

Classification of data from Charge Coupled Devices

Asbjørn Bonefeld Preuss

Daniel Lomholt Christensen

Jens Kinch

Sejr Sebastian Haahr Bergmann

KØBENHAVNS UNIVERSITET



Outline

- 1 Data
- 2 Approaches
- 3 Supervised classification
- 4 Self-supervised classification
- 5 Conclusion and evaluation
- 6 Appendix

Data

Approaches

Supervised
classification

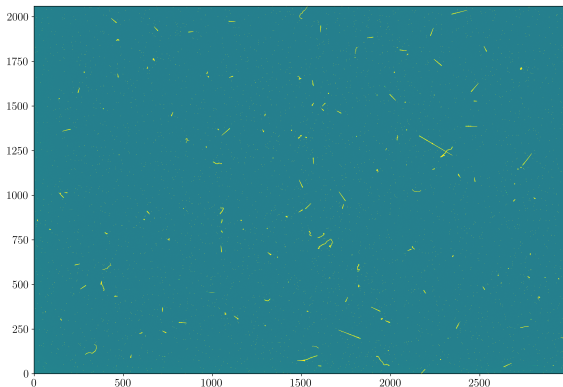
Self-supervised
classification

Conclusion
and evaluation

Appendix

Data

Data are energy pictures, in 2D showing traces of particles from space:



Data

Approaches

Supervised
classification

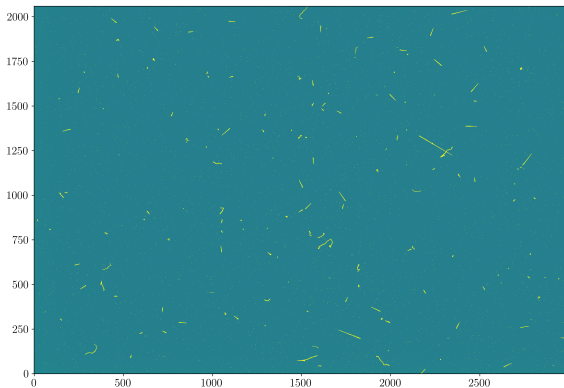
Self-supervised
classification

Conclusion
and evaluation

Appendix

Data

Data are energy pictures, in 2D showing traces of particles from space: Project Purpose: To classify the particles faster than humans can, and possibly classify more difficult particles than we as humans can.



Data

Approaches

Supervised
classification

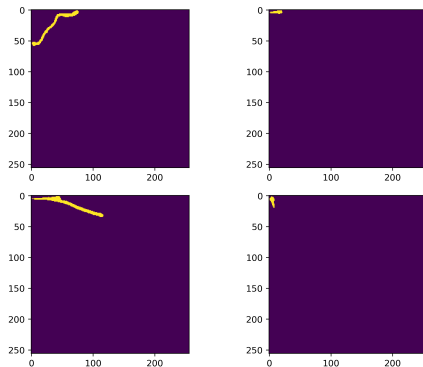
Self-supervised
classification

Conclusion
and evaluation

Appendix

Data

We can perform binary clustering and padding, to get all clusters as binary 256 x 256 images:



Data

Approaches

Supervised
classification

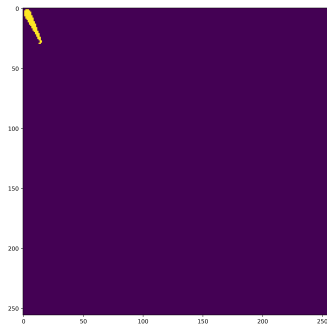
Self-supervised
classification

Conclusion
and evaluation

Appendix

Thoughts on data processing:

- Images are quite large (256 X 256 pixels) where most images are nothing



Data

Approaches

Supervised
classification

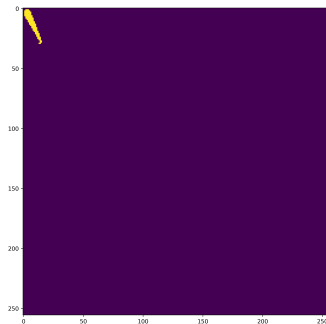
Self-supervised
classification

Conclusion
and evaluation

Appendix

Thoughts on data processing:

- Images are quite large (256 X 256 pixels) where most images are nothing
- Images are binary



Data

Approaches

Supervised
classification

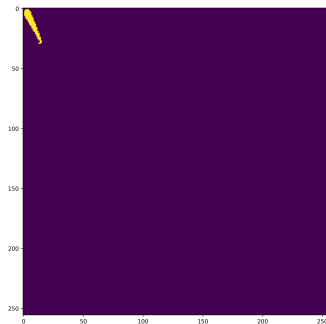
Self-supervised
classification

Conclusion
and evaluation

Appendix

Thoughts on data processing:

- Images are quite large (256 X 256 pixels) where most images are nothing
- Images are binary
- The labels we can give, are (perhaps) quite poor



Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Thoughts on data processing:

- Images are quite large (256 X 256 pixels) where most images are nothing
- Images are binary
- The labels we can give, are (perhaps) quite poor

Solutions:



Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

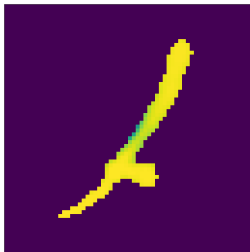
Appendix

Thoughts on data processing:

- Images are quite large (256 X 256 pixels) where most images are nothing
- Images are binary
- The labels we can give, are (perhaps) quite poor

Solutions:

- Downscale the large images to 64 X 64 pixels^a, pad the small images up to 64 X 64
- Re-insert energy
- Use self-supervised learning?



^aSee appendix for explanation of downscaling

Approaches

Supervised approach

- Use CNN to classify the data

Self-Supervised approach

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Approaches

Supervised approach

- Use CNN to classify the data
- This requires labels. These are made manually by us (≈ 500)

Self-Supervised approach

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Approaches

Supervised approach

- Use CNN to classify the data
- This requires labels. These are made manually by us (≈ 500)
- Will have a hard time classifying better than us.

Self-Supervised approach

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Approaches

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Supervised approach

- Use CNN to classify the data
- This requires labels. These are made manually by us (≈ 500)
- Will have a hard time classifying better than us.

Self-Supervised approach

- Use variational auto-encoder and clustering on latent space

Approaches

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Supervised approach

- Use CNN to classify the data
- This requires labels. These are made manually by us (≈ 500)
- Will have a hard time classifying better than us.

Self-Supervised approach

- Use variational auto-encoder and clustering on latent space
- Could show us more than we can predict

Approaches

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Supervised approach

- Use CNN to classify the data
- This requires labels. These are made manually by us (≈ 500)
- Will have a hard time classifying better than us.

Self-Supervised approach

- Use variational auto-encoder and clustering on latent space
- Could show us more than we can predict
- Does not need labels

Approaches

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Supervised approach

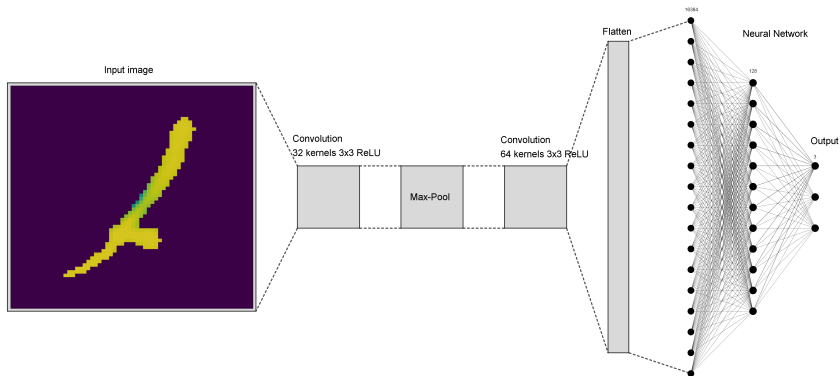
- Use CNN to classify the data
- This requires labels. These are made manually by us (≈ 500)
- Will have a hard time classifying better than us.

Self-Supervised approach

- Use variational auto-encoder and clustering on latent space
- Could show us more than we can predict
- Does not need labels
- Experimental \Leftrightarrow Might not work

Supervised classification

Use Convolutional Neural Network, with a categorical cross entropy loss function.
Quick run down of CNNs:



Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Our CNN

- Built in Tensorflow
- Hyper parameters optimized with optuna (Bayesian optimization)



TensorFlow

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Results from Supervised Classification

Data

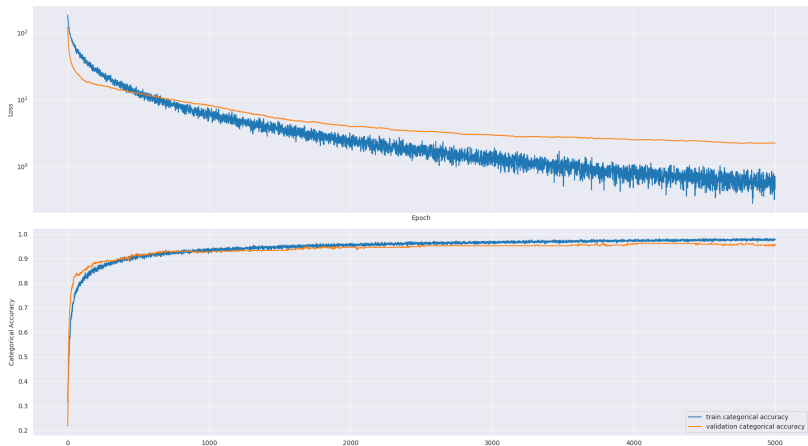
Approaches

Supervised
classification

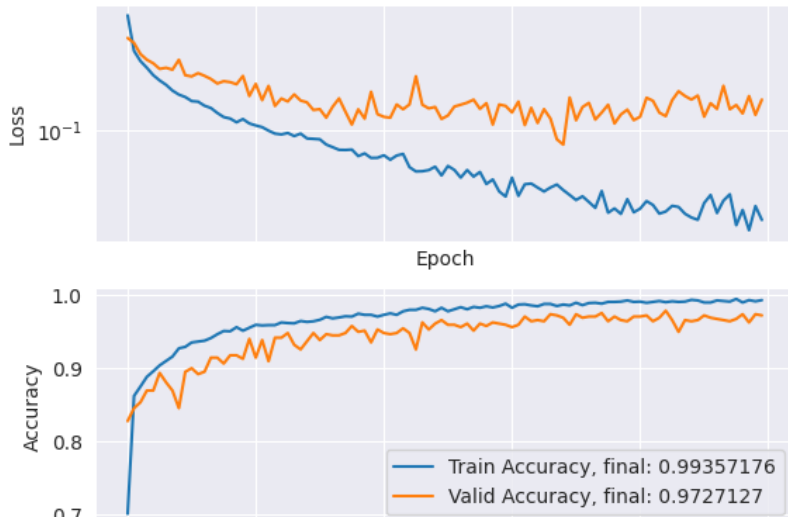
Self-supervised
classification

Conclusion
and evaluation

Appendix



Results from Supervised Classification



Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Results from Supervised Classification

Data

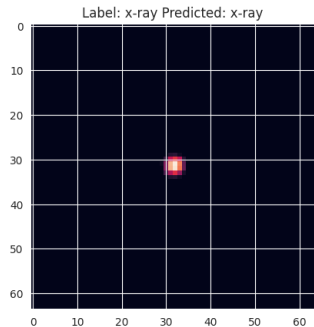
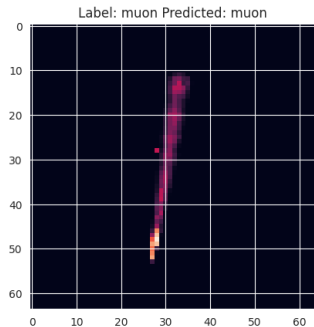
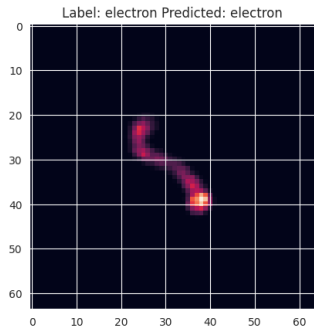
Approaches

**Supervised
classification**

Self-supervised
classification

Conclusion
and evaluation

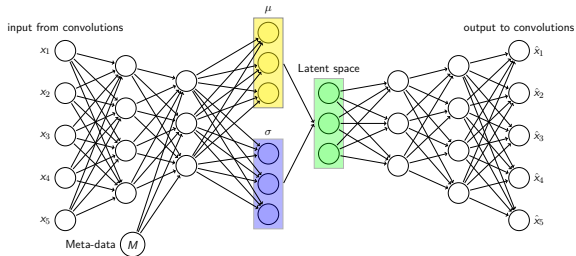
Appendix



Self-supervised classification

Use variational auto-encoder via a convolutional neural network and finally umap cluster on the latent space:

- Parameters for convolutions initially based on promising results from the CNN
- Meta data (Width, Height, Number of pixels, sum of those pixels) fed into the second layer of the neural network



Variational layers

Pseudo-code:

- Encode data to latent space
- Pass latent space to linear layer that returns array of means
- pass latent space to linear layer that array of variances
- return an array of means with Gaussian noise based on variances

Code (PyTorch):

```
### Variational part
self.variational_mean = nn.Linear(encoded_space_dim, encoded_space_dim)
self.variational_var = nn.Linear(encoded_space_dim, encoded_space_dim)

def reparameterization(self, mean, var): # Stolen from MNIST example, device
    epsilon = torch.randn_like(var).cpu()
    z = mean + var * epsilon
    return z

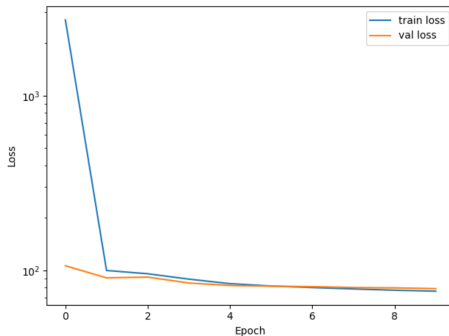
def forward(self, x):
    #CNN and removal of meta data:
    meta = x[:4]
    x = x[:4].unflatten(0, (1, 1, 64, 64))
    x, indices = self.encoder_cnn(x) # Capture indices from MaxPool2d
    x = self.flatten(x)
    x = self.encoder_lin_pic(x)

    #Add meta data and encode:
    y = torch.cat((x, meta), dim=0)
    x = self.encoder_lin_w_meta(y)
    mean = self.variational_mean(x)
    log_var = self.variational_var(x)
    z = self.reparameterization(mean, torch.exp(0.5 * log_var))
    return z, mean, log_var, indices
```


Loss

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{BCE}} + \lambda \mathcal{L}_{\text{KL}}, \quad \mathcal{L}_{\text{KL}} = -\frac{1}{2} \sum_n (1 + \log \sigma_n^2 - \mu_n^2 - \sigma_n^2)$$

second term is the Kullback-Leibler (KL) divergence, that enforces a Gaussian latent space



Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Results from our variational auto-encoder

Data

Approaches

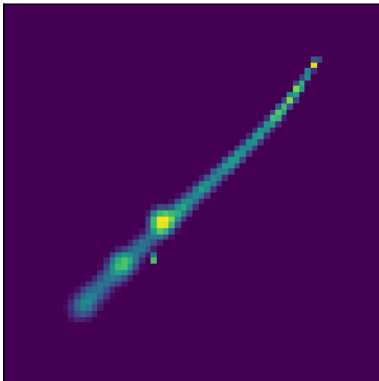
Supervised
classification

Self-supervised
classification

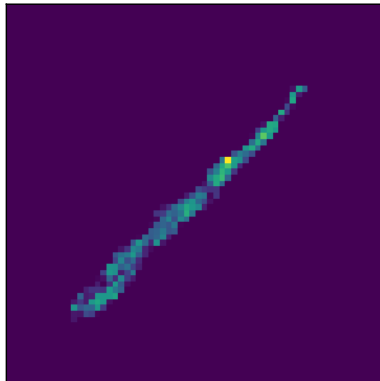
Conclusion
and evaluation

Appendix

Original image



Reconstructed image



Dimensionality reduction and clustering of the latent space

Data

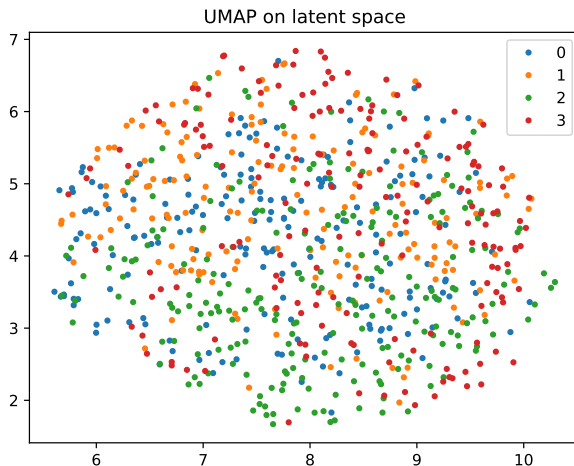
Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix



Dimensionality reduction and clustering of the latent space

Data

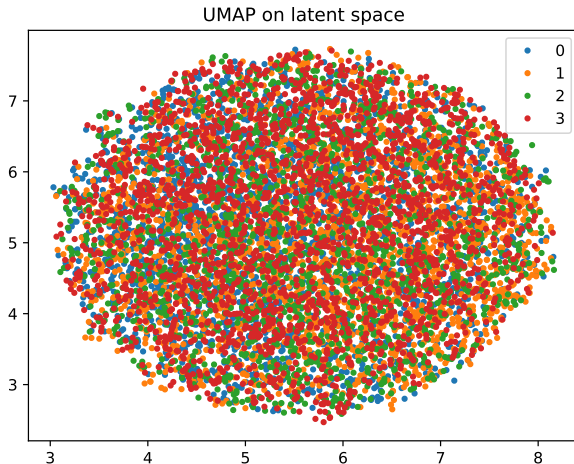
Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix



Result from single cluster

Data

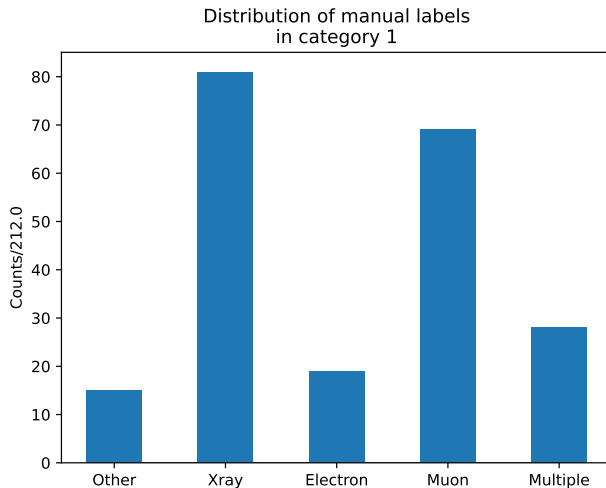
Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix



Conclusion

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

CNN accurately categorises the samples, but only as well as we could.

Conclusion

CNN accurately categorises the samples, but only as well as we could.

VAE via CNN can cluster the latent space, and finds clusters that are not that great.

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Evaluation of models

CNN:

Started working quite fast

Reached a high accuracy

Does not provide any new information/
can only predict what we saw

VAE via CNN:

Took a long time to get working (i.e many
issues with preprocessing and long training
time)

Clustered based on directions at best.

Might get better with more training.

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Appendix

Project statement: All authors contributed equally.

Github: <https://github.com/AsbjornPreuss/AppMLFinalProject>

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Hyperparameters CNN

Hyperparameter name	Parameter space	Final value
N_conv_layers	1-3	2
N_dense_layers	1-3	3
dropout_rate	0-0.5	0.4981
conv_filter_exp0	3-7	6
conv_kernel_size0	2-5	2
conv_filter_exp1	3-7	4
conv_kernel_size1	2-5	5
dense_units_exp0	3-5	4
dense_units_exp1	3-5	3
adam_learning_rate	1e-7-1e-1	0.0003676
adam_decay_steps	100-10000	3866
adam_decay_rate	0.8-1.0	0.9122

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Hyperparameters CNN

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

N_conv_layers	3
N_dense_layers	3
dropout_rate	0.427
filters_layer0	41
filter_size_layer0	6
stride_layer0	1
padding_layer0	3
filters_layer1	49
filter_size_layer1	3
stride_layer1	1
padding_layer1	2
filters_layer2	9
filter_size_layer2	7
stride_layer2	2
padding_layer2	2

Hyperparameters VAE via CNN

No hyperparameter optimization was reached. This is the hyperparameters we manually chose.

Encoder:

Hyperparameter name	Final value
Learning rate	0.001
Latent space dim	50
Conv_1	12X12
Conv_2	6X6
Dense_layer_1	43264X200
Dense_layer_2	204X50

Decoder:

Hyperparameter name	Final value
Learning rate	0.001
Latent space dim	50
Conv_2	12X12
Conv_1	6X6
Dense_layer_2	200X43264
Dense_layer_1	50X200

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Preprocessing algorithms

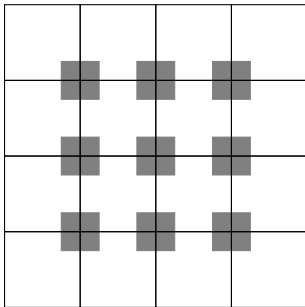
Depth first algorithm:

- ① Make data binary by threshold
- ② Scan for values, starting from upper left corner
- ③ When a value is found, search neighbours to see if they have a value (larger than background noise)
- ④ If the neighbour has a value, add them to the cluster
- ⑤ Go to each neighbour in the cluster, and scan their neighbours
- ⑥ Repeat until no more neighbours with values are found
- ⑦ Go back and keep on scanning the image, whilst ignoring the ones that are already assigned a cluster.
- ⑧ Save height, width, size (in pixels) and summed value of the pixels in the original image (an estimate for total energy) for meta-data

Downscaling

The downscaling of images was done by averaging four nearest neighbours, thus reducing the image sizes with 1 along both axes.

This is done recursively, allowing for the images to be scaled to the required size.



Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix

Some thoughts on Convolutions in PyTorch

Debugging within training is difficult in PyTorch. The functions within an encoder can be accessed any time, this is hopefully a helpful snippet for debugging:

```
# encoder:
A = nn.Conv2d(1, 32, 3, stride=1, padding=1)
B = nn.Conv2d(32, 64, 3, stride=1, padding=1)
C = nn.MaxPool2d(2)

# flatten
D = nn.Flatten()

# linear
E = nn.Linear(65536, 128)
F = nn.Linear(128, 10)

print(f'original data: {x.shape}')

print('encoder:')

print(f'Conv1: {A(x).shape}')
print(f'Conv2: {B(A(x)).shape}')
print(f'MaxPool: {C(B(A(x))).shape}')
print(f'Flatten: {D(C(B(A(x))))}.shape}')
print(f'Linear1: {E(D(C(B(A(x))))).shape}')
print(f'Linear2: {F(E(D(C(B(A(x))))).shape)')

```

Data

Approaches

Supervised
classification

Self-supervised
classification

Conclusion
and evaluation

Appendix