RUTGERS

THE STATE UNIVERSITY
OF NEW JERSEY

# Searching for New Physics with Deep Autoencoders
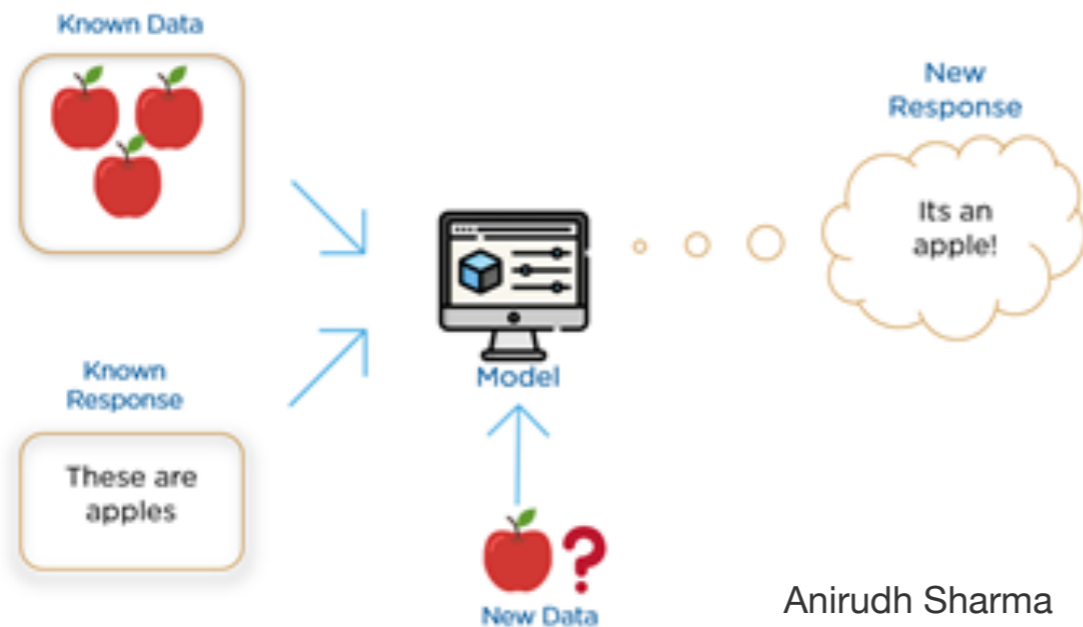
**Yuichiro Nakai (Rutgers)**

**Based on M. Farina, YN and D. Shih, arXiv:1808.08992 [hep-ph].**

**Beyond Standard Model: Where do we go from here?**
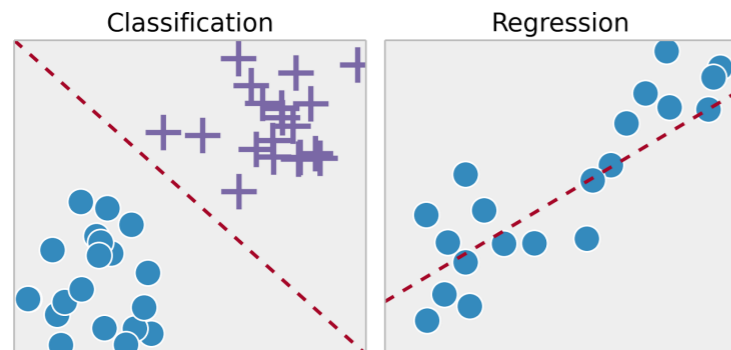
# Supervised or Unsupervised

Machine learning algorithms can be classified into:

## Supervised learning

Learn from <u>labeled</u> data



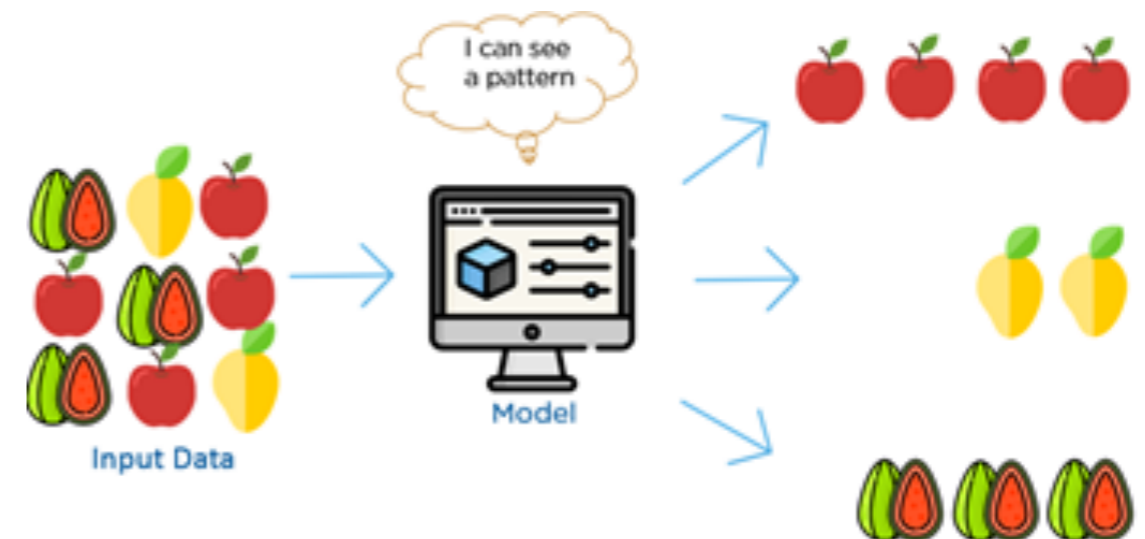Anirudh Sharma

Applications )



## Unsupervised learning

Learn from <u>unlabeled</u> data



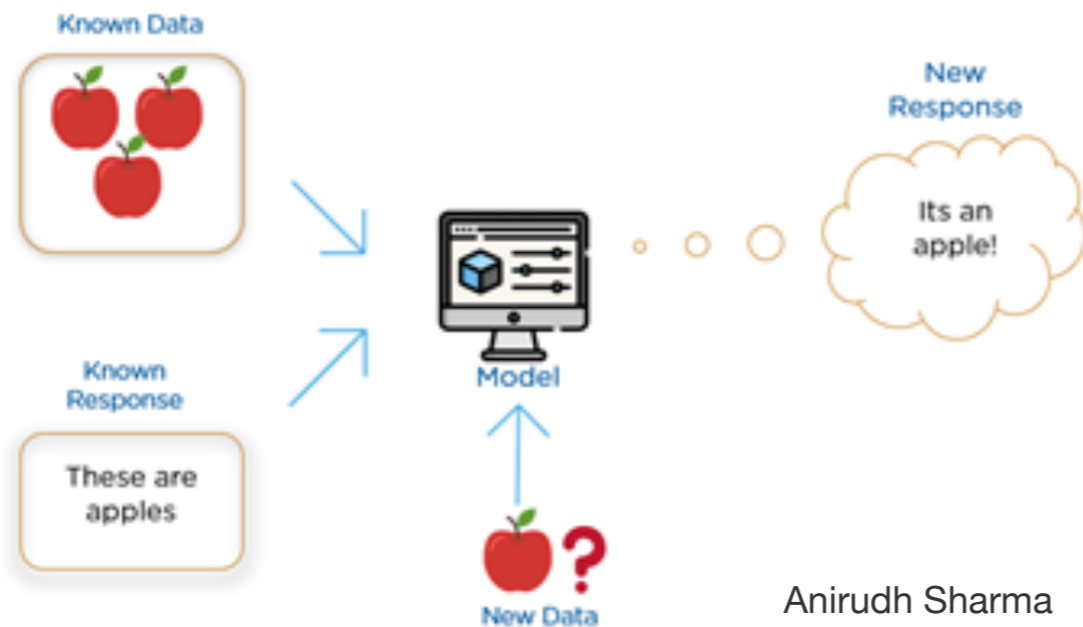The system looks for patterns and extracts features in data.

Applications )  Clustering

Anomaly detection
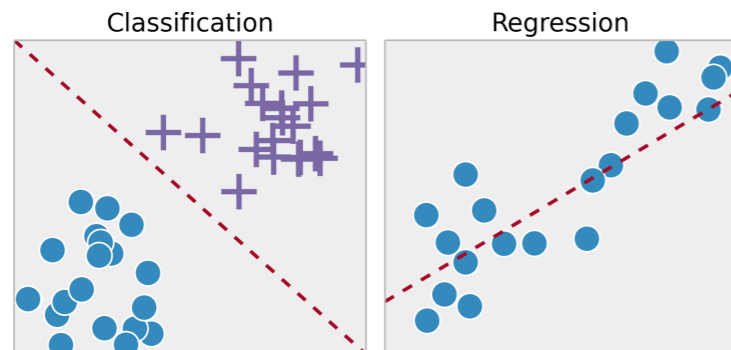
# Supervised or Unsupervised

Machine learning algorithms can be classified into:

## Supervised learning

Learn from <u>labeled</u> data
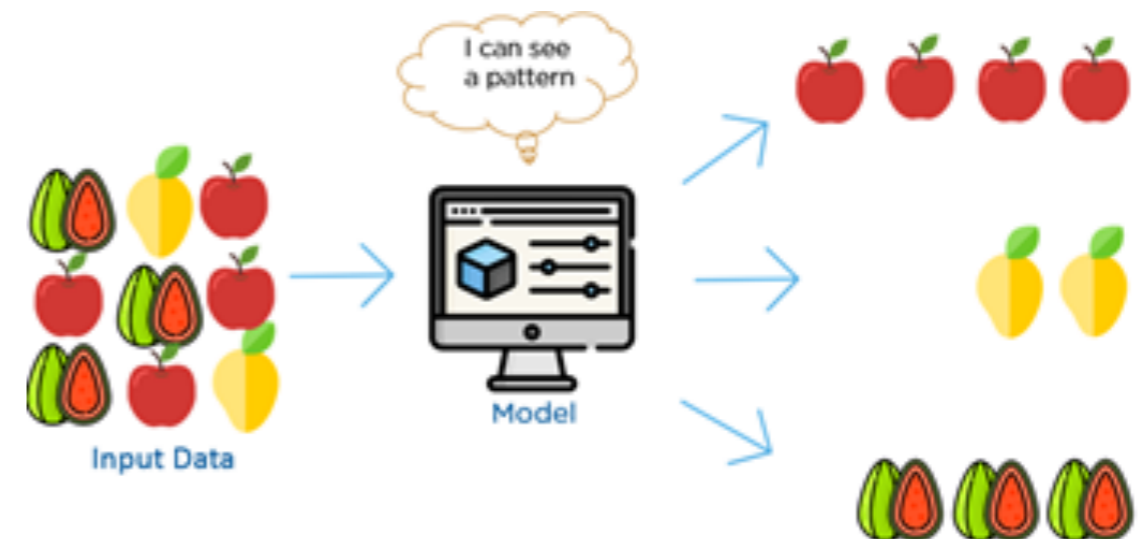


Anirudh Sharma

Applications )

## Unsupervised learning

Learn from <u>unlabeled</u> data



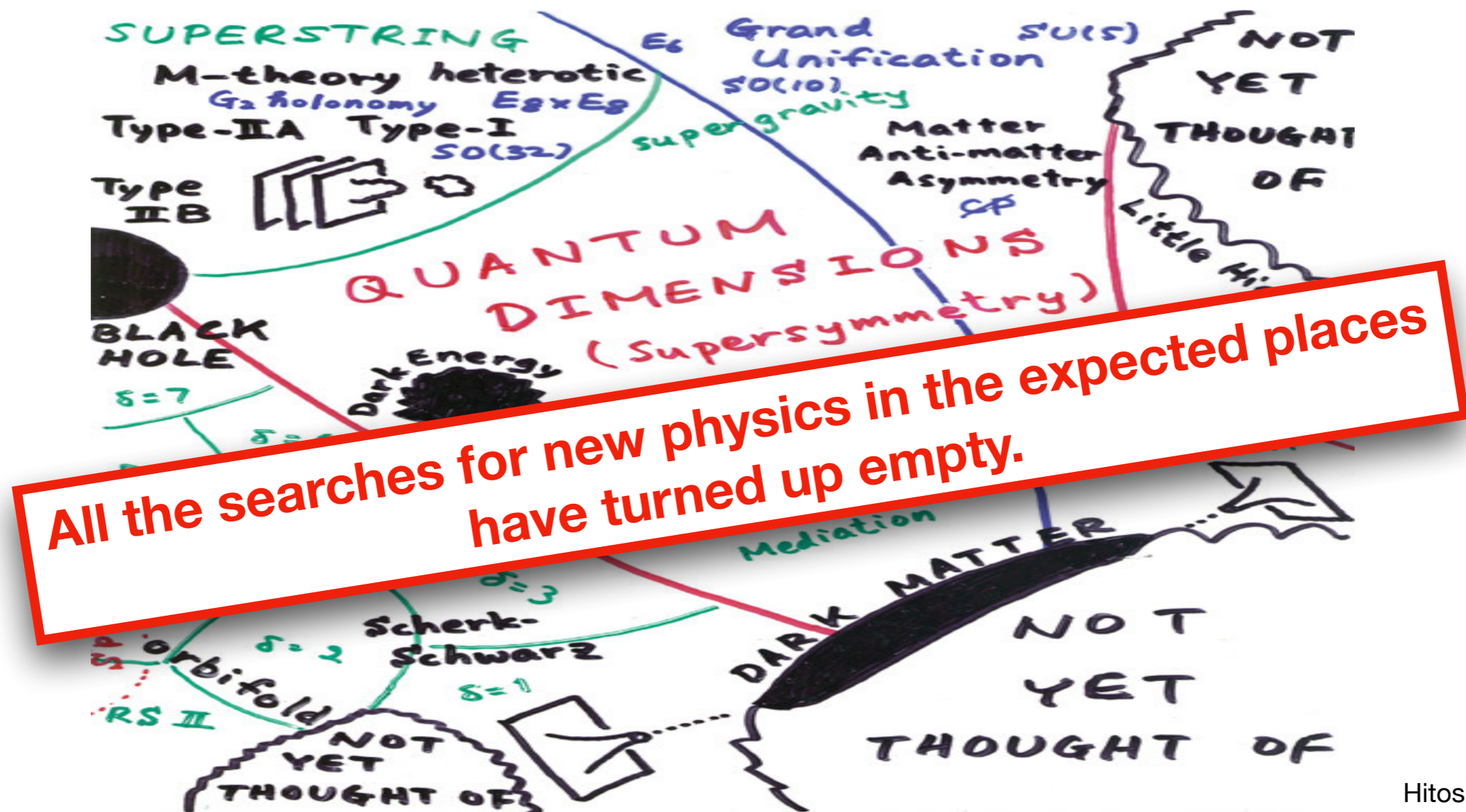The system looks for patterns and extracts features in data.

Applications )  Clustering

Anomaly detection

# Anomaly Detection

We have considered many possibilities of BSM physics with top-down theory prejudice (supersymmetry, extra dimension, …)



Hitoshi Murayama

We need more ways to discover the unexpected at the LHC, and here is where unsupervised machine learning comes into play.

# Autoencoder

Autoencoder is an unsupervised learning algorithm that maps an input to <u>a latent compressed representation</u> and then back to itself.



The Keras Blog

## Anomaly detection with autoencoder

- Autoencoder learns to map background events back to themselves.

- It <u>fails to</u> reconstruct anomalous events that it has never encountered.

➡ Signal the existence of anomaly !

# Sample Generation

*The idea is general, but concentrate on detection of anomalous jets.*

Generate jet samples by using PYTHIA for hadronization and Delphes for detector simulation.

Background : <u>QCD jets</u>     $p_T \in [800, \ 900] \ \mathrm{GeV}$     $|\eta| < 1$

Signal jets: <u>top jets</u>, <u>RPV gluino jets</u>   $m_{\tilde{g}} = 400 \ \mathrm{GeV}$

(decay to 3 light quark jets)

Match requirement : heavy resonance is within the fat jet, $\Delta R < 0.6$

Merge requirement : the partonic daughters of heavy resonance is within the fat jet, $\Delta R < 0.6$

We use sample sizes of 100k events for training and testing.

(The performance seems to saturate.)

# Jet Images

Concentrate on jet images ( 2D of eta and phi ) whose pixel intensities correspond to total pT.
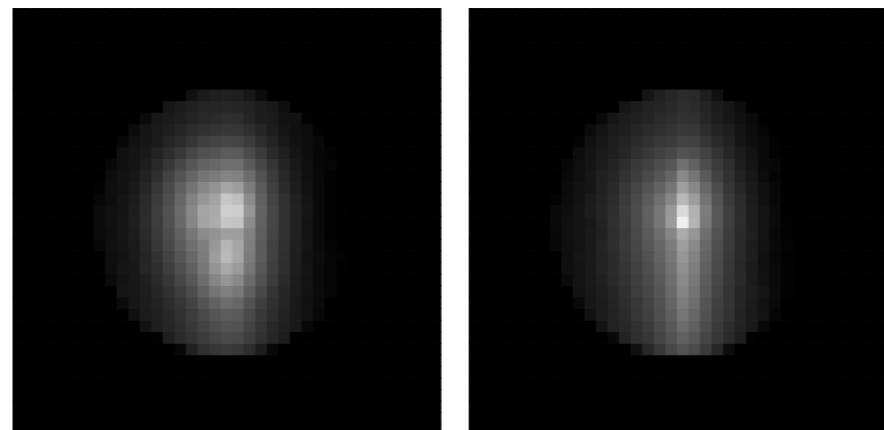
## Image pre-processing

1. Shift an image so that the centroid is at the origin

2. Rotate the image so that the major principal axis is vertical

3. Flip the image so that the maximum intensity is in the upper right region

4. Normalize the image to unit total intensity

5. Pixelate the image ( 37 x 37 pixels )

**Average images**

Left : top jets

Right : QCD jets



Macaluso, Shih (2018)

# Autoencoder Architectures

**Reconstruction error** : a measure for how well autoencoder performs.

$$L(x, \hat{x}) = \frac{1}{n} \sum_{i=1}^{n} \left| x_i - \hat{x}_i \right|^2 \qquad \begin{aligned} & x : \text{inputs} \\ & \hat{x} : \text{outputs} \end{aligned}$$

Train autoencoder to <u>minimize the reconstruction error</u> on background events.



**Architectures we consider** :

✓ Principal Component Analysis (PCA)

✓ Simple (dense) autoencoder

✓ Convolutional (CNN) autoencoder

# Principal Component Analysis

PCA is a technique to drop the least important variables by focusing on <u>variance</u> of data.

Original data       First PC       Reconstruction

Find the axis and project data to the axis

**"Encoder"**       **"Decoder"**

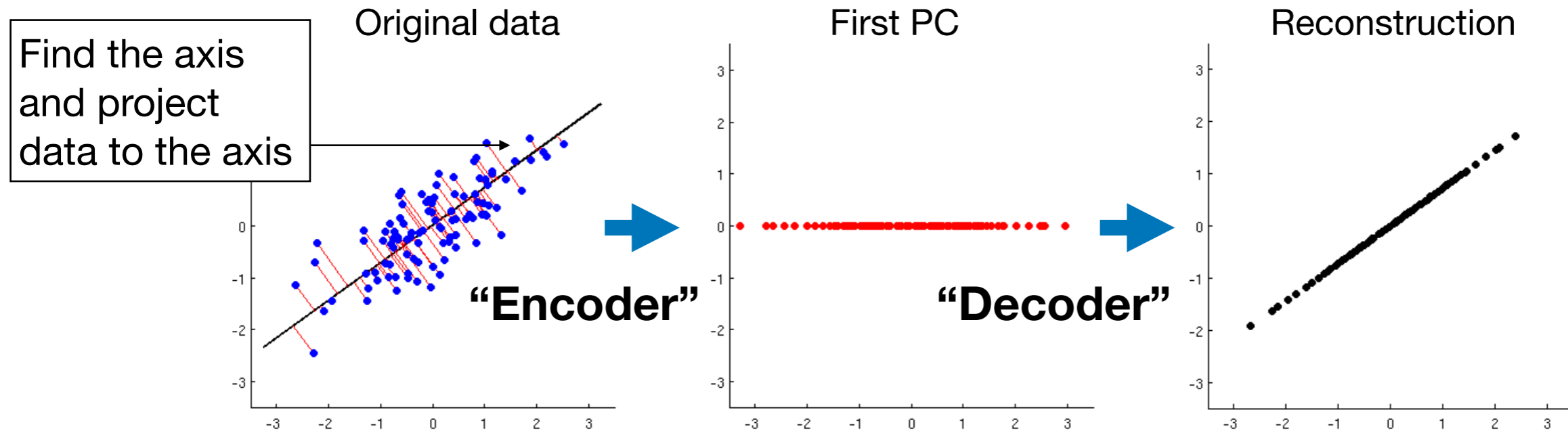Eigenvectors of covariance matrix of $\boldsymbol{x}_n - \boldsymbol{c}_0$ $\left( \boldsymbol{c}_0 = \sum_n \boldsymbol{x}_n / N \right)$ give desired axes.

$$\Gamma = (\boldsymbol{e}_1 \; \boldsymbol{e}_2 \; ... \; \boldsymbol{e}_d)$$    $d$ : the number of principal components ( $d < D$ )

**"PCA autoencoder"**

"Encoder" : $\tilde{\boldsymbol{x}}_n = (\boldsymbol{x}_n - \boldsymbol{c}_0)\Gamma$       "Decoder" : $\boldsymbol{x}'_n = \tilde{\boldsymbol{x}}_n \Gamma^T + \boldsymbol{c}_0$

# Simple Autoencoder

Autoencoder with a single <u>dense (fully-connected)</u> layer
as encoder and as decoder.

✓ Encoder and decoder are <u>symmetric</u>.

✓ The number of neurons in a hidden layer = 32.

✓ Flatten a jet image into <u>a single column vector</u>.

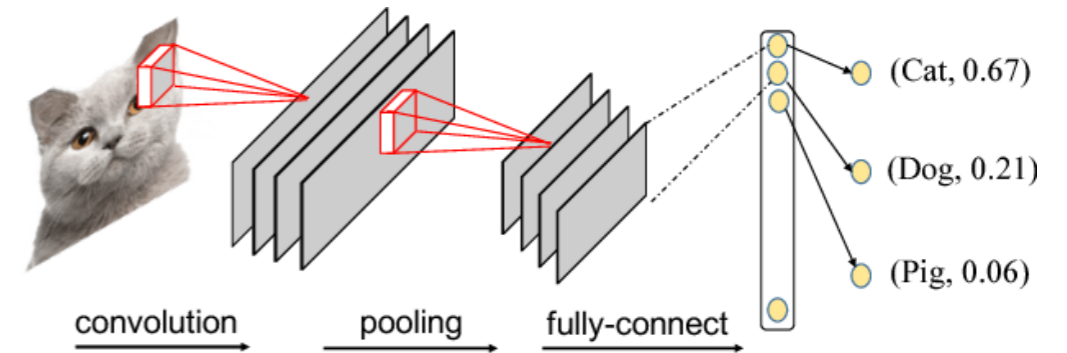✓ We use Keras with Tensorflow backend for implementation.

**Training details**

✦ The default Adam algorithm for optimizer.

✦ Minibatch size of 1024 ⟵ The number of images fed into the network at one time

~100 iterations of optimization in one epoch

✦ Early stopping : threshold = 0 and patience = 5 ⟵ To avoid overtraining

# Convolutional Autoencoder

## Convolutional Neural Network (CNN)

✓ Show high performance for <u>image recognitions</u>

✓ Maintain the <u>spacial information</u> of images



(Cat, 0.67)
(Dog, 0.21)
(Pig, 0.06)

convolution     pooling     fully-connect

arXiv:1712.01670

## Convolutional layer

$4\times1+9\times0+2\times(-1)$
$+5\times1+6\times0+2\times(-1)$
$+2\times1+4\times0+5\times(-1)=2$

Weights          Feature maps



## Max pooling

Reduce the image size



Single depth slice

max pool with 2x2 filters and stride 2

**Up sampling (pooling)**
also exists in autoencoder.

# Convolutional Autoencoder

Autoencoder architecture :



M. Ke, C. Lin, Q. Huang (2017)

**Encoder**          **Latent space**          **Decoder**

128C3-MP2-128C3-MP2-128C3-32N-6N-32N-12800N-128C3-US2-128C3-US2-1C3

128C3 : 128 filters with
    a 3x3 kernel

MP2 : max pooling with
    a 2x2 reduction factor

32N : a fully-connected layer
    with 32 neurons

US2 : up sampling with
    a 2x2 expansion factor

# Weakly-supervised mode

Weakly-supervised case with <u>pure</u> background events for training.

**Convolutional autoencoder**



**Average images**

QCD | Top | Gluino

Inputs

Outputs

Pixel-wise squared error

Autoencoder learns to reconstruct the QCD background.

More error

Autoencoder fails to reconstruct the signals.

**Reconstruction error is used as an anomaly threshold.**

# Autoencoder Performance

Performance measure :

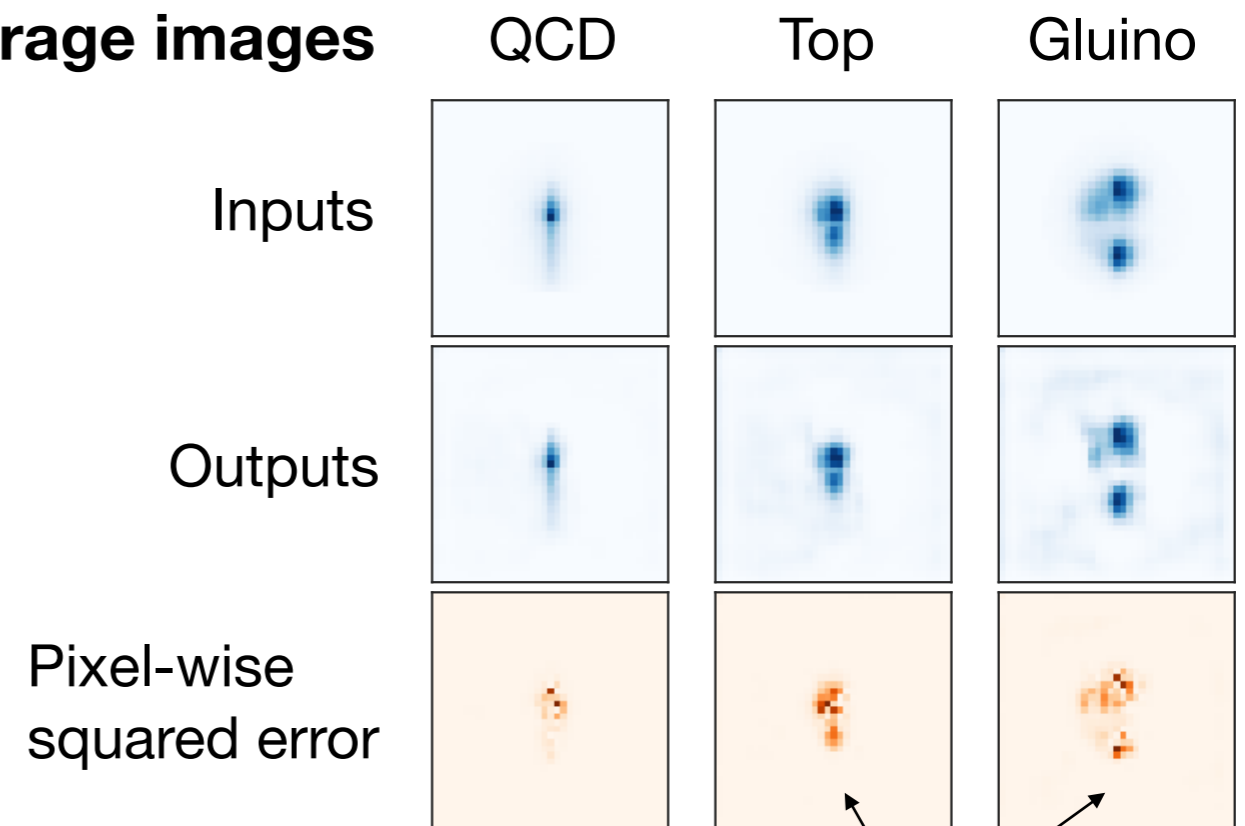$$\varepsilon_S = \frac{\text{(Correctly classified into signals)}}{\text{(Total number of signal jets)}}$$

$$\varepsilon_B = \frac{\text{(Misclassified into signals)}}{\text{(Total number of backgrounds)}}$$



Larger $\varepsilon_S$
Larger $\varepsilon_B$

Smaller $\varepsilon_S$
Smaller $\varepsilon_B$



**Top jets**

CNN outperforms the others.

**Gluino jets**

PCA outperforms CNN.

Jet mass as anomaly threshold

For gluino jets, PCA ROC curve approaches jet mass ROC curve, suggesting PCA reconstruction error is highly correlated with jet mass.

# Choosing the Latent Dimension k

Too small k  ➡️  Autoencoder cannot capture all the features.

Too large k  ➡️  Autoencoder approaches trivial representation.

Optimizing the latent dimension <u>using various signals</u> is **NOT** a good idea.

Instead, we use <u>the number of principal components</u> in PCA and <u>reconstruction error</u>.

Amount of variance ("scree plot") :

Reconstruction error :

**Choose k close to the "elbow"**

or

**Consider cumulative % of total variance**

Similar behavior as scree plot.

We choose <u>k = 6</u>.

# Choosing the Latent Dimension k

Let's examine our choice by looking at the top signal.

$E_{10,100}$ : the signal efficiency at 90% and 99% background rejection



Each dot corresponds to the average of 5 independent training runs.

**Autoencoder performance plateaus around k = 6.**

# Robustness with Other Monte Carlo

*Autoencoder really does not learn artifacts special to a Monte Carlo?*

One possible check :

Evaluate autoencoder (trained on PYTHIA samples) on <u>jet samples produced with HERWIG</u>.

**Comparison of reconstruction error (top jets, CNN)**

**The differences are small.**

Separation between background and anomaly is preserved.



Autoencoder probably learns fundamental jet features.

# Unsupervised mode

A much more exciting possibility is to train autoencoder on <u>actual data</u> (which may contain some amount of signals).

Train autoencoder on a sample of backgrounds <u>contaminated by</u> <u>a small fraction of signal events</u>.

➡ **Autoencoder performance is remarkably stable against signal contamination.**

Top jets for anomalous events

Reduction is not dramatic !

# Correlation with Jet Mass

*In actual new physics searches, we look for subtle signals …*

It's more powerful to combine autoencoder with another variable such as <u>jet mass</u>.

Cut hard on reconstruction error to clean out the QCD background and look for a bump in jet mass distribution.

**Reconstruction error should not be correlated with jet mass.**

Mean jet mass in bins of reco error for the QCD background

For **PCA and dense**, reco error is <u>correlated</u> with jet mass.

Jet mass distribution is <u>stable</u> against cutting on **CNN loss**.



PCA
Dense
CNN

Mean Jet Mass [GeV]

Reconstruction Error × $10^6$

# Correlation with Jet Mass

Jet mass distributions after cuts on CNN loss



Reduce the QCD background by a factor of 10, 100 and 1000.

**Convolutional autoencoder is useful for a bump hunt in jet mass above 300 GeV.**

Jet mass histograms normalized to LO gluino and QCD cross sections



**Before the cut**

$$S / B \approx 4\%$$

**After the cut**

$$S / B \approx 25\%$$

# Comments on "QCD or What?"

T. Heimel, G. Kasieczka, T. Plehn, J. Thompson, arXiv:1808.08979 [hep-ph].

They also consider anomaly detection through autoencoder.

Signal jets : top jets, scalar decay to jets, dark showers

Performance is comparable.

$$pp \to (\phi \to aa \to c\bar{c}\ c\bar{c}) + \text{jets}$$

$$m_a = 4 \text{ GeV} \qquad m_\phi = m_t = 175 \text{ GeV}$$



**Top jets**

Non-Adversarial
Bottleneck 32
AUC 0.90

Adversarial
Bottleneck 32
AUC 0.60

# Comments on "QCD or What?"

## Correlation with jet mass

They take an alternative approach using <u>adversarial networks</u>.

Additional adversary tries to extract jet mass from autoencoder output. ⟷ Autoencoder wants the adversary to be as unsuccessful as possible.

Autoencoder will avoid all information on jet mass.



Non-adversarial

Adversarial

Fake peak

Flatten

# Summary

✓ Autoencoder learns to map background events back to themselves but <u>fails to reconstruct signals</u> that it has never encountered before.

✓ Reconstruction error is used as an anomaly threshold.

✓ Autoencoder performance is <u>stable against signal contamination</u> which enables us to <u>train autoencoder on actual data</u>.

✓ Jet mass distribution is stable against cutting on CNN loss and <u>convolutional autoencoder is useful for a bump hunt in jet mass</u>.

✓ Thresholding on reco error gives <u>a significant improvement of S/B</u>.

# Future directions

✓ Testing out autoencoder on other signals.
( Other numbers of subjets, non-resonant particles, … )

✓ Training autoencoder to flag entire events as anomalous,
instead of just individual fat jets.

✓ Trying other autoencoder architectures on the market to improve
the performance.

✓ Understanding what the latent space actually learns.
( Jet mass? N-subjettiness? )                                    …

**Autoencoder is a powerful new method to search for
any signal of new physics without prejudice !**

*Thank you.*

# Backup Material

# What is Machine Learning?

**Machine learning** : technique to give computer systems the ability
to learn with data <u>without being explicitly programmed</u>.

Machine can learn the feature of data which human has not realized !

## Neural Networks

✓ Powerful machine learning-based techniques
used to solve many real-world problems

✓ Modeled loosely after the human brain

✓ Containing <u>weights</u> between neurons
that are tuned by learning from data

Input layer    Hidden layer    Output layer

Weights

Weights

Output

Networks contain multiple hidden layers    **Deep learning**

# What is Machine Learning?

The goal of training is to minimize <u>loss function</u> :

$$L = \sum_i f(p(\theta, x_i), y_i)$$

$p(\theta, x_i)$ : Prediction     $\theta$ : Weights

$x_i$ : Input     $y_i$ : Target value of example $i$

Mean squared error (MSE) :   $f(p, y) = (p - y)^2$

Cross entropy :   $f(p, y) = -(y \log p + (1 - y) \log(1 - p))$

**Loss function** $L$



Initial weights

Minimum

**Weights** $\theta$

Weights are updated according to derivative of loss function :

$$\Delta\theta = -\eta \nabla L$$

Learning rate

# Keras Codes

- **Simple autoencoder**

```
 1  input_img = Input(shape=(37*37,))
 2  layer = Dense(32, activation='relu')(input_img)
 3  encoded = Dense(6, activation='relu')(layer)
 4
 5  layer = Dense(32, activation='relu')(encoded)
 6  layer = Dense(37*37, activation='relu')(layer)
 7  decoded=Activation('softmax')(layer)
 8
 9  autoencoder=Model(input_img,decoded)
10  autoencoder.compile(loss=keras.losses.mean_squared_error,
11              optimizer=keras.optimizers.Adam())
```

# Keras Codes

- **Convolutional autoencoder**

```
1  input_img=Input(shape= (40, 40, 1))
2
3  layer=input_img
4  layer=Conv2D(128, kernel_size=(3, 3),
5              activation='relu',padding='same')(layer)
6  layer=MaxPooling2D(pool_size=(2, 2),padding='same')(layer)
7  layer=Conv2D(128, kernel_size=(3, 3),
8              activation='relu',padding='same')(layer)
9  layer=MaxPooling2D(pool_size=(2, 2),padding='same')(layer)
10 layer=Conv2D(128, kernel_size=(3, 3),
11              activation='relu',padding='same')(layer)
12 layer=Flatten()(layer)
13 layer=Dense(32, activation='relu')(layer)
14 layer=Dense(6)(layer)
15 encoded=layer
16
17 layer=Dense(32, activation='relu')(encoded)
18 layer=Dense(12800, activation='relu')(layer)
19 layer=Reshape((10,10,128))(layer)
20 layer=Conv2D(128, kernel_size=(3, 3),
21              activation='relu',padding='same')(layer)
22 layer=UpSampling2D((2,2))(layer)
23 layer=Conv2D(128, kernel_size=(3, 3),
24              activation='relu',padding='same')(layer)
25 layer=UpSampling2D((2,2))(layer)
26 layer=Conv2D(1, kernel_size=(3, 3),padding='same')(layer)
27 layer=Reshape((1,1600))(layer)
```

# CWoLa Hunting

J. Collins, K. Howe, B. Nachman, arXiv:1805.02664 [hep-ph].

Another approach to anomaly detection to extend bump hunt with machine learning.

Mass distribution    target

Auxiliary information



signal + bg

bg only



CWoLa labels

sideband

**Toy model**

$$Y = (x, y)$$

Background : $\begin{array}{l} -0.5 < x < 0.5 \\ -0.5 < y < 0.5 \end{array}$

Signal : $\begin{array}{l} -w/2 < x < w/2 \\ -w/2 < y < w/2 \end{array}$

Mixed Sample 1    Mixed Sample 2



0    1

Classifier

**Classification without labels (CWoLa)**

A classifier is trained to distinguish statistical mixtures of classes.

Metodiev, Nachman, Thaler