

# PART V: A VERY BRIEF INTRODUCTION TO DEEP UNSUPERVISED LEARNING

\*elements taken from D. Kirkby lectures at KSPA19

# WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for

**SUPERVISED**

LEARNS A MAP FROM  
X [FEATURES] TO Y  
[LABELS]

$$P(X|Y)$$

the machine is NOT told what to look for

**UN-SUPERVISED**

NO LABELS - DISCOVER  
PATTERNS

$$P(X)$$

IN THIS LAST PART WE ARE GOING TO BRIEFLY INTRODUCE CURRENT TECHNIQUES OF UNSUPERVISED LEARNING WITH NEURAL NETWORKS

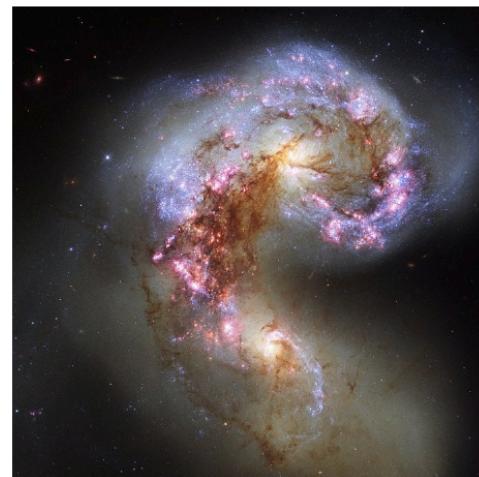
THERE ARE 3 MAJOR APPLICATIONS TO ASTRONOMY (THAT I CAN THINK OF):

- **DIMENSIONALITY REDUCTION:** HOW CAN I REPRESENT MY COMPLEX DATA MORE EFFICIENTLY TO GET NEW INSIGHTS INTO ITS STRUCTURE?
- **GENERATIVE MODELING:** HOW CAN I INTERPOLATE / EXTRAPOLATE A (SMALL, SPARSE) DATASET TO GENERATE NEW DATA SAMPLED FROM THE SAME (UNKNOWN) DISTRIBUTION?
- **PROBABILISTIC MODELING:** WHAT IS THE PROBABILITY THAT A NEW OBSERVATION IS DRAWN FROM THE SAME (UNKNOWN) DISTRIBUTION AS SOME REFERENCE (SMALL, SPARSE) DATASET?

# DIMENSIONALITY REDUCTION

NEED OF REPRESENTING THE DATA IN A SMALLER  
DIMENSIONALITY SPACE

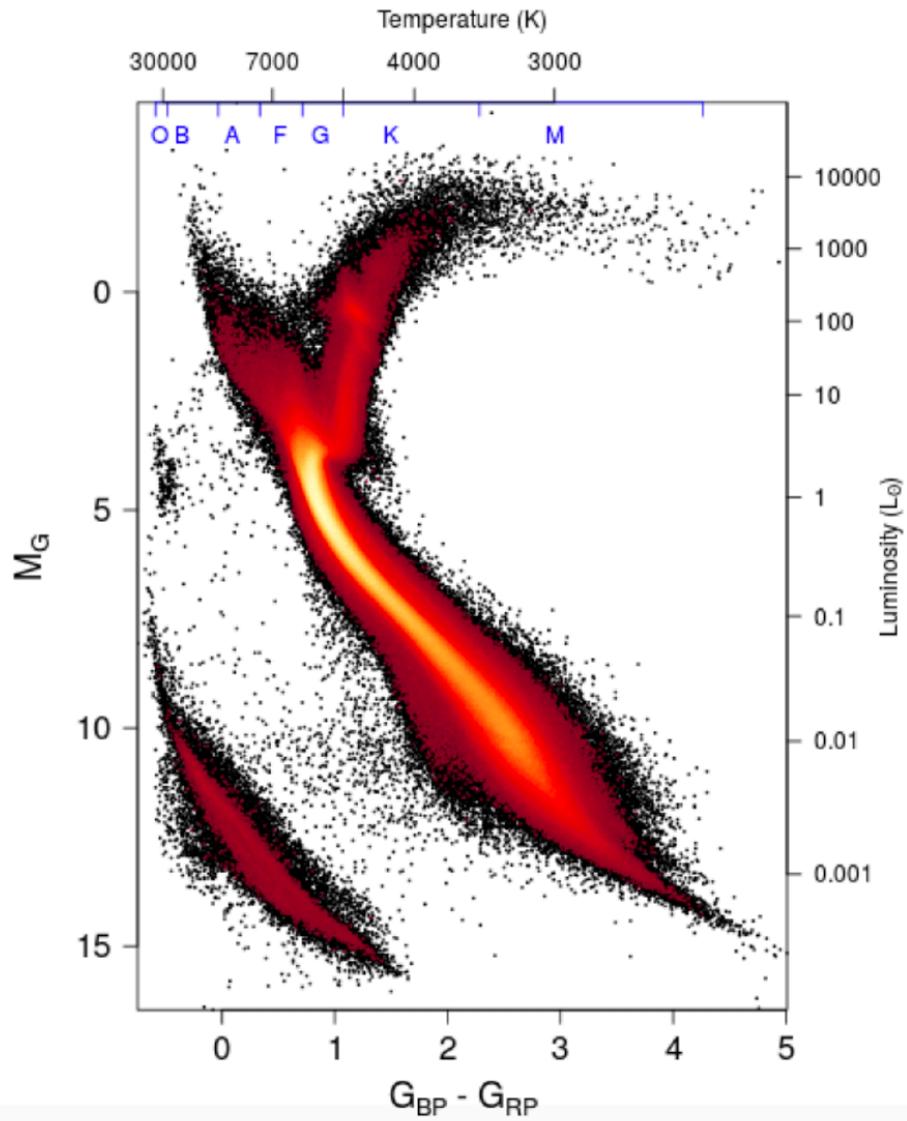
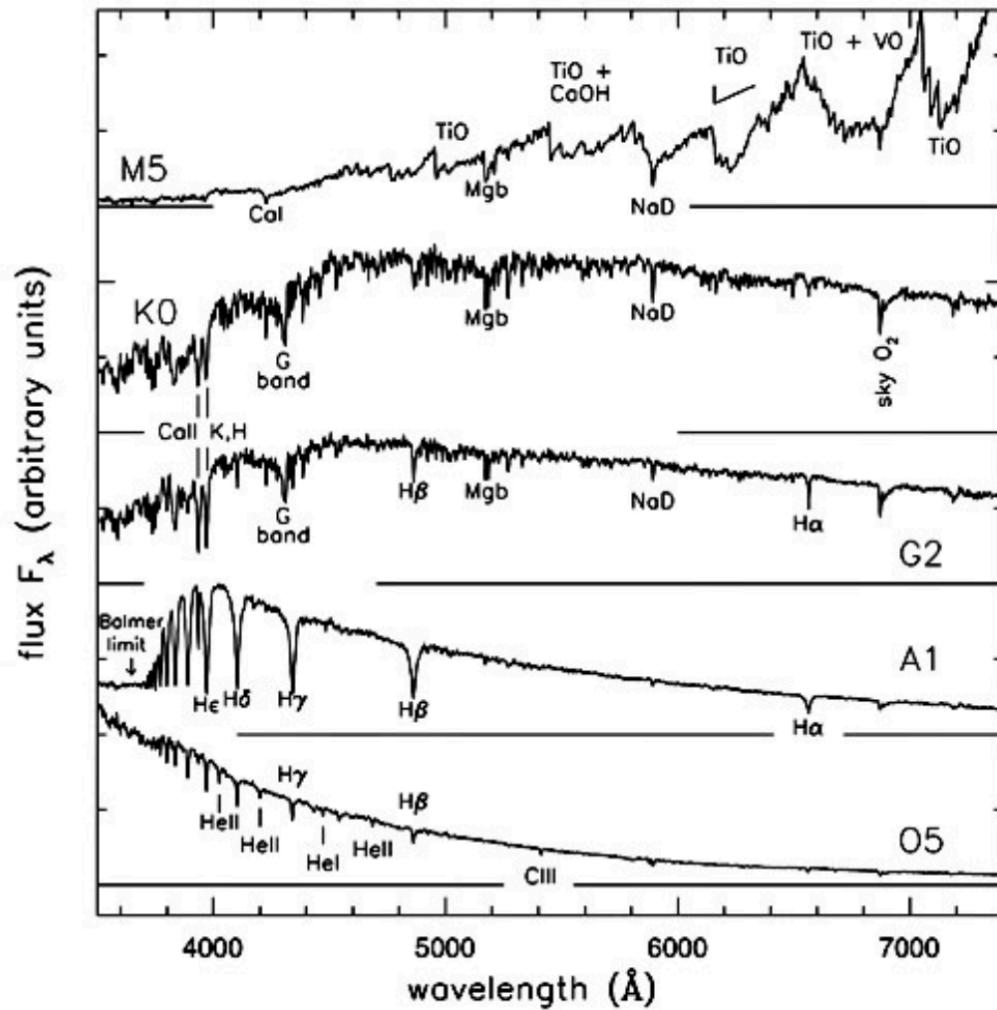
$$G_{\mu\nu} = \frac{8\pi G}{c^4} T_{\mu\nu}$$



@D. BARON

MOST OF THE SCIENTIFIC WORK IS ABOUT COMPRESSION OF  
NATURE IN SOME EXPLANATORY EQUATIONS

From: Gaia Collaboration et al. 2018



@D. BARON

**THE STELLAR MAIN SEQUENCE IS A DIMENSIONALITY  
REDUCTION**

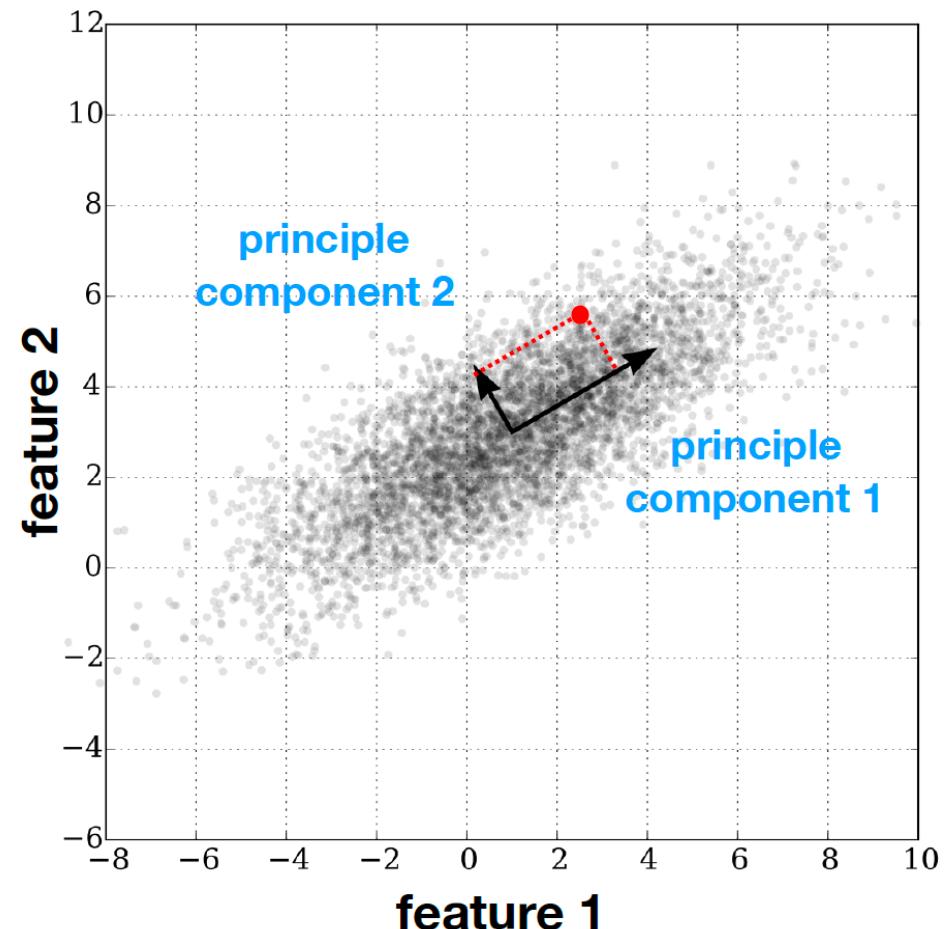
$$\begin{matrix} D \text{ dimensions (columns = features)} \\ N \text{ samples (rows)} \end{matrix} \boxed{\mathbf{X} - \boldsymbol{\mu}} \quad \approx \quad \begin{matrix} d \text{ latent variables} \\ \mathbf{Y} \\ \mathbf{M} \end{matrix}$$

CLASSICAL METHODS FOR DIMENSIONALITY REDUCTION  
SEEK A LINEAR DECOMPOSITION THAT BEST EXPLAINS  
THE OBSERVATIONS  $\mathbf{X}$  IN TERMS OF **LATENT VARIABLES  $\mathbf{Y}$**

# PRINCIPAL COMPONENT ANALYSIS

PCA CONVERT A SET OF  
(CORRELATED) VARIABLES INTO A  
SET  
OF VALUES LINEARLY  
UNCORRELATED

1. THE FIRST PRINCIPLE COMPONENT (“PROTOTYPE”), HAS THE LARGEST POSSIBLE VARIANCE
2. THE FOLLOWING COMPONENTS HAVE THE LARGEST VARIANCES WITH THE ADDITIONAL CONSTRAINT THAT THEY ARE ORTHOGONAL TO THE PRECEDING COMPONENTS

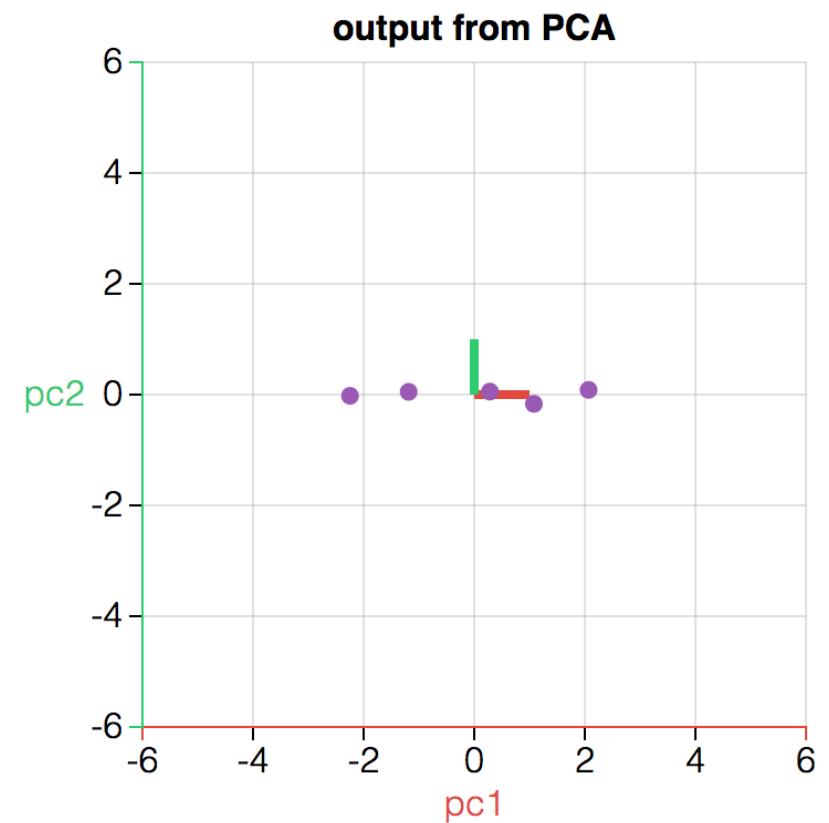
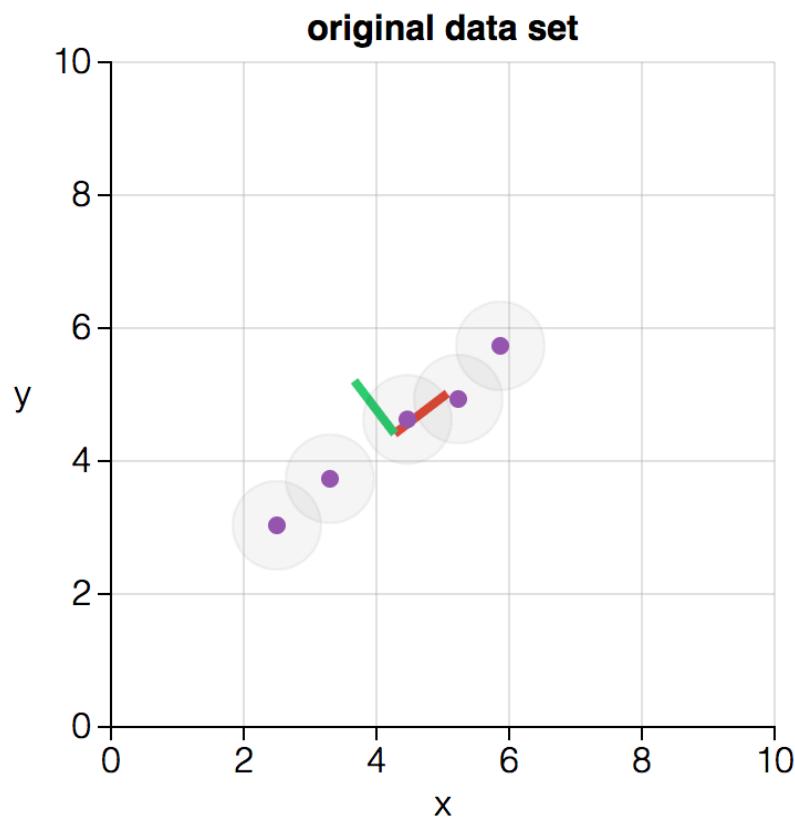


# PRINCIPAL COMPONENT ANALYSIS

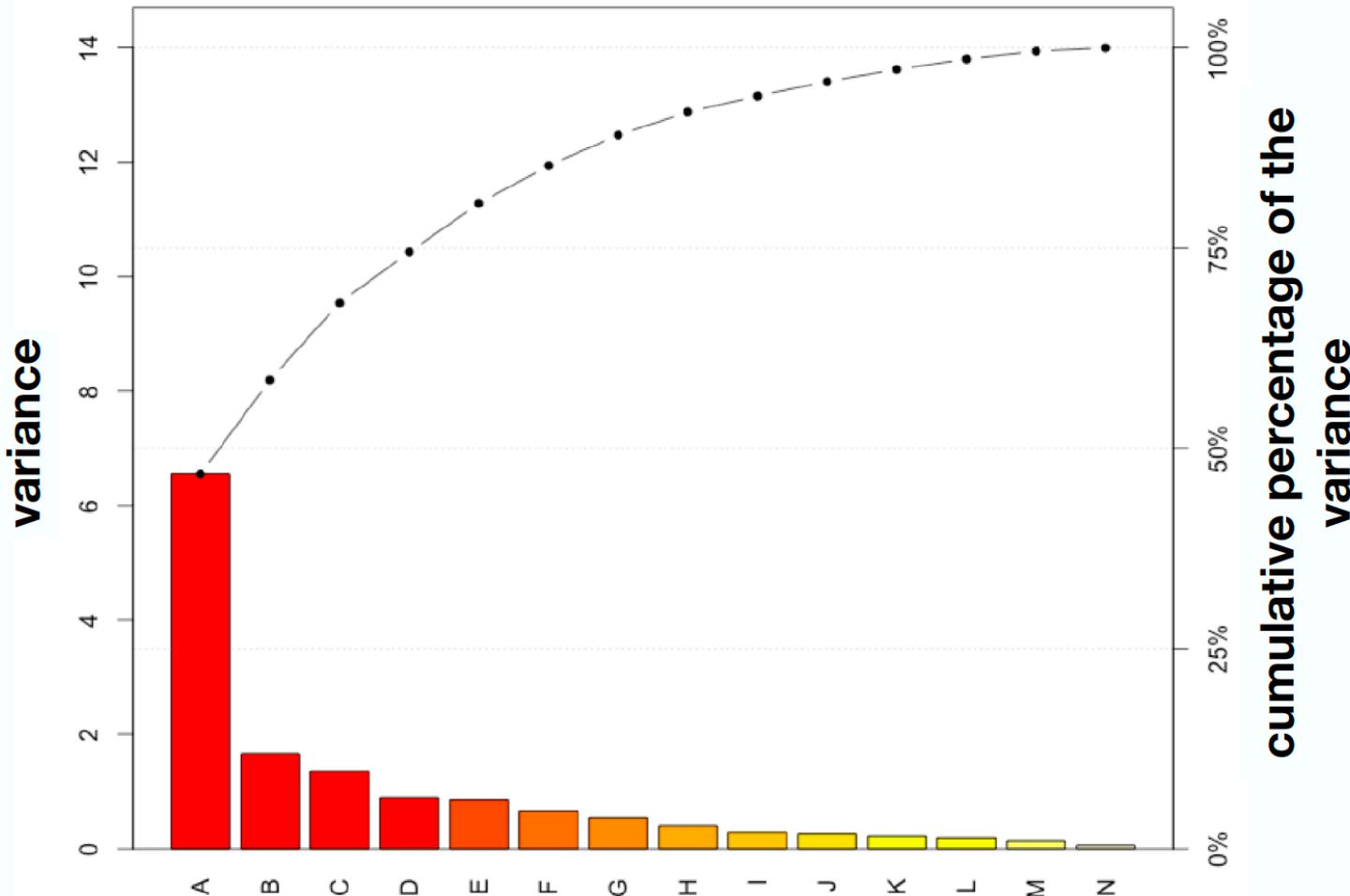
1. COMPUTE THE **COVARIANCE MATRIX** OF YOUR DATA
2. COMPUTE THE **EIGENVALUES AND EIGENVECTORS** OF THE **COVARAINCE MATRIX**
3. THE **EIGENVECTORS** PROVIDE THE DIRECTION OF MAXIMUM VARIANCE AND THE **EIGENVALUES** THE 'IMPORTANCE' OF THAT PARTICULAR FEATURE

**REMEMBER, THE COVARIANCE MATRIX IS NOTHING ELSE THAN THAT:**

$$K_{XX} = \begin{bmatrix} E[(X_1 - E[X_1])(X_1 - E[X_1])] & E[(X_1 - E[X_1])(X_2 - E[X_2])] & \cdots & E[(X_1 - E[X_1])(X_n - E[X_n])] \\ E[(X_2 - E[X_2])(X_1 - E[X_1])] & E[(X_2 - E[X_2])(X_2 - E[X_2])] & \cdots & E[(X_2 - E[X_2])(X_n - E[X_n])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - E[X_n])(X_1 - E[X_1])] & E[(X_n - E[X_n])(X_2 - E[X_2])] & \cdots & E[(X_n - E[X_n])(X_n - E[X_n])] \end{bmatrix}$$

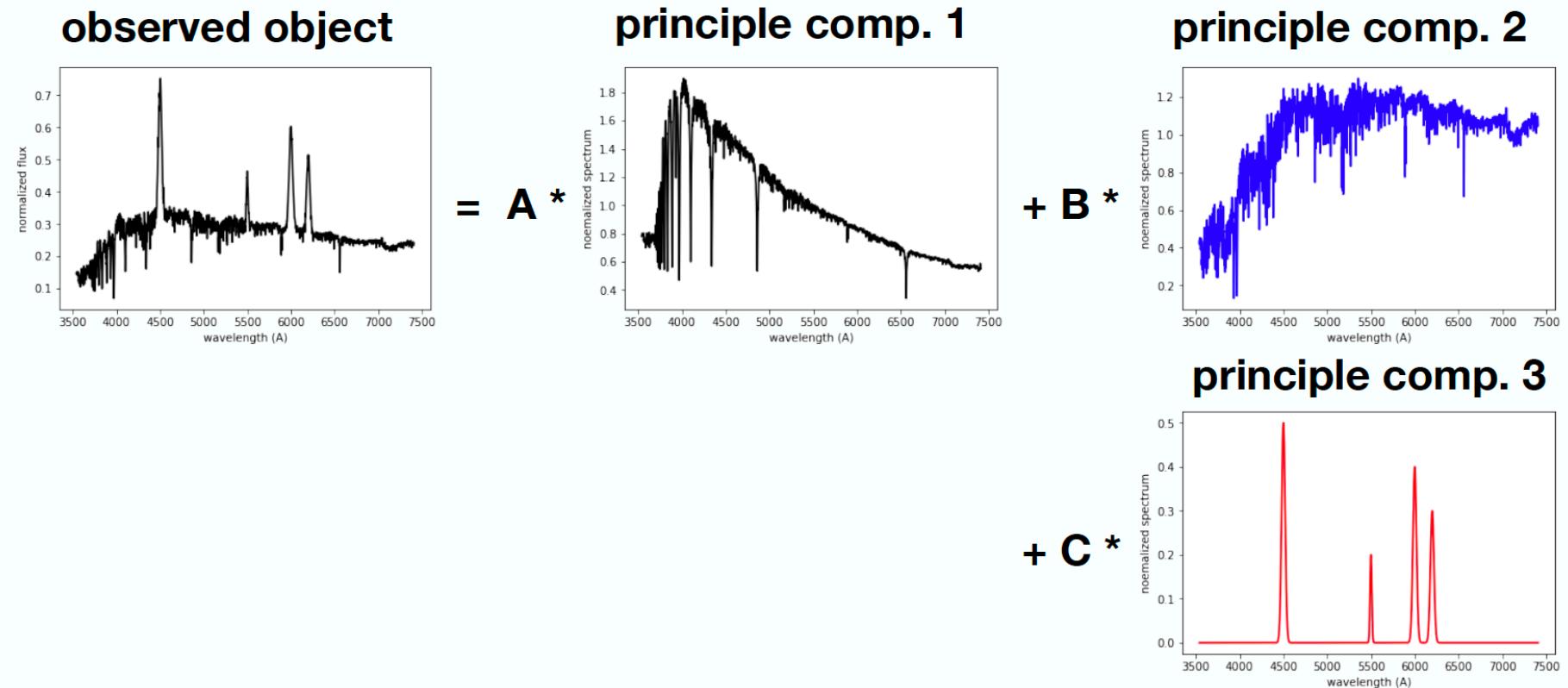


# PRINCIPAL COMPONENT ANALYSIS



IT RESULTS IN DATA COMPRESSION,  
BY REPRESENTING EACH OBJECT AS A PROJECTION OF THE FIRST PRINCIPLE COMPONENTS

# PRINCIPAL COMPONENT ANALYSIS



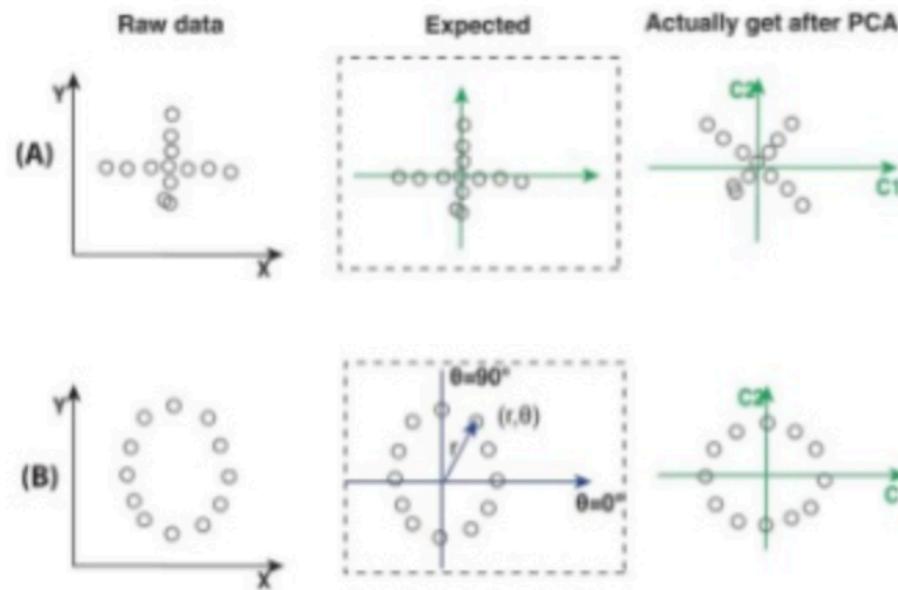
OBJECTS ARE DIVIDED INTO MAJOR PROTOTYPES AND ALL OBJECTS CAN BE OBTAINED AS LINEAR COMBINATIONS

@D. BARON

# LIMITATIONS OF PCA

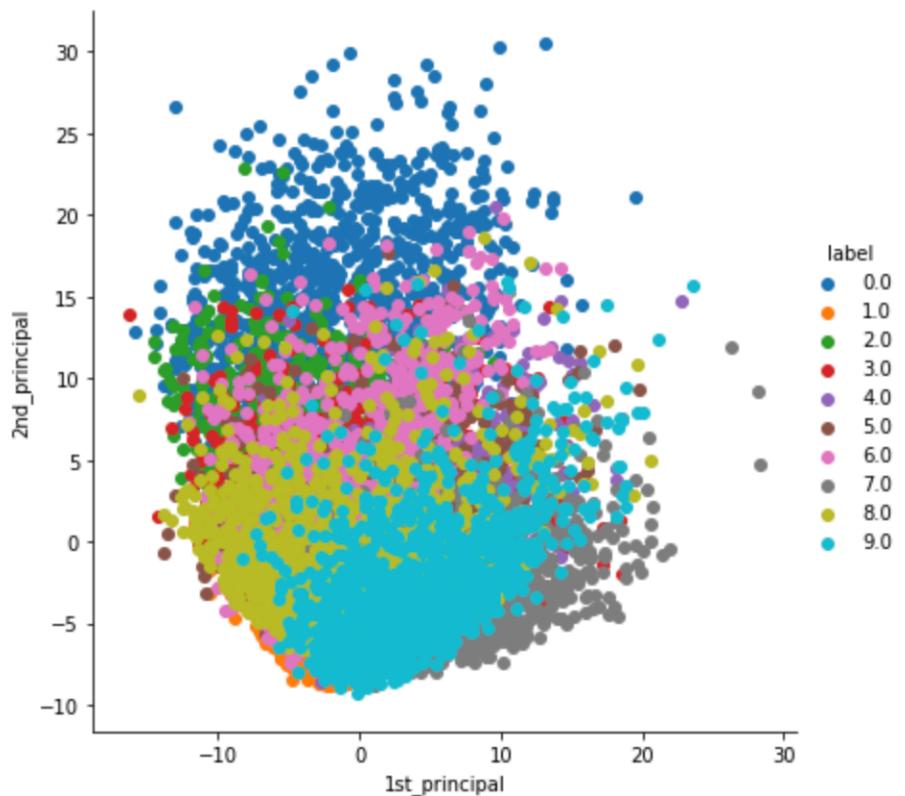
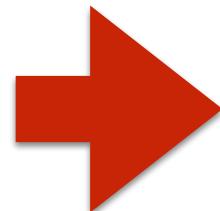
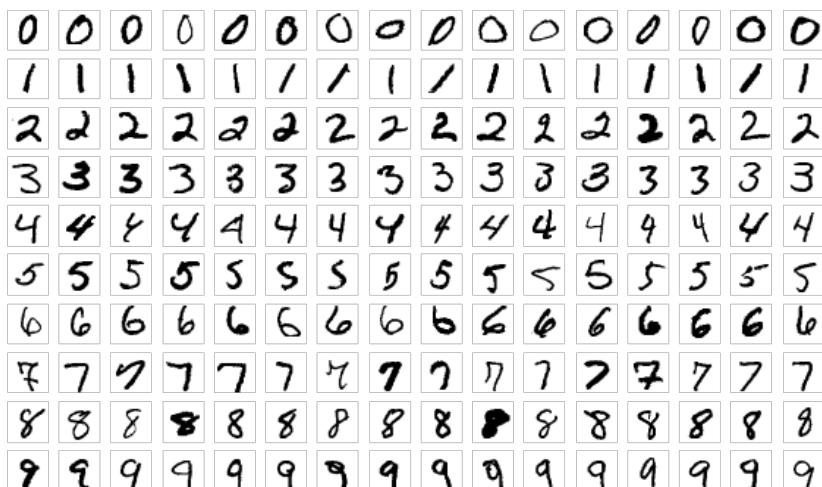
PCA APPLY LINEAR TRANSFORMATIONS

SINCE WE USE THE COVARIANCE MATRIX, IT ASSUMES THAT THE DATA FOLLOWS A **MULTIDIMENSIONAL NORMAL DISTRIBUTION**



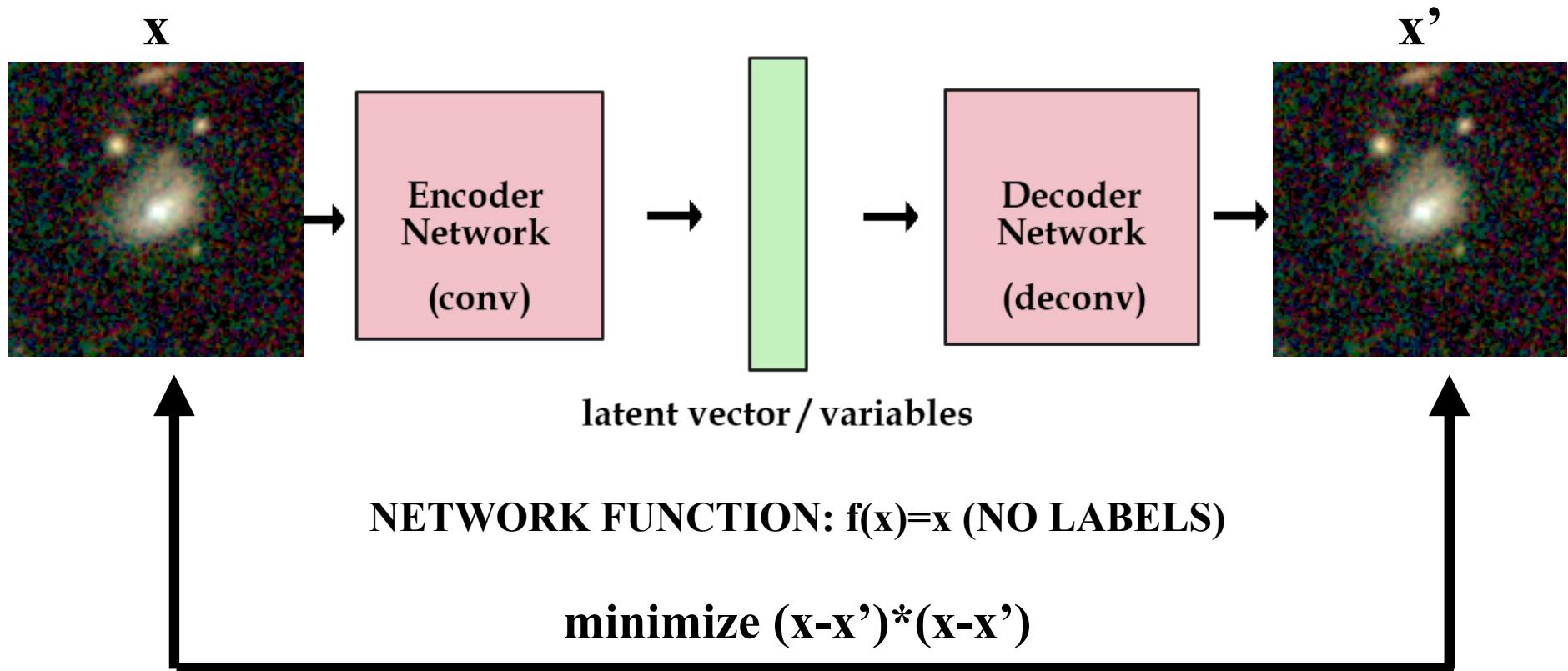
# LIMITATIONS OF PCA

AND DATA IS NOT ALWAYS GAUSSIAN....



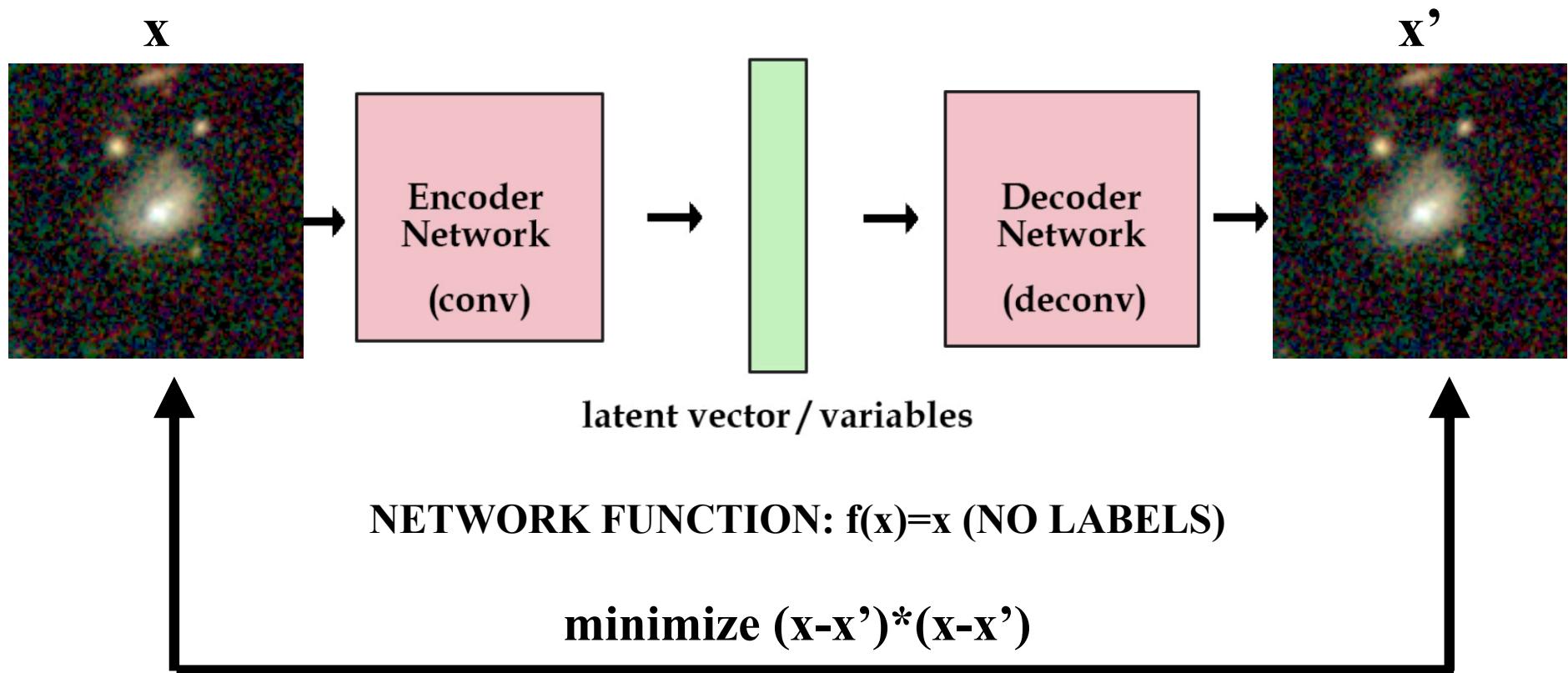
**CAN WE GENERALIZE THAT?**

# CONVOLUTIONAL AUTO-ENCODER



AN AUTO-ENCODER IS ANY NETWORK WITH IDENTICAL INPUT AND OUTPUT

# CONVOLUTIONAL AUTO-ENCODER



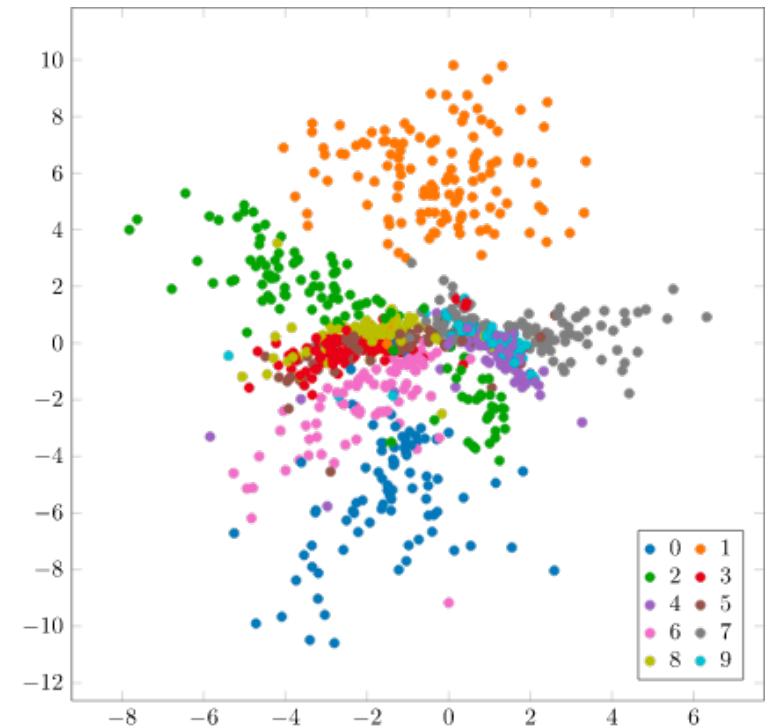
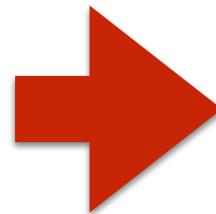
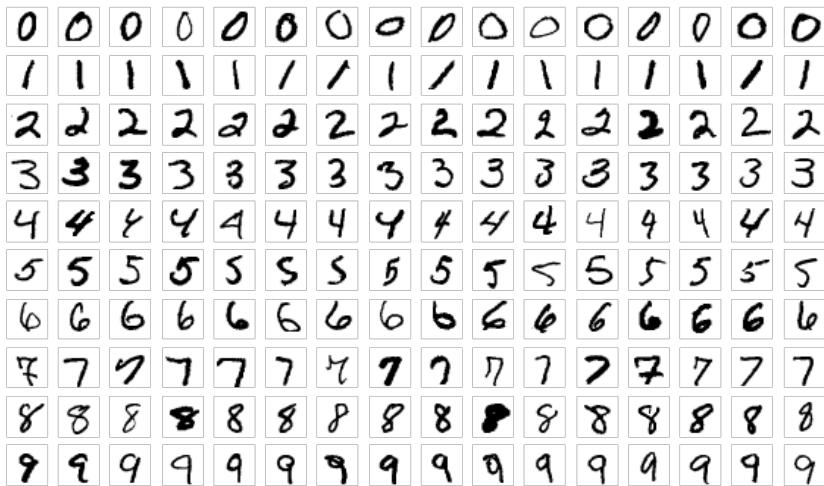
BY REDUCING THE DIMENSIONALITY IN THE LATENT SPACE WE FORCE THE NETWORK TO LEARN A REPRESENTATION OF THE INPUT DATA IN A LOWER DIMENSIONALITY SPACE

\* NO NEED TO BE CONVOLUTIONAL - ANY NEURAL NETWORK WITH A BOTTLENECK WILL DO THE JOB

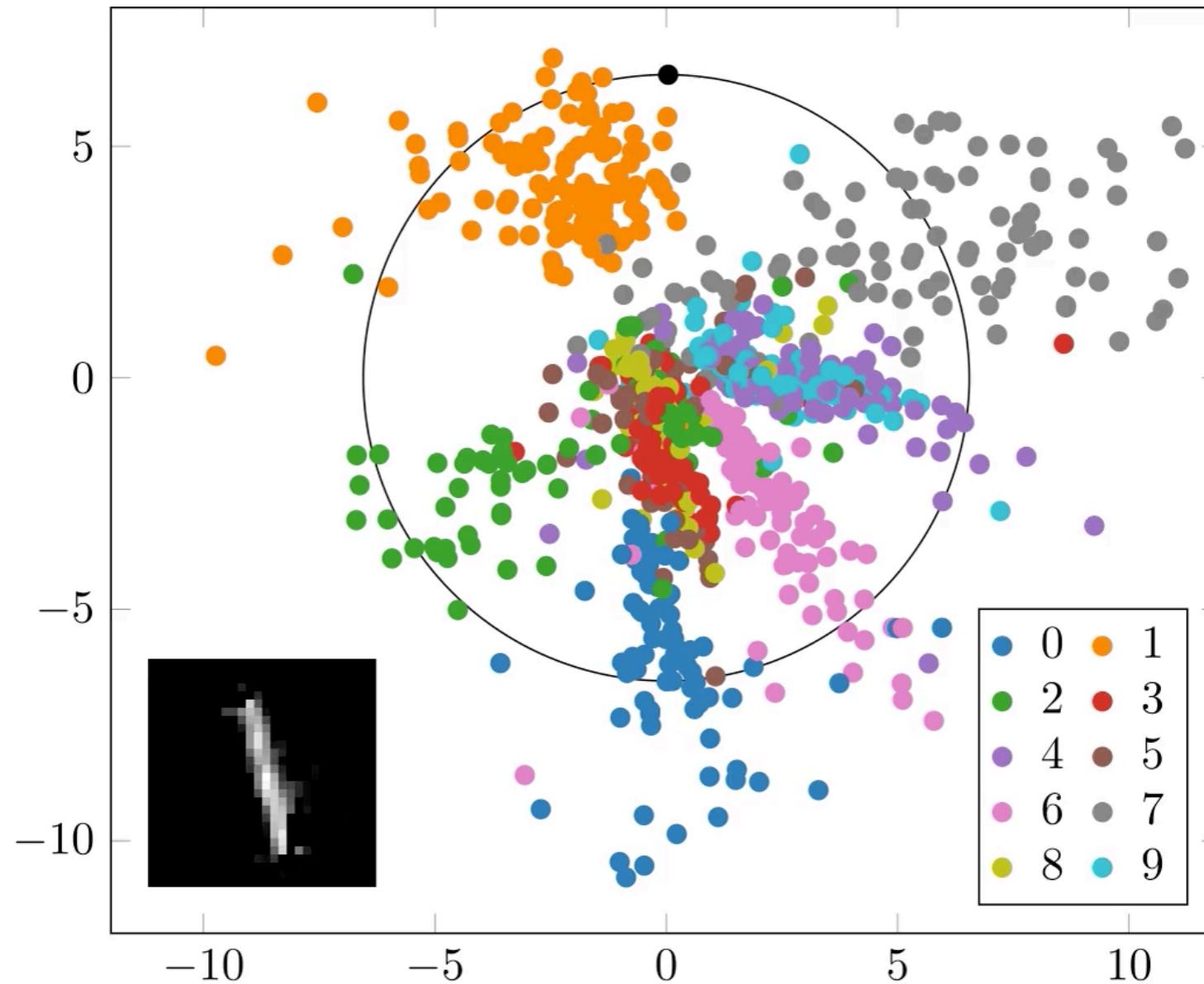
\* **QUESTION:** WHAT WOULD HAPPEN IF WE SET AN AUTOENCODER WITH NO ACTIVATION FUNCTIONS?

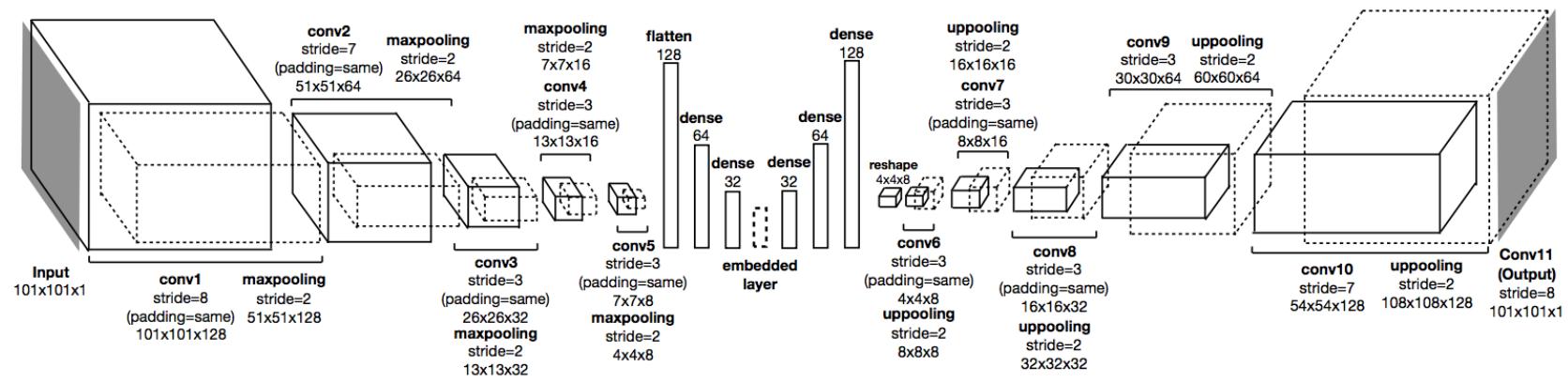
SEE EXAMPLE FROM TUTORIALS FOR STAR-GALAXY SEPARATION

# AUTOENCODER REPRESENTATION OF MNIST

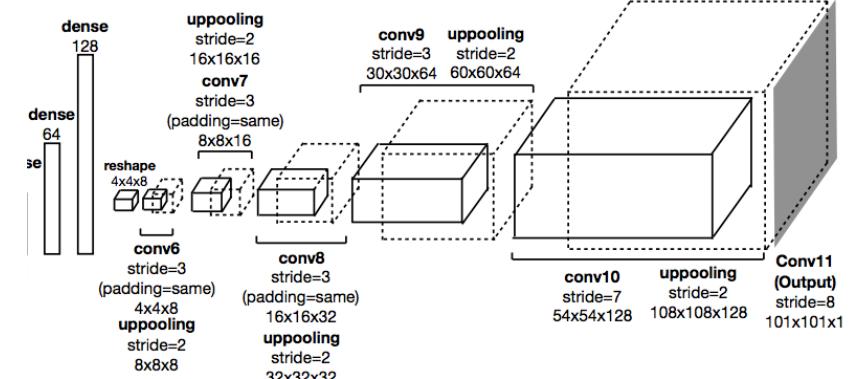
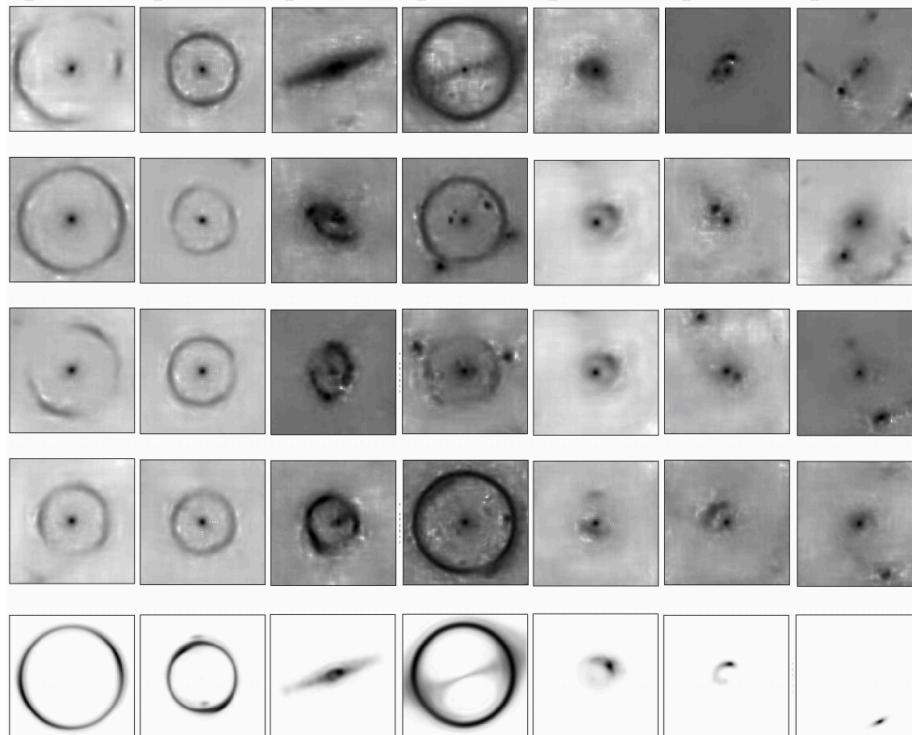


# AUTOENCODER REPRESENTATION OF MNIST

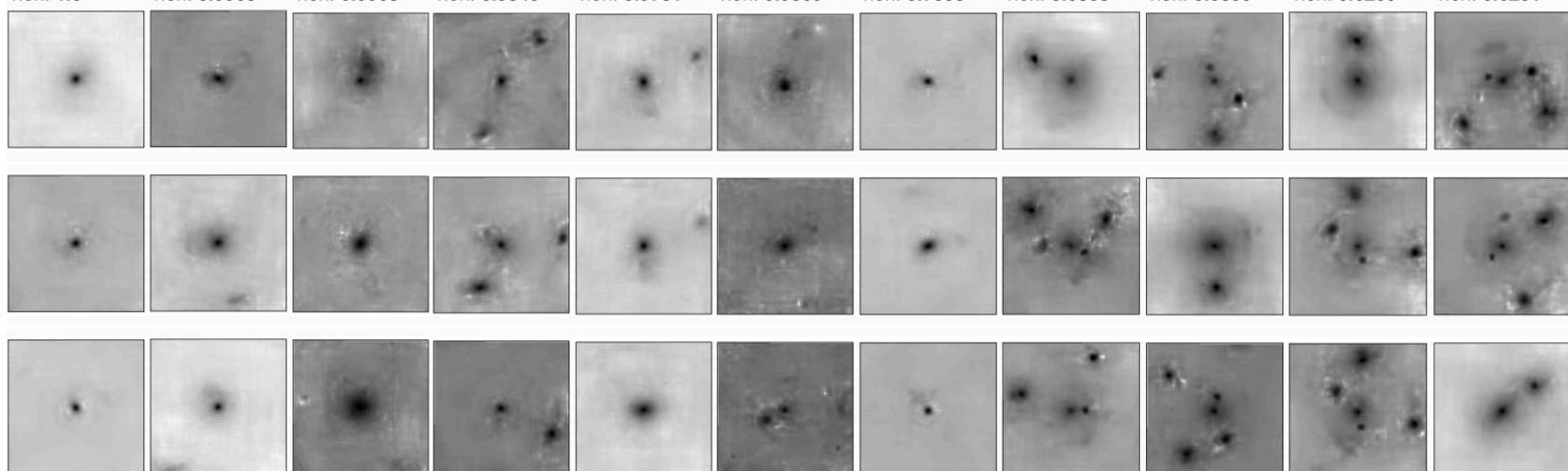




|   |   |  |  |  |   |  |
|---|---|--|--|--|---|--|
| <b>Cluster 17:</b><br>lensing: 0.9873<br>non: 0.0127<br>F_len: 0.0914 | <b>Cluster 21:</b><br>lensing: 0.9448<br>non: 0.0552<br>F_len: 0.0695 | <b>Cluster 1:</b><br>lensing: 0.9159<br>non: 0.0841<br>F_len: 0.0731 | <b>Cluster 6:</b><br>lensing: 0.8997<br>non: 0.1003<br>F_len: 0.0729 | <b>Cluster 2:</b><br>lensing: 0.803<br>non: 0.197<br>F_len: 0.1945 | <b>Cluster 20:</b><br>lensing: 0.6206<br>non: 0.3794<br>F_len: 0.0575 | <b>Cluster 5:</b><br>lensing: 0.6170<br>non: 0.3830<br>F_len: 0.0734 |
|---|---|--|--|--|---|--|



|  |   |   |  |  |   |   |  |  |  |   |
|--|---|---|--|--|---|---|--|--|--|---|
| <b>Cluster 14:</b><br>lensing: 0.0<br>non: 1.0 | <b>Cluster 3:</b><br>lensing: 0.0037<br>non: 0.9963 | <b>Cluster 8:</b><br>lensing: 0.0037<br>non: 0.9963 | <b>Cluster 22:</b><br>lensing: 0.0154<br>non: 0.9846 | <b>Cluster 16:</b><br>lensing: 0.0219<br>non: 0.9781 | <b>Cluster 4:</b><br>lensing: 0.0431<br>non: 0.9569 | <b>Cluster 0:</b><br>lensing: 0.2642<br>non: 0.7358 | <b>Cluster 18:</b><br>lensing: 0.3132<br>non: 0.6868 | <b>Cluster 19:</b><br>lensing: 0.3601<br>non: 0.6399 | <b>Cluster 13:</b><br>lensing: 0.3731<br>non: 0.6269 | <b>Cluster 9:</b><br>lensing: 0.3769<br>non: 0.6231 |
|--|---|---|--|--|---|---|--|--|--|---|



Cheng+20

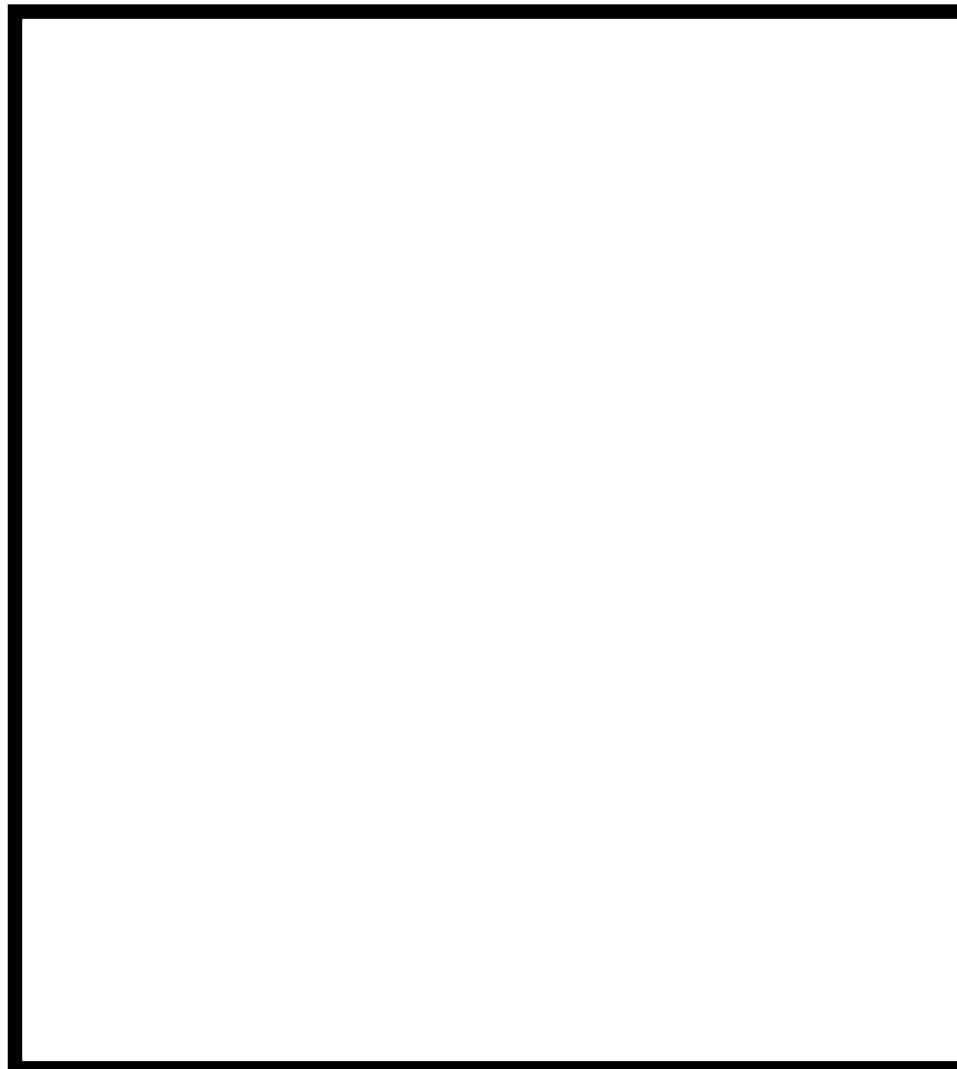
# THIS IS OF COURSE NOT THE ONLY WAY TO OBTAIN USEFUL REPRESENTATIONS OF YOUR DATA....

IT IS AN ACTIVE FIELD OF RESEARCH BECAUSE LABELLING IS ONE OF THE MAIN BOTTLENECKS OF MACHINE LEARNING

- ▶ “Pure” Reinforcement Learning (**cherry**)
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ A few bits for some samples
- ▶ Supervised Learning (**icing**)
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ 10→10,000 bits per sample
- ▶ Self-Supervised Learning (**cake génoise**)
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ Millions of bits per sample

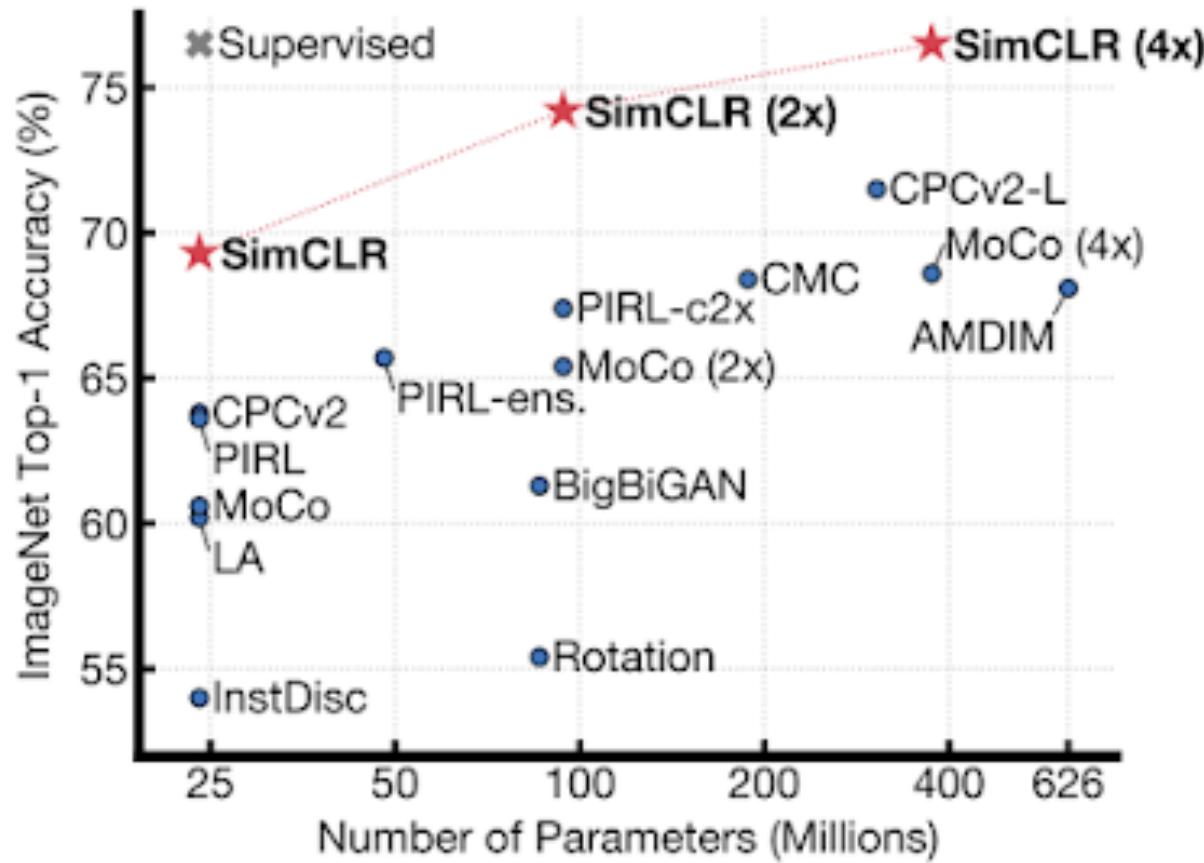


# **LEARNING REPRESENTATIONS THROUGH SELF-SUPERVISED LEARNING**



Chen+20

# SELF-SUPERVISED LEARNING REACHES COMPARABLE ACCURACY TO FULLY SUPERVISED APPROACHES...



# **UMAP AND T-SNE ARE ALSO TWO USEFUL TECHNIQUES TO VISUALIZE LARGE DIMENSIONAL SPACES**

INSTEAD OF MATRIX FACTORIZATION SUCH AS PCA, THEY  
ARE BASED ON GRAPH LAYOUT

1. WE BUILD A GRAPH OF OUR HIGH-DIMENSIONALITY DATA
2. WE OPTIMIZE TO GET THE MOST SIMILAR GRAPH IN TWO DIMENSIONS

[https://  
lvdmaaten.github.io/tsne/](https://lvdmaaten.github.io/tsne/)

[https://umap-learn.readthedocs.io/en/latest/how\\_umap\\_works.html](https://umap-learn.readthedocs.io/en/latest/how_umap_works.html)

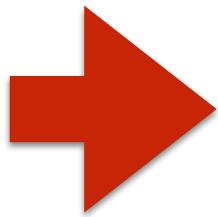
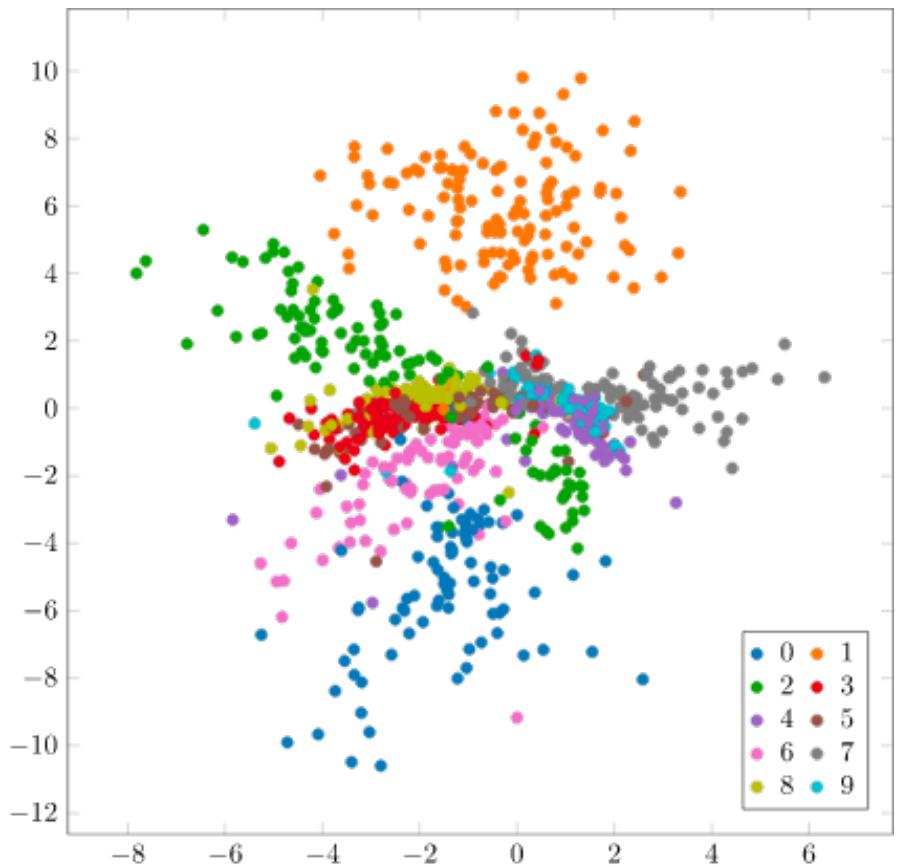
# **UMAP EXAMPLES AND EXPLANATIONS;**

<https://pair-code.github.io/understanding-umap/>  
#:~:text=In%20the%20simplest%20sense%2C%20UMAP,behind%20them%  
20is%20remarkably%20simple.

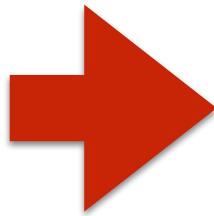
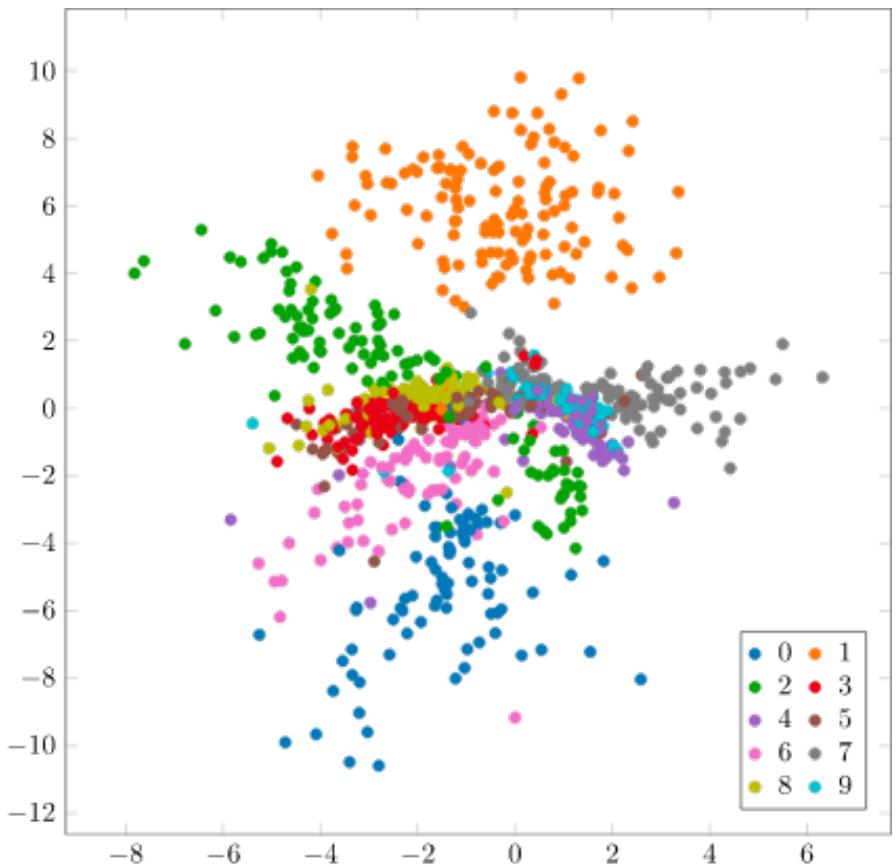
HOW MIGHT YOU SOLVE THESE  
RELATED PROBLEMS ONCE YOU HAVE  
THE LATENT SPACE COORDINATES FOR  
YOUR TRAINING SAMPLE?

GENERATE A RANDOM SAMPLE DRAWN FROM THE INPUT  
DISTRIBUTION ("**GENERATIVE MODEL**")

ESTIMATE THE PROBABILITY DENSITY OF AN ARBITRARY  
INPUT, RELATIVE TO THE INPUT DISTRIBUTION  
 ("**PROBABILISTIC MODEL**")



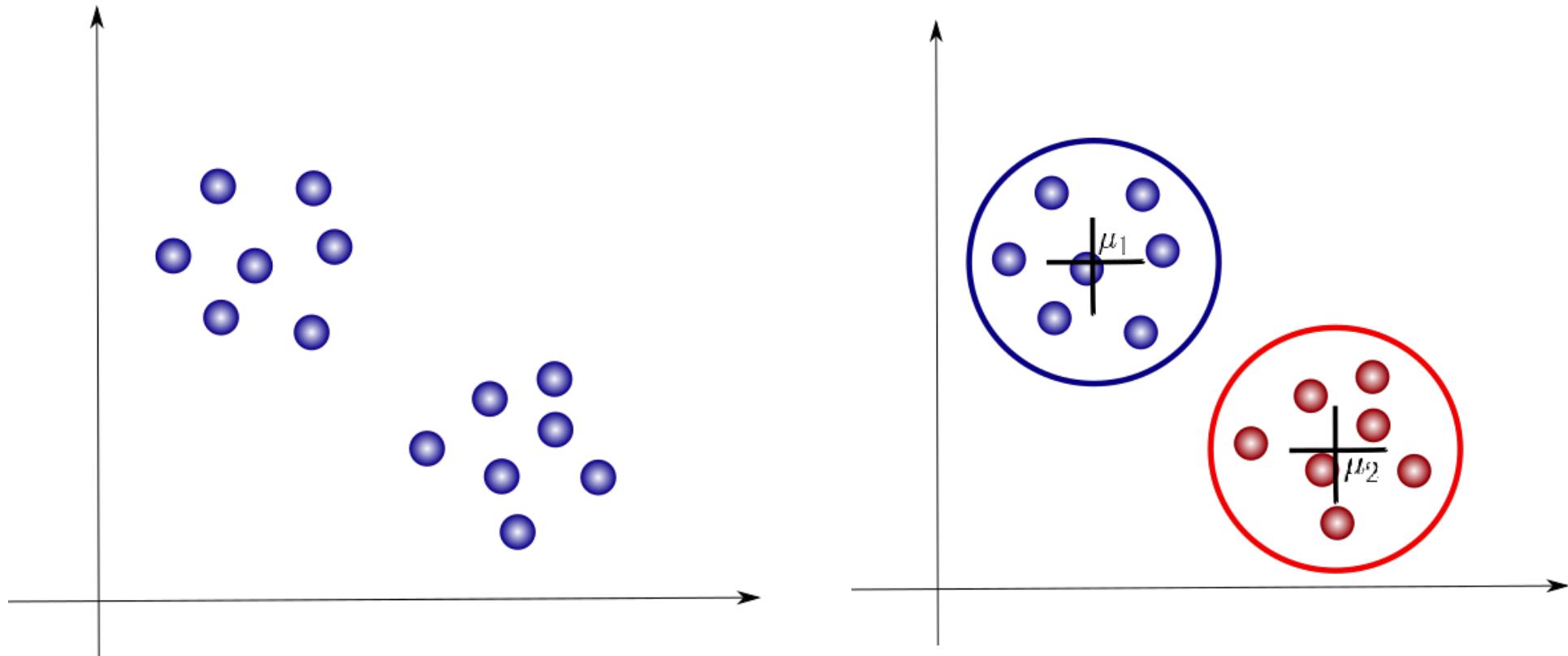
**HOW CAN I ESTIMATE  
 $P(X)$ ?**



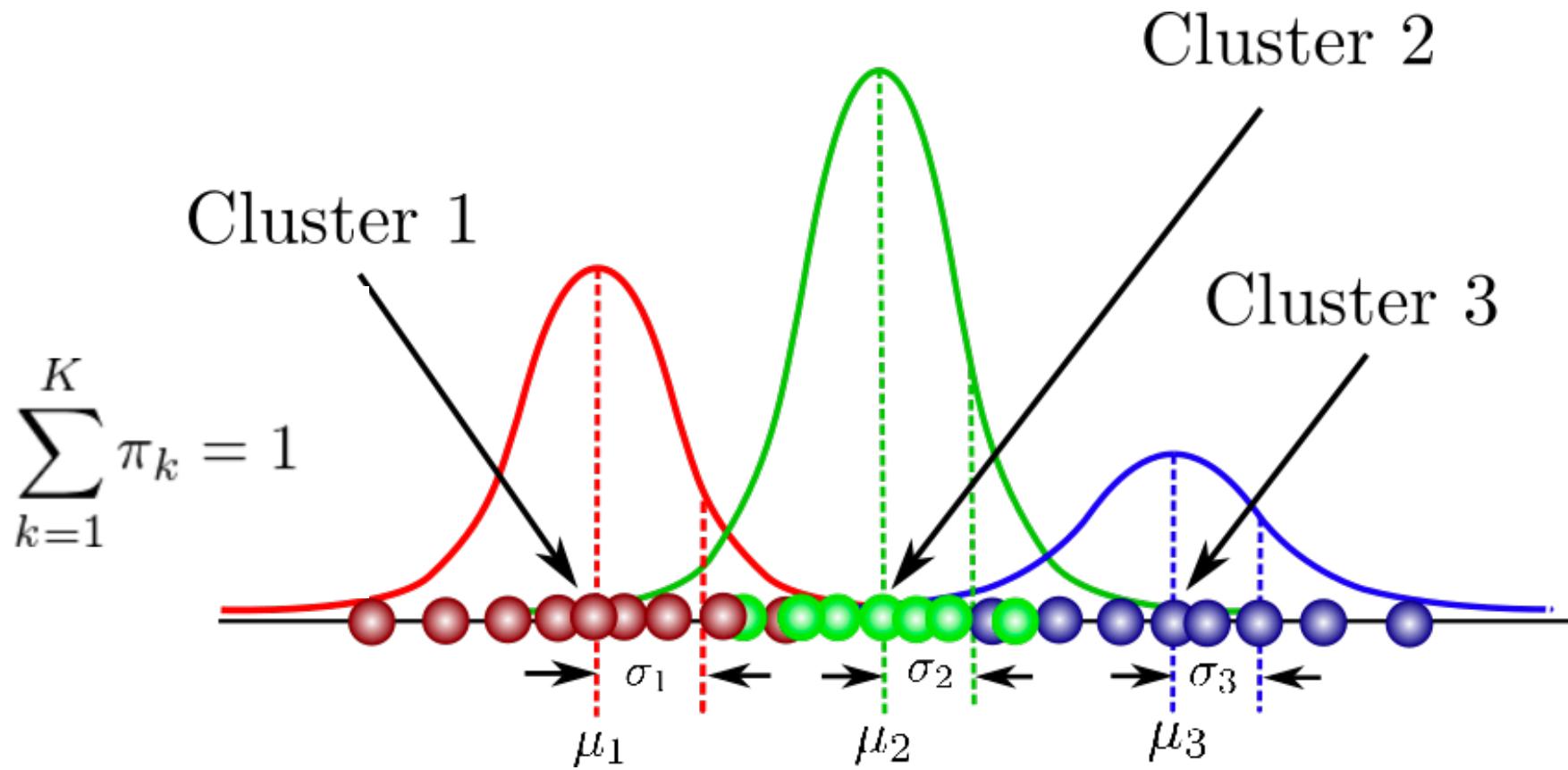
**HOW CAN I ESTIMATE  
 $P(X)$ ?**

**WHEN YOU DON'T KNOW, ASSUME IT IS GAUSSIAN....**

# GAUSSIAN MIXTURE MODELS (GMMs) ARE DENSITY ESTIMATOR METHODS THAT FIT MULTIPLE GAUSSIANS TO THE REPRESENTATION

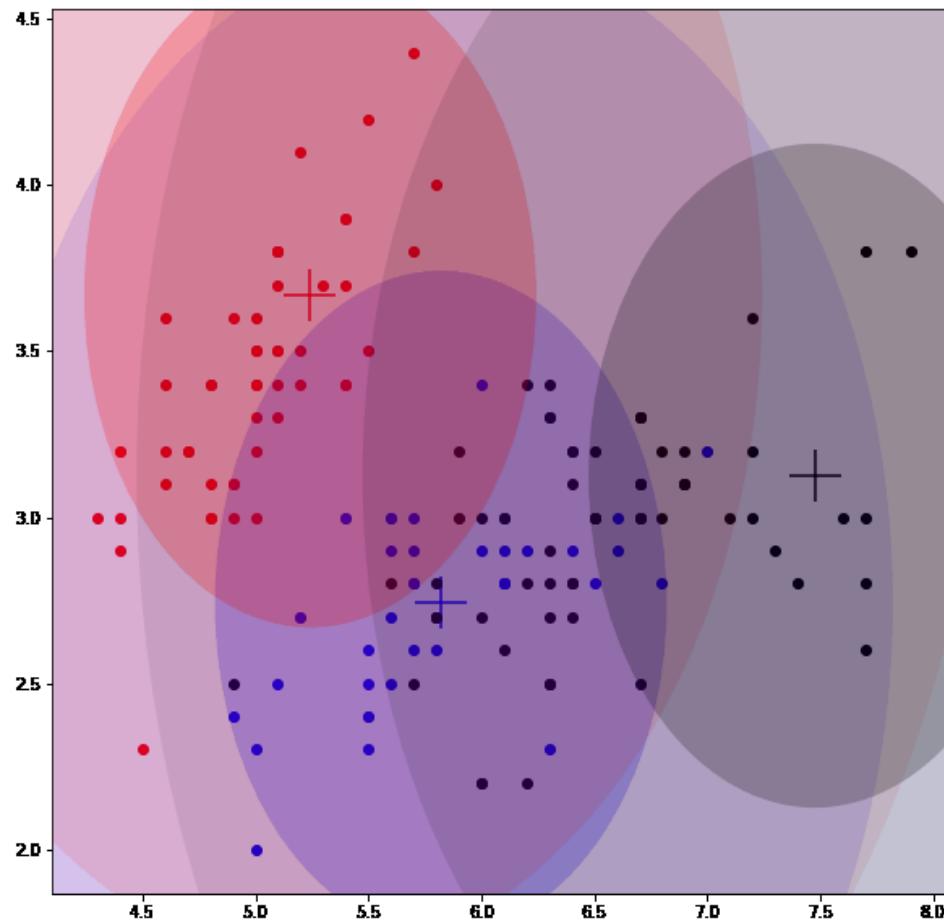


# GAUSSIAN MIXTURE MODELS (GMMs) ARE DENSITY ESTIMATOR METHODS THAT FIT MULTIPLE GAUSSIANS TO THE REPRESENTATION

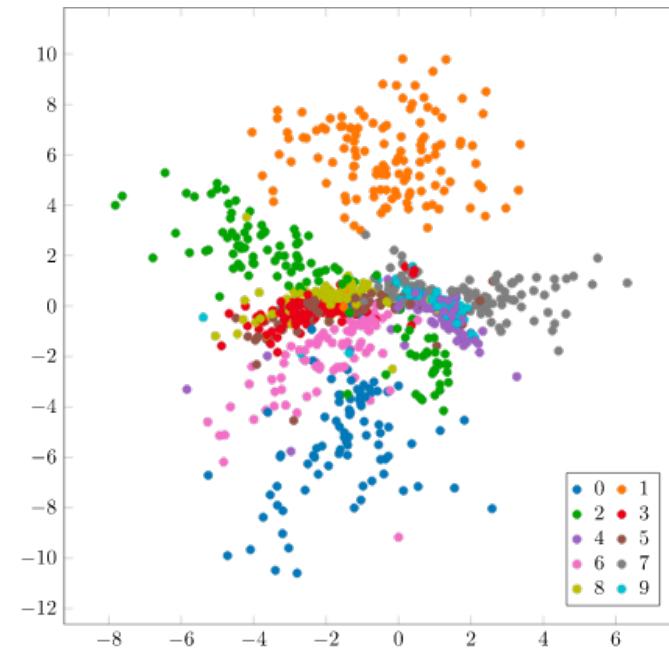
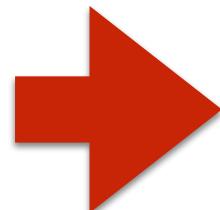
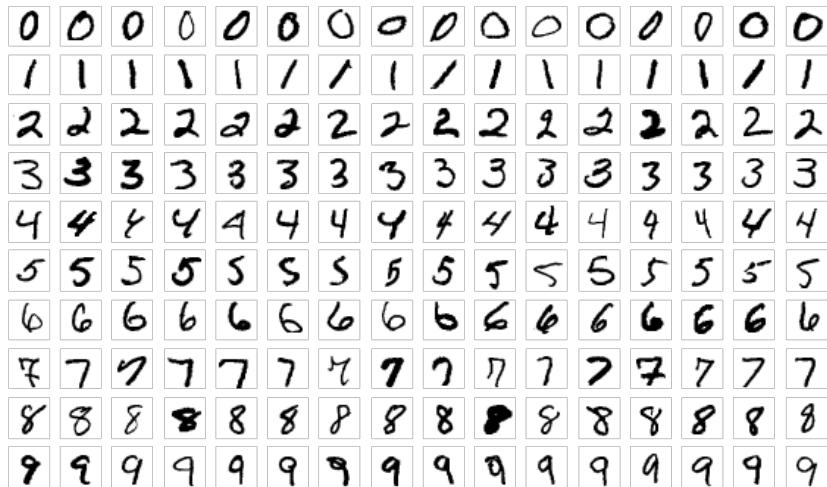


**means, sigmas and scale factors of each gaussian are free parameters**

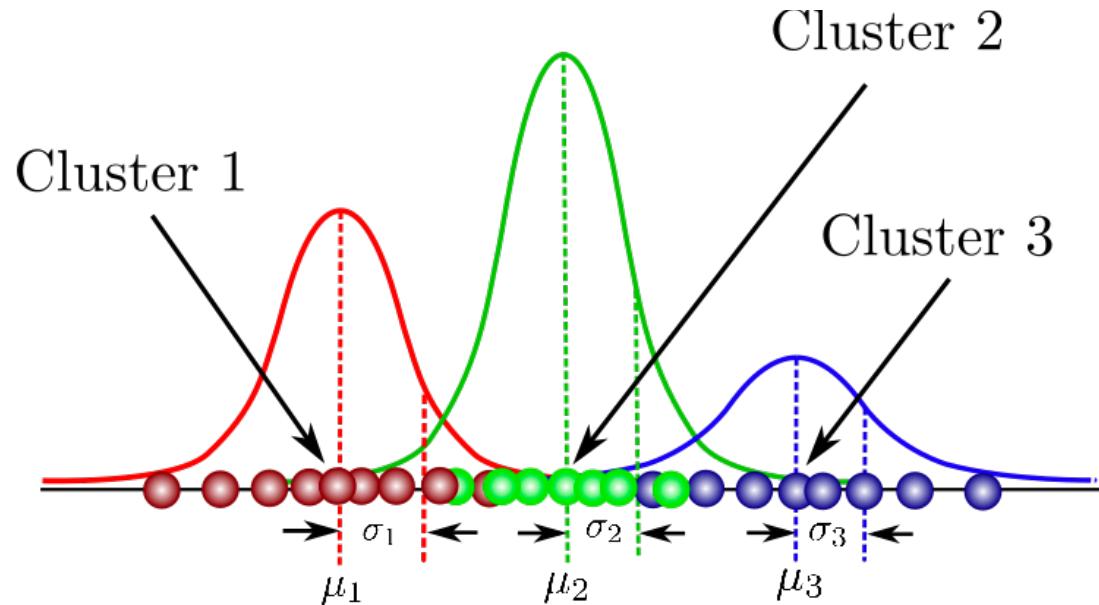
# GAUSSIAN MIXTURE MODELS (GMMs) ARE DENSITY ESTIMATOR METHODS THAT FIT MULTIPLE GAUSSIANS TO THE REPRESENTATION



# DATA



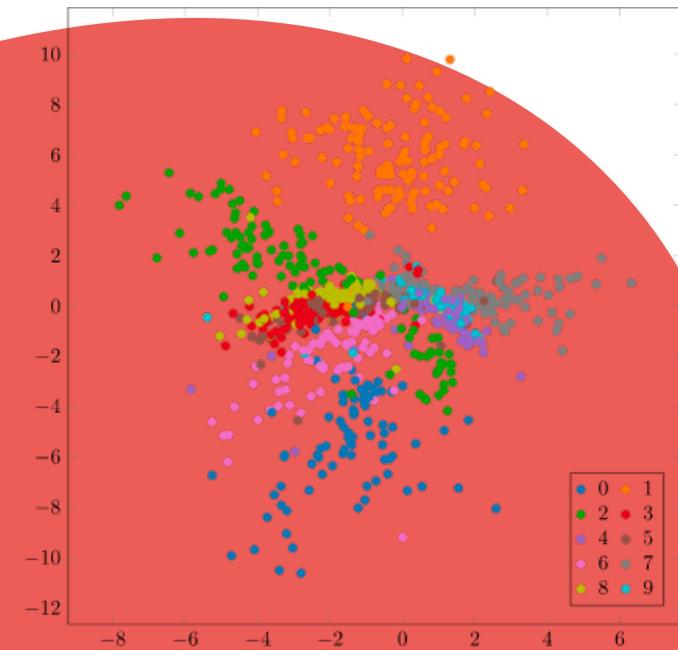
REPRESENTATION  
(AUTOENCODER)



MODELING OF  $P(\mathbf{x})$  WITH GMMs

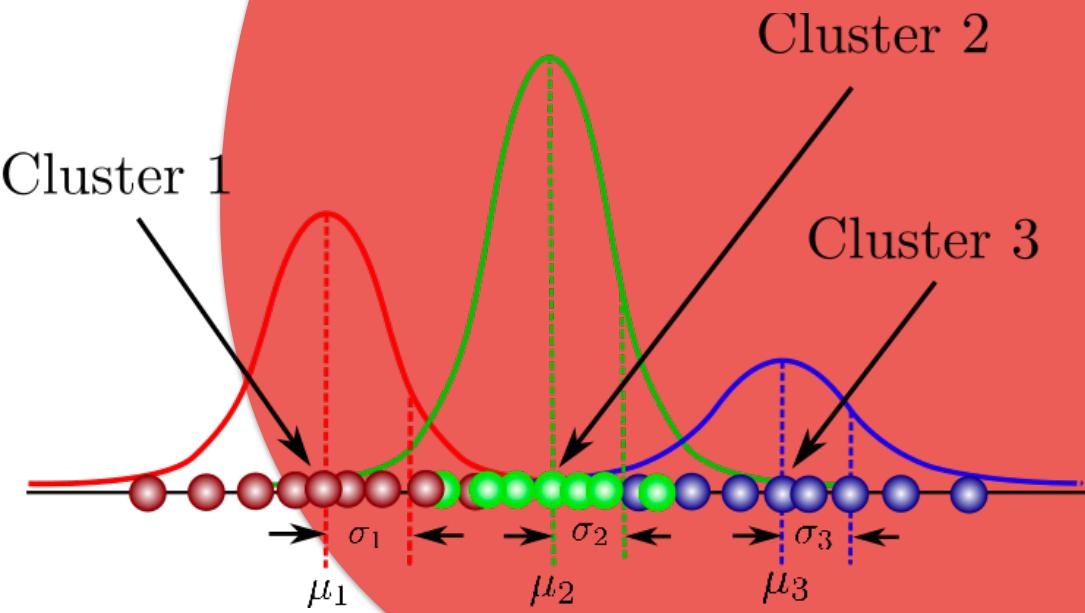
# DATA

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 |
| 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 |
| 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 |



# REPRESENTATION (AUTOENCODER)

CAN WE COMBINE THESE 2 STEPS?



**MODELING OF  $P(X)$  WITH GMMs**