# Target-Guided AutoEncoder

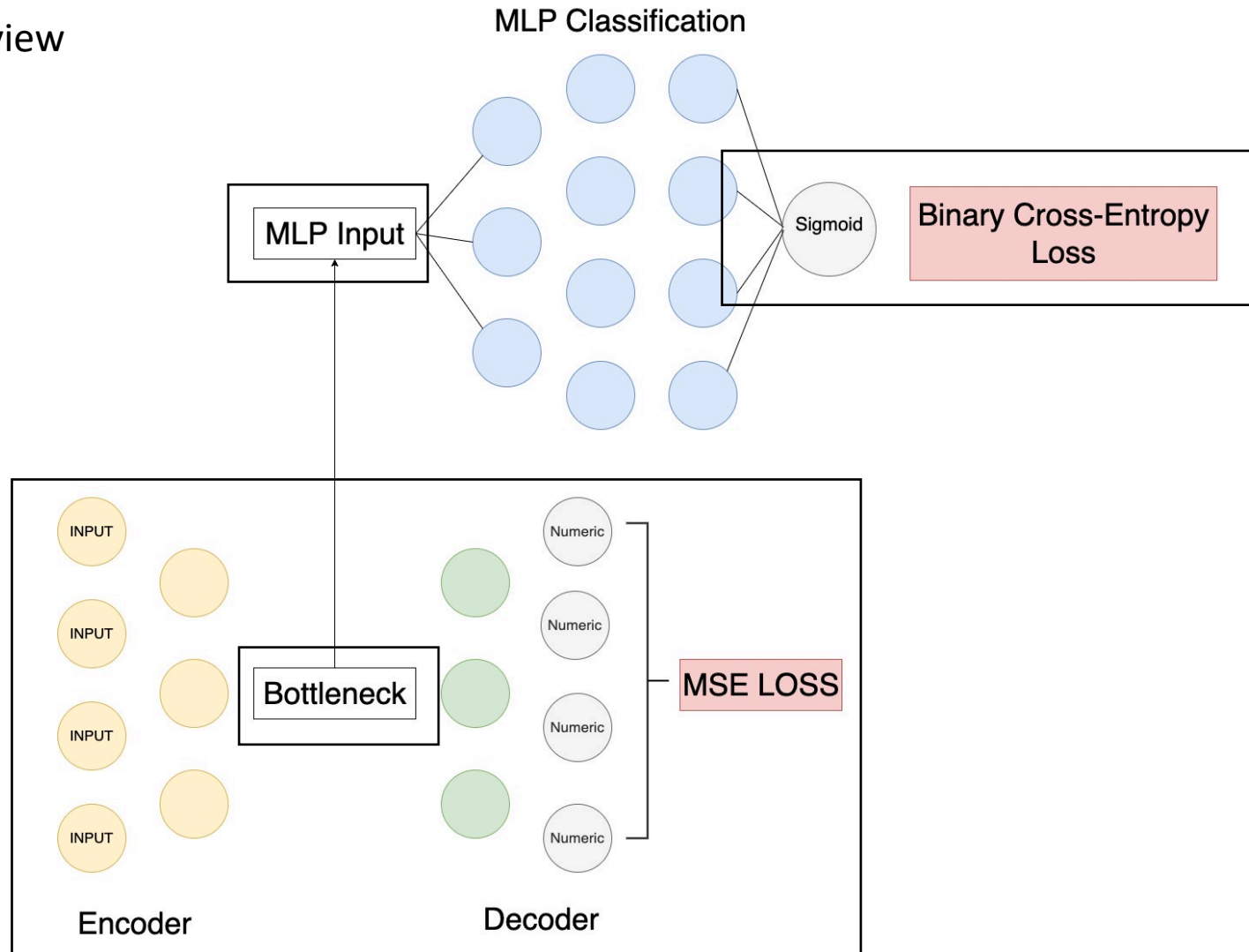Niclas Wölner-Hanssen, VT 2023, STAN47 – Deep Learning

# Targeted-Guided AutoEncoder

- Feature extraction/ dimension reduction technique that is guided to construct latent variables that are informative.

- Often times the indirect purpose of dimension reduction techniques such as PCA and AE is to reduce the number of dimensions to make better predictions (overcome multicollinearity, overfitting etc.)

- However, these transformed new features, the latent variables, are not necessarily predictive.

- No information about the target, the thing we try to predict, is introduced in the transformation of the features into the latent variables when using PCA or Simple AutoEncoder.
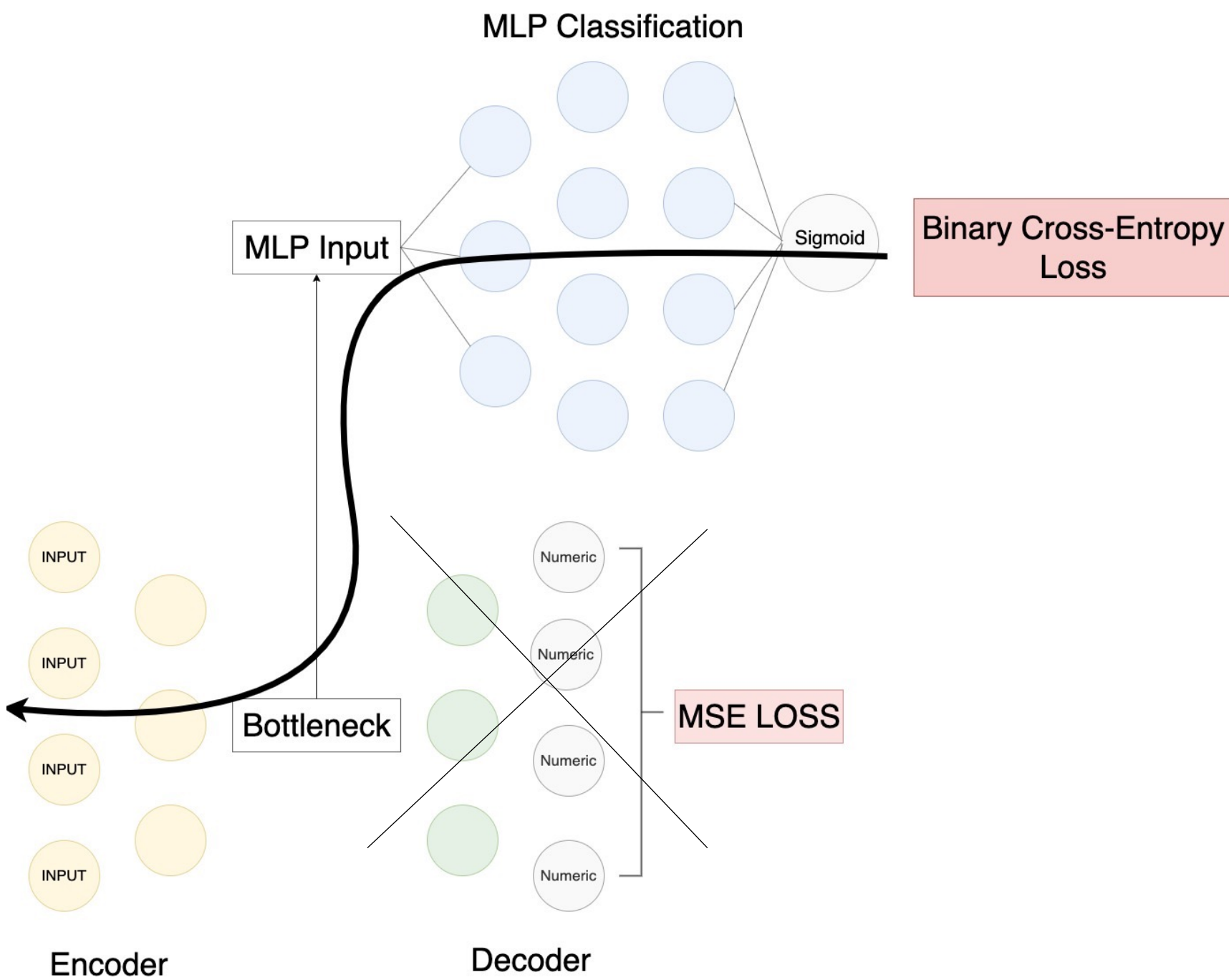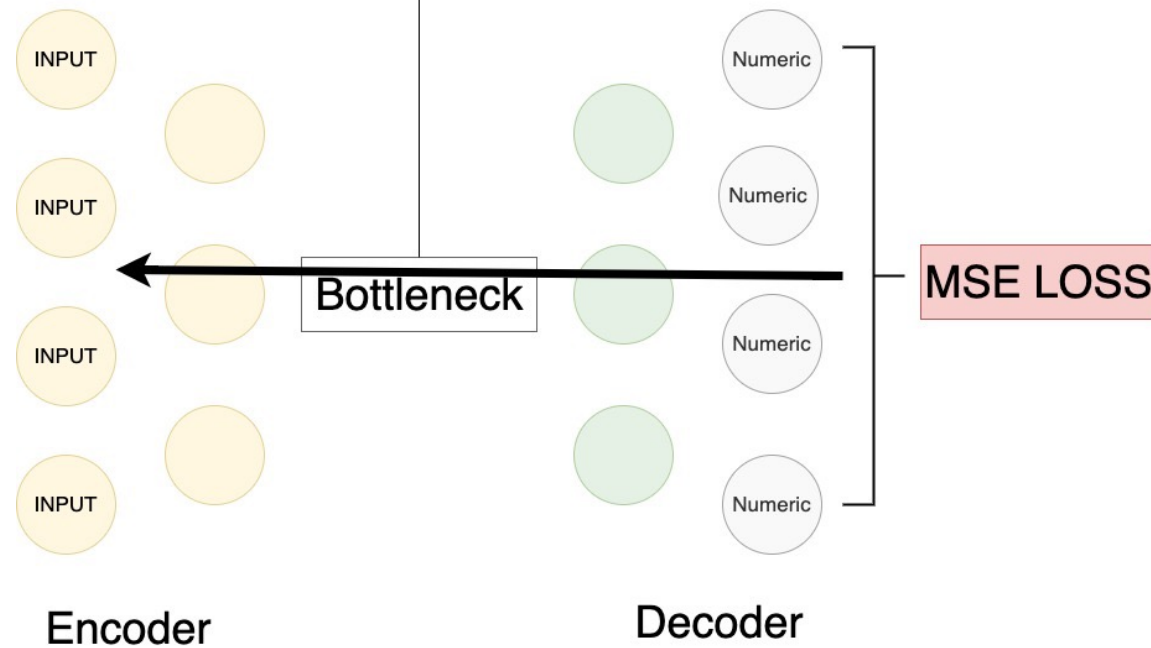
# Targeted-Guided AutoEncoder

Model Overview

# Target-Guided AutoEncoder

- By including a seperate MLP branch in the overall model architecture, the target information (i.e y_train) is introduced to the shared bottleneck layer of the autoencoder.

- During training, the gradients of the MLP loss with respect to the weights of the MLP are propagated backwards through the shared bottleneck layer of the autoencoder, allowing the model to update the bottleneck representation in a way that is informed by the classification task.

- Bottomline is that the inclusion of the MLP branch enables the model to learn the bottleneck representation that is not only optimized for reconstruction of the input data but also takes the target information into account.
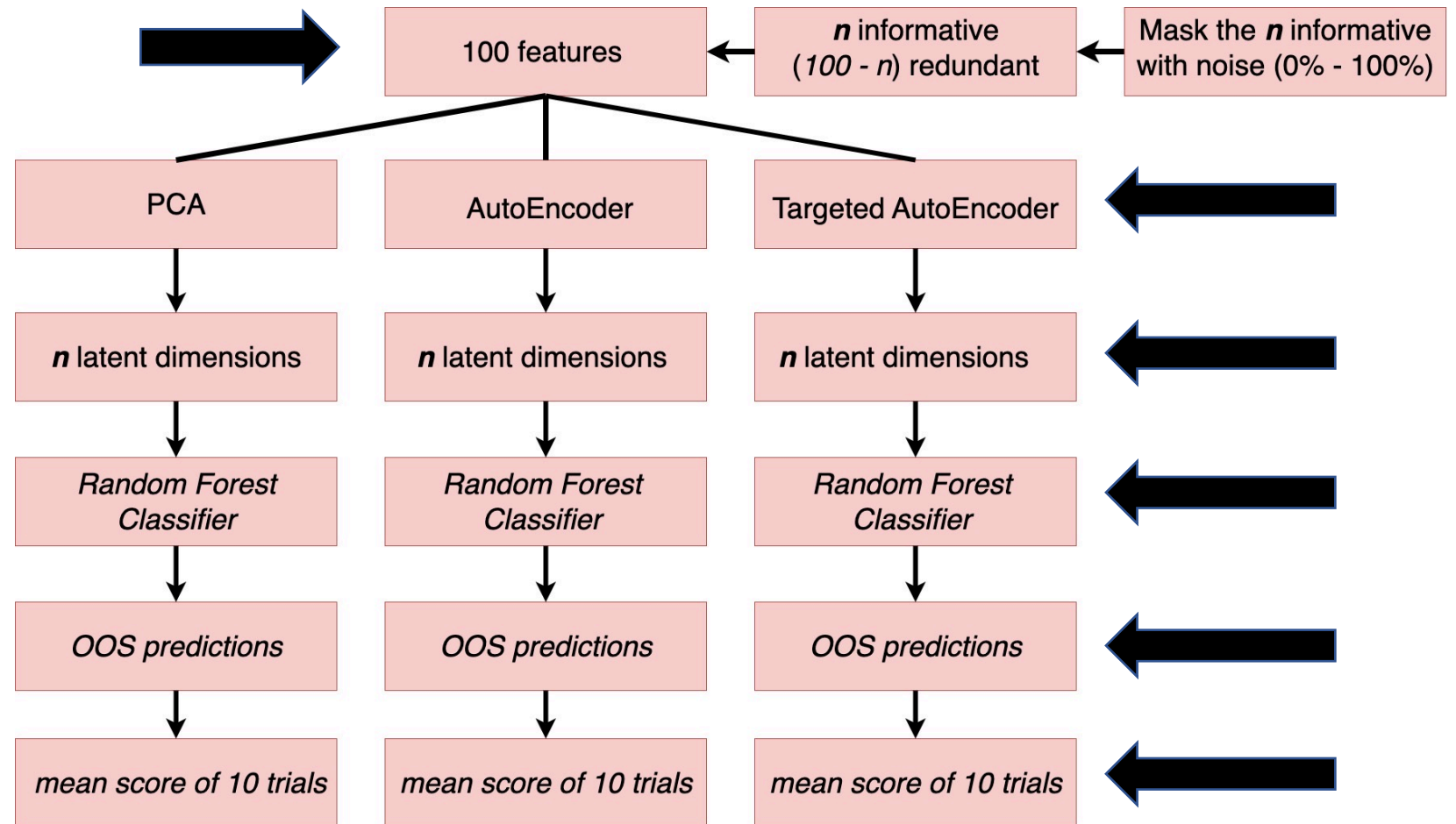
Testing

```
                          ┌──────────────┐        ┌──────────────────┐        ┌─────────────────────┐
        ═══════▶          │ 100 features │ ◀───── │   n informative   │ ◀───── │ Mask the n informative│
                          └──────────────┘        │ (100 - n) redundant│        │ with noise (0% - 100%)│
                                                   └──────────────────┘        └─────────────────────┘
```

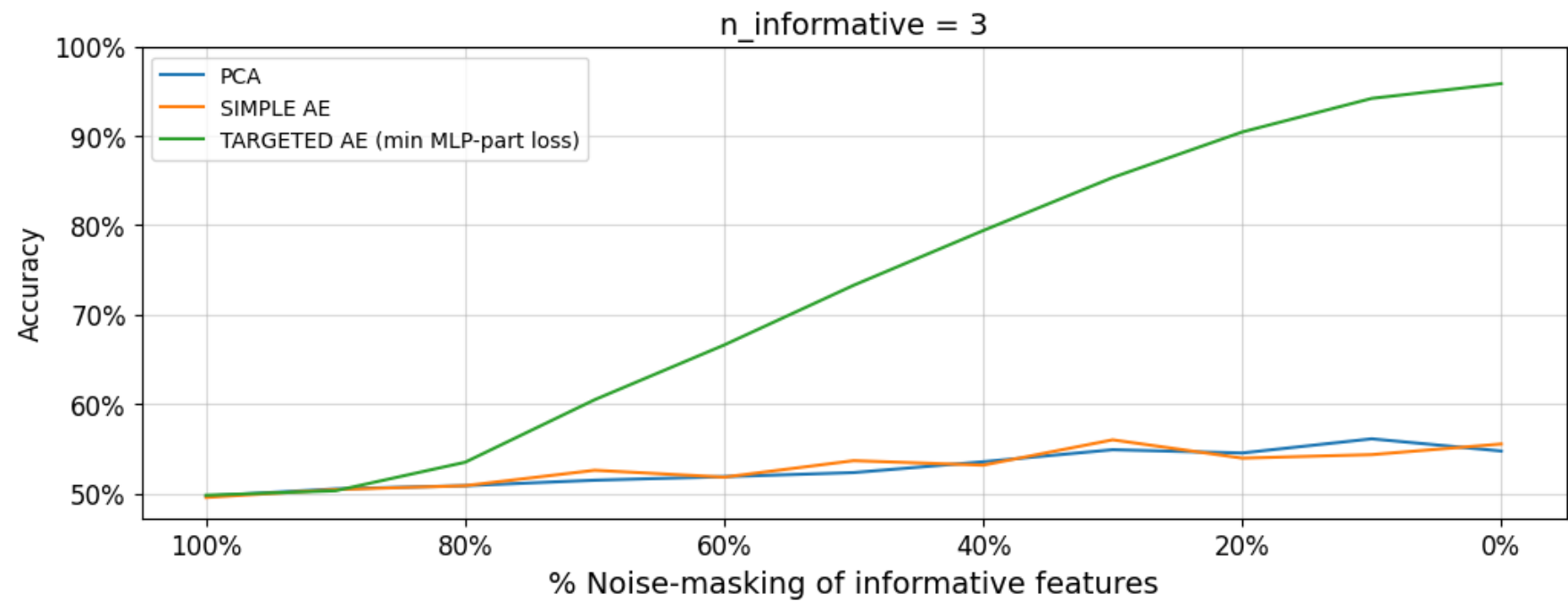| PCA | AutoEncoder | Targeted AutoEncoder |
|---|---|---|
| $n$ latent dimensions | $n$ latent dimensions | $n$ latent dimensions |
| *Random Forest Classifier* | *Random Forest Classifier* | *Random Forest Classifier* |
| *OOS predictions* | *OOS predictions* | *OOS predictions* |
| *mean score of 10 trials* | *mean score of 10 trials* | *mean score of 10 trials* |

*RF: n_estimators =500, max_depth = 5, max_samples = 100

n_informative = 3
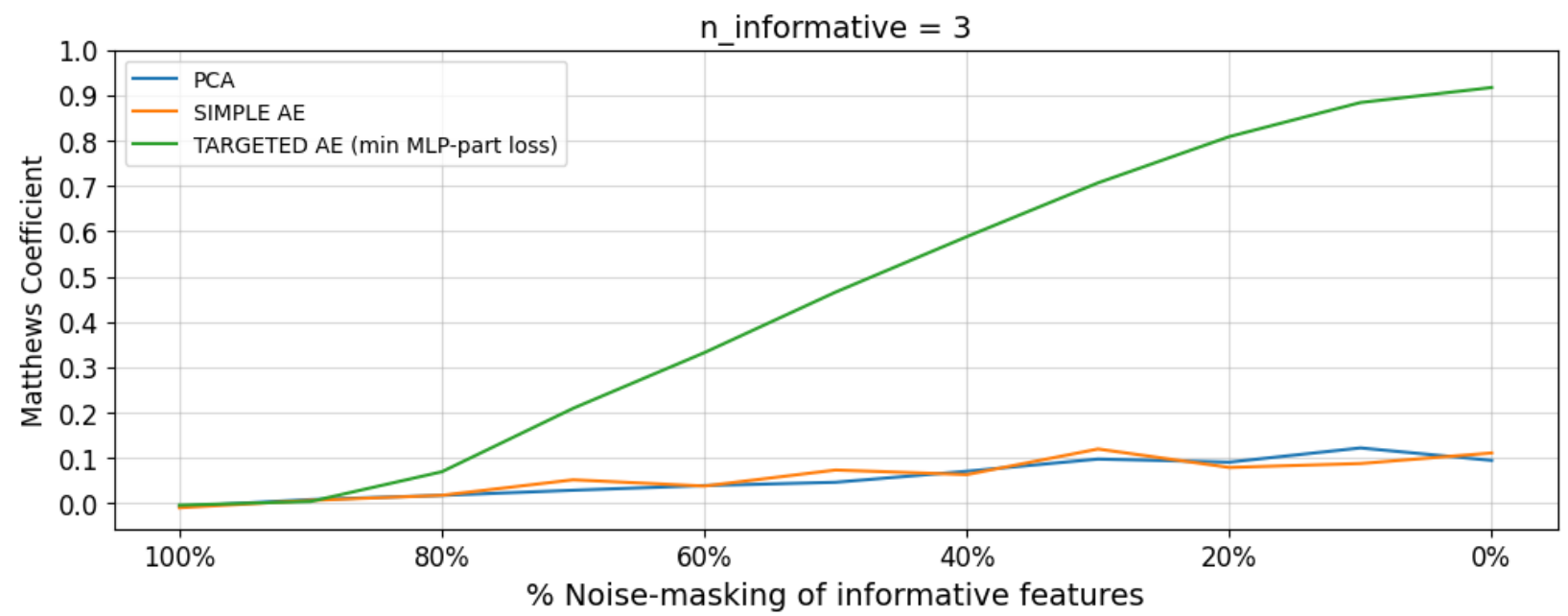
Testing

- 100 features (IID Std Normal, 10K samples)
- The sum of the first 3 features makes up the binary target:
  - y>0 = 1
    y<0 = 0
- The other features are thus just noise.
- In the plot the 3 informative features are after the construction of the target masked with noise. Going from complete noise (100%) to no-noise (0%)(x-axis).
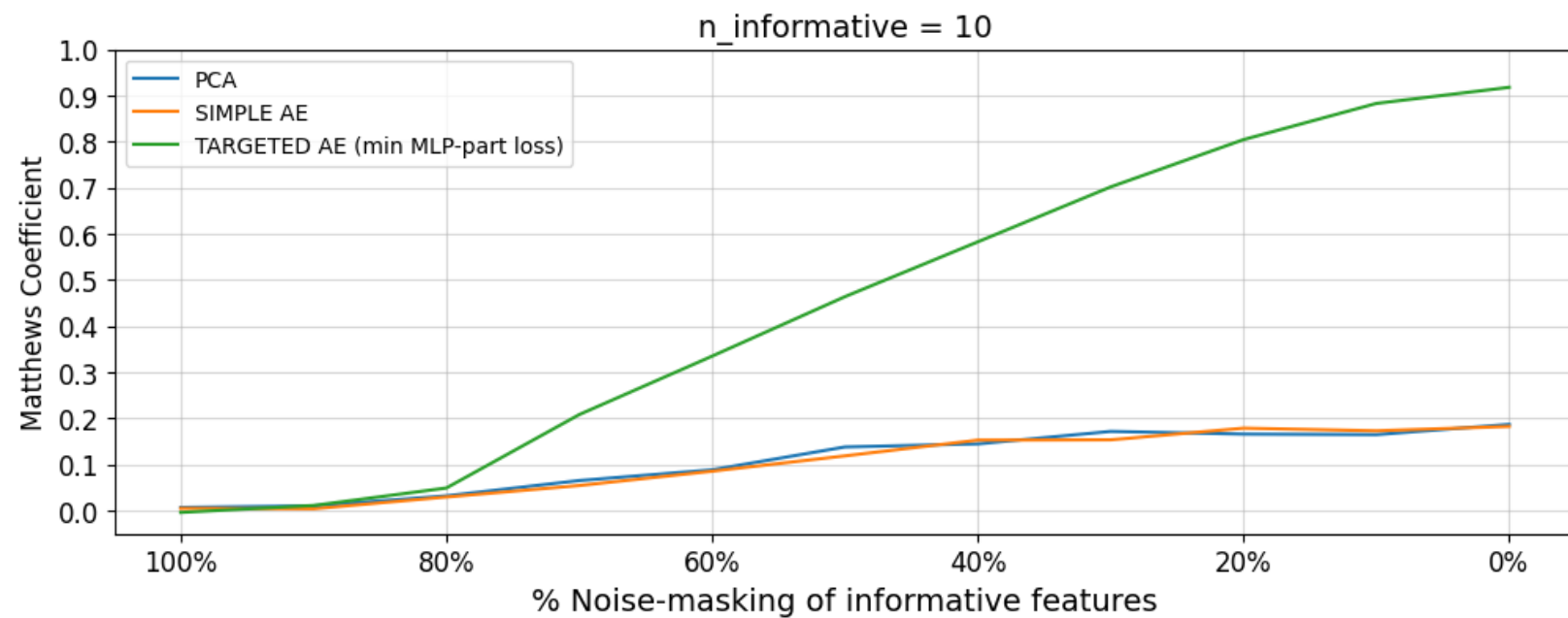
Testing

- Matthews correlation coefficient

Testing

**n_informative = 10**

Accuracy vs % Noise-masking of informative features

- PCA
- SIMPLE AE
- TARGETED AE (min MLP-part loss)

**n_informative = 10**

Matthews Coefficient vs % Noise-masking of informative features

- PCA
- SIMPLE AE
- TARGETED AE (min MLP-part loss)

Testing

n_informative = 50

PCA
SIMPLE AE
TARGETED AE (min MLP-part loss)

Accuracy
% Noise-masking of informative features

n_informative = 50

PCA
SIMPLE AE
TARGETED AE (min MLP-part loss)

Matthews Coefficient
% Noise-masking of informative features

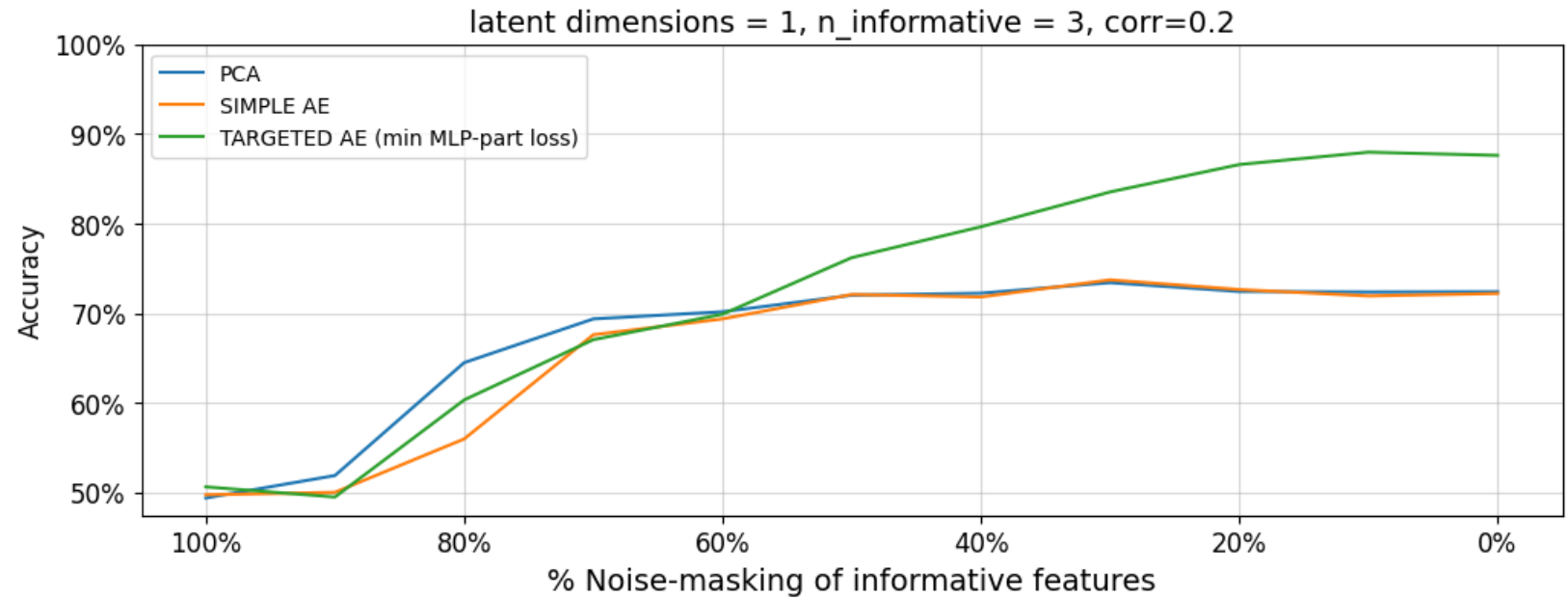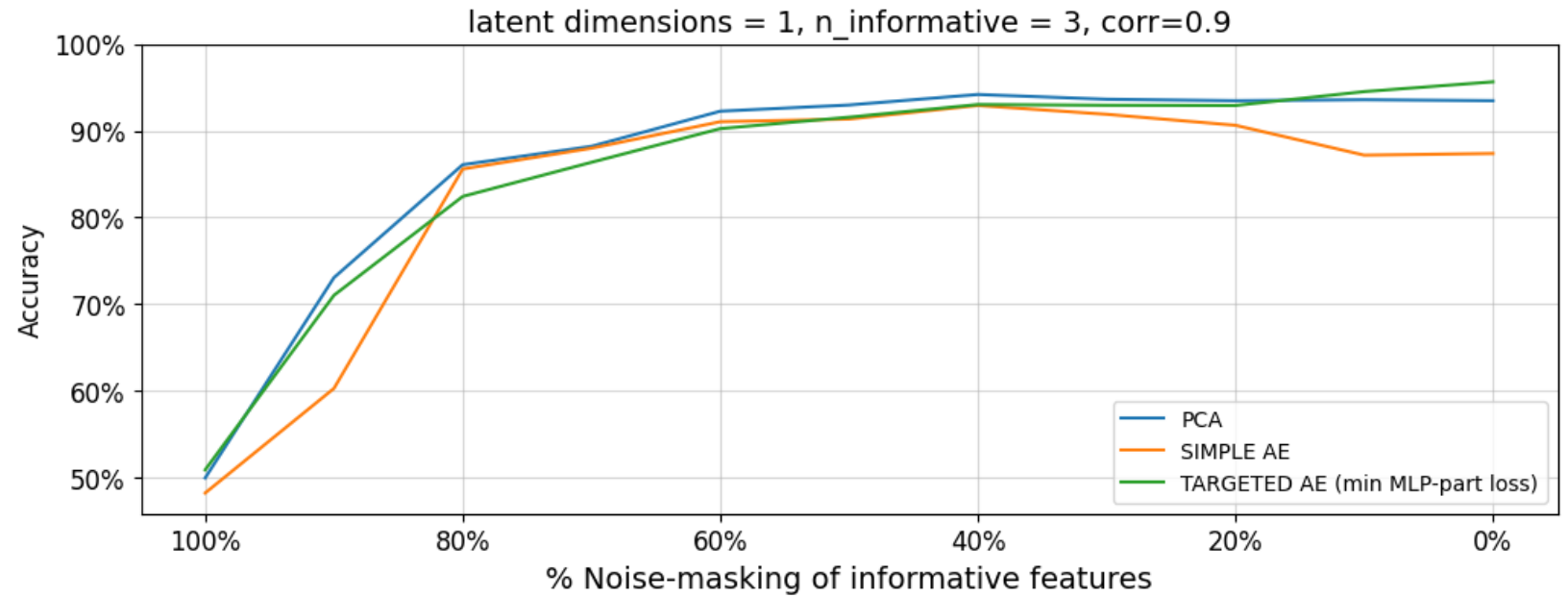# Target-Guided AutoEncoder

- So far, the samples has been IID. However, PCA is based on the covariance/correlation matrix and generally the higher the correlation between the features the more information (i.e. variance) of these it can store in fewer dimensions.

- For this reason, lets introduce some correlation between the features.

Testing

- The target is constructed as the sign of the sum of the 3 informative features, the other 97 features are "redundant".
- However, the correlation between all features is approx 0.2 when noise = 0%.
- PCA: First principal component
- AE & Targeted AE: 1 latent dimension

latent dimensions = 1, n_informative = 3, corr=0.9
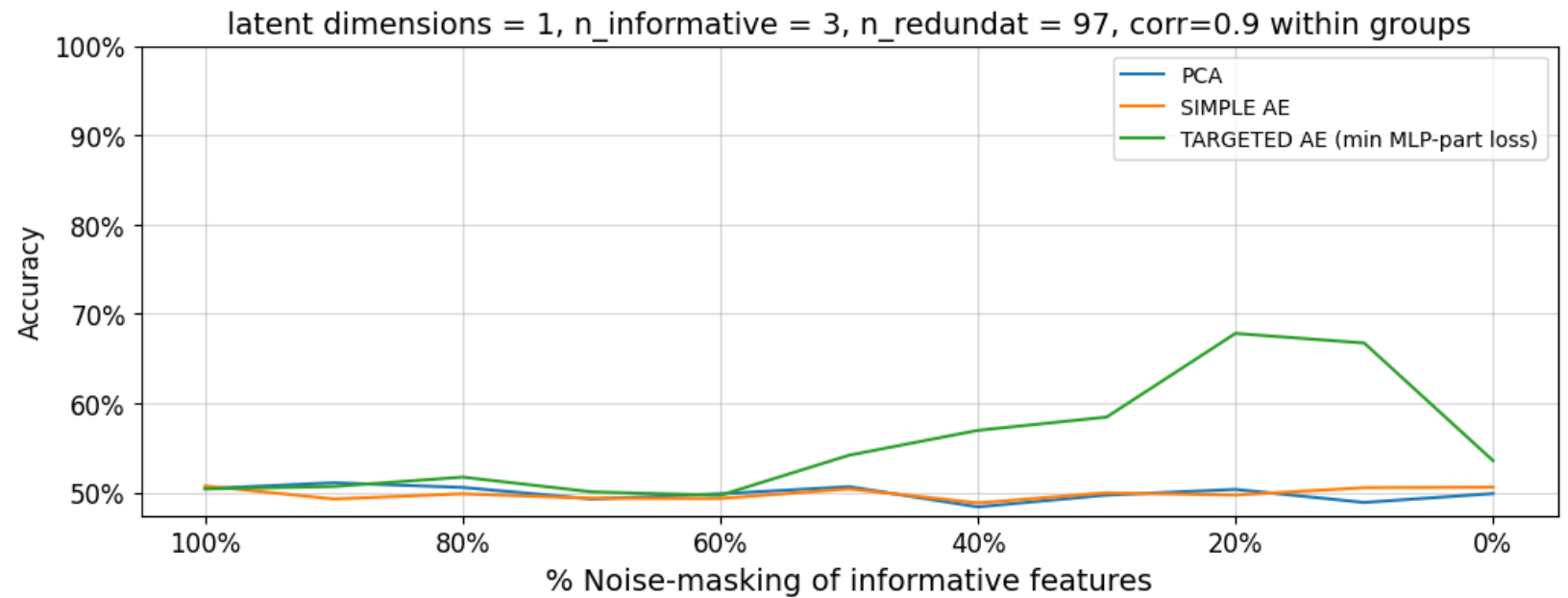
**Testing**

- The target is constructed as the sign of the sum of the 3 informative features, the other 97 features are "redundant".
- However, the correlation between all features is approx 0.9 when noise = 0%.
- PCA: First principal component
- AE & Targeted AE: 1 latent dimension

latest dimensions = 1, n_informative = 3, n_redundat = 97, corr=0.9 within groups

Testing

- Within the 3 informative features the correlation is 0.9 when noise = 0%.
- Within the 97 redundant features the correlation is 0.9 when noise = 0%.
- Between these groups the correlation is approx 0.
- PCA: First principal component
- AE & Targeted AE: 1 latent dimension

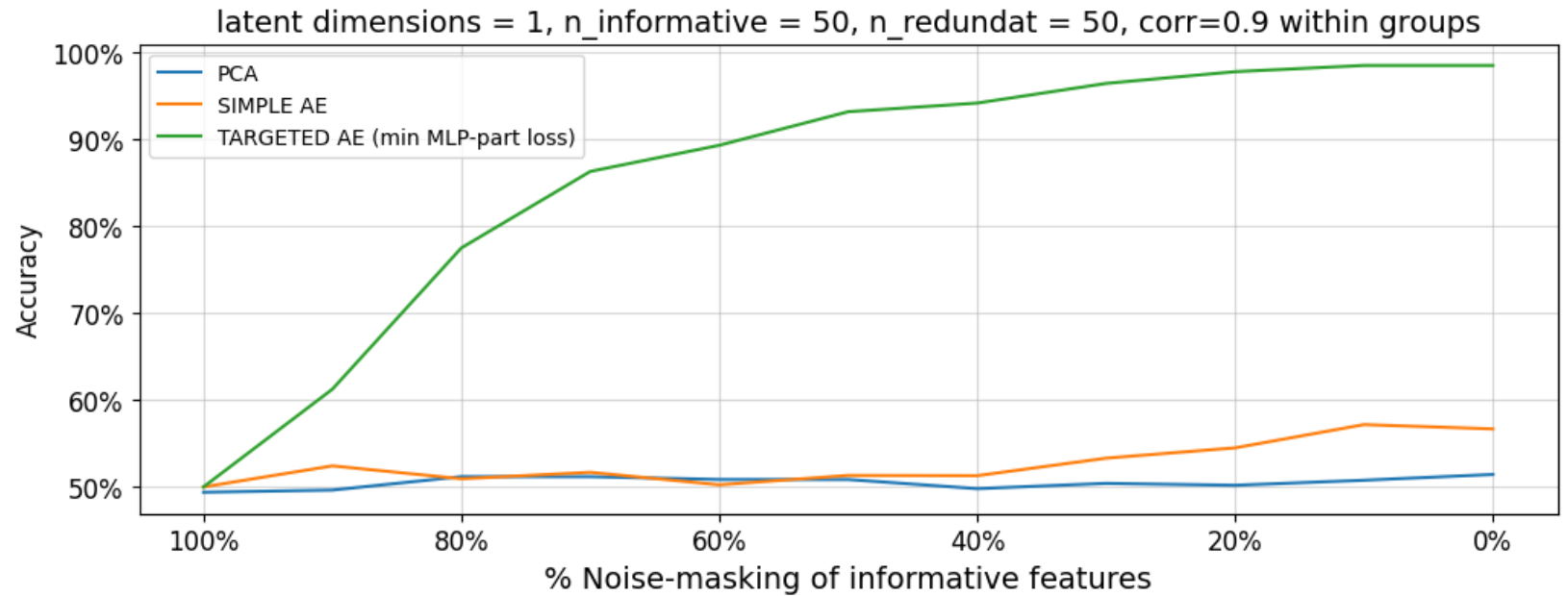latent dimensions = 1, n_informative = 50, n_redundat = 50, corr=0.9 within groups

Testing

- Within the 50 informative features the correlation is 0.9 when noise = 0%.
- Within the 50 redundant features the correlation is 0.9 when noise = 0%.
- Between these groups the correlation is approx 0.
- PCA: First principal component
- AE & Targeted AE: 1 latent dimension

# Conclusion:

Testing

- The Target-Guided AutoEncoder tend to at least never underperform PCA and simple AE, and most often it outperforms in terms of informative quality in constructing latent variables.

- The results also shows that PCA and Simple AutoEncoder seems to from time to time miss the informative information in the features, whereas the Target-Guided AutoEncoder finds these with realtive ease most of the time.

```python
In [94]:   1  # input regularization
           2  input_x = Input(shape=(x_inputs,))
           3  input_x = GaussianNoise(0.03)(input_x)
           4  input_x = BatchNormalization()(input_x)
           5  # encoder level 1
           6  e = Dense(x_inputs, activation='tanh',kernel_initializer=RandomNormal(mean=0.0, stddev=0.01))(input_x) #layer 1
           7  e = BatchNormalization()(e)
           8  e = Dropout(dropout_rate)(e)
           9  # encoder level 2
          10  e = Dense(int(x_inputs/2),activation='tanh')(e) #layer 2
          11  e = BatchNormalization()(e)
          12  e = Dropout(dropout_rate)(e)
          13  # bottleneck
          14  bottleneck = Dense(latent_dimensions)(e) #Bottleneck
          15
          16  # decoder, level 1
          17  d = Dense(int(x_inputs/2),activation='tanh')(bottleneck) #layer 1
          18  d = BatchNormalization()(d)
          19  d = Dropout(dropout_rate)(d)
          20  # decoder level 2
          21  d = Dense(x_inputs,activation='tanh')(d)  #layer 2
          22  d = BatchNormalization()(d)
          23  d = Dropout(dropout_rate)(d)
          24  # output layer
          25  decoder_output = Dense(x_inputs, activation='linear',name = 'decoder_output')(d)
          26
          27  merged = concatenate([bottleneck])
          28

          29  ####MLP Classification
          30
          31  mlp = Dense(merged.shape[1], activation='tanh')(merged) #Input Layer MLP
          32  mlp = BatchNormalization()(mlp)
          33  mlp = Dropout(dropout_rate)(mlp)
          34
          35  mlp = Dense(x_inputs*3, activation='tanh')(mlp) #layer 1
          36  mlp = BatchNormalization()(mlp)
          37  mlp = Dropout(dropout_rate)(mlp)
          38
          39  mlp = Dense(x_inputs*3, activation='tanh')(mlp) #layer 2
          40  mlp = BatchNormalization()(mlp)
          41  mlp = Dropout(dropout_rate)(mlp)
          42
          43  mlp = Dense(x_inputs*2, activation='tanh')(mlp) #layer 3
          44  mlp = BatchNormalization()(mlp)
          45  mlp = Dropout(dropout_rate)(mlp)
          46
          47  mlp_output = Dense(1, activation = 'sigmoid', name = 'mlp_output')(mlp) #Output Layer MLP
          48
          49  # define autoencoder model
          50  model   = Model(inputs= input_x,     outputs=[decoder_output,mlp_output])
          51  encoder = Model(inputs= input_x   , outputs=bottleneck)
          52  decoder = Model(inputs= bottleneck, outputs=decoder_output)
          53  MLP     = Model(inputs= merged    , outputs=mlp_output)
          54
          55  # compile model
          56  model.compile(optimizer = tf.keras.optimizers.Adam(learning_rate = 1e-3),
          57                loss=   {'decoder_output': 'mse', 'mlp_output': 'binary_crossentropy'},
          58                metrics={'decoder_output': 'mse', 'mlp_output': AUC()},
          59                sample_weight_mode={'decoder_output': None, 'mlp_output': 'temporal'},
          60                weighted_metrics={'mlp_output': 'binary_crossentropy'}
          61               )
```

Code

Encoder

Decoder

Code

```python
es = EarlyStopping(monitor = 'val_mlp_output_loss',min_delta = 1e-4, patience = 20,
                          baseline = None, restore_best_weights = True, verbose = 0)



history = model.fit(

    x = X_train,
    y = {'decoder_output': X_train, 'mlp_output': y_train},
    epochs=200,
    batch_size=100,
    verbose=0,
    shuffle=True,
    validation_split = 0.2,
    callbacks = [ckp, es])
```

Thank you!