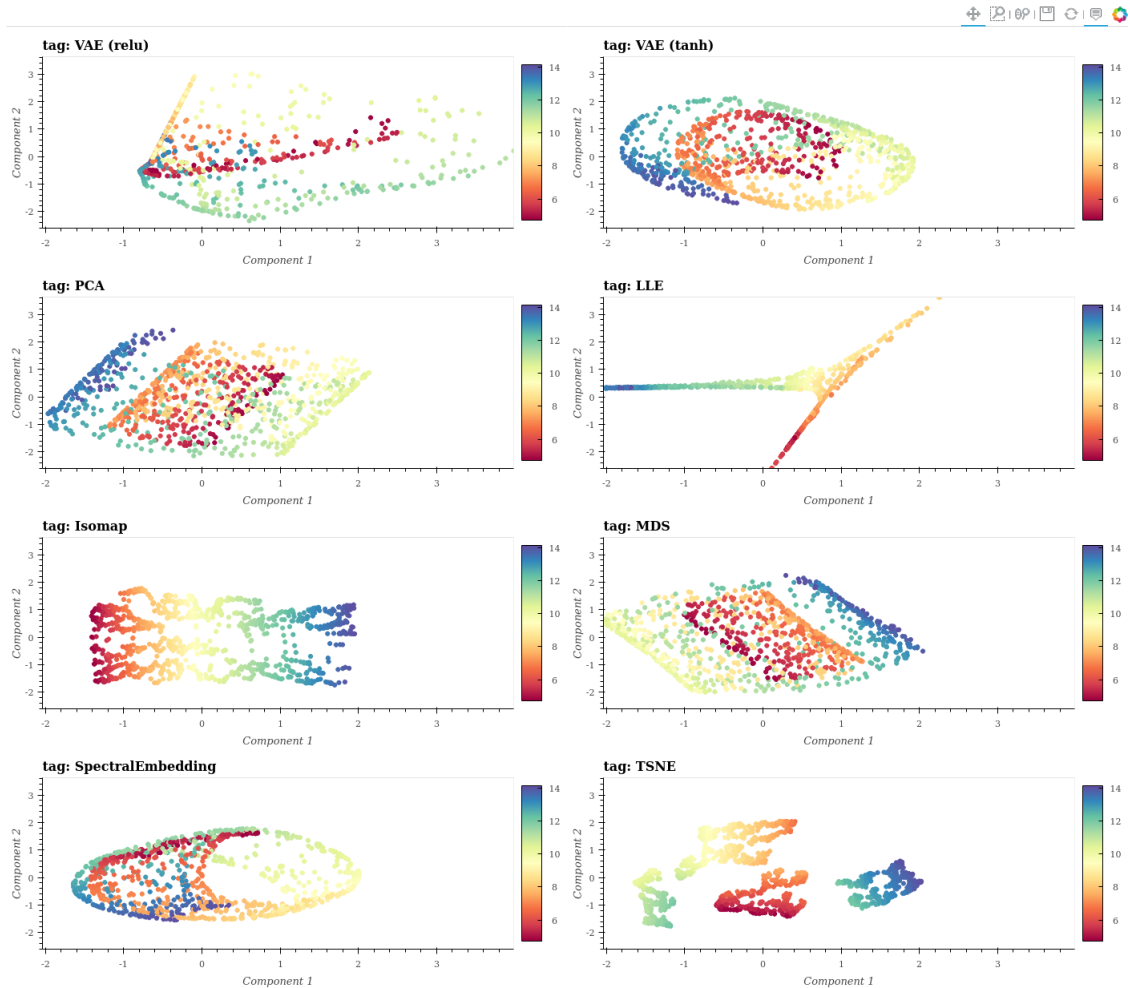


Figure 7: Swiss Roll manifold learning

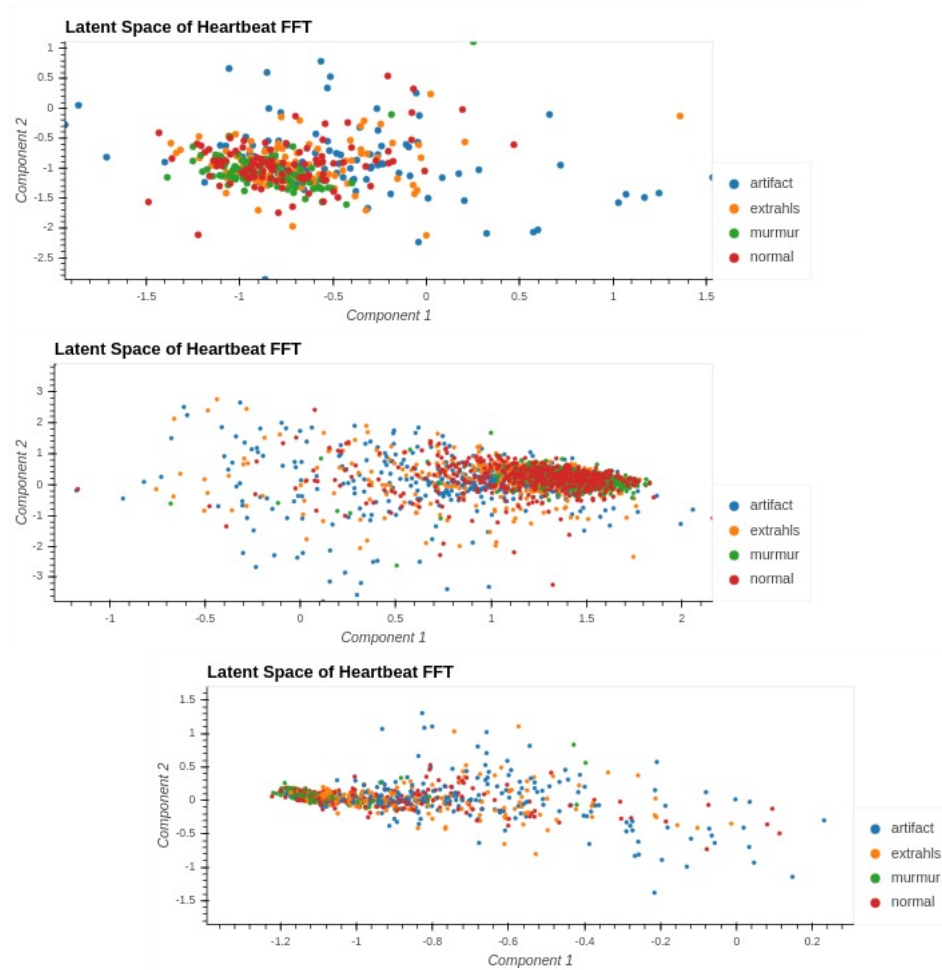


Figure 8: Heartbeat Latent Space Representation

```
31 from sklearn.datasets import load_iris
32 import pandas as pd
33 import numpy as np
34 import holoviews as hv
35 import hvplot.pandas
36 from toolz.curried import *
37 from sklearn.decomposition import PCA
38 from sklearn.manifold import MDS, TSNE
39 from sklearn.preprocessing import StandardScaler
40 from sklearn.pipeline import make_pipeline
41 from sklearn.model_selection import train_test_split, KFold
42
43
44 sys.path.append("../")
```



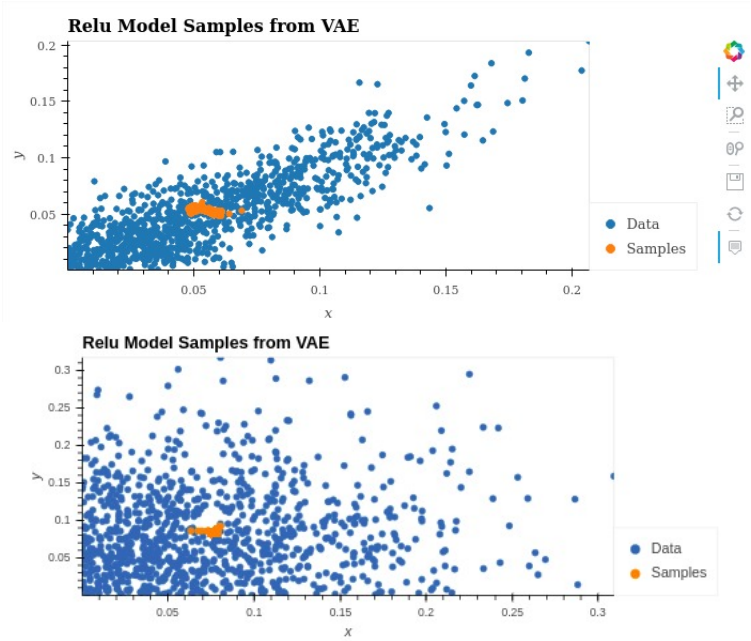


Figure 9: Gaussian Generative Modeling example

```
45 hv.extension("bokeh")
```

```
46
```

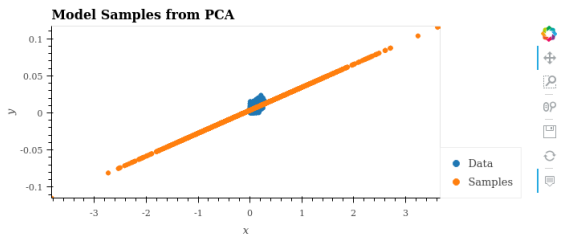


Figure 10: PCA on Bivariate Gaussian Generative Modeling example

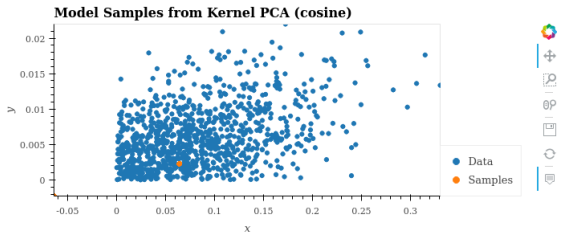


Figure 11: Cosine Kernel PCA on Bivariate Gaussian Generative Modeling

```
47 # %% {"slideshow": {"slide_type": "skip"}}
48 from super_spirals.neural_network import VAE
49
50 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
51 # # Preliminary Analysis
52
53 # %% {"slideshow": {"slide_type": "skip"}}
54 X = load_iris()
55
56 # %% {"slideshow": {"slide_type": "skip"}}
57 pipeline = make_pipeline(
58     StandardScaler(),
59     VAE(
60         hidden_layer_sizes=(5, 2), activation="tanh", divergence_weight
61         =5, max_iter=500
62     ),
63 )
64 # %% {"slideshow": {"slide_type": "skip"}}
65 pipeline.fit(X=X.data)
66
67 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
68 # __Inspect model__
69
70 # %% [markdown] {"slideshow": {"slide_type": "-"}}
71 # Encoder
72
73 # %% {"slideshow": {"slide_type": "-"}}
74 pipeline.named_steps["vae"].encoder.summary()
75
76 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
77 # __Decoder__
78
79 # %% {"slideshow": {"slide_type": "-"}}
80 pipeline.named_steps["vae"].decoder.summary()
81
82 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
83 # __Denoising__
84
85 # %% [markdown] {"slideshow": {"slide_type": "-"}}
86 # Original data
87
88 # %% {"slideshow": {"slide_type": "-"}}
89 # original data
90 sample = X.data[:10, :]
91 pipe(sample, partial(pd.DataFrame, columns=X.feature_names))
92
93 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
94 # 'Denoised' Data
```

```
95
96 # %% {"slideshow": {"slide_type": "-"}}
97 # denoised
98 pipe(
99     sample,
100     pipeline.transform,
101     pipeline.inverse_transform,
102     partial(pd.DataFrame, columns=X.feature_names),
103 )
104
105 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
106 # __Generate new samples__
107
108 # %% {"slideshow": {"slide_type": "-"}}
109 pipe(
110     np.random.multivariate_normal(np.zeros(2), np.diag(np.ones(2))),
111     size=10),
112     pipeline.inverse_transform,
113 )
114
115 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
116 # __Dimensionality Reduction__
117
118 # %% {"slideshow": {"slide_type": "-"}}
119 # %%output filename='../media/01-iris-latent' fig='png'
120 (
121     pipe(
122         X.data,
123         pipeline.transform,
124         partial(pd.DataFrame, columns=["Component 1", "Component 2"]),
125     )
126     .assign(label=X.target)
127     .assign(label=lambda d: d.label.replace(dict(enumerate(X.
128 target_names))))
129     .hvplot.scatter(
130         x="Component 1",
131         y="Component 2",
132         color="label",
133         label="Variational Autoencoder Latent Encoding",
134     )
135 )
136
137 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
138 # __Anomaly Detection__
139
140 # %% [markdown] {"slideshow": {"slide_type": "-"}}
141 # Assume normal, use Z-scores to filter outliers
```

### Dataset Benchmark Code

```
1 # ——
2 # jupyter:
3 #   jupyter:
4 #     text_representation:
5 #       extension: .py
6 #       format_name: percent
7 #       format_version: '1.2'
8 #       jupyter_version: 1.2.4
9 #   kernelspec:
10 #     display_name: Python 3
11 #     language: python
12 #     name: python3
13 # ——
14
15 # %% {"slideshow": {"slide_type": "skip"}, "language": "html"}
16 # <style>
17 #   div.input {
18 #     display: none;
19 #   }
20 # </style>
21
22 # %% {"slideshow": {"slide_type": "skip"}, "language": "html"}
23 # <style>
24 #   div.input {
25 #     display: contents;
26 #   }
27 # </style>
28
29 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
30 # # Dimensionality Reduction
31 # 1. Compressed representation which capture all relevant structure in
32 #    the data
33 # 2. Representation useful in supervised and unsupervised learning
34 #    tasks
35
36 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
37 # #### Models
38 # 1. VAE (Tanh)
39 # 2. VAE (ReLU)
40 # 3. PCA
41 # 4. ICA
42 # 5. Kernel PCA (RBF)
43 # 6. Kernel PCA (Cosine)
44
45 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
46 # #### Data
47 # 1. Fisher's Iris data
```

```
46 # 2. UCI ML Wine Data
47 # 3. UCI ML Breast Cancer Wisconsin (Diagnostic)
48 #
49
50 # %% {"slideshow": {"slide_type": "skip"}}
51 import sys
52 from sklearn.datasets import (
53     load_iris,
54     load_wine,
55     load_breast_cancer,
56     make_multilabel_classification,
57 )
58 import pandas as pd
59 import numpy as np
60 import holoviews as hv
61 import hvplot.pandas
62 from typing import Dict
63 from toolz.curried import *
64 from sklearn.decomposition import PCA, FastICA, KernelPCA
65 from sklearn.preprocessing import StandardScaler
66 from sklearn.pipeline import make_pipeline
67
68
69 sys.path.append("../")
70 hv.extension("bokeh")
71
72 # %% {"slideshow": {"slide_type": "skip"}}
73 from super_spirals.metrics import reconstruction_benchmark
74 from super_spirals.neural_network import VAE
75
76
77 # %% {"slideshow": {"slide_type": "skip"}}
78 def get_models():
79     return {
80         "VAE (relu)": make_pipeline(
81             StandardScaler(),
82             VAE(activation="relu", max_iter=300, hidden_layer_sizes=(4,
83 2)),
84         ),
85         "VAE (tanh)": make_pipeline(
86             StandardScaler(),
87             VAE(activation="tanh", max_iter=300, hidden_layer_sizes=(4,
88 2)),
89         ),
90         "PCA": make_pipeline(StandardScaler(), PCA(n_components=2)),
91         "ICA": make_pipeline(StandardScaler(), FastICA(n_components=2))
92     },
93     "KPCA (rbf)": make_pipeline(
94         StandardScaler(),
```

```
92         KernelPCA(n_components=2, kernel="rbf",
93             fit_inverse_transform=True),
94         "KPCA (cosine)": make_pipeline(
95             StandardScaler(),
96             KernelPCA(n_components=2, kernel="cosine",
97                 fit_inverse_transform=True),
98         )
99     }
100
101 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
102 # # Iris
103
104 # %% {"slideshow": {"slide_type": "skip"}}
105 iris_models = get_models()
106 iris_df, iris_reconstruction, iris_silhouette, iris_plot =
107     reconstruction_benchmark(
108         load_iris(), iris_models, "Iris Dataset"
109     )
110
111 # %% {"slideshow": {"slide_type": "slide"}}
112 # %%output filename='../media/02-iris-loss' fig='png'
113 iris_reconstruction
114
115 # %% {"slideshow": {"slide_type": "slide"}}
116 # %%output filename='../media/02-iris-silhouette' fig='png'
117 iris_silhouette
118
119 # %% {"slideshow": {"slide_type": "slide"}}
120 # %%output filename='../media/02-iris-latent' fig='png'
121 iris_plot
122
123 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
124 # # Wine
125
126 # %% {"slideshow": {"slide_type": "skip"}}
127 wine_models = get_models()
128 wine_df, wine_reconstruction, wine_silhouette, wine_plot =
129     reconstruction_benchmark(
130         load_wine(), wine_models, "Wine Dataset"
131     )
132
133 # %% {"slideshow": {"slide_type": "slide"}}
134 # %%output filename='../media/02-wine-loss' fig='png'
135 wine_reconstruction
136
137 # %% {"slideshow": {"slide_type": "slide"}}
138 # %%output filename='../media/02-wine-silhouette' fig='png'
```

```
137 wine_silhouette
138
139 # %% {"slideshow": {"slide_type": "slide"}}
140 # %%output filename='../media/02-wine-latent' fig='png'
141 wine_plot
142
143 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
144 # # Cancer
145
146 # %% {"slideshow": {"slide_type": "skip"}}
147 cancer_models = get_models()
148 cancer_df, cancer_reconstruction, cancer_silhouette, cancer_plot =
    reconstruction_benchmark(
149     load_breast_cancer(), cancer_models, "Cancer Dataset"
150 )
151
152 # %% {"slideshow": {"slide_type": "slide"}}
153 # %%output filename='../media/02-cancer-loss' fig='png'
154 cancer_reconstruction
155
156 # %% {"slideshow": {"slide_type": "slide"}}
157 # %%output filename='../media/02-cancer-silhouette' fig='png'
158 cancer_silhouette
159
160 # %% {"slideshow": {"slide_type": "slide"}}
161 # %%output filename='../media/02-cancer-latent' fig='png'
162 cancer_plot
163
164
165 # %% {"slideshow": {"slide_type": "skip"}}
166 class get_random:
167     random = make_multilabel_classification(n_samples=1000, n_features
    =20, n_classes=5)
168     data = random[0]
169     target = pd.DataFrame(random[1]).idxmax(1).to_numpy()
170
171
172 # %% {"slideshow": {"slide_type": "skip"}}
173 random_models = get_models()
174 random_df, random_reconstruction, random_silhouette, random_plot =
    reconstruction_benchmark(
175     get_random(), random_models, "Random Dataset"
176 )
177
178 # %% {"slideshow": {"slide_type": "skip"}}
179 # %%output filename='../media/02-random-loss' fig='png'
180 random_reconstruction
181
182 # %% {"slideshow": {"slide_type": "skip"}}
```

```
183 # %%output filename='../media/02-random-silhouette' fig='png'
184 random_silhouette
185
186 # %% {"slideshow": {"slide_type": "skip"}}
187 # %%output filename='../media/02-random-latent' fig='png'
188 random_plot
189
190 # %% {"slideshow": {"slide_type": "skip"}}
```

### Manifold Learning Comparison

```
1 # —
2 # jupyter:
3 #   jupyter:
4 #     text_representation:
5 #       extension: .py
6 #       format_name: percent
7 #       format_version: '1.2'
8 #       jupyter_text_version: 1.2.4
9 #   kernelspec:
10 #     display_name: Python 3
11 #     language: python
12 #     name: python3
13 # —
14
15 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
16 # # Manifold Learning
17 # 1. Data Exists in Manifold within the high dimensional space
18 # 2. Often non-linear surface
19
20 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
21 # ### Models
22 # - VAE (relu)
23 # - VAE (tanh)
24 # - PCA
25 # - LLE
26 # - Isomap
27 # - MDS
28 # - SpectralEmbedding
29 # - TSNE
30
31 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
32 # ### Data
33 # - S-curve
34 # - Swiss Roll
35
36 # %% {"slideshow": {"slide_type": "skip"}, "language": "html"}
37 # <style>
38 #   div.input {
```



```
39 #         display:none;
40 #     }
41 # </style>
42
43 # %% {"slideshow": {"slide_type": "skip"}, "language": "html"}
44 # <style>
45 #     div.input {
46 #         display:contents;
47 #     }
48 # </style>
49
50 # %% {"slideshow": {"slide_type": "skip"}}
51 import sys
52 from sklearn.datasets import samples_generator
53 import pandas as pd
54 import numpy as np
55 import holoviews as hv
56 import hvplot.pandas
57 from toolz.curried import *
58 from sklearn import manifold
59 from sklearn.decomposition import PCA
60 from sklearn.preprocessing import StandardScaler
61 from sklearn.pipeline import make_pipeline
62 from sklearn.model_selection import train_test_split, KFold
63
64
65 sys.path.append("../")
66 hv.extension("bokeh")
67
68 # %% {"slideshow": {"slide_type": "skip"}}
69 from super_spirals.neural_network import VAE
70
71 # %% {"slideshow": {"slide_type": "skip"}}
72 n_points = 1000
73
74
75 # %% {"slideshow": {"slide_type": "skip"}}
76 def get_models():
77     return {
78         "VAE (relu)": make_pipeline(
79             StandardScaler(),
80             VAE(activation="relu", max_iter=300, hidden_layer_sizes=(4,
81                 5, 2)),
82         "VAE (tanh)": make_pipeline(
83             StandardScaler(),
84             VAE(activation="tanh", max_iter=300, hidden_layer_sizes=(4,
85                 5, 2)),
86     )
```

```

86         "PCA": make_pipeline(StandardScaler(), PCA(n_components=2)),
87         "LLE": make_pipeline(
88             StandardScaler(), manifold.LocallyLinearEmbedding(
89                 n_components=2
90             ),
91         "Isomap": make_pipeline(StandardScaler(), manifold.Isomap(
92             n_components=2)),
93         "MDS": make_pipeline(StandardScaler(), manifold.MDS(
94             n_components=2)),
95         "SpectralEmbedding": make_pipeline(
96             StandardScaler(), manifold.SpectralEmbedding(n_components
97                 =2)
98         ),
99         "TSNE": make_pipeline(StandardScaler(), manifold.TSNE(
100             n_components=2)),
101     }
102
103     # %% {"slideshow": {"slide_type": "skip"}}
104     def get_components(model, X, y, tag):
105
106         latent = pipe(
107             X,
108             model.fit_transform,
109             StandardScaler().fit_transform,
110             PCA(whiten=True).fit_transform,
111             partial(pd.DataFrame, columns=["Component 1", "Component 2"]),
112         )
113
114         return latent.assign(y=y).assign(tag=tag)
115
116     # %% {"slideshow": {"slide_type": "skip"}}
117     s_curve_models = get_models()
118     s_curve_X, s_curve_color = samples_generator.make_s_curve(n_points,
119         random_state=0)
120
121     s_curve_components = pd.concat(
122         [get_components(m, s_curve_X, s_curve_color, t) for t, m in
123          s_curve_models.items()]
124     )
125
126     # %% {"slideshow": {"slide_type": "slide"}}
127     # %%output filename='../media/03-scurve-latent' fig='png'
128     (
129         s_curve_components.hvplot.scatter(
130             x="Component 1", y="Component 2", color="y", groupby="tag",
131             cmap="spectral"
132         )
133     )

```

```
127     .layout()
128     .opts(title="S-Curve Manifold", shared_axes=False)
129     .cols(2)
130 )
131
132 # %% {"slideshow": {"slide_type": "skip"}}
133 swissroll_models = get_models()
134 swissroll_X, swissroll_color = samples_generator.make_swiss_roll(
135     n_points, random_state=0
136 )
137
138 swissroll_components = pd.concat(
139     [
140         get_components(m, swissroll_X, swissroll_color, t)
141         for t, m in swissroll_models.items()
142     ]
143 )
144
145 # %% {"slideshow": {"slide_type": "slide"}}
146 # %%output filename='../media/03-swissroll-latent' fig='png'
147 (
148     swissroll_components.hvplot.scatter(
149         x="Component 1", y="Component 2", color="y", groupby="tag",
150         cmap="spectral"
151     )
152     .layout()
153     .opts(title="Swill-roll Manifold", shared_axes=False)
154     .cols(2)
155 )
```

## Latent Space Sampling

```
1 # —
2 # jupyter:
3 #   jupyter:
4 #     text_representation:
5 #       extension: .py
6 #       format_name: percent
7 #       format_version: '1.2'
8 #       jupyter_version: 1.2.4
9 #     kernelspec:
10 #       display_name: Python 3
11 #       language: python
12 #       name: python3
13 # —
14
15 # %% [markdown] {"slideshow": {"slide_type": "skip"}}
16 # # Sampling
17
```

## Multivariate Report

---

```
18 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
19 # ____Aim____
20 # Truncated Bivariate Gaussian  $\rightarrow$  1D Representation  $\rightarrow$  Truncated Bivariate Gaussian
21
22 # %% {"slideshow": {"slide_type": "skip"}}
23 import sys
24 import pandas as pd
25 import numpy as np
26 from sklearn.decomposition import PCA, KernelPCA
27 import holoviews as hv
28 import hvplot.pandas
29 from toolz.curried import *
30
31
32 sys.path.append("../")
33 hv.extension("bokeh")
34
35 # %% {"slideshow": {"slide_type": "skip"}}
36 from super_spirals.neural_network import VAE
37
38 # %% {"slideshow": {"slide_type": "skip"}}
39 X = np.random.multivariate_normal(
40     mean=np.zeros(2), cov=np.diag(np.ones(2)), size=100000
41 )
42
43 beta = np.random.uniform(-0.1, 0.1, size=(2, 2))
44 data = pd.DataFrame(X @ beta).where(lambda d: d > 0).dropna(how="any").
45     to_numpy()
46
47 # %% {"slideshow": {"slide_type": "skip"}}
48 vae = VAE(
49     hidden_layer_sizes=(6, 3, 1),
50     max_iter=500,
51     activation="relu",
52     alpha=0.0005,
53     divergence_weight=0,
54     batch_size=100000 / 10,
55 )
56
57 # %% {"slideshow": {"slide_type": "skip"}}
58 vae.fit(x=data)
59
60 # %% {"slideshow": {"slide_type": "slide"}}
61 # %%output filename='../media/04-generative-samples' fig='png'
62 (
63     (
64         pd.DataFrame(data, columns=["x", "y"])
65         .sample(1000)
```

```
65         .hvplot.scatter(x="x", y="y", label="Data")
66     )
67     * (
68         pd.DataFrame(vae.sample(1000), columns=["x", "y"]).hvplot.
scatter(
69         x="x", y="y", label="Samples"
70     )
71 )
72 ).opts(title="Relu Model Samples from VAE", tools=[])
73
74 # %%
75 kpca = KernelPCA(1, kernel="rbf")
76 kpca.fit(X=data)
77
78 # %% {"slideshow": {"slide_type": "skip"}}
79 # %%output filename='../media/04-generative-pca' fig='png'
80 (
81     (
82         pd.DataFrame(data, columns=["x", "y"])
83         .sample(1000)
84         .hvplot.scatter(x="x", y="y", label="Data")
85     )
86     * (
87         pd.DataFrame(
88             pca.inverse_transform(np.random.normal(size=(1000,)).
reshape(-1, 1)),
89             columns=["x", "y"],
90         ).hvplot.scatter(x="x", y="y", label="Samples")
91     )
92 ).opts(title="Relu Model Samples from PCA", tools=[])
93
94 # %%
95 pca = PCA(1)
96 pca.fit(X=data)
97
98 # %%
99 # %%output filename='../media/04-generative-pca' fig='png'
100 (
101     (
102         pd.DataFrame(data, columns=["x", "y"])
103         .sample(1000)
104         .hvplot.scatter(x="x", y="y", label="Data")
105     )
106     * (
107         pd.DataFrame(
108             pca.inverse_transform(np.random.normal(size=(1000,)).
reshape(-1, 1)),
109             columns=["x", "y"],
110
```

```
111         ).hvplot.scatter(x="x", y="y", label="Samples")
112     )
113 ).opts(title="Relu Model Samples from PCA", tools=[])
```

### Heartbeats Dataset

```
1 # —
2 # jupyter:
3 #   jupyter:
4 #     text_representation:
5 #       extension: .py
6 #       format_name: percent
7 #       format_version: '1.2'
8 #       jupyter_text_version: 1.2.4
9 #   kernelspec:
10 #     display_name: Python 3
11 #     language: python
12 #     name: python3
13 # —
14
15 # %% {"language": "html"}
16 # <style>
17 #   div.input {
18 #     display:none;
19 #   }
20 # </style>
21
22 # %% {"language": "html"}
23 # <style>
24 #   div.input {
25 #     display:contents;
26 #   }
27 # </style>
28
29 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
30 # # Application
31
32 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
33 # ____Data____
34 # 
35
36 # %% {"slideshow": {"slide_type": "skip"}}
37 # # ! kaggle datasets download -d kinguistics/heartbeat-sounds -p ../
38 #   data/raw
39 # %% {"slideshow": {"slide_type": "skip"}}
40 import os
41 import glob
42 import zipfile
```

```
43 import sys
44 from sklearn.datasets import load_iris
45 import pandas as pd
46 import numpy as np
47 import holoviews as hv
48 import hvplot.pandas
49 from toolz.curried import *
50 from sklearn.decomposition import PCA
51 from sklearn.preprocessing import StandardScaler
52 from sklearn.pipeline import make_pipeline
53 from sklearn.model_selection import train_test_split, KFold
54 from scipy.io import wavfile # get the api
55 from scipy.fftpack import fft, irfft
56 from scipy.signal import find_peaks, resample_poly
57 from scipy.spatial import procrustes
58 from scipy.stats import iqr
59 import panel as pn
60 import param
61 from random import sample
62
63
64 sys.path.append("../")
65 hv.extension("bokeh")
66
67 # %% {"slideshow": {"slide_type": "skip"}}
68 from super_spirals.neural_network import VAE
69
70 # %% [markdown] {"slideshow": {"slide_type": "skip"}}
71 # Load Data
72
73 # %% {"slideshow": {"slide_type": "skip"}}
74 data_path = os.path.join("../", "data", "raw", "heartbeat-sounds")
75
76 if not os.path.exists(data_path):
77     with zipfile.ZipFile(
78         os.path.join("../", "data", "raw", "heartbeat-sounds.zip"), "r"
79     ) as zip_ref:
80         zip_ref.extractall(data_path)
81
82 # %% {"slideshow": {"slide_type": "skip"}}
83 heartbeat_path = os.path.join(data_path, "set_b")
84
85 # %% {"slideshow": {"slide_type": "skip"}}
86 files = pipe(
87     os.listdir(heartbeat_path),
88     map(str),
89     map(lambda f: os.path.join(heartbeat_path, f)),
90     list,
91 )
```

```

92
93 # %% {"slideshow": {"slide_type": "skip"}}
94 set_a = pipe(data_path, lambda f: os.path.join(f, "set_a.csv"), pd.
    read_csv)
95
96 # %% {"slideshow": {"slide_type": "skip"}}
97 set_b = pipe(data_path, lambda f: os.path.join(f, "set_b.csv"), pd.
    read_csv)
98
99 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
100 # __Preprocessing__
101 # Data $\rightarrow$ Standardize $\rightarrow$ Clip $\rightarrow$ Split
    into single beats $\rightarrow$ Resample resolution $\rightarrow$
    Fourier Transform
102
103 # %% [markdown] {"slideshow": {"slide_type": "skip"}}
104 # High pass filter
105
106 # %% {"slideshow": {"slide_type": "skip"}}
107 filter_signal = lambda a: (
108     pipe(
109         a,
110         partial(find_peaks, distance=1000),
111         get(0),
112         lambda x: pipe(a[x], np.median),
113         lambda x: np.clip(a, -x, x),
114     )
115 )
116
117 # %% {"slideshow": {"slide_type": "skip"}}
118 get_signal = lambda f: pipe(f, wavfile.read, get(1))
119
120 # %% {"slideshow": {"slide_type": "slide"}}
121 a = pipe(files[0], get_signal)
122
123 b = pipe(files[0], get_signal, filter_signal)
124
125 c = pipe(
126     b,
127     partial(find_peaks, distance=1000),
128     get(0),
129     lambda x: np.vstack((np.arange(b.shape[0]), b)).T[x, :],
130 )
131
132 (
133     hv.Curve(a) * hv.Curve(b).opts(color="orange") * hv.Scatter(c).opts
    (color="green")
134 ).opts(width=600)
135

```



```

136
137 # %% [markdown] {"slideshow": {"slide_type": "skip"}}
138 # split into individual heartbeats, resample to equal length and get
    Fourier Transform
139
140 # %% {"slideshow": {"slide_type": "skip"}}
141 def explode(b, components=10):
142     return pipe(
143         b,
144         partial(find_peaks, distance=1000),
145         get(0),
146         sliding_window(2),
147         map(lambda x: b[x[0] : x[1]]),
148         map(lambda x: (x) / (np.quantile(np.abs(x), 0.9))),
149         map(
150             lambda x: x[
151                 np.round(np.linspace(0, x.shape[0] - 1, num=1000)).
152                 astype(np.int)
153             ],
154         map(partial(fft, n=components)),
155         map(lambda x: np.hstack((np.real(x)).reshape(-1))),
156         list,
157     )
158
159
160 # %% {"slideshow": {"slide_type": "skip"}}
161 def get_explotion(f, components=25):
162     try:
163         return pipe(f, wavfile.read, get(1), partial(explode,
164             components=components))
165     except:
166         return [np.zeros(int(round(components))) * np.nan]
167
168
169 # %% {"slideshow": {"slide_type": "skip"}}
170 frequencies_a = set_a.fname.apply(lambda f: os.path.join(data_path, f))
171     .apply(
172         get_explotion
173     )
174
175 # %% {"slideshow": {"slide_type": "skip"}}
176 frequencies_b = pipe(files, pd.Series).apply(get_explotion)
177
178 # %% [markdown] {"slideshow": {"slide_type": "skip"}}
179 # Merge data with frequencies
180
181 # %% {"slideshow": {"slide_type": "skip"}}

```

```

181 set_a_freq = (
182     set_a.assign(frequencies=frequencies_a)
183     .explode("frequencies")
184     .reset_index(drop=True)
185     .assign(label=lambda d: d.label.fillna("None"))
186     .where(lambda d: ~d.label.str.startswith("None"))
187     .dropna(how="all")
188 )
189
190 # %% {"slideshow": {"slide_type": "skip"}}
191 X_a = set_a_freq.frequencies.apply(pd.Series)
192
193 # %% {"slideshow": {"slide_type": "skip"}}
194 set_a_filtered = set_a_freq.loc[~X_a.isna().all(axis=1), :]
195
196 # %% {"slideshow": {"slide_type": "skip"}}
197 X_b = (
198     pd.DataFrame({"frequencies": frequencies_b})
199     .explode("frequencies")
200     .reset_index(drop=True)
201     .frequencies.apply(pd.Series)
202 )
203
204 # %% [markdown] {"slideshow": {"slide_type": "skip"}}
205 # #### Train
206
207 # %% {"slideshow": {"slide_type": "skip"}}
208 pipeline = make_pipeline(
209     StandardScaler(),
210     PCA(whiten=True),
211     VAE(
212         hidden_layer_sizes=(15, 10, 2),
213         max_iter=1000,
214         divergence_weight=10,
215         activation="tanh",
216     ),
217 )
218
219 # %% {"slideshow": {"slide_type": "skip"}}
220 pipeline.fit(X_b.dropna())
221
222 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
223 # __Model__
224
225 # %% {"slideshow": {"slide_type": "-"}}
226 pipeline.named_steps["vae"].encoder.summary()
227
228 # %% [markdown] {"slideshow": {"slide_type": "slide"}}
229 # #### Dashboard

```

```

230
231 # %% {"slideshow": {"slide_type": "skip"}}
232 latent_a = pipeline.transform(X_a.dropna()) # .loc[~X_a.isna().all(
      axis=1),:]
233
234 # %% {"slideshow": {"slide_type": "skip"}}
235 latent_a_df = (
236     pd.DataFrame(latent_a, columns=["Component 1", "Component 2"])
237     .assign(label=set_a_filtered.label.fillna("None"))
238     .reset_index()
239     .groupby("label")
240     .apply(lambda d: d.sample(250, replace=False))
241     .reset_index(drop=True)
242     .set_index("index")
243 )
244
245 # %% {"slideshow": {"slide_type": "skip"}}
246 latent_a_df.head()
247
248 # %% {"slideshow": {"slide_type": "skip"}}
249 print(
250     latent_a_df.loc[:, ["Component 1", "label"]]
251     .groupby("label")
252     .describe()
253     .iloc[:, 1:3]
254     .to_latex()
255 )
256
257 # %% {"slideshow": {"slide_type": "skip"}}
258 print(
259     latent_a_df.loc[:, ["Component 2", "label"]]
260     .groupby("label")
261     .describe()
262     .iloc[:, 1:3]
263     .to_latex()
264 )
265
266 # %% {"slideshow": {"slide_type": "skip"}}
267 clips = set_a_filtered.loc[latent_a_df.index, "fname"].to_list()
268
269
270 # %% {"slideshow": {"slide_type": "skip"}}
271 class Dashboard(param.Parameterized):
272     files = pn.widgets.Select(name="Audio Clip", value=clips[0],
      options=clips)
273
274     @pn.depends("files.value")
275     def update(self, index):
276         if index:

```

```
277         self.files.value = clips[index[0]]
278         wav_file = pipe(self.files.value, lambda f: os.path.join(
data_path, f))
279
280         data = pipe(wav_file, wavfile.read, get(1))
281
282         time = pipe(data, lambda x: x[:, :400] / np.max(np.abs(x)), hv.
Curve).opts(
283             width=400, xlabel="time", ylabel="waveform", height=300
284         )
285
286         frequency = pipe(
287             data,
288             partial(fft, n=1000),
289             np.real,
290             lambda x: x / np.max(np.abs(x)),
291             hv.Curve,
292         ).opts(xlabel="frequency", ylabel="aplitude", width=400, height
=300)
293
294         return time + frequency
295
296     @pn.depends("files.value")
297     def view(self):
298
299         latent = latent_a_df.hvplot.scatter(
300             x="Component 1",
301             y="Component 2",
302             color="label",
303             title="Latent Space of Heartbeat FFT",
304             width=800,
305             size=10,
306             height=300,
307             tools=["tap"],
308         )
309
310         stream = hv.streams.Selection1D(source=latent)
311
312         reg = hv.DynamicMap(self.update, kdims=[], streams=[stream])
313
314         audio = pn.widgets.Audio(
315             name="Audio",
316             value=pipe(self.files.value, lambda f: os.path.join(
data_path, f)),
317         )
318
319         return pn.Column(latent, reg, audio)
320
321
```

```
322 # %% {"slideshow": {"slide_type": "skip"}}
323 d = Dashboard()
324
325 # %% {"slideshow": {"slide_type": "-"}}
326 pn.Column(d.files, d.view)
327
328 # %% {"slideshow": {"slide_type": "skip"}}
```

### Library Code

```
1 # wavfile.py (Enhanced)
2 #
3 # Mod by X-Raym
4 # Date: 20181906_1222
5 # * corrected loops
6 # * unsupported chunk read and write
7 # * LIST-INFO support
8 # * renamed variables to avoid conflict with python native functions
9 # * correct bytes error
10 # * correct write function
11 #
12 # URL: https://gist.github.com/josephernest/3f22c5ed5dabf1815f16efa8fa53d476
13 # Source: scipy/io/wavfile.py
14 #
15 # Mod by Joseph Basquin
16 # Date: 20180430_2335
17 # * read: also returns bitrate, cue markers + cue marker labels (sorted
18 #   ), loops, pitch
19 # * read: 24 bit & 32 bit IEEE files support (inspired from
20 #   wavio_weckesser.py from Warren Weckesser)
21 # * read: added normalized (default False) that returns everything as
22 #   float in [-1, 1]
23 # * read: added forcastereo that returns a 2-dimensional array even if
24 #   input is mono
25 #
26 # * write: can write cue markers, cue marker labels, loops, pitch
27 # * write: 24 bit support
28 # * write: can write from a float normalized in [-1, 1]
29 # * write: 20180430_2335: bug fixed when size of data chunk is odd (
30 #   previously, metadata could become unreadable because of this)
31 #
32 # * removed RIFX support (big-endian) (never seen one in 10+ years of
33 #   audio production/audio programming), only RIFF (little-endian) are
34 #   supported
35 # * removed read(..., mmap)
36 #
37 #
38 # Test:
```

```
32 # ..\wav\_____wavfile_demo.py
33
34
35 """
36 Module to read / write wav files using numpy arrays
37
38 Functions
39 _____
40 'read': Return the sample rate (in samples/sec) and data from a WAV
41       file.
42
43 'write': Write a numpy array as a WAV file.
44
45 """
46 # from __future__ import division, print_function, absolute_import
47
48 import numpy
49 import struct
50 import warnings
51 import collections
52
53 # from operator import itemgetter
54
55 class WavFileWarning(UserWarning):
56     pass
57
58
59 _ieee = False
60
61 # assumes file pointer is immediately
62 # after the 'fmt' id
63 def _read_fmt_chunk(fid):
64     res = struct.unpack("<ihIIHH", fid.read(20))
65     size, comp, noc, rate, sbytes, ba, bits = res
66     if comp != 1 or size > 16:
67         if comp == 3:
68             global _ieee
69             _ieee = True
70             # warnings.warn("IEEE format not supported", WavFileWarning)
71         else:
72             warnings.warn("Unfamiliar format bytes", WavFileWarning)
73     if size > 16:
74         fid.read(size - 16)
75     return size, comp, noc, rate, sbytes, ba, bits
76
77
78 # assumes file pointer is immediately
```

```

79 # after the 'data' id
80 def _read_data_chunk(fid, noc, bits, normalized=False):
81     size = struct.unpack("<i", fid.read(4))[0]
82
83     if bits == 8 or bits == 24:
84         dtype = "u1"
85         bytes_val = 1
86     else:
87         bytes_val = bits // 8
88         dtype = "<i%d" % bytes_val
89
90     if bits == 32 and _ieee:
91         dtype = "float32"
92
93     data = numpy.fromfile(fid, dtype=dtype, count=size // bytes_val)
94
95     if bits == 24:
96         a = numpy.empty((len(data) // 3, 4), dtype="u1")
97         a[:, :3] = data.reshape((-1, 3))
98         a[:, 3:] = (a[:, 3 - 1 : 3] >> 7) * 255
99         data = a.view("<i4").reshape(a.shape[: -1])
100
101     if noc > 1:
102         data = data.reshape(-1, noc)
103
104     if bool(
105         size & 1
106     ): # if odd number of bytes, move 1 byte further (data chunk is
word-aligned)
107         fid.seek(1, 1)
108
109     if normalized:
110         if bits == 8 or bits == 16 or bits == 24:
111             normfactor = 2 ** (bits - 1)
112             data = numpy.float32(data) * 1.0 / normfactor
113
114     return data
115
116
117 def _skip_unknown_chunk(fid):
118     data = fid.read(4)
119     size = struct.unpack("<i", data)[0]
120     if bool(
121         size & 1
122     ): # if odd number of bytes, move 1 byte further (data chunk is
word-aligned)
123         size += 1
124     fid.seek(size, 1)
125

```

```
126
127 def _read_unknown_chunk(fid , name):
128     data = fid.read(4)
129     size = struct.unpack("<i", data)[0]
130     # string = fid.read(size).rstrip(bytes('\x00', 'UTF-8')).decode("
utf-8")
131     offset = 0
132     if bool(
133         size & 1
134     ): # if odd number of bytes, move 1 byte further (data chunk is
word-aligned)
135         offset = 1
136     string = fid.read(size)
137     fid.seek(offset , 1)
138     return string
139
140
141 def _read_riff_chunk(fid):
142     str1 = fid.read(4)
143     if str1 != b"RIFF":
144         raise ValueError("Not a WAV file.")
145     fsize = struct.unpack("<I", fid.read(4))[0] + 8
146     str2 = fid.read(4)
147     if str2 != b"WAVE":
148         raise ValueError("Not a WAV file.")
149     return fsize
150
151
152 def read_wav(
153     file ,
154     readmarkers=False ,
155     readmarkerlabels=False ,
156     readmarkerslist=False ,
157     readloops=False ,
158     readpitch=False ,
159     normalized=False ,
160     forcestereo=False ,
161     log=True ,
162     readlistinfo=True ,
163     readunsupported=True ,
164 ):
165     """
166     Return the sample rate (in samples/sec) and data from a WAV file
167
168     Parameters
169     -----
170     file : file
171         Input wav file.
172
```



```
173 Returns
174 -----
175 rate : int
176     Sample rate of wav file
177 data : numpy array
178     Data read from wav file
179
180 Notes
181 -----
182
183 * The file can be an open file or a filename.
184
185 * The returned sample rate is a Python integer
186 * The data is returned as a numpy array with a
187   data-type determined from the file.
188
189 """
190 if hasattr(file, "read"):
191     fid = file
192 else:
193     fid = open(file, "rb")
194
195 fsize = _read_riff_chunk(fid)
196 noc = 1
197 bits = 8
198 # _cue = []
199 # _cuelabels = []
200 _markersdict = collections.defaultdict(lambda: {"position": -1, "
label": ""})
201 unsupported = {}
202 loops = []
203 list_info_index = [
204     "IARL",
205     "IART",
206     "ICMS",
207     "ICMT",
208     "ICOP",
209     "ICRD",
210     "IENG",
211     "IGNR",
212     "IKEY",
213     "IMED",
214     "INAM",
215     "IPRD",
216     "ISBJ",
217     "ISFT",
218     "ISRC",
219     "ISRF",
220     "ITCH",
```

```

221 ]
222 info = {}
223 pitch = 0.0
224 while fid.tell() < fsize:
225     # read the next chunk
226     chunk_id = fid.read(4)
227     chunk_id_str = chunk_id.decode("utf-8")
228     if chunk_id == b"fmt ":
229         size, comp, noc, rate, sbytes, ba, bits = _read_fmt_chunk(
fid)
230     elif chunk_id == b"data":
231         data = _read_data_chunk(fid, noc, bits, normalized)
232     elif chunk_id == b"cue ":
233         str1 = fid.read(8)
234         size, numcue = struct.unpack("<i", str1)
235         for c in range(numcue):
236             str1 = fid.read(24)
237             idx, position, datachunkid, chunkstart, blockstart,
sampleoffset = struct.unpack(
238                 "<iiiiii", str1
239             )
240             # _cue.append(position)
241             _markersdict[idx][
242                 "position"
243             ] = position # needed to match labels and markers
244
245     elif chunk_id == b"LIST":
246         str1 = fid.read(8)
247         size, datatype = struct.unpack("<i", str1)
248     elif (
249         chunk_id_str in list_info_index
250 ): # see http://www.pjb.com.au/midi/sfspec21.html#i5
251         s = _read_unknown_chunk(fid, chunk_id_str)
252         info[chunk_id_str] = s.decode("UTF-8")
253     elif chunk_id == b"label":
254         str1 = fid.read(8)
255         size, idx = struct.unpack("<i", str1)
256         size = size + (
257             size % 2
258         ) # the size should be even, see WAV specification, e.g.
16=>16, 23=>24
259         label = fid.read(size - 4).rstrip(
260             bytes("\x00", "UTF-8")
261         ) # remove the trailing null characters
262         # _cuelabels.append(label)
263         _markersdict[idx]["label"] = label # needed to match
labels and markers
264
265     elif chunk_id == b"smp1":

```

```

266         str1 = fid.read(40)
267         size, manuf, prod, sampleperiod, midiunitynote,
midipitchfraction, smpteformat, smpteoffs, numsampleloops, samplerdata
= struct.unpack(
268             "<iiiiIiiii", str1
269         )
270         cents = midipitchfraction * 1.0 / (2 ** 32 - 1)
271         pitch = 440.0 * 2 ** ((midiunitynote + cents - 69.0) / 12)
272         for i in range(numsampleloops):
273             str1 = fid.read(24)
274             cuepointid, datatype, start, end, fraction, playcount =
struct.unpack(
275                 "<iiiiii", str1
276             )
277             loops.append(
278                 {
279                     "cuepointid": cuepointid,
280                     "datatype": datatype,
281                     "start": start,
282                     "end": end,
283                     "fraction": fraction,
284                     "playcount": playcount,
285                 }
286             )
287         else:
288             if log:
289                 warnings.warn("Chunk " + str(chunk_id) + " skipped",
WavFileWarning)
290             if readunsupported:
291                 # print( chunk_id.decode("utf-8") + " unsupported")
292                 unsupported[chunk_id] = _read_unknown_chunk(fid,
chunk_id_str)
293             else:
294                 _skip_unknown_chunk(fid)
295         fid.close()
296
297         if data.ndim == 1 and forcestereo:
298             data = numpy.column_stack((data, data))
299
300         _markerslist = sorted(
301             [_markersdict[l] for l in _markersdict], key=lambda k: k["
position"])
302         ) # sort by position
303         _cue = [m["position"] for m in _markerslist]
304         _cuelabels = [m["label"] for m in _markerslist]
305
306         return (
307             (rate, data, bits)
308             + ((_cue,) if readmarkers else ()))

```

```

309         + ((_cuelabels,) if readmarkerlabels else ())
310         + ((_markerslist,) if readmarkerslist else ())
311         + ((_loops,) if readloops else ())
312         + ((_pitch,) if readpitch else ())
313         + ((_info,) if readlistinfo else ())
314         + ((_unsupported,) if readunsupported else ())
315     )
316
317
318 def write_wav(
319     filename,
320     rate,
321     data,
322     bitrate=None,
323     markers=None,
324     loops=None,
325     pitch=None,
326     normalized=False,
327     infos=None,
328     unsupported=None,
329 ):
330     """
331     Write a numpy array as a WAV file
332
333     Parameters
334     -----
335     filename : file
336         The name of the file to write (will be over-written).
337     rate : int
338         The sample rate (in samples/sec).
339     data : ndarray
340         A 1-D or 2-D numpy array of integer data-type.
341
342     Notes
343     -----
344     * Writes a simple uncompressed WAV file.
345     * The bits-per-sample will be determined by the data-type.
346     * To write multiple-channels, use a 2-D array of shape
347       (Nsamples, Nchannels).
348
349     """
350
351     # normalization and 24-bit handling
352     if (
353         bitrate == 24
354     ): # special handling of 24 bit wav, because there is no numpy.
355         int24...
356         if normalized:
357             data[data > 1.0] = 1.0

```

```

357         data[data < -1.0] = -1.0
358         a32 = numpy.asarray(data * (2 ** 23 - 1), dtype=numpy.int32
359     )
360     else:
361         a32 = numpy.asarray(data, dtype=numpy.int32)
362         if a32.ndim == 1:
363             a32.shape = a32.shape + (1,) # Convert to a 2D array with
364             a single column.
365             a8 = (
366                 a32.reshape(a32.shape + (1,)) >> numpy.array([0, 8, 16])
367             ) & 255 # By shifting first 0 bits, then 8, then 16, the
368             resulting output is 24 bit little-endian.
369             data = a8.astype(numpy.uint8)
370         else:
371             if normalized: # default to 32 bit int
372                 data[data > 1.0] = 1.0
373                 data[data < -1.0] = -1.0
374                 data = numpy.asarray(data * (2 ** 31 - 1), dtype=numpy.
375 int32)
376
377     fid = open(filename, "wb")
378     fid.write(b"RIFF")
379     fid.write(b"\x00\x00\x00\x00")
380     fid.write(b"WAVE")
381
382     # fmt chunk
383     fid.write(b"fmt ")
384     if data.ndim == 1:
385         noc = 1
386     else:
387         noc = data.shape[1]
388     bits = data.dtype.itemsize * 8 if bitrate != 24 else 24
389     sbytes = rate * (bits // 8) * noc
390     ba = noc * (bits // 8)
391     fid.write(struct.pack("<iHHHH", 16, 1, noc, rate, sbytes, ba,
392 bits))
393
394     if unsupported:
395         for key, val in unsupported.items():
396             if len(key) % 2 == 1:
397                 key += b"\x00"
398             if len(val) % 2 == 1:
399                 val += b"\x00"
400             info = key
401             size = len(val) # because \x00
402             info = struct.pack("<i", size)
403             info += val
404             fid.write(key)
405             size = len(info)

```

```

401         fid.write(info)
402
403     # cue chunk
404     if markers: # != None and != []
405         if isinstance(
406             markers[0], dict
407         ): # then we have [{'position': 100, 'label': 'marker1'}, ...]
408             labels = [m["label"] for m in markers]
409             markers = [m["position"] for m in markers]
410         else:
411             labels = [""] for m in markers]
412
413     fid.write(b"cue ")
414     size = 4 + len(markers) * 24
415     fid.write(struct.pack("<i", size, len(markers)))
416     for i, c in enumerate(markers):
417         s = struct.pack(
418             "<iiiiii", i + 1, c, 1635017060, 0, 0, c
419         ) # 1635017060 is struct.unpack('<i',b'data')
420         fid.write(s)
421
422     lbls = b""
423     for i, lbl in enumerate(labels):
424         lbls += b"label"
425         label = lbl + (b"\x00" if len(lbl) % 2 == 1 else b"\x00\x00
426 ")
427         size = len(lbl) + 1 + 4 # because \x00
428         lbls += struct.pack("<i", size, i + 1)
429         lbls += label
430
431     fid.write(b"LIST")
432     size = len(lbls) + 4
433     fid.write(struct.pack("<i", size))
434     fid.write(
435         b"adt1"
436     ) # https://web.archive.org/web/20141226210234/http://www.
437     sonicspot.com/guide/wavefiles.html#list
438     fid.write(lbls)
439
440     # smpl chunk
441     if loops or pitch:
442         if not loops:
443             loops = []
444         if pitch:
445             midiunitynote = 12 * numpy.log2(pitch * 1.0 / 440.0) + 69
446             midipitchfraction = int(
447                 (midiunitynote - int(midiunitynote)) * (2 ** 32 - 1)
448             )
449             midiunitynote = int(midiunitynote)

```

```

448         # print(midipitchfraction, midiunitynote)
449     else:
450         midiunitynote = 0
451         midipitchfraction = 0
452     fid.write(b"smp1")
453     size = 36 + len(loops) * 24
454     sampleperiod = int(1000000000.0 / rate)
455
456     fid.write(
457         struct.pack(
458             "<iiiiIiiii",
459             size,
460             0,
461             0,
462             sampleperiod,
463             midiunitynote,
464             midipitchfraction,
465             0,
466             0,
467             len(loops),
468             0,
469         )
470     )
471     for i, loop in enumerate(loops):
472         fid.write(
473             struct.pack(
474                 "<iiiiii",
475                 loop["cuepointid"],
476                 loop["datatype"],
477                 loop["start"],
478                 loop["end"],
479                 loop["fraction"],
480                 loop["playcount"],
481             )
482         )
483
484     # data chunks
485     fid.write(b"data")
486     fid.write(struct.pack("<i", data.nbytes))
487     import sys
488
489     if data.dtype.byteorder == ">" or (
490         data.dtype.byteorder == "=" and sys.byteorder == "big"
491     ):
492         data = data.byteswap()
493
494     data.tofile(fid)
495
496     if (

```

```
497     data.nbytes % 2 == 1
498 ): # add an extra padding byte if data.nbytes is odd: https://web.
archive.org/web/20141226210234/http://www.sonicspot.com/guide/
wavefiles.html#data
499     fid.write("\x00")
500
501 # This need to be made modular !
502 if infos:
503     info = b""
504     for key, val in infos.items():
505         key = bytes(key, "UTF-8")
506         val = bytes(val, "UTF-8")
507         # val += b'\x00' # Note: Fix windows display error. Is this
valid ?
508         size = len(val) # because \x00
509         if len(val) % 2 == 1:
510             val += b"\x00"
511         info += key
512         info += struct.pack("<i", size)
513         info += val
514     # info += b'\x00'
515     if len(info) % 2 == 1:
516         info += b"\x00"
517     fid.write(b"LIST")
518     size = len(info) + 4
519     fid.write(struct.pack("<i", size))
520     fid.write(
521         b"INFO"
522     ) # https://web.archive.org/web/20141226210234/http://www.
sonicspot.com/guide/wavefiles.html#list
523     fid.write(info)
524
525 # Determine file size and place it in correct
526 # position at start of the file.
527 size = fid.tell()
528 fid.seek(4)
529 fid.write(struct.pack("<i", size - 8))
530 fid.close()
531 return "success"
```

```
1 import pandas as pd
2 import numpy as np
3 import holoviews as hv
4 import hvplot.pandas
5 from typing import Dict
6 from toolz.curried import *
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import silhouette_score
9 from sklearn.base import TransformerMixin
```



```
10
11
12 def get_latent_space(
13     model, X_train, X_test, y_train, y_test, tag, labels, feature_names
14 ):
15
16     latent = pipe(
17         X_test,
18         model.transform,
19         partial(pd.DataFrame, columns=["Component 1", "Component 2"]),
20     )
21
22     return (
23         pd.concat([latent], axis=1)
24         .assign(label=y_test)
25         .assign(label=lambda d: d.label.replace(labels))
26         .assign(tag=tag)
27     )
28
29
30 def reconstruction_benchmark(
31     dataset: Dict[str, np.ndarray], models: Dict[str, TransformerMixin
32 ]), label: str
33 ):
34     """
35     """
36     data = dataset
37
38     if hasattr(data, "target_names"):
39         labels = dict(enumerate(data.target_names))
40     else:
41         labels = pipe(data.target, np.unique, map(str), list, np.array)
42         .astype(str)
43
44     if hasattr(data, "feature_names"):
45         names = dict(enumerate(data.feature_names))
46     else:
47         names = pipe(range(data.data.shape[1]), map(str), list, np.
48 array).astype(str)
49
50     X_train, X_test, y_train, y_test = train_test_split(
51         data.data, data.target, test_size=0.33, random_state=42
52     )
53
54     for m in models.values():
55         m.fit(X_train)
```

```

56         get_latent_space(m, X_train, X_test, y_train, y_test, t,
labels, names)
57         for t, m in models.items()
58     ]
59 )
60
61 reconstruction_loss = pd.DataFrame(
62     {
63         t: pipe(
64             X_train,
65             m.transform,
66             m.inverse_transform,
67             lambda x: np.subtract(x, X_train),
68             lambda x: np.power(x, 2),
69             lambda x: np.array(x).flatten(),
70             np.mean,
71         )
72         for t, m in models.items()
73     },
74     index=["Reconstruction Loss"],
75 ).T
76
77 silhouette = pd.DataFrame(
78     {
79         t: [silhouette_score(X=pipe(X_test, m.transform), labels=
y_test)]
80         for t, m in models.items()
81     },
82     index=["Silhouette"],
83 ).T
84
85 return (
86     latent,
87     reconstruction_loss.hvplot.bar(title=f"{label}: Reconstruction
Loss"),
88     silhouette.hvplot.bar(title=f"{label}: Silhouette Scores"),
89     latent.hvplot.scatter(
90         x="Component 1", y="Component 2", color="label", groupby="
tag", label=label
91     )
92     .layout()
93     .cols(2),
94 )

```

```

1 from sklearn.base import BaseEstimator, TransformerMixin
2 import numpy as np
3 from toolz.curried import *
4 from keras.layers import Lambda, Input, Dense
5 from keras.models import Model

```

```
6 from keras.datasets import mnist
7 from keras.losses import mse, binary_crossentropy
8 from keras.utils import plot_model
9 from keras import backend as K
10 from keras.regularizers import l2
11
12
13 def sampling(args):
14     """Reparameterization trick by sampling from an isotropic unit
15     Gaussian.
16
17     # Arguments
18         args (tensor): mean and log of variance of Q(z|X)
19
20     # Returns
21         z (tensor): sampled latent vector
22     """
23     z_mean, z_log_var = args
24     batch = K.shape(z_mean)[0]
25     dim = K.int_shape(z_mean)[1]
26     # by default, random_normal has mean = 0 and std = 1.0
27     epsilon = K.random_normal(shape=(batch, dim))
28     return z_mean + K.exp(0.5 * z_log_var) * epsilon
29
30
31 class VAE(BaseEstimator, TransformerMixin):
32     """Transform data using vae"""
33
34     def __init__(
35         self,
36         hidden_layer_sizes: tuple = (25, 2),
37         activation: str = "relu",
38         solver: str = "adam",
39         divergence_weight: float = 1,
40         alpha: float = 0.0001,
41         batch_size: str = "auto",
42         learning_rate: str = "constant",
43         learning_rate_init: float = 0.001,
44         power_t: float = 0.5,
45         max_iter: int = 200,
46         shuffle: bool = True,
47         random_state=None,
48         tol: float = 0.0001,
49         verbose: bool = False,
50         warm_start: bool = False,
51         momentum: float = 0.9,
52         nesterovs_momentum: bool = True,
53         early_stopping: bool = False,
```

```

54     validation_fraction: float = 0.1,
55     beta_1: float = 0.9,
56     beta_2: float = 0.999,
57     epsilon: float = 1e-08,
58     n_iter_no_change=10,
59 ):
60     """
61     """
62     self.alpha = alpha
63     self.regularizer = l2(alpha)
64
65     self.hidden_layer_sizes = hidden_layer_sizes
66     self.activation = activation
67     self.max_iter = max_iter
68     if solver == "auto":
69         self.solver = "adam"
70     else:
71         self.solver = solver
72     self.divergence_weight = divergence_weight
73     self.model = None
74     self.validation_fraction = validation_fraction
75
76 def build_encoder_(self, layers: tuple):
77     # VAE model = encoder + decoder
78     # build encoder model
79     input_shape, *encoder_shape, latent_dim = layers
80
81     inputs = Input(shape=(input_shape,), name="encoder_input")
82     transformations = pipe(
83         encoder_shape,
84         map(
85             lambda d: Dense(
86                 units=d,
87                 kernel_regularizer=self.regularizer,
88                 activation=self.activation,
89             )
90         ),
91         lambda f: compose_left(*f),
92     )
93
94     x = pipe(inputs, transformations)
95
96     z_mean = Dense(latent_dim, name="z_mean")(x)
97     z_log_var = Dense(latent_dim, name="z_log_var")(x)
98
99     # use reparameterization trick to push the sampling out as
100     input
    z = Lambda(sampling, output_shape=(latent_dim,), name="z")([
        z_mean, z_log_var])

```

```

101
102     # note that "output_shape" isn't necessary with the TensorFlow
backend
103     # instantiate encoder model
104     encoder = Model(inputs, [z_mean, z_log_var, z], name="encoder")
105
106     return inputs, encoder, z_mean, z_log_var
107
108 def build_decoder_(self, layers: tuple):
109     # build decoder model
110     latent_shape, *decoder_layers, original_dim = layers
111
112     latent_inputs = Input(shape=(latent_shape,), name="z_sampling")
113
114     transformations = pipe(
115         decoder_layers,
116         map(
117             lambda d: Dense(
118                 units=d,
119                 kernel_regularizer=self.regularizer,
120                 activation=self.activation,
121             ),
122         ),
123         lambda f: compose_left(*f),
124     )
125     final_layer = Dense(original_dim, name="original_dim")
126
127     outputs = pipe(latent_inputs, transformations, final_layer)
128
129     # instantiate decoder model
130     decoder = Model(latent_inputs, outputs, name="decoder")
131     return decoder
132
133 def build_model_(self, layers: tuple):
134     """ """
135     inputs, self.encoder, z_mean, z_log_var = self.build_encoder_(
layers)
136     self.decoder = self.build_decoder_(reversed(layers))
137
138     outputs = pipe(inputs, self.encoder, get(2), self.decoder)
139     vae = Model(inputs, outputs, name="vae_mlp")
140
141     #         if args.mse:
142     reconstruction_loss = mse(inputs, outputs)
143     #         else:
144     #             reconstruction_loss = binary_crossentropy(inputs,
145     #                                                         outputs
146
)

```

```

147     reconstruction_loss *= layers[0]
148     kl_loss = 1 + z_log_var - K.square(z_mean) - K.exp(z_log_var)
149     kl_loss = K.sum(kl_loss, axis=-1)
150     kl_loss *= -0.5
151     vae_loss = K.mean(reconstruction_loss + self.divergence_weight
* kl_loss)
152
153     vae.add_loss(vae_loss)
154     vae.compile(optimizer=self.solver)
155
156     return vae
157
158 def fit(self, x: np.ndarray, y: np.ndarray = None):
159     """
160     """
161     layers = x.shape[1], *self.hidden_layer_sizes
162
163     #         if self.model is None:
164     self.model = self.build_model_(layers)
165
166     n_samples = x.shape[0]
167     if self.solver == self.solver:
168         self.batch_size = n_samples
169     elif self.batch_size == self.solver:
170         self.batch_size = min(200, n_samples)
171     else:
172         if self.batch_size < 1 or self.batch_size > n_samples:
173             warnings.warn(
174                 "Got 'batch_size' less than 1 or larger than "
175                 "sample size. It is going to be clipped"
176             )
177         self.batch_size = np.clip(self.batch_size, 1, n_samples)
178
179     self.model.fit(
180         x,
181         epochs=self.max_iter,
182         batch_size=self.batch_size,
183         validation_split=self.validation_fraction,
184     )
185
186     #         vae.save_weights('vae_mlp_mnist.h5')
187
188     return self
189
190 def transform(self, X: np.ndarray):
191     """
192     """
193     return pipe(self.encoder.predict(X), get(0))
194

```

```
195     def sample(self, n: int):
196         """
197         """
198         dim = self.hidden_layer_sizes[-1]
199         N = np.random.multivariate_normal(np.zeros(dim), np.diag(np.
ones(dim)), size=n)
200         return self.decoder.predict(N)
201
202     def inverse_transform(self, Xt: np.ndarray):
203         """
204         """
205         return self.decoder.predict(Xt)
206
207     def predict(self, X: np.ndarray):
208         """
209         """
210         return self.model.predict(X)
211
212     def score(self, X: np.ndarray, y: np.ndarray = None):
213         """
214         """
215         return self.model.evaluate(X)
```

### Environment

```
1 name: super-spirals
2 channels:
3   - bioconda
4   - conda-forge
5   - intel
6   - defaults
7 dependencies:
8   - _libgcc_mutex=0.1=main
9   - absl-py=0.7.1=py36_0
10  - appdirs=1.4.3=py_1
11  - asn1crypto=0.24.0=py36_3
12  - astor=0.8.0=py36_0
13  - attrs=19.1.0=py_0
14  - backcall=0.1.0=py36_2
15  - backports=1.0=py36_9
16  - black=19.3b0=py_0
17  - bleach=2.1.3=py36_2
18  - bokeh=1.3.4=py36_0
19  - c-ares=1.15.0=h7b6447c_1001
20  - ca-certificates=2019.9.11=hecc5488_0
21  - certifi=2019.9.11=py36_0
22  - cffi=1.11.5=py36_3
23  - chardet=3.0.4=py36_3
24  - click=7.0=py_0
25  - cryptography=2.3=py36_2
26  - cycler=0.10.0=py36_7
27  - daal=2019.5=intel_281
28  - daal4py=2019.5=py36ha68da19_2
29  - decorator=4.3.0=py36_3
30  - entrypoints=0.2.3=py36_2
31  - fastdtw=0.2.0=py_1
32  - fontconfig=2.13.1=h86ecdb6_1001
33  - freetype=2.9.1=h8a8886c_1
34  - gast=0.2.2=py36_0
35  - get_terminal_size=1.0.0=py36_7
36  - glob2=0.7=py_0
37  - google-pasta=0.1.7=py_0
38  - grpcio=1.23.0=py36he9aelf9_0
39  - h5py=2.8.0=py36h989c5e5_3
40  - hdf5=1.10.2=2
41  - holoviews=1.12.3=py_2
42  - html5lib=1.0.1=py36_4
43  - hvplot=0.4.0=py_1
44  - icc_rt=2019.5=intel_281
45  - icu=64.2=he1b5a44_1
46  - idna=2.6=py36_3
47  - impi_rt=2019.5=intel_281
```



```
48 - intel-openmp=2019.5=intel_281
49 - intelpython=2019.5=0
50 - ipykernel=4.6.1=py36_2
51 - ipython=6.3.1=py36_3
52 - ipython_genutils=0.2.0=py36_2
53 - jedi=0.12.0=py36_2
54 - jinja2=2.10.1=py_0
55 - joblib=0.13.2=py36_1
56 - jpeg=9b=h024ee3a_2
57 - json5=0.8.5=py_0
58 - jsonschema=2.6.0=py36_2
59 - jupyter_client=5.1.0=py36_5
60 - jupyter_core=4.4.0=py36_6
61 - jupyterlab=1.1.4=py_0
62 - jupyterlab_server=1.0.0=py_0
63 - jupyter_text=1.2.4=0
64 - kaggle=1.5.6=py36_0
65 - keras=2.2.4=0
66 - keras-applications=1.0.8=py_0
67 - keras-base=2.2.4=py36_0
68 - keras-preprocessing=1.1.0=py_1
69 - kiwisolver=1.0.1=py36_2
70 - libffi=3.2.1=11
71 - libgcc-ng=9.1.0=hdf63c60_0
72 - libiconv=1.15=h516909a_1005
73 - libpng=1.6.36=2
74 - libprotobuf=3.8.0=hd408876_0
75 - libsodium=1.0.16=3
76 - libstdcxx-ng=9.1.0=hdf63c60_0
77 - libtiff=4.0.10=h2733197_2
78 - libuuid=2.32.1=h14c3975_1000
79 - libxml2=2.9.9=hee79883_5
80 - markdown=3.1.1=py36_0
81 - markupsafe=1.0=py36_3
82 - matplotlib=3.1.1=py36_2
83 - mistune=0.8.3=py36_2
84 - mkl=2019.5=intel_281
85 - mkl-service=2.3.0=py36_0
86 - mkl_fft=1.0.14=py36ha68da19_1
87 - mkl_random=1.0.4=py36ha68da19_2
88 - nbconvert=5.2.1=py36_2
89 - nbformat=4.4.0=py36_2
90 - ncurses=6.1=he6710b0_1
91 - nodejs=10.13.0=he6710b0_0
92 - notebook=5.2.2=py36_1
93 - numexpr=2.6.9=py36_0
94 - numpy=1.17.0=py36ha68da19_13
95 - numpy-base=1.17.0=py36_13
96 - olefile=0.46=py36_0
```

```
97 - openssl=1.1.1c=h516909a_0
98 - packaging=19.1=py36_0
99 - pandas=0.25.0=py36_5
100 - pandocfilters=1.4.1=py36_2
101 - panel=0.6.2=h39e3cac_0
102 - param=1.9.1=py_0
103 - parso=0.2.0=py36_2
104 - path.py=11.0.1=py36_2
105 - pexpect=4.2.1=py36_4
106 - phantomjs=2.1.1=1
107 - pickleshare=0.7.4=py36_3
108 - pillow=6.1.0=py36h34e0f95_0
109 - pip=19.1.1=py36_0
110 - prompt_toolkit=1.0.15=py36_2
111 - protobuf=3.8.0=py36he6710b0_0
112 - ptyprocess=0.5.2=py36_2
113 - pycparser=2.18=py36_2
114 - pyct=0.4.6=py36_0
115 - pygments=2.2.0=py36_5
116 - pyopenssl=17.5.0=py36_2
117 - pyparsing=2.2.0=py36_2
118 - pysocks=1.6.7=py36_1
119 - python=3.6.8=h0371630_0
120 - python-dateutil=2.8.0=py36_0
121 - python-slugify=3.0.3=py_0
122 - pytz=2019.1=py36_0
123 - pyviz_comms=0.7.2=py_0
124 - pyyaml=5.1.1=py36_0
125 - pyzmq=16.0.2=py36_6
126 - readline=7.0=h7b6447c_5
127 - requests=2.20.1=py36_1
128 - scikit-learn=0.21.3=py36ha68da19_4
129 - scipy=1.3.1=py36ha68da19_2
130 - selenium=3.141.0=py36h7b6447c_0
131 - setuptools=41.0.1=py36_0
132 - simplegeneric=0.8.1=py36_7
133 - six=1.12.0=py36_0
134 - sqlite=3.28.0=0
135 - tbb=2019.8=intel_281
136 - tbb4py=2019.8=py36_intel_0
137 - tcl=8.6.4=24
138 - tensorboard=1.14.0=py36hf484d3e_0
139 - tensorflow=1.14.0=py36_0
140 - tensorflow-base=1.14.0=0
141 - tensorflow-estimator=1.14.0=py_0
142 - termcolor=1.1.0=py36_1
143 - terminado=0.8.1=py36_2
144 - testpath=0.3.1=py36_2
145 - text-unidecode=1.2=py_0
```

```
146 - tk=8.6.8=hbc83047_0
147 - toml=0.10.0=py_0
148 - toolz=0.10.0=py_0
149 - tornado=4.5.2=py36_5
150 - tqdm=4.36.1=py_0
151 - traitlets=4.3.2=py36_3
152 - unicode=1.1.1=py_0
153 - urllib3=1.24.1=py36_2
154 - wcwidth=0.1.7=py36_6
155 - webencodings=0.5.1=py36_0
156 - werkzeug=0.14.1=py36_0
157 - wheel=0.31.0=py36_3
158 - wrapt=1.11.2=py36h7b6447c_0
159 - xz=5.2.4=5
160 - yaml=0.1.7=2
161 - zeromq=4.2.3=2
162 - zip=3.0=0
163 - zlib=1.2.11=5
164 - zstd=1.3.7=h0b5b093_0
165 - pip:
166   - intel-tensorflow==1.14.0
167 prefix: /home/marcusskky/.conda/envs/super-spirals
```

## References

- Aeberhard, Stefan, Danny Coomans, and Olivier De Vel (1994), “Comparative analysis of statistical pattern recognition methods in high dimensional settings.” *Pattern Recognition*, 27, 1065–1077.
- An, Jinwon and Sungzoon Cho (2015), “SNU Data Mining Center 2015-2 Special Lecture on IE Variational Autoencoder based Anomaly Detection using Reconstruction Probability.”
- Bengio, Yoshua, Aaron Courville, and Pascal Vincent (2013), “Representation learning: A review and new perspectives.” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35, 1798–1828.
- Bengio, Yoshua, Pascal Lamblin, Dan Popovici, and Hugo Larochelle (2007), “Greedy Layer-Wise Training of Deep Networks Yoshua.” *In Advances in neural information processing system*.
- Bentley, P., G. Nordehn, M. Coimbra, and S. Mannor (????), “The PASCAL Classifying Heart Sounds Challenge 2011 (CHSC2011) Results.” <http://www.peterjbentley.com/heartchallenge/index.html>.

- Fisher, R. A. (1936), “The use of multiple measurements in taxonomic problems.” *Annals of Eugenics*, 7, 179–188.
- Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio (2014), “Generative Adversarial Nets.” iii, URL <https://papers.nips.cc/paper/5423-generative-adversarial-nets>{\%}0Ahttp://doi.wiley.com/10.1002/9781118472507.fmatter{\%}0Ahttp://linkinghub.elsevier.com/retrieve/pii/B9780408001090500018.
- Kingma, Diederik P and Max Welling (2013), “Auto-Encoding Variational Bayes.” *arXiv e-prints*, arXiv:1312.6114.
- Street, W. N., W. H. Wolberg, and O. L. Mangasarian (1993). *Biomedical Image Processing and Biomedical Visualization*.
- Tolstikhin, Ilya, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf (2017), “Wasserstein auto-encoders.” *arXiv preprint arXiv:1711.01558*.
- Zhao, Shengjia, Jiaming Song, and Stefano Ermon (2017), “Infovae: Information maximizing variational autoencoders.” *arXiv preprint arXiv:1706.02262*.