

# Guided Tour of Machine Learning in Finance

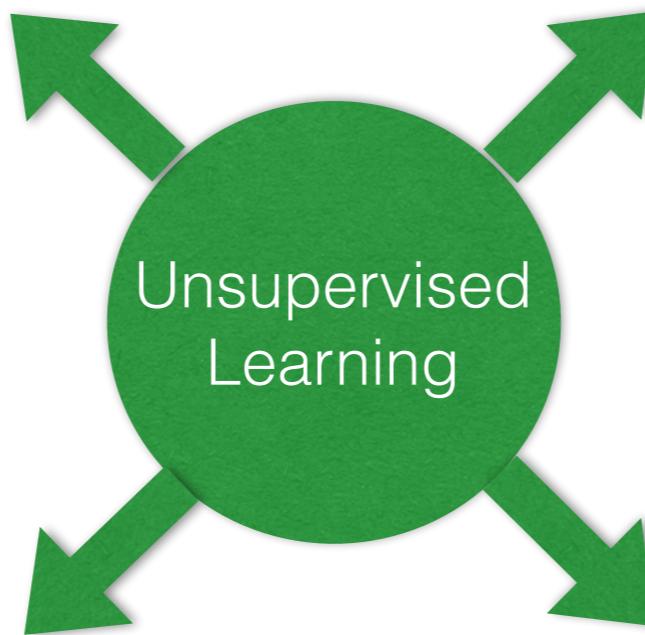
**Week 3: Unsupervised Learning**

**Core concept of UL**

Igor Halperin

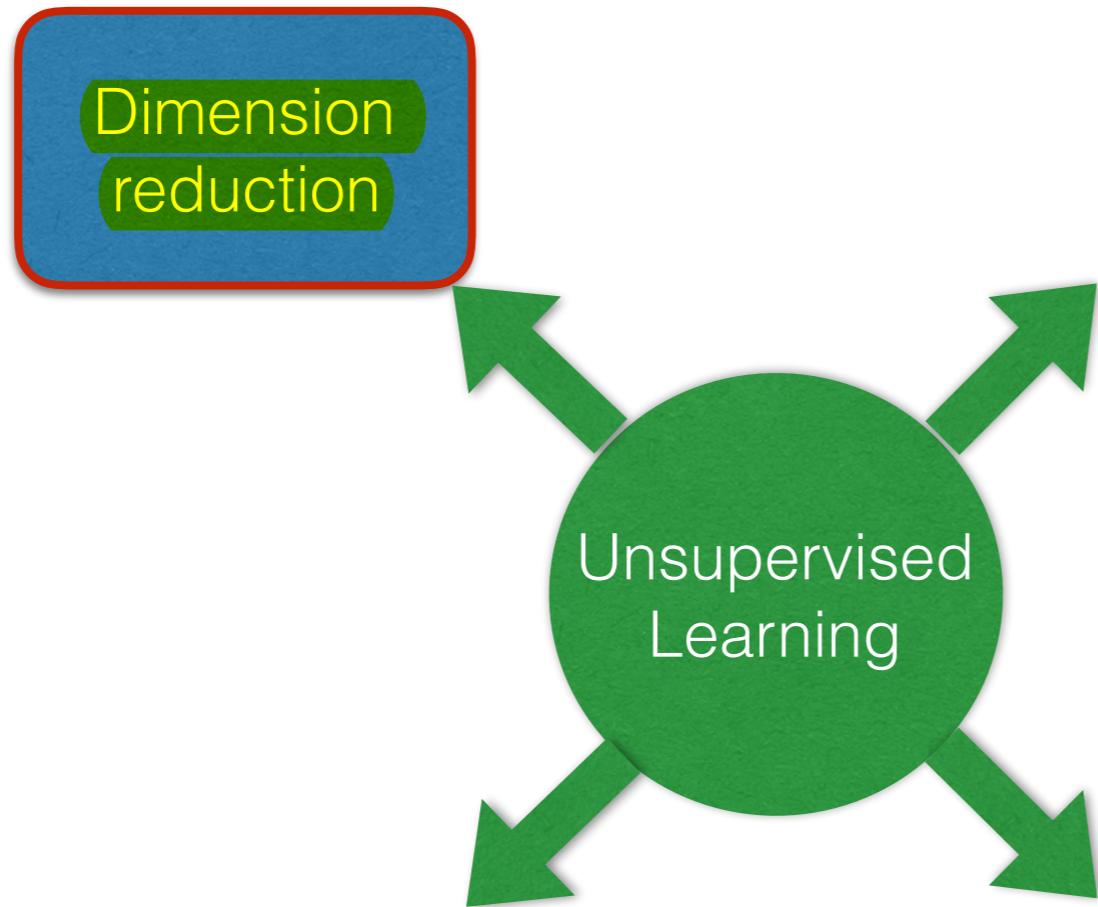
NYU Tandon School of Engineering, 2017

# Unsupervised Learning



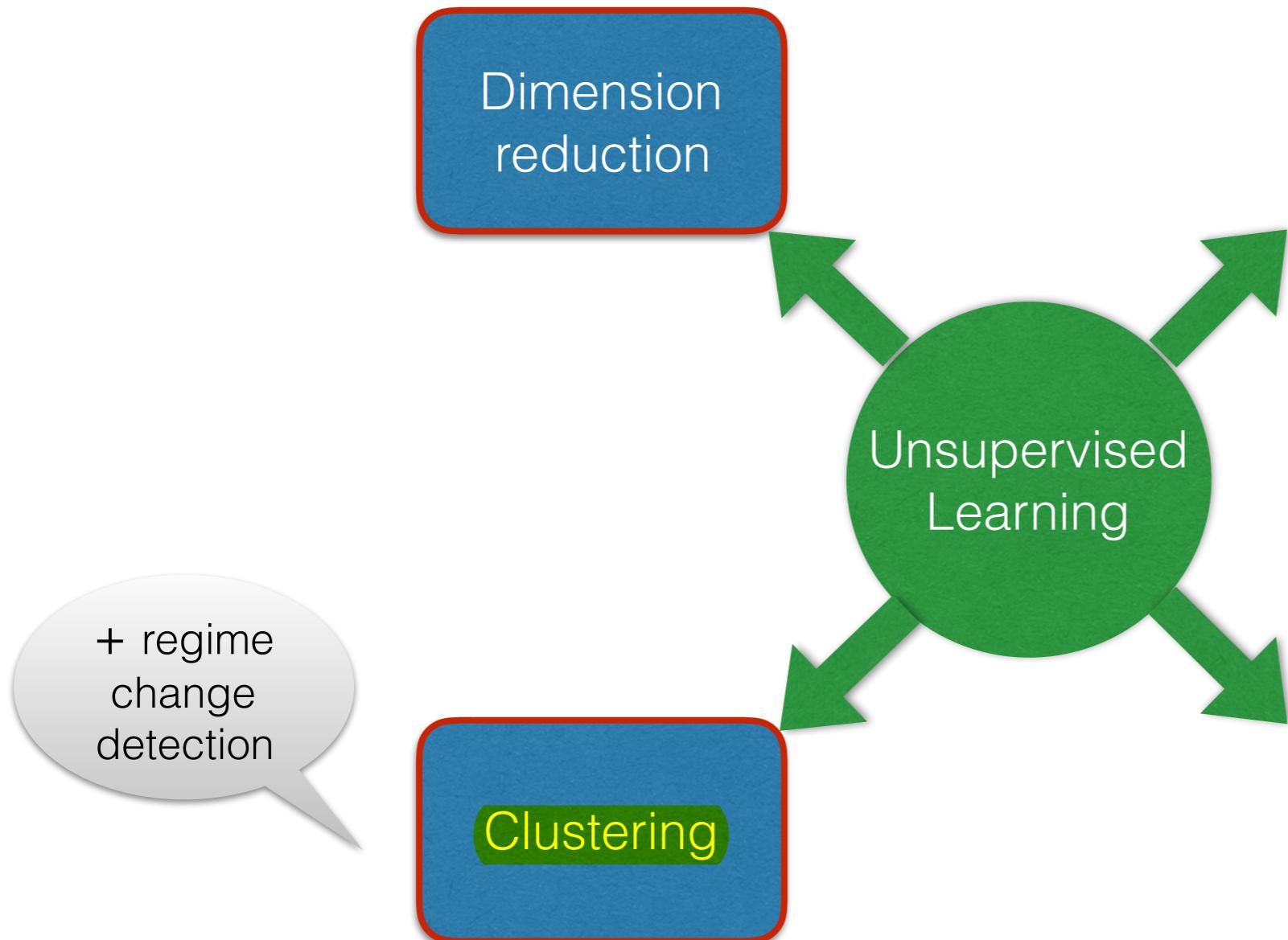
Learning = ability to generalize  
(for UL) = ability to learn the data representation

# Unsupervised Learning



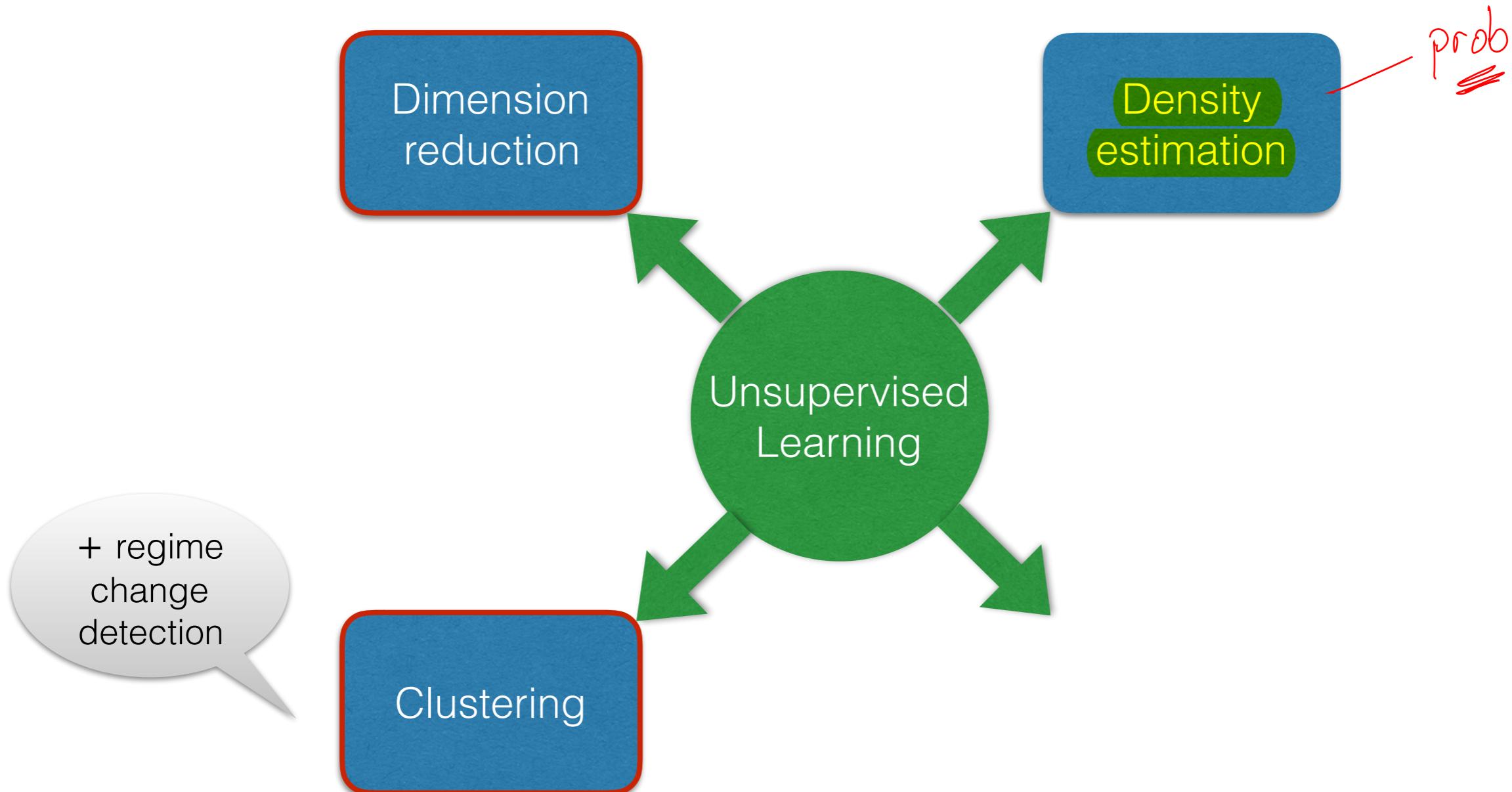
**Learning = ability to generalize**  
**(for UL) = ability to learn the data representation**

# Unsupervised Learning



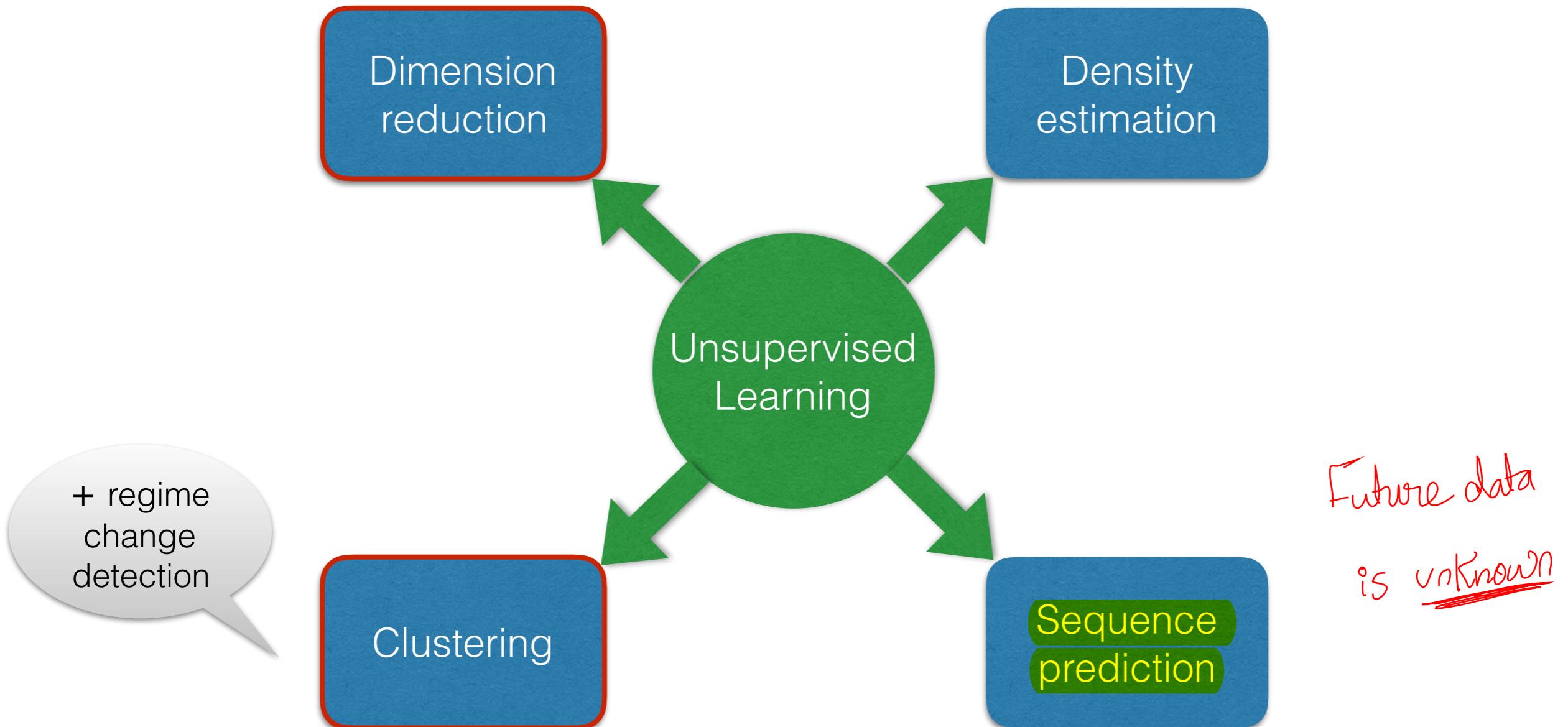
**Learning = ability to generalize**  
**(for UL) = ability to learn the data representation**

# Unsupervised Learning



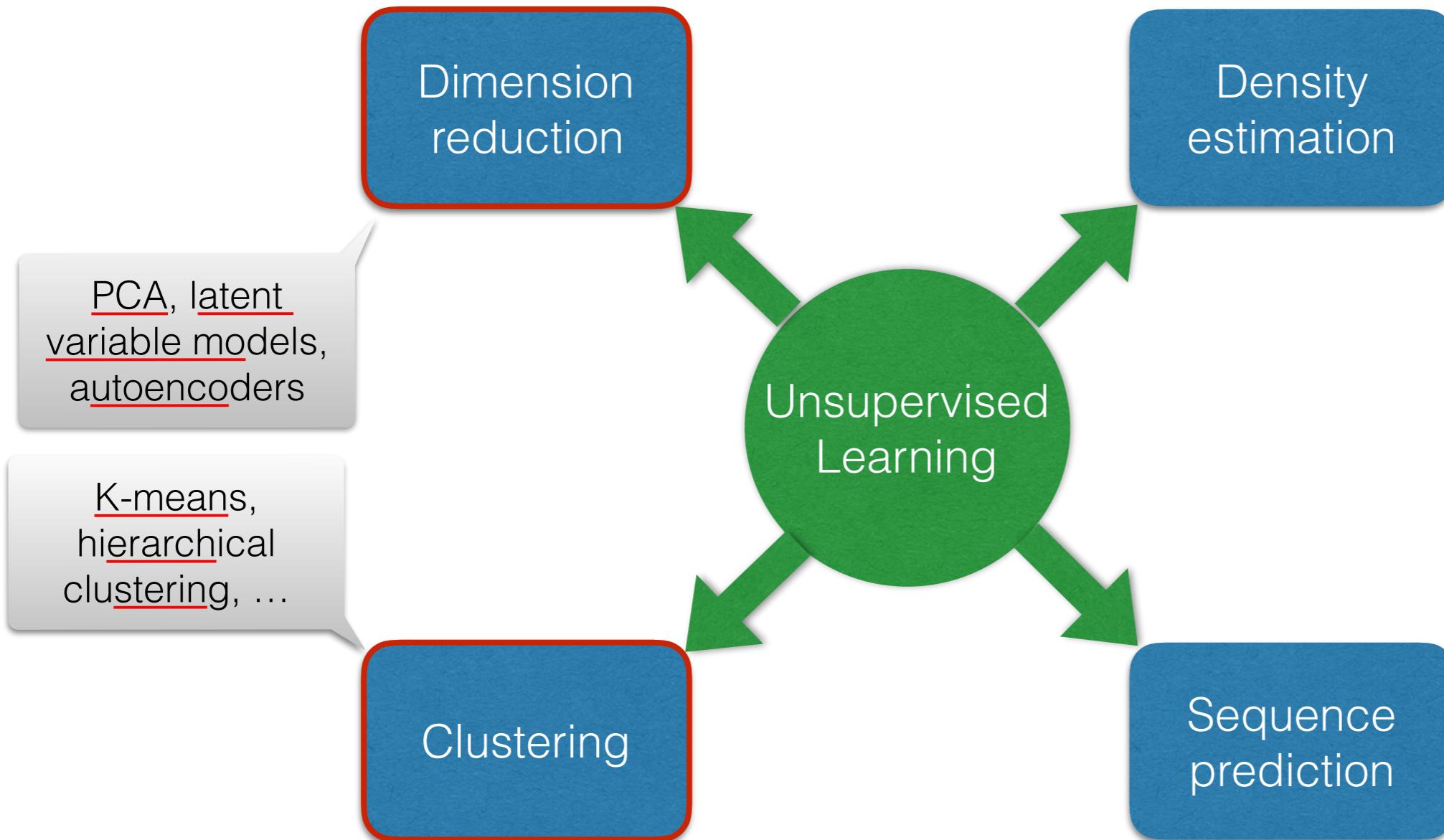
**Learning = ability to generalize**  
**(for UL) = ability to learn the data representation**

# Unsupervised Learning



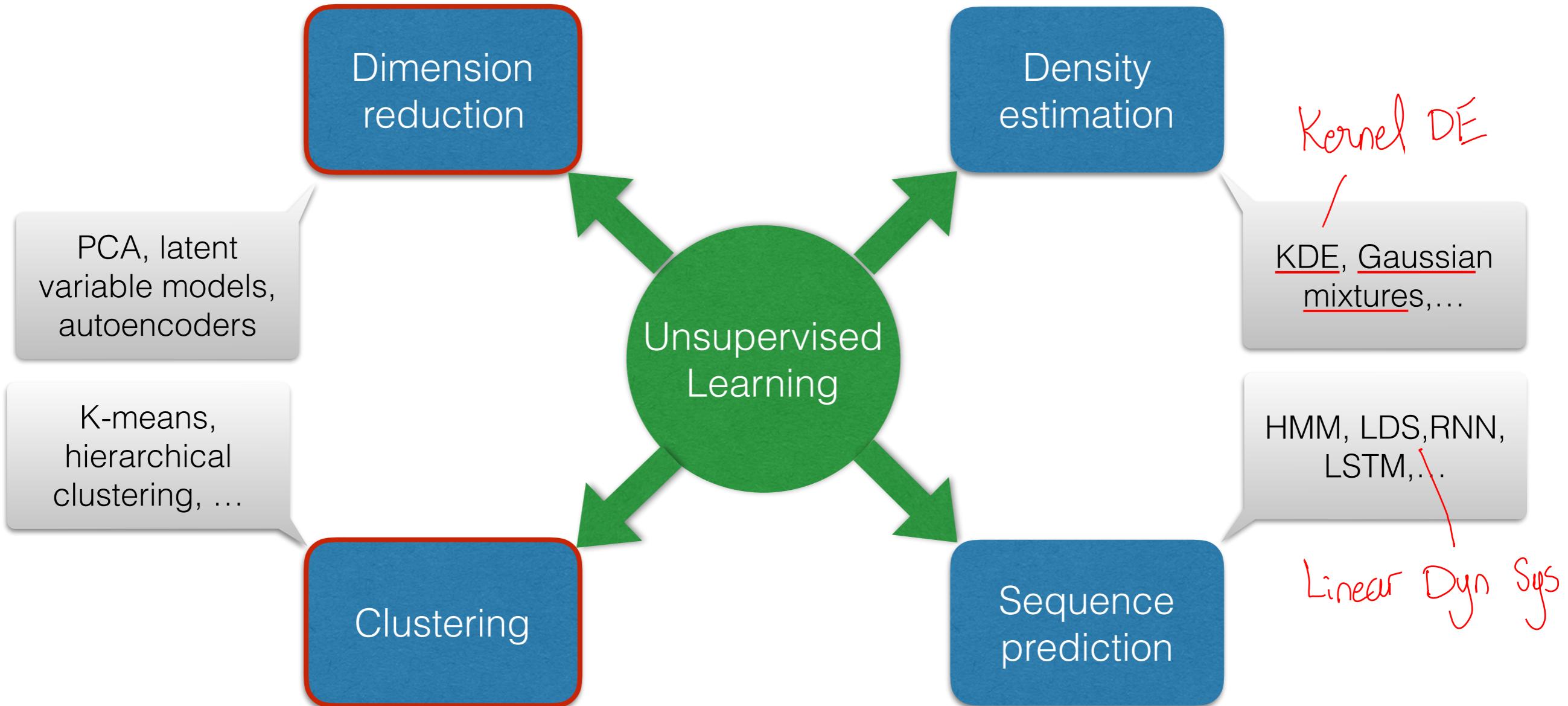
**Learning = ability to generalize**  
**(for UL) = ability to learn the data representation**

# Unsupervised Learning



**Learning = ability to generalize**  
**(for UL) = ability to learn the data representation**

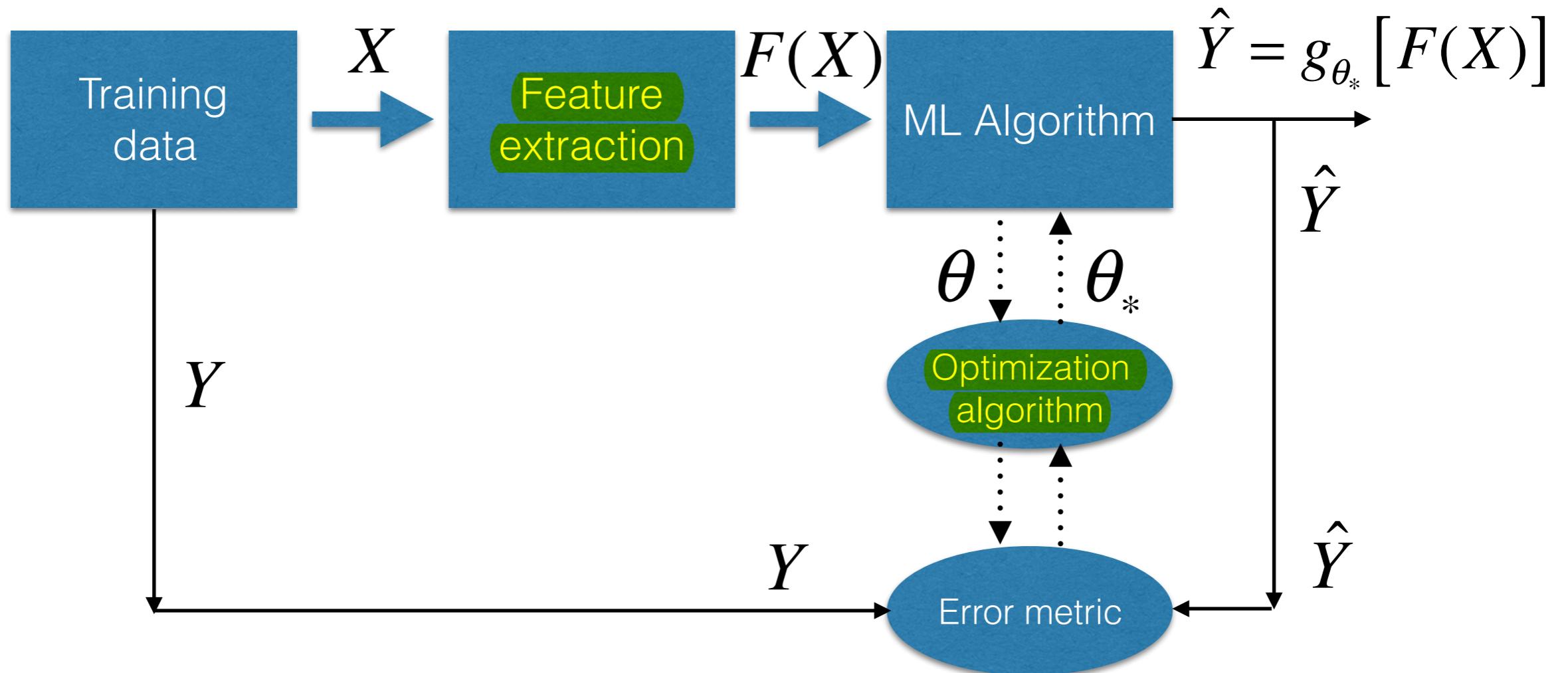
# Unsupervised Learning



**Learning = ability to generalize**  
**(for UL) = ability to learn the data representation**

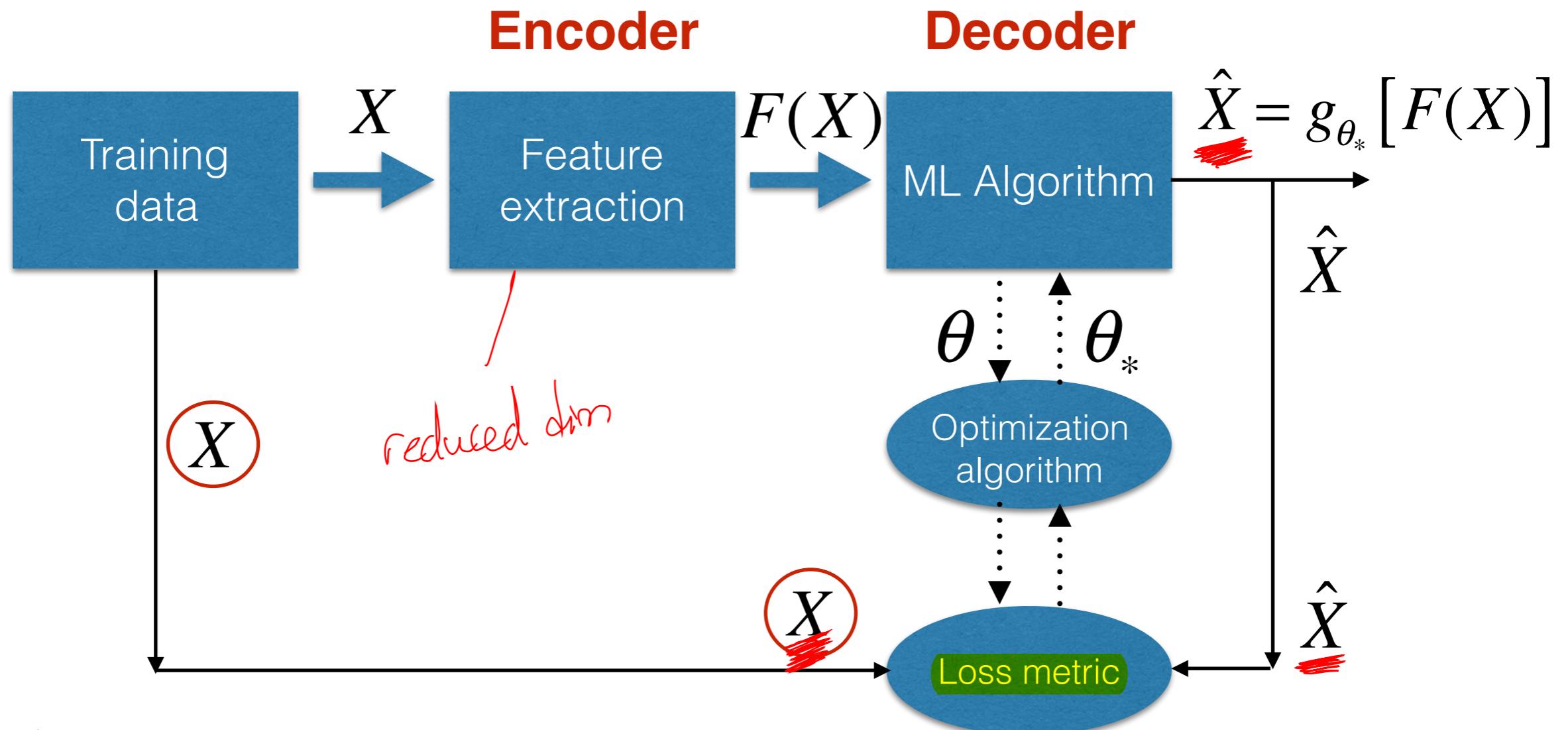
# Supervised learning diagram

The general process for training supervised learning algorithms



# Unsupervised learning diagram: Autoencoder

**Autoencoder:** replace  $Y$  with  $X$ !

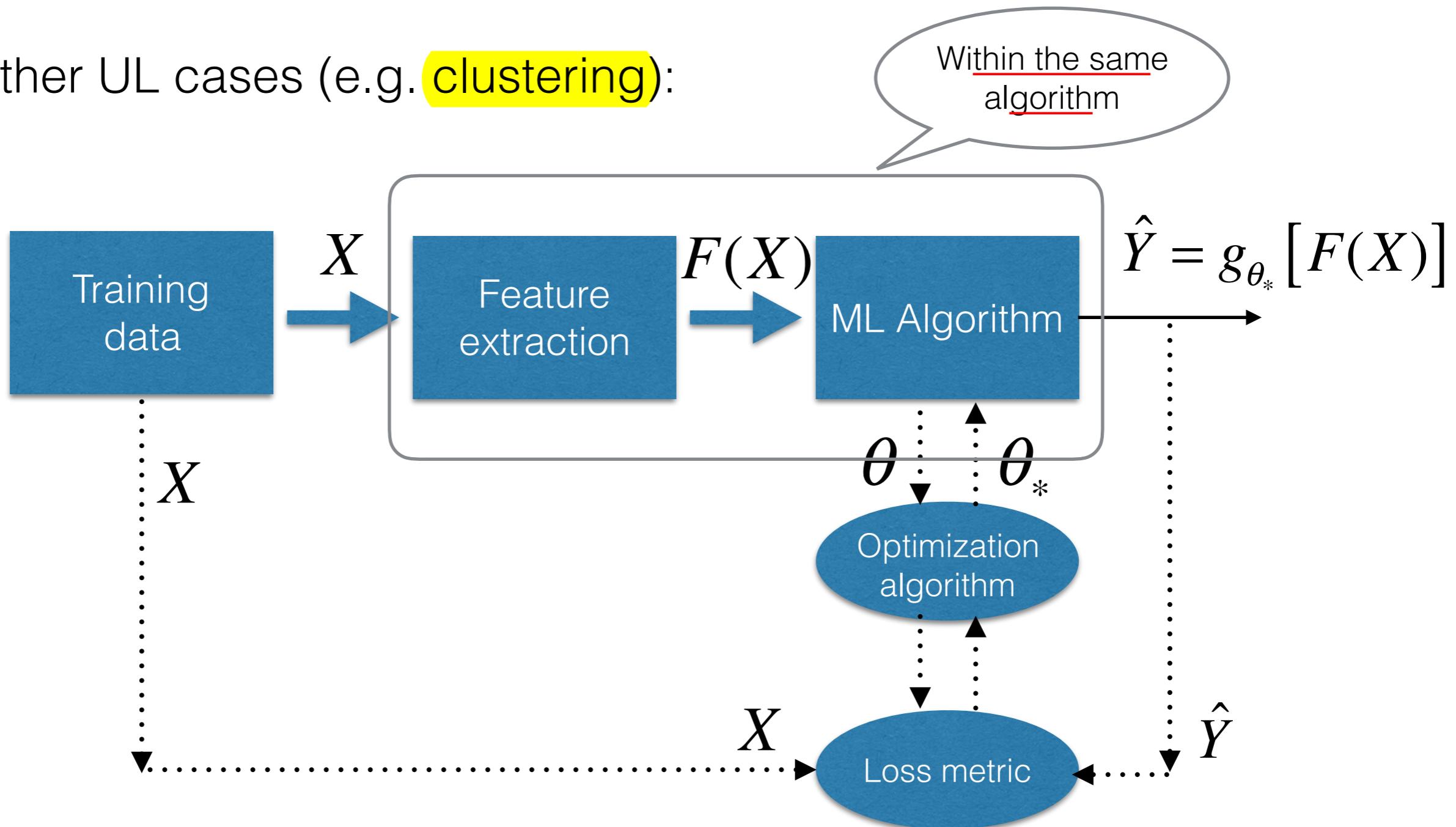


\* Model Architecture prevents the model from memorizing the input.

This is done because of the reduced dim of encoded input

# Unsupervised learning diagram

Other UL cases (e.g. clustering):



# Control question

Q: Select all correct statements:

1. As there is no “teacher” in Unsupervised Learning, it can not have a well-defined and practically useful objective.
2. The general objective of Unsupervised Learning is to learn useful representations of data.
3. A loss function for Unsupervised Learning can be obtained by comparing true labels for available data with labels produced by the model.
4. A loss function for Unsupervised Learning can be obtained by comparing an actual data distribution with a recovered data distribution as obtained from the model.



Correct answers: 2, 4.

# **Guided Tour of Machine Learning in Finance**

**Week 3-1-2-1: Unsupervised Learning**

**Principal Component Analysis (PCA)**

Igor Halperin

NYU Tandon School of Engineering, 2017

# Stock returns

Analyze returns of 30 stocks from the DJI index:

$$\underline{r_i(t)} = \log\left(\frac{S_i(t)}{S_i(t - \Delta t)}\right), \quad i = 1, \dots, N_{stocks}$$

$$\underline{r_i(t)} = \alpha_i + \beta_i r_M(t) + \varepsilon_{it}$$

Here:

$r_i(t)$  : the log-return for stock  $i$

$r_M(t)$  : the log-return for the “market”  $M$  (the DJI index)

$\beta_i$  : the “beta” - measure of correlation with the market

$\varepsilon_{it}$  : the “noise” term (residual)  $\Leftrightarrow \mathbb{E}[\varepsilon_{it}] = 0$

$\Delta t$  : the time step (1m, 1h, 1d, 1w, 1m, 1y, etc.)

a simple linear reg of  
all log-returns on a single  
predictor

Market Factor

This is the regression equation of the Capital Asset Pricing Model (CAPM),  
see e.g. Cochrane, “Asset Pricing” 2005).

# Eigen-portfolios as a UL tasks

Trading signals for statistical arbitrage from the ML perspective

In statistical arbitrage strategies:

- Trading signals are systematic or rule-based
- Trading book has zero beta with the market (i.e. it is market-neutral)
- Use diversification to produce low volatility portfolios uncorrelated with the market

trying to filter out impact of the market from stock return

i.e. market is  
too random  
to be predicted

$$\hat{r}_i(t) = \hat{\alpha}_i + \hat{\beta}_i r_M(t)$$

$$\hat{\varepsilon}_i(t) = r_i(t) - \hat{r}_i(t) = r_i(t) - \hat{\alpha}_i - \hat{\beta}_i r_M(t)$$

\* We subtract market impact & focus on the residuals, which are parts of the equity return that are explained by the market returns.

\* Residuals are used as features to produce trading signals.

# Eigen-portfolios as a UL tasks

Trading signals for statistical arbitrage from the ML perspective

In statistical arbitrage strategies:

- Trading signals are systematic or rule-based
- Trading book has zero beta with the market (i.e. it is market-neutral)
- Use diversification to produce low volatility portfolios uncorrelated with the market
- See Avellaneda and Lee (2008) for details
- Multi-factor model: individual returns  $R_i$  ( $i = 1, \dots, N$ ) are driven by  $m$  factors  $F_j$  ( $j = 1, \dots, m$ ) and idiosyncratic components  $\tilde{R}_i$  with  $\mathbb{E}[\tilde{R}_i] = 0$

$$R_i = \alpha_i + \sum_{j=1}^m \beta_{ij} F_j + \tilde{R}_i$$

# Eigen-portfolios as a UL tasks

Trading signals for statistical arbitrage from the ML perspective

In statistical arbitrage strategies:

- Trading signals are systematic or rule-based
- Trading book has zero beta with the market (i.e. it is market-neutral)
- Use diversification to produce low volatility portfolios uncorrelated with the market
- See Avellaneda and Lee (2008) for details
- Multi-factor model: individual returns  $R_i$  ( $i = 1, \dots, N$ ) are driven by  $m$  factors  $F_j$  ( $j = 1, \dots, m$ ) and idiosyncratic components  $\tilde{R}_i$  with  $\mathbb{E}[\tilde{R}_i] = 0$

$$R_i = \alpha_i + \sum_{j=1}^m \beta_{ij} F_j + \tilde{R}_i$$

- Factors  $F_j$  can be thought of as returns of “benchmark” portfolios representing **systematic factors** (e.g. industry and geography factors)

\* We focus on a more data-driven approach.

↑  
Classic Finance

# Control question

Select all correct answers

1. In the factor model, all individual returns are only correlated with “market”  $M$  with correlations  $\beta_i$ , but remain uncorrelated between themselves.
2. Correlation of two assets  $i$  and  $j$  in the factor model is given by  $\beta_i \beta_j$ .
3. Statistical arbitrage deals with “factor investing” to find factors  $F_j$  with highest returns.
4. Statistical arbitrage use residuals of stock returns “unexplained” by factors to extract trading signals.

**Correct answer:** 2, 4

# **Guided Tour of Machine Learning in Finance**

**Week 3-1-2-2: Unsupervised Learning**

**Principal Component Analysis (PCA)**

Igor Halperin

NYU Tandon School of Engineering, 2017

# PCA approach: eigen-portfolios

Try PCA as a way to **extract factors** directly from the data

**Data:** history of  $N + 1$  daily prices of  $M$  stocks  $S_n^{(i)} \equiv S_{t_n}^{(i)}$  for  $i = 1, \dots, M$   
measured at times  $t = [t_0, t_1, \dots, t_N]$

1. Compute **daily returns**:

$$\underline{R}_{ni} = \frac{S_n^{(i)} - S_{n-1}^{(i)}}{S_{n-1}^{(i)}} \simeq \log \frac{S_n^{(i)}}{S_{n-1}^{(i)}}, \quad n = 1, \dots, N, \quad i = 1, \dots, M$$

2. **Standardized returns** (data normalization):

$$\underline{X}_{ni} = \frac{\underline{R}_{ni} - \bar{R}_i}{\bar{\sigma}_i}, \quad \bar{R}_i = \frac{1}{N} \sum_{n=1}^N R_{ni}, \quad \bar{\sigma}_i^2 = \frac{1}{N-1} \sum_{n=1}^N (\underline{R}_{ni} - \bar{R}_i)^2$$

3. The empirical **correlation matrix** is the covariance matrix of standardized returns:

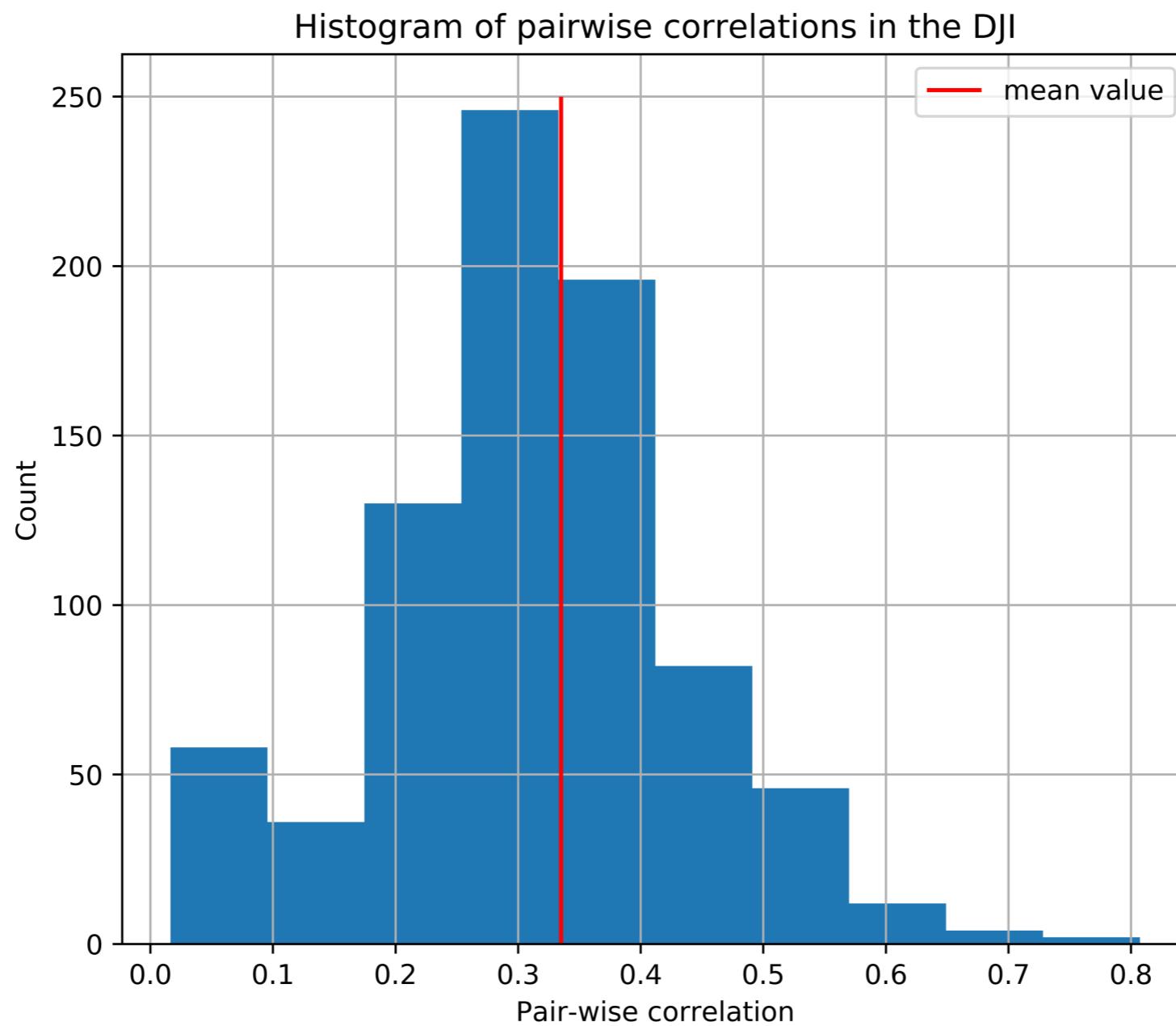
$$\underline{C}_{ij} = \frac{1}{N-1} \sum_{n=1}^N X_{ni} X_{nj} = \frac{1}{N-1} (\mathbf{X}^T \mathbf{X})_{ij}$$



This matrix is not diagonal!

# Distribution of pairwise correlations in the DJI

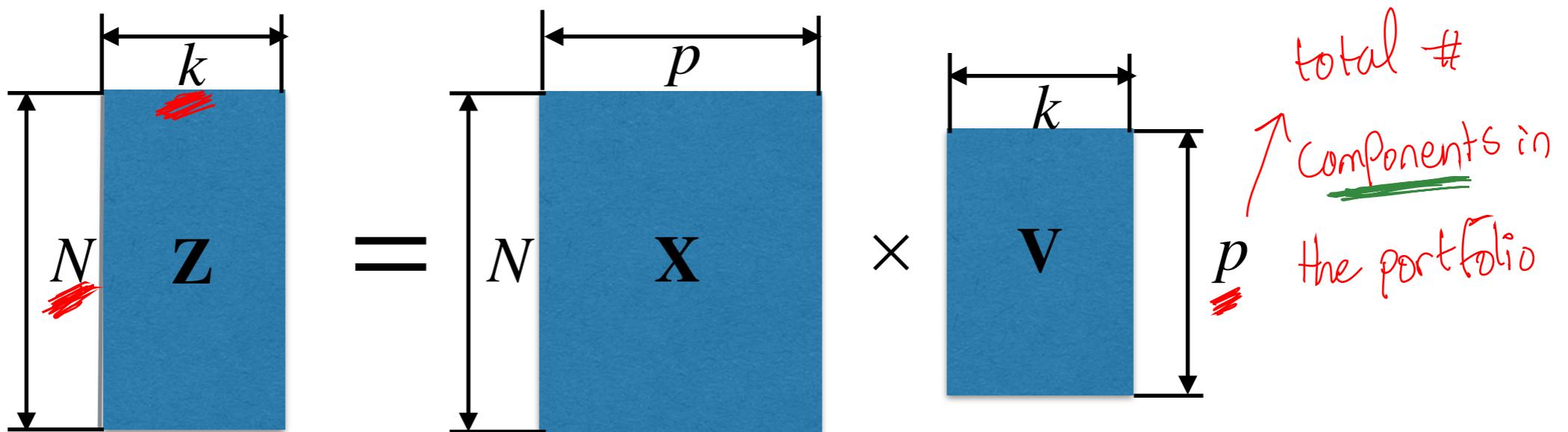
Histogram of pairwise return correlations for stocks in the DJI index:



# PCA as a coordinate transform

How to make correlation matrix  $\mathbf{C}$  diagonal?

Introduce a linear transform (linear encoder) of the data  $\mathbf{Z} = \mathbf{X}\mathbf{V}$  parametrized by a  $p \times k$  orthogonal matrix  $\mathbf{V}$  with  $\mathbf{V}\mathbf{V}^T = 1$



$\mathbf{V}$  is an orthogonal matrix that stores

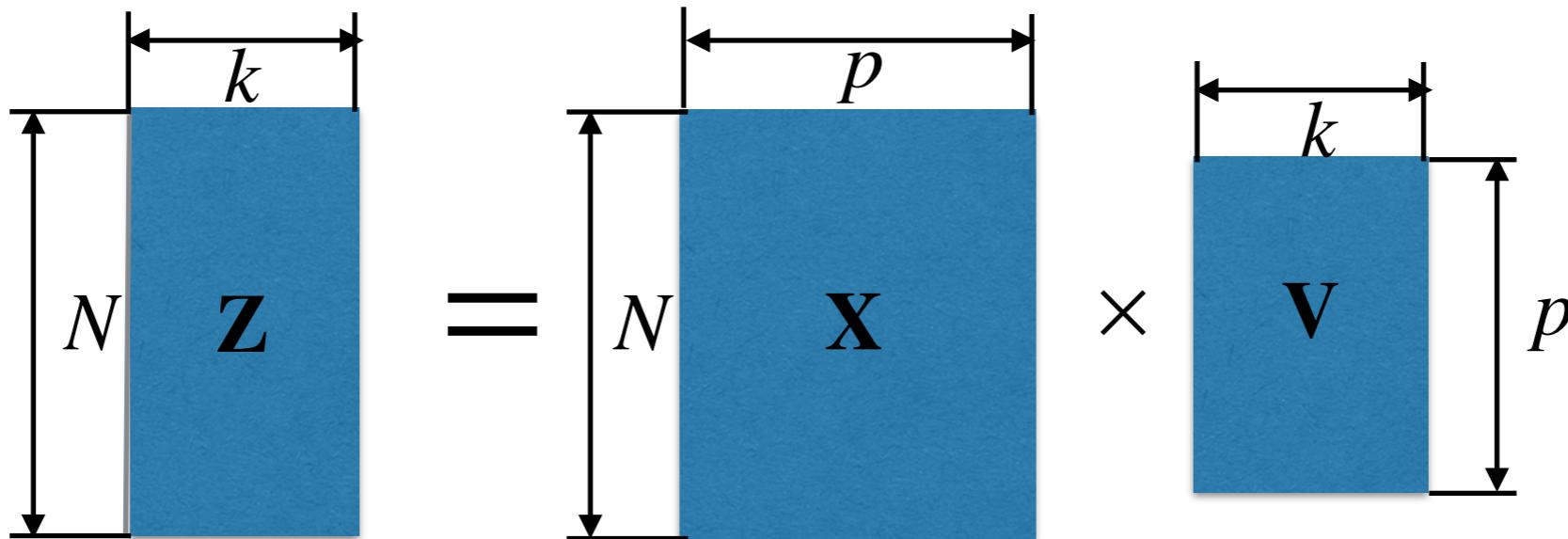
eigenvectors column-wise

$p \times k$  s.t.  $p \geq k$

# PCA as a coordinate transform

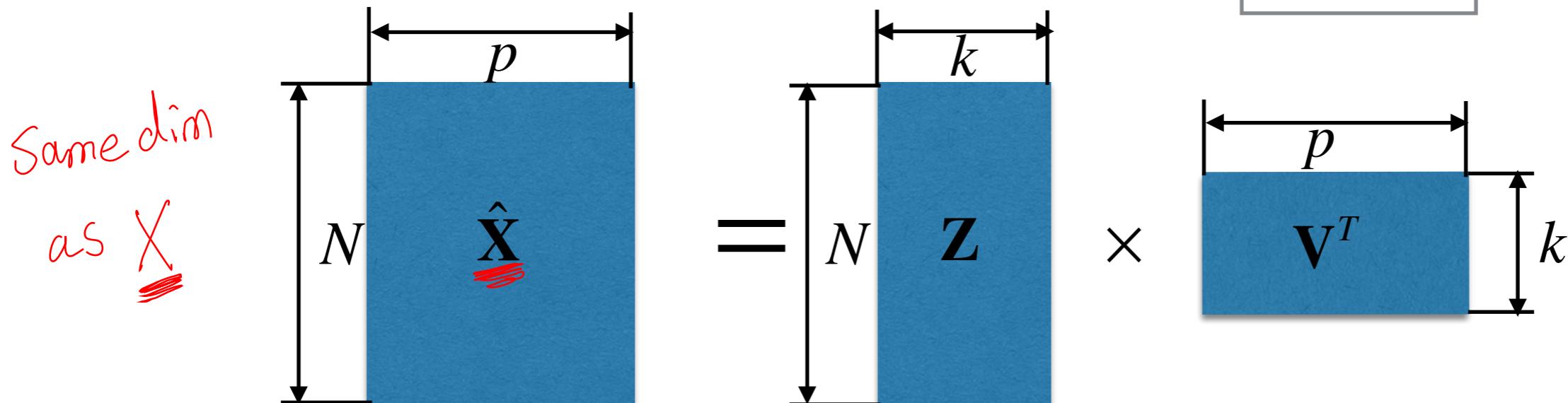
How to make correlation matrix  $\mathbf{C}$  diagonal?

Introduce a linear transform (linear encoder) of the data  $\mathbf{Z} = \mathbf{X}\mathbf{V}$  parametrized by a  $p \times k$  orthogonal matrix  $\mathbf{V}$  with  $\mathbf{V}\mathbf{V}^T = 1$



A **decoded signal** (inverse transform) is obtained as

$$\hat{\mathbf{X}} = \mathbf{Z}\mathbf{V}^T$$



# PCA: Eigenvalue decomposition

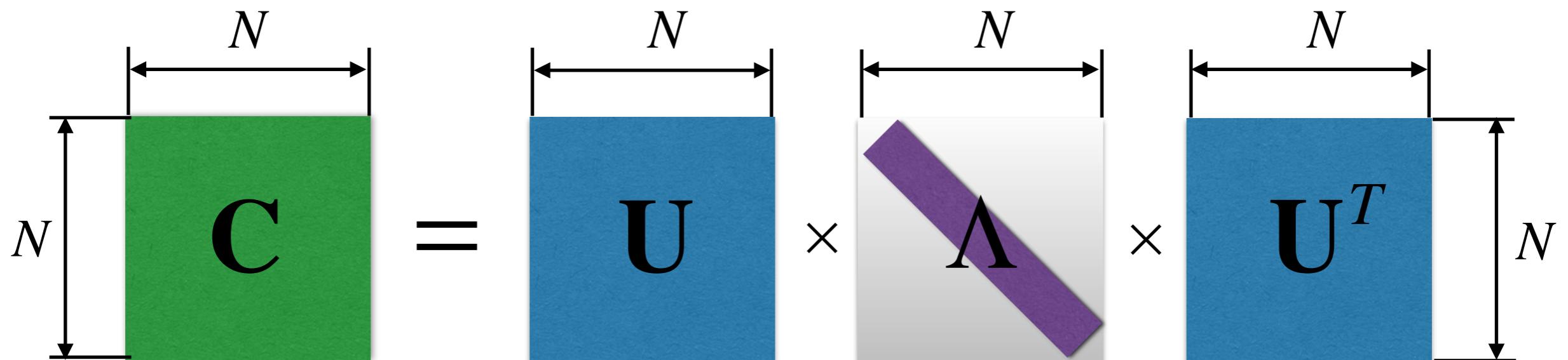
**Eigenvalue decomposition** of the correlation matrix (note that  $\mathbf{C}$  is non-negative definite!)

$$\rightarrow \mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^T$$

where (assuming that  $N > p$ ):

$\Lambda$  =  $diag(\lambda_1, \dots, \lambda_p)$  ( $\lambda_1 \geq \dots \geq \lambda_p$ ) is a diagonal matrix of ordered eigenvalues

$\mathbf{U}$  is a  $N \times N$  orthogonal matrix ( $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ ) that stores eigenvectors column-wise



# PCA: Eigenvalue decomposition

Eigenvalue decomposition of the correlation matrix (note that  $\mathbf{C}$  is non-negative definite!)

$$\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^T$$

where (assuming that  $N > p$ ):

$\Lambda = diag(\lambda_1, \dots, \lambda_p)$  ( $\lambda_1 \geq \dots \geq \lambda_p$ ) is a diagonal matrix of ordered eigenvalues

$\mathbf{U}$  is a  $N \times N$  orthogonal matrix ( $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ ) that stores eigenvectors column-wise

Use this to compute the covariance matrix of  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

$$Cov[\mathbf{Z}] = \frac{1}{N-1} \mathbf{Z}^T \mathbf{Z} = \frac{1}{N-1} \mathbf{V}^T \mathbf{X}^T \mathbf{X} \mathbf{V} = \mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{V}^T \mathbf{U} \Lambda (\mathbf{V}^T \mathbf{U})^T$$

# PCA: Eigenvalue decomposition

Eigenvalue decomposition of the correlation matrix (note that  $\mathbf{C}$  is non-negative definite!)

$$\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^T$$

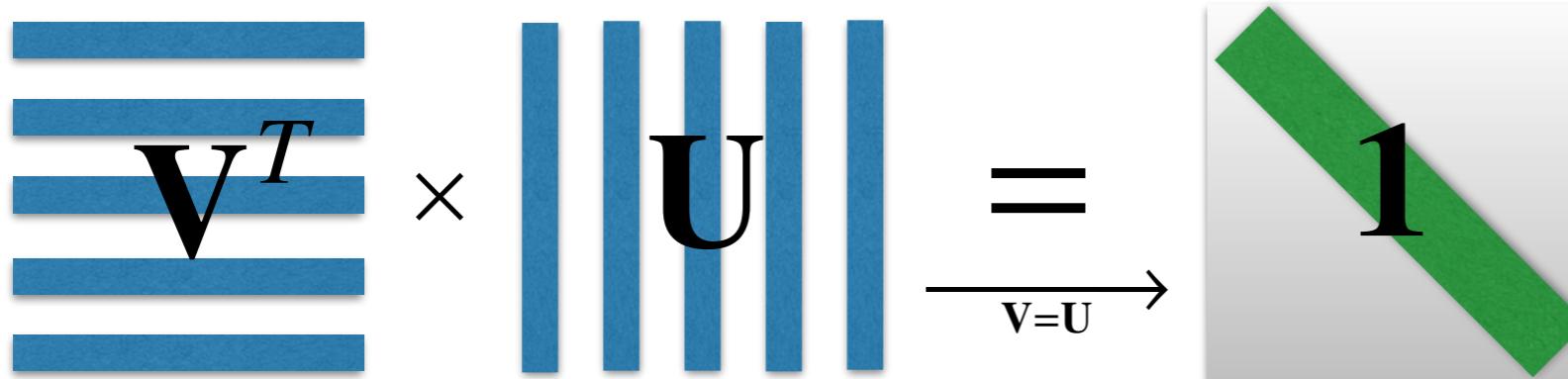
where (assuming that  $N > p$ ):

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  ( $\lambda_1 \geq \dots \geq \lambda_p$ ) is a diagonal matrix of ordered eigenvalues

$\mathbf{U}$  is a  $N \times N$  orthogonal matrix ( $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ ) that stores eigenvectors column-wise

Use this to compute the covariance matrix of  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

$$\cancel{\mathbf{Cov}[\mathbf{Z}] = \frac{1}{N-1}\mathbf{Z}^T\mathbf{Z}} = \frac{1}{N-1}\mathbf{V}^T\mathbf{X}^T\mathbf{X}\mathbf{V} = \mathbf{V}^T\mathbf{C}\mathbf{V} = \mathbf{V}^T\mathbf{U}\Lambda(\mathbf{V}^T\mathbf{U})^T \xrightarrow[\mathbf{V}=\mathbf{U}]{} \Lambda$$


$$\mathbf{V}^T \times \mathbf{U} = \mathbf{1}$$

$\mathbf{Z} \rightarrow$  made of uncorrelated  
Components with a diagonal  
Correlation matrix

# PCA as a coordinate transform

Eigenvalue decomposition of the correlation matrix (note that  $\mathbf{C}$  is non-negative definite!)

$$\mathbf{C} = \mathbf{U}\Lambda\mathbf{U}^T$$

where (assuming that  $N > p$ ):

$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$  ( $\lambda_1 \geq \dots \geq \lambda_p$ ) is a diagonal matrix of ordered eigenvalues

$\mathbf{U}$  is a  $N \times N$  orthogonal matrix ( $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ ) that stores eigenvectors column-wise

Use this to compute the covariance matrix of  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

$$\text{Cov}[\mathbf{Z}] = \frac{1}{N-1} \mathbf{Z}^T \mathbf{Z} = \frac{1}{N-1} \mathbf{V}^T \mathbf{X}^T \mathbf{X} \mathbf{V} = \mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{V}^T \mathbf{U} \Lambda (\mathbf{V}^T \mathbf{U})^T \xrightarrow{\mathbf{V}=\mathbf{U}} \Lambda$$

Therefore, when  $\mathbf{V} = \mathbf{U}$  (and  $k = p$ ), encoding  $\mathbf{Z} = \mathbf{X}\mathbf{V}$  preserves the total variation of data

↑

$$\text{TotVar}[\mathbf{X}] = \frac{1}{N-1} \text{Tr}[\mathbf{X}^T \mathbf{X}] = \text{Tr}[\mathbf{C}] = \text{Tr}[\mathbf{U} \Lambda \mathbf{U}^T] = \text{Tr}[\Lambda \mathbf{U}^T \mathbf{U}] = \text{Tr}[\Lambda]$$

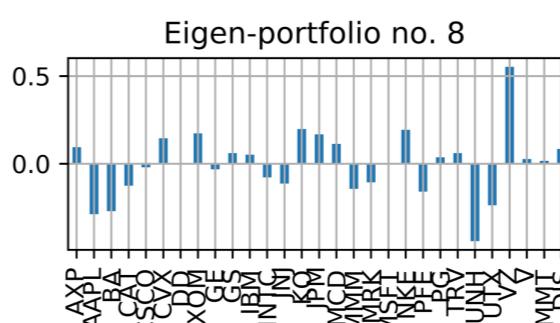
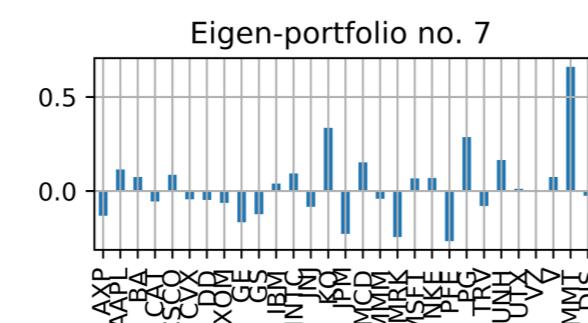
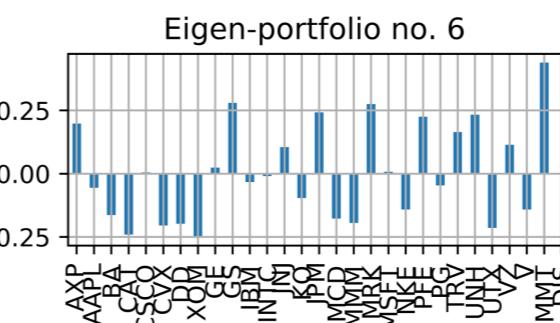
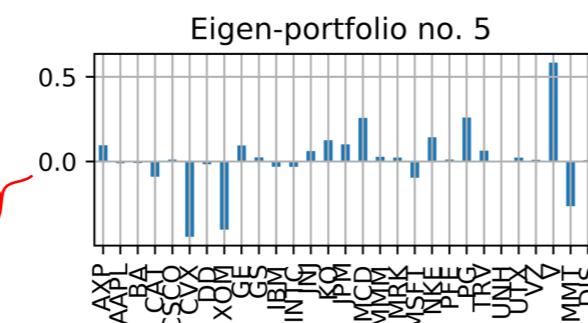
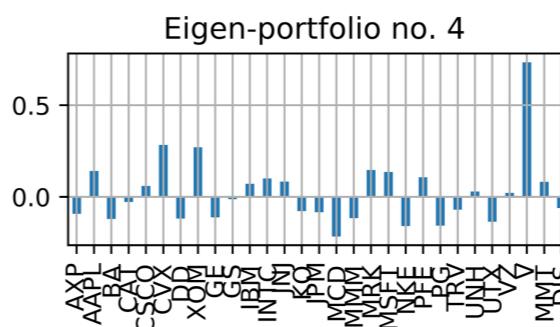
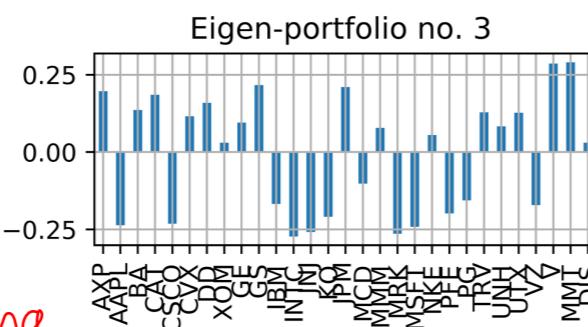
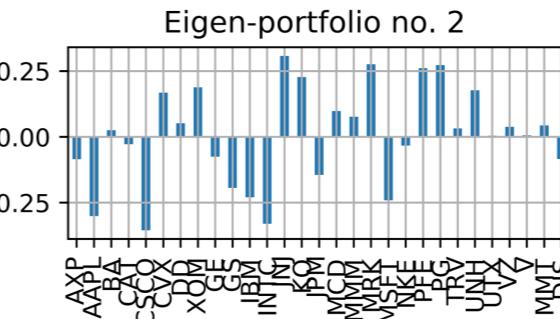
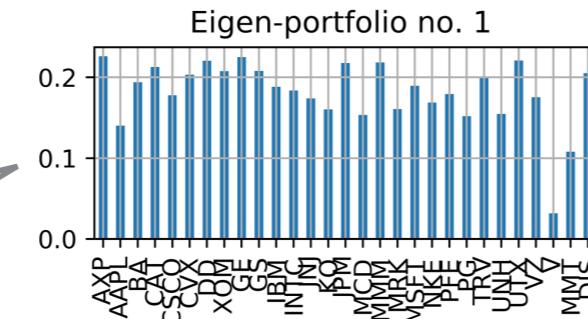
$$\text{TotVar}[\mathbf{Z}] = \frac{1}{N-1} \text{Tr}[\mathbf{Z}^T \mathbf{Z}] = \text{Tr}[\mathbf{V}^T \mathbf{U} \Lambda (\mathbf{V}^T \mathbf{U})^T] \xrightarrow{\mathbf{V}=\mathbf{U}} \text{Tr}[\Lambda]$$

\* When  $K=p$ , we preserve all variance in data by making a linear encoding  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

# PCA for stocks in DJI: eigenvectors

Eigenvectors = weights of stocks in “eigen-portfolios”

Approximates  
the market (DJI)



Orthogonal to  
the market (DJI)

\* We reduce our problem  
of optimal investment among  
universe of  $P$  correlated  
stocks → into a simpler  
one of  $P$  uncorrelated  
eigenportfolios

i.e. uncorrelated  
to market, and  
uncorrelated to the  
rest

uncorrelated  
≠  
independent

# Control question

Select all correct answers

- 1. The PCA implements a Linear Encoder that converts correlated inputs  $\mathbf{X}$  into uncorrelated features  $\mathbf{Z}$  by a linear transform  $\mathbf{Z} = \mathbf{X}\mathbf{V}$
- 2. If the orthogonal matrix  $\mathbf{V}$  has dimension  $p \times p$ , i.e. it keeps all eigenvectors of correlation matrix of  $\mathbf{X}$ , then  $\mathbf{Z}$  preserves the total variation of  $\mathbf{X}$ .
- 3. When the data is noisy or non-stationary,  $\mathbf{Z}$  can have a higher variance than  $\mathbf{X}$ .
- 4. The PCA is a probabilistic method that enables simulating from data.

**Correct answer:** 1, 2

# **Guided Tour of Machine Learning in Finance**

## **Week 3: Unsupervised Learning**

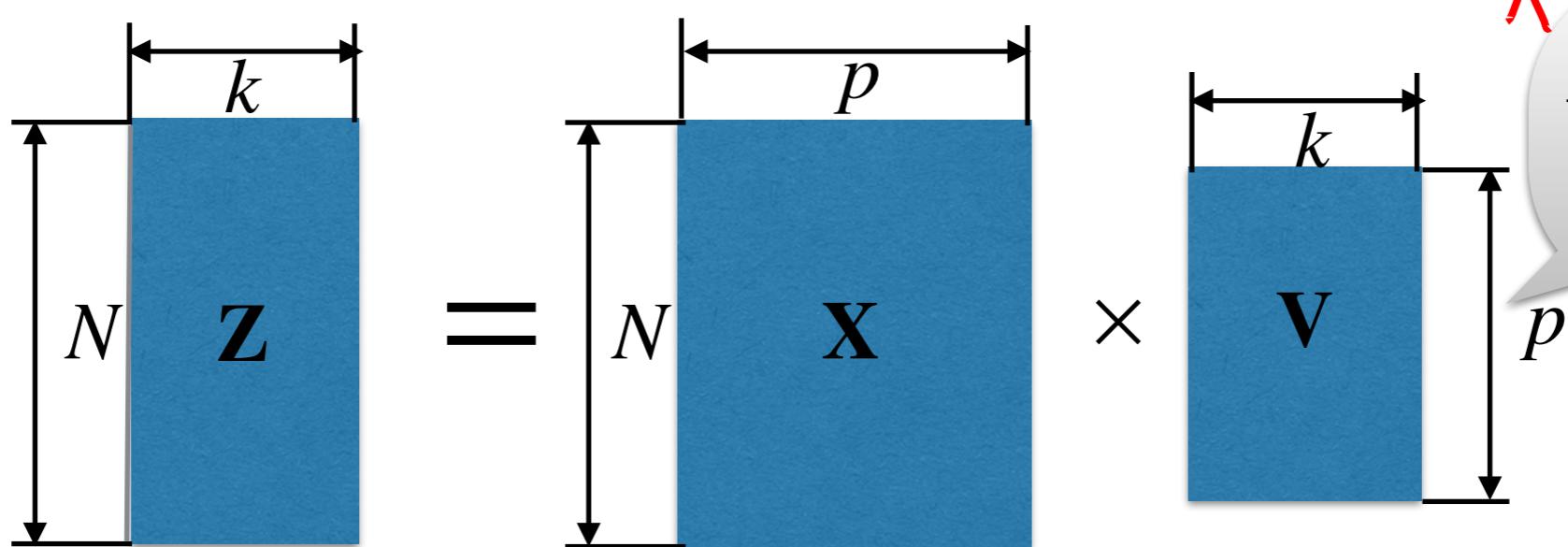
### **Dimension reduction with the PCA**

Igor Halperin

NYU Tandon School of Engineering, 2017

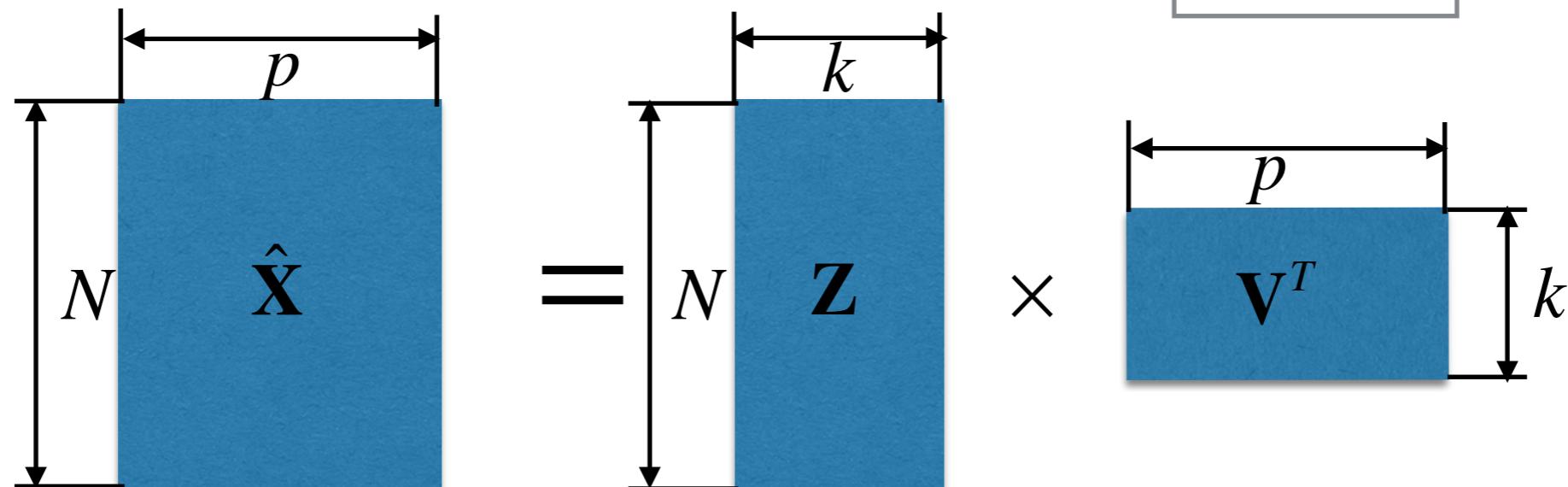
# PCA as a dimension reduction method

We had a linear transform (linear encoder) of the data  $\mathbf{Z} = \mathbf{X}\mathbf{V}$  parametrized by a  $p \times k$  orthogonal matrix  $\mathbf{V}$  with  $\mathbf{V}\mathbf{V}^T = 1$



A decoded signal (inverse transform) is obtained as

$$\hat{\mathbf{X}} = \mathbf{Z}\mathbf{V}^T$$



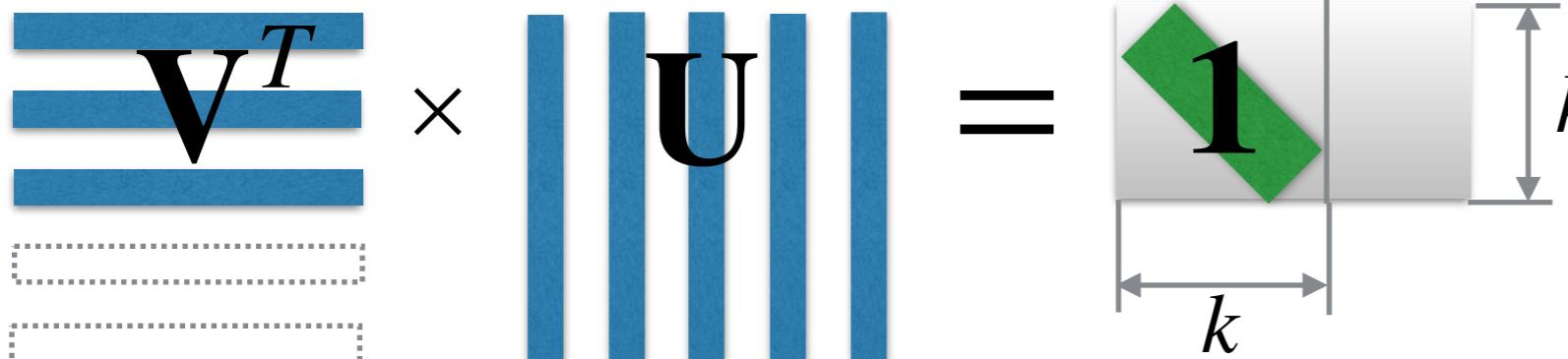
# Dimension reduction with PCA

Dimension reduction: make the projection matrix  $\mathbf{V}$  of the first  $k \leq p$  eigenvectors of  $\mathbf{C}$

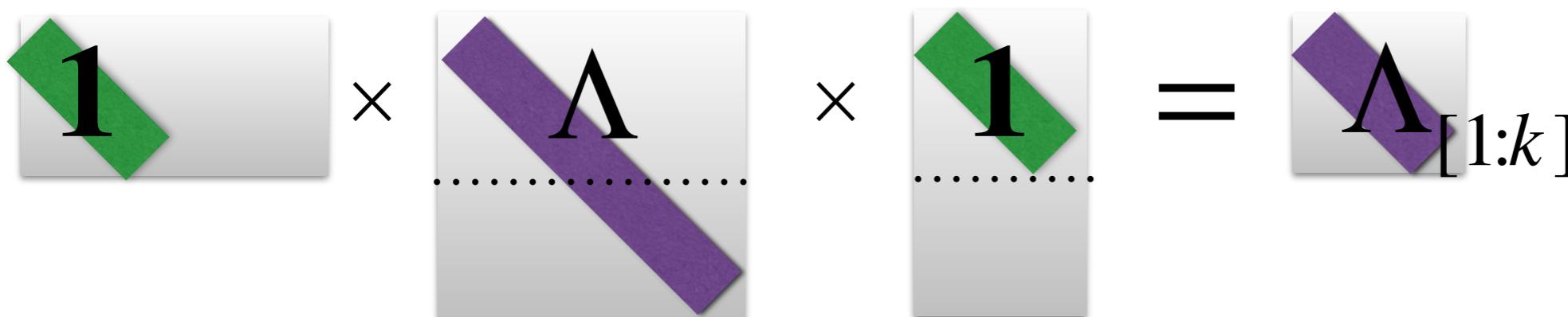
$$\mathbf{V} = \mathbf{U}^{[1:k]}$$

Use this to compute the covariance matrix of  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

$$Cov[\mathbf{Z}] = \frac{1}{N-1} \mathbf{Z}^T \mathbf{Z} = \frac{1}{N-1} \mathbf{V}^T \mathbf{X}^T \mathbf{X} \mathbf{V} = \mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{V}^T \mathbf{U} \boldsymbol{\Lambda} (\mathbf{V}^T \mathbf{U})^T \xrightarrow[\mathbf{V}=\mathbf{U}^{[1:k]}]{} \boldsymbol{\Lambda}_{[1:k]}$$

  $\mathbf{V}^T \times \mathbf{U} = \begin{matrix} & & 1 & & \\ & & \vdots & & \\ & & 1 & & \\ & & \vdots & & \\ & & 1 & & \end{matrix}$

$\iff \mathbf{W}_{ij} = \mathbf{U}_i^T \mathbf{U}_j = \delta_{ij},$   
 $i = 1, \dots, k, j = 1, \dots, p$



# Dimension reduction with PCA

Dimension reduction: make the projection matrix  $\mathbf{V}$  of the first  $k \leq p$  eigenvectors of  $\mathbf{C}$

$$\mathbf{V} = \mathbf{U}^{[1:k]}$$

Use this to compute the covariance matrix of  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

$$Cov[\mathbf{Z}] = \frac{1}{N-1} \mathbf{Z}^T \mathbf{Z} = \frac{1}{N-1} \mathbf{V}^T \mathbf{X}^T \mathbf{X} \mathbf{V} = \mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{V}^T \mathbf{U} \boldsymbol{\Lambda} (\mathbf{V}^T \mathbf{U})^T \xrightarrow[\mathbf{V}=\mathbf{U}^{[1:k]}]{} \boldsymbol{\Lambda}_{[1:k]}$$

$$V^T \times U = \begin{matrix} \text{1} \\ \text{k} \end{matrix}$$

Only a part of total variation of  $\mathbf{X}$  is preserved:

$$TotVar[\mathbf{Z}] = \text{Tr}[\Lambda_{[1:k]}] = \sum_{i=1}^k \lambda_i$$

$$\begin{matrix} \text{1} \\ \times & \Lambda \\ \dots & \dots \end{matrix} = \begin{matrix} \text{1} \\ \Lambda_{[1:k]} \end{matrix}$$

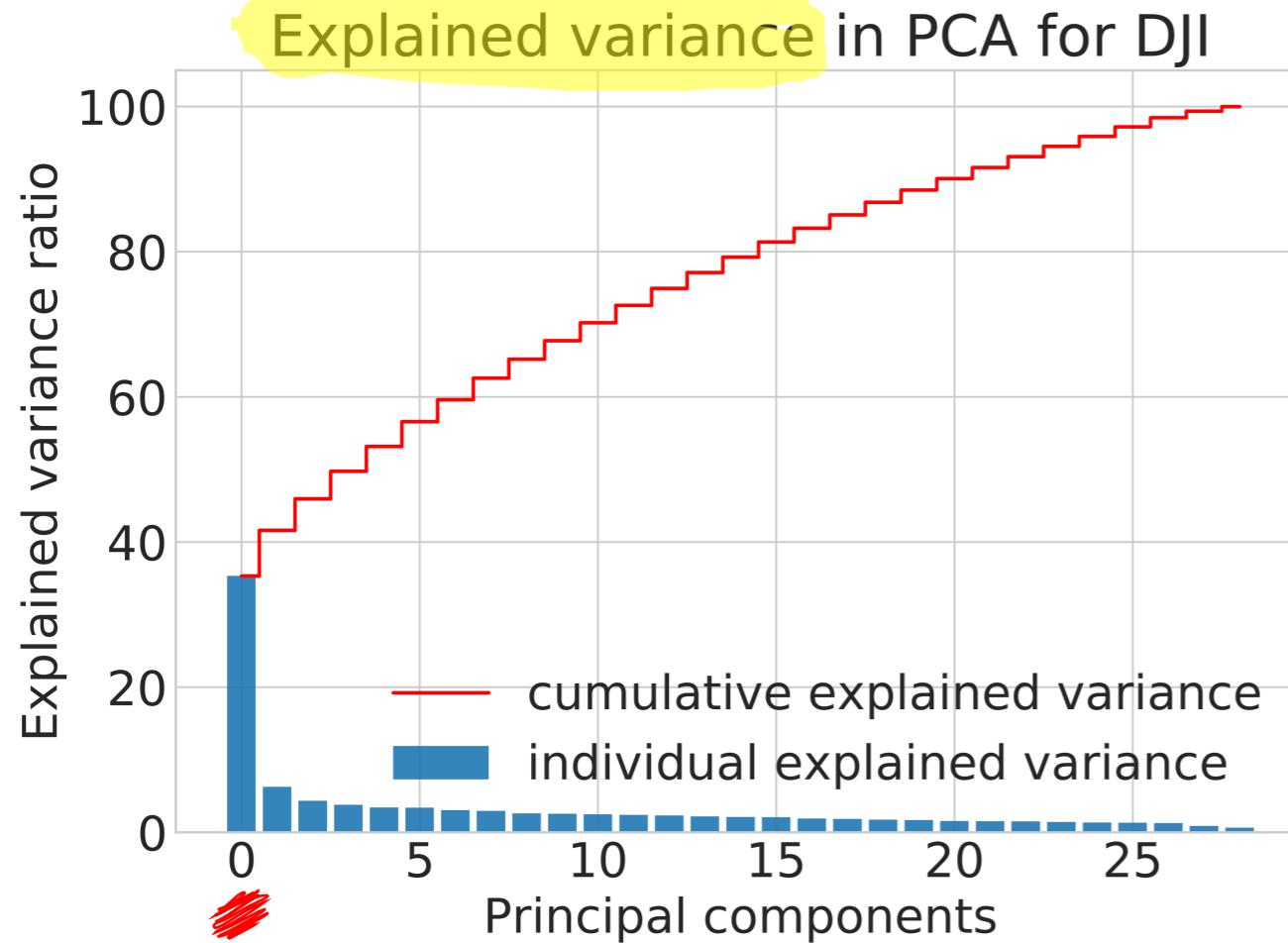
# Variance explained by PCA

Dimension reduction: make the projection matrix  $\mathbf{V}$  of the first  $k \leq p$  eigenvectors of  $\mathbf{C}$

$$\mathbf{V} = \mathbf{U}^{[1:k]}$$

Use this to compute the covariance matrix of  $\mathbf{Z} = \mathbf{X}\mathbf{V}$

$$Cov[\mathbf{Z}] = \frac{1}{N-1} \mathbf{Z}^T \mathbf{Z} = \frac{1}{N-1} \mathbf{V}^T \mathbf{X}^T \mathbf{X} \mathbf{V} = \mathbf{V}^T \mathbf{C} \mathbf{V} = \mathbf{V}^T \mathbf{U} \boldsymbol{\Lambda} (\mathbf{V}^T \mathbf{U})^T \xrightarrow{\mathbf{V}=\mathbf{U}^{[1:k]}} \boldsymbol{\Lambda}_{[1:k]}$$



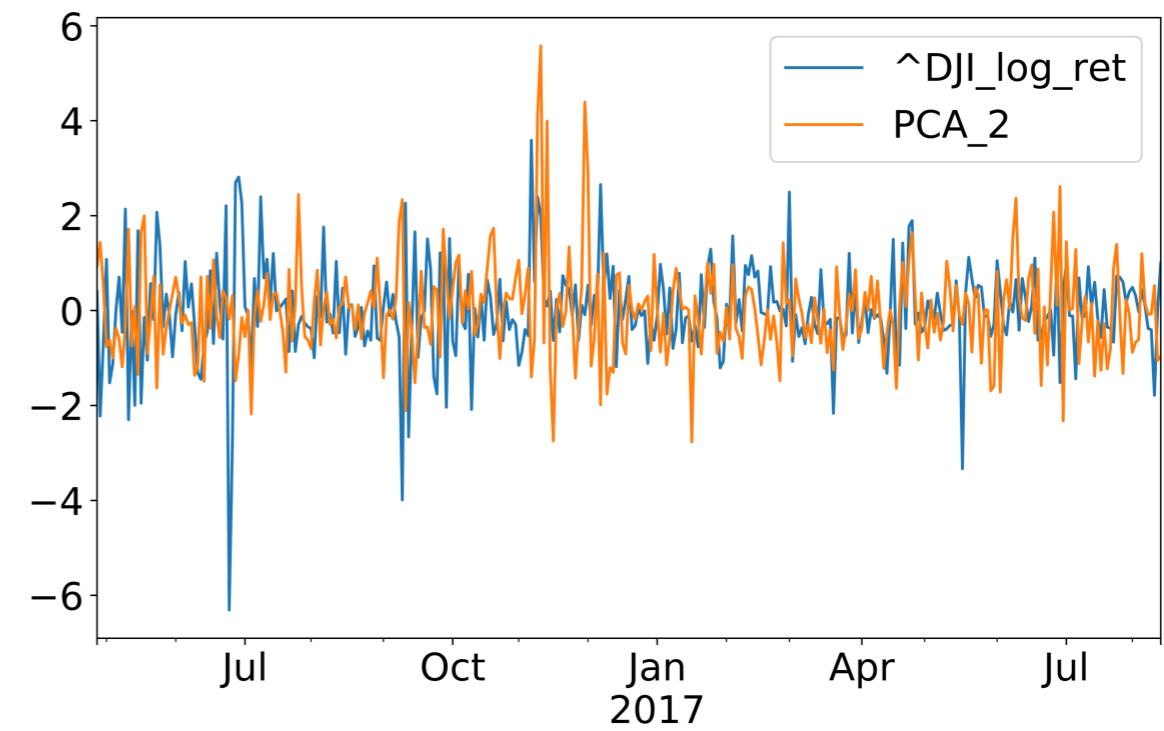
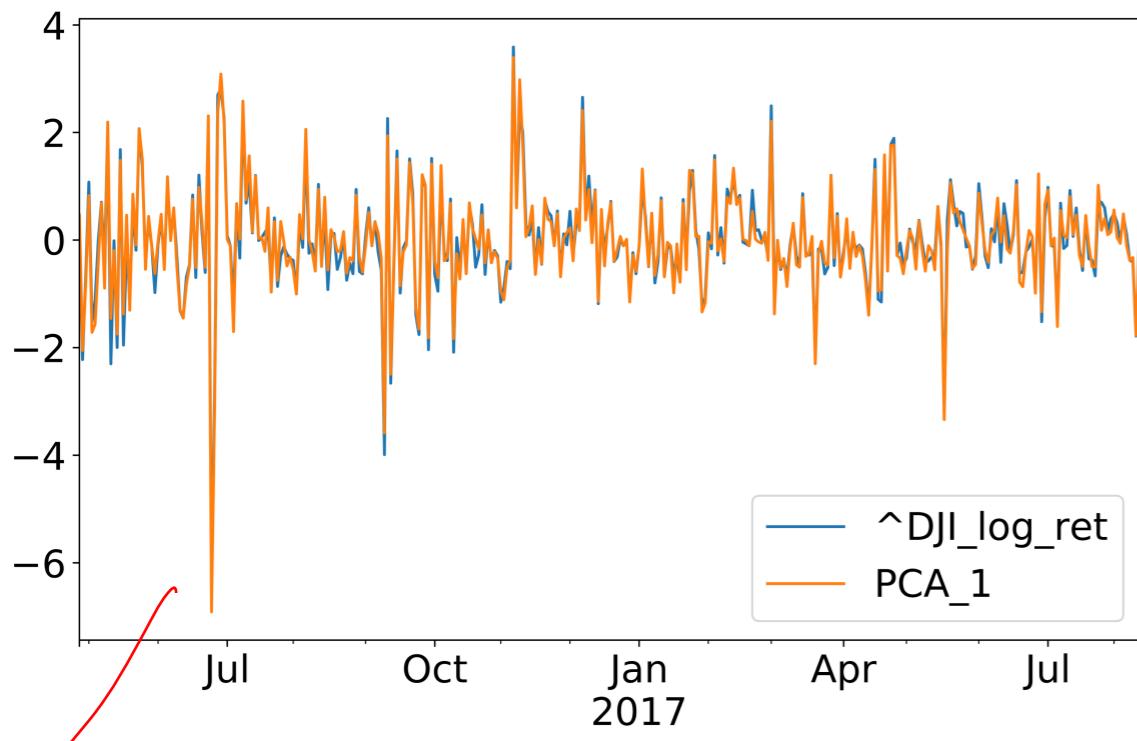
Only a part of total variation of  $\mathbf{X}$  is preserved:

$$TotVar[\mathbf{Z}] = \text{Tr}[\boldsymbol{\Lambda}_{[1:k]}] = \sum_{i=1}^k \lambda_i$$

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^N \lambda_i} \%$$

# PCA projections vs the market

Compare the actual returns on the DJI index with returns of the two leading eigen-portfolios



large values  
of returns exhibit  
obvious mismatch

PC 1

Correlated  
with  
the market

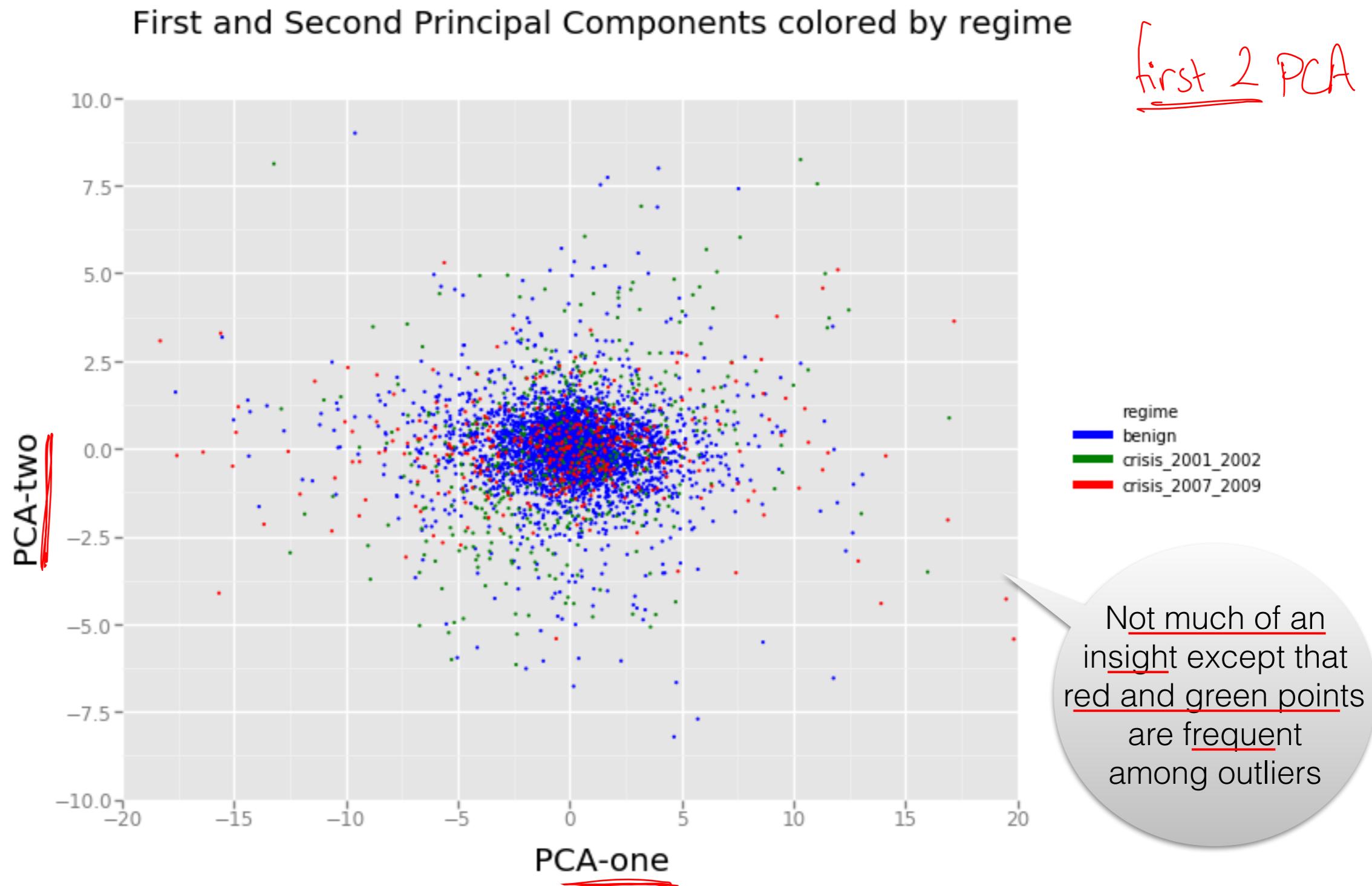
PC 2

Uncorrelated  
with the market

much more mismatch

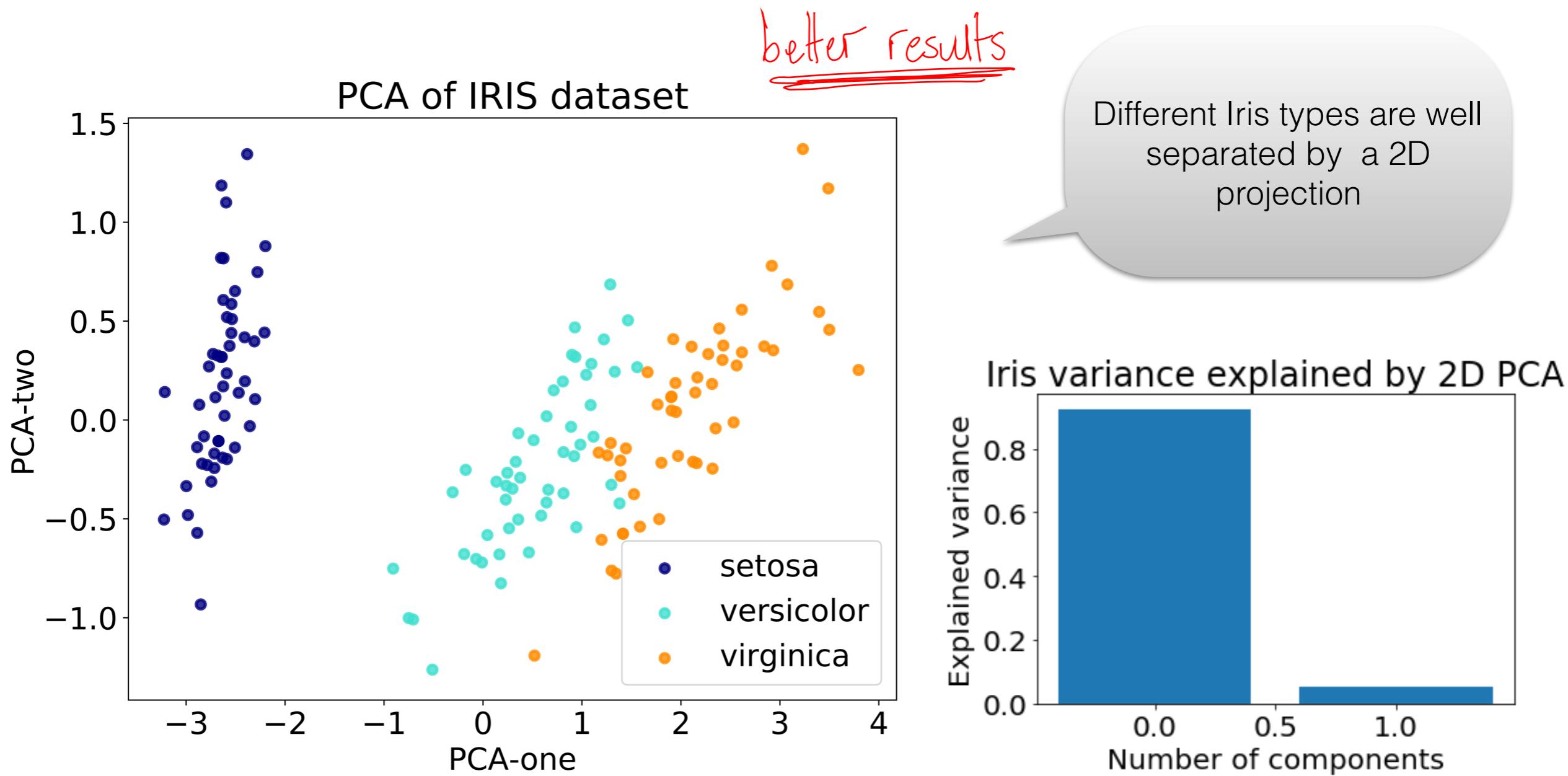
# PCA projections: 2D visualization

Compare the actual returns on the DJI index with returns of the two leading eigen-portfolios



# PCA 2D visualization for the Iris dataset

Compare with the 2D PCA projection for the Iris dataset: 3 types of Iris, 4 features: Sepal Length, Sepal Width, Petal Length and Petal Width, 150 examples, see [http://scikit-learn.org/stable/auto\\_examples/datasets/plot\\_iris\\_dataset.html](http://scikit-learn.org/stable/auto_examples/datasets/plot_iris_dataset.html).



# Quizz: Why does the 2D PCA perform worse for stocks than for Iris?

1. Because stock prices are produced by human interactions, while variability of the Iris species is produced by nature
2. The Iris data is stationary, and is more homogeneous/has less outliers than the stock data
3. Smaller dimension of the feature space
4. Stronger correlations between features
5. The labels are real for Iris (they are Iris' types), but “noisy” for the “regime” label in the stock analysis

## Multiple choice answers:

1. Answer 1
2. Answer 3
3. Answers 2 and 4
4. Answer 2 and 5
5. Answers 3, 4, 5
6. Answers 2, 3, 4, 5
7. All of the above

\* PCA is linear method, any non-linearity  
will not be captured

## Correct answer: 6



## 2D PCA for stocks vs Iris

Why does the 2D PCA perform much worse for stocks than for Iris?

- Smaller dimension of the feature space
- Stronger correlations between the features (the first PC explains 92% of total variance!)
- The Iris data is stationary, and is more homogeneous/has less outliers than the stock data
- The labels are real for Iris (they are Iris' types), but “noisy” for the “regime” label in the stock analysis



PCA is expected to work well for problems that are largely determined

by lowest moments of distributions i.e mean & covariance

Higher moments like Skewness & Kurtosis won't.

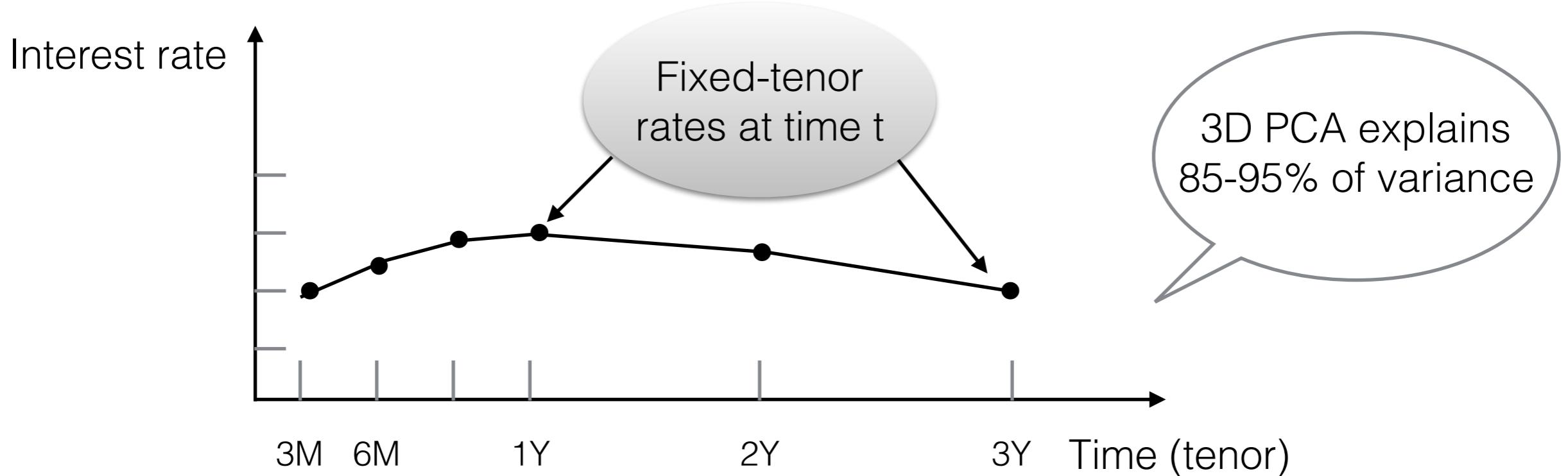


Such app rely on analysis of the total variance in the data versus variance explained by first Principal Components.

# Financial applications of the PCA

Numerous financial applications of the PCA:

- Systematic trading: construction of market-neutral portfolios with low volatility, extraction of trading signals
- Systemic risk detection methods (e.g. M. Kritzman et. al. “Principal Components as a Measure of Systemic Risk”, 2010)
- Risk management: sector risk exposure management of trading book or banking book portfolios, de-noising of empirical correlation matrices
- 2D PCA for data visualization
- PCA for yield curve modeling (fixed income, commodities)



# **Guided Tour of Machine Learning in Finance**

## **Week 3: Unsupervised Learning**

**Dimension reduction and data visualization with t-SNE**

Igor Halperin

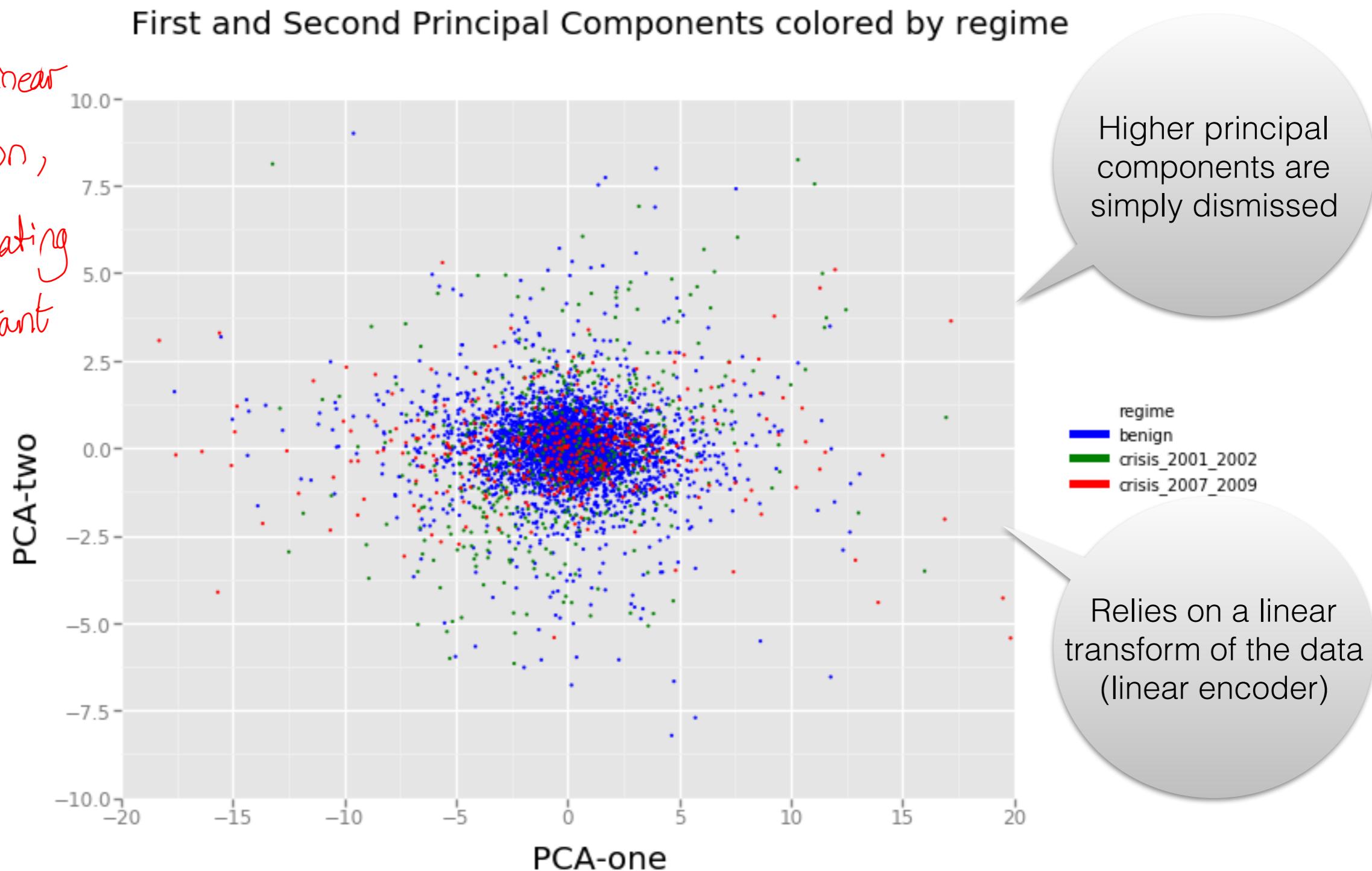
NYU Tandon School of Engineering, 2017

# 2D visualization: what can be improved?

Can we find a better 2D visualization of the DJI stock returns than the 2D PCA?

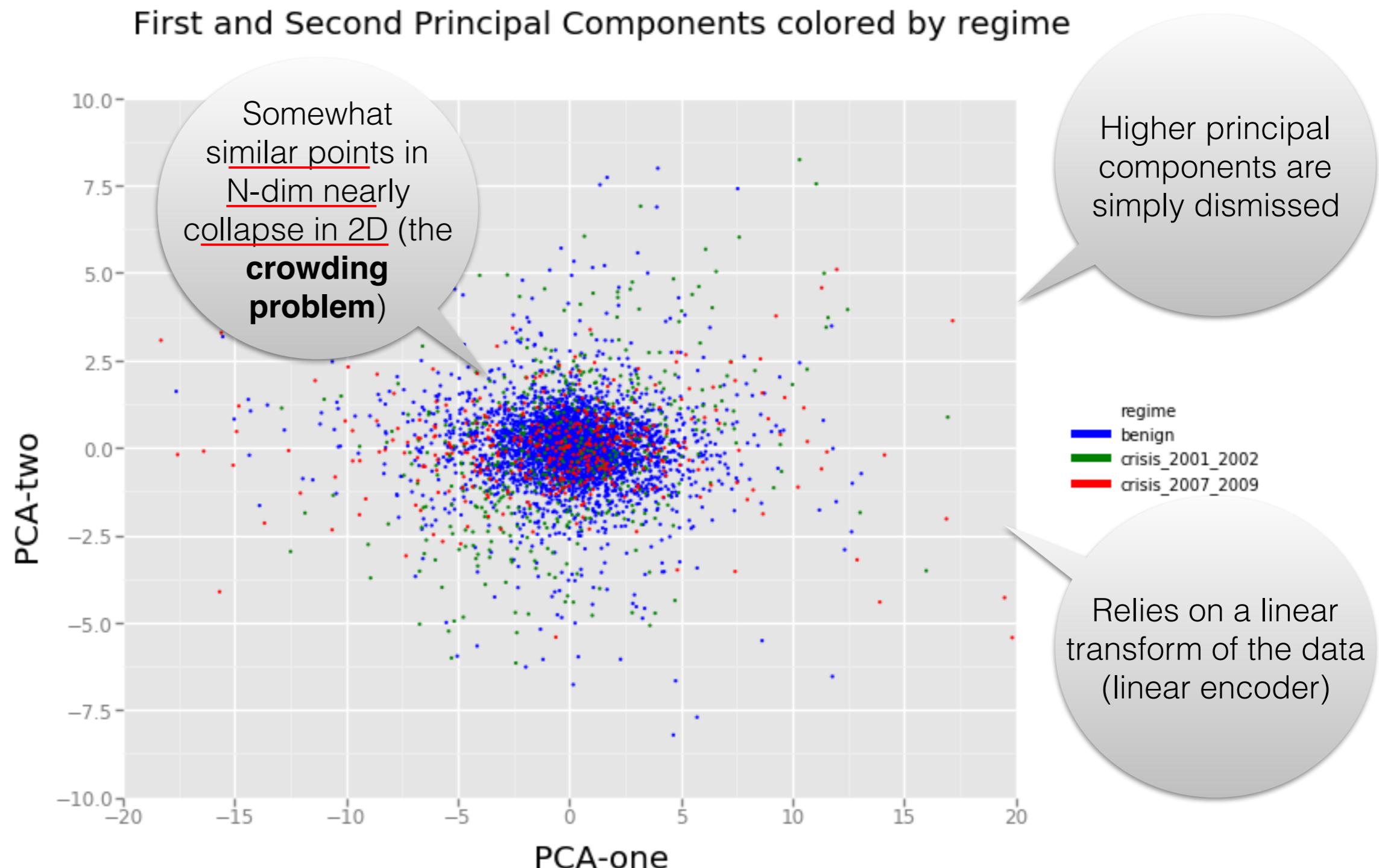
PCA

relies on linear  
transformation,  
while truncating  
less important  
dim



# 2D visualization: the “crowding problem”

Can we find a better 2D visualization of the DJI stock returns than the 2D PCA?



# Quiz: What do you think is the origin of the crowding problem?

1. The crowding problem is specific to stock analyses because of high level of noise in financial data
2. Two dimensions are just not enough to explain the amount of variation in multi-dimensional data
3. When you project multi-dimensional data on a plane, the need to display outliers forces non-outliers to crowd together
4. Because in higher dimensions there are typically more outliers than in lower dimensions, the crowding problem becomes worse with increased dimensions
5. A uniform distribution in an  $N$ -dimensional hyper-cube of size  $r$  has the number of points proportional to the volume of the hyper-cube  $V \sim r^N$ , therefore a *fixed* number of uniformly distributed observed points in  $N$  dimensions will generally translate into a *non-uniform* distributions of their projections onto a 2D plan.

## Multiple choice answers:

1. Answer 1
2. Answer 2
3. Answers 3 and 4
4. Answer 5
5. Answers 3, 4, 5
6. All of the above



## Quiz: What do you think is the origin of the crowding problem?

1. The crowding problem is specific to stock analyses because of high level of noise in financial data
2. Two dimensions are just not enough to explain the amount of variation in multi-dimensional data
3. When you project multi-dimensional data on a plane, the need to display outliers forces non-outliers to crowd together
4. Because in higher dimensions there are typically more outliers than in lower dimensions, the crowding problem becomes worse with increased dimensions
5. A uniform distribution in an  $N$ -dimensional hyper-cube of size  $r$  has the number of points proportional to the volume of the hyper-cube  $V \sim r^N$ , therefore a fixed number of uniformly distributed observed points in  $N$  dimensions will generally translate into a non-uniform distributions of their projections onto a 2D plan.

### Multiple choice answers:

1. Answer 1
2. Answer 2
3. Answers 3 and 4
4. Answer 5
5. Answers 3, 4, 5
6. All of the above

**Correct answer: 5**

# tSNE vs PCA

**t-SNE:** t-distributed Stochastic Neighbor Embedding (van der Maaten and Hinton 2008)

	PCA	t-SNE
Type of algorithm	Deterministic	Stochastic
Projection onto a low dimensional space	Linear	Non-linear
Global or local approach?	Global	Local/global
Handling discarded dimensions	Hard truncation of extra PCA dimensions	No truncation of dimensions
Unique solution?	Yes	No
Interpretability of results	Straightforward (PCA is just a rotation of axes...)	Subjective

# Probabilistic dimension reduction with t-SNE

matching 2 distributions  
data projection

1. Define **similarity** of point  $\mathbf{x}_j \in \mathbb{R}^D$  to point  $\mathbf{x}_i \in \mathbb{R}^D$  by assuming that it is generated by a **Gaussian distribution** centered at  $\mathbf{x}_i$

$$p_{j|i} = \frac{\exp\left(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\sigma_i^2)\right)}{\sum_{k \neq i} \exp\left(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / (2\sigma_i^2)\right)}$$

and define the **joint probability** as  $p_{ij} = (p_{j|i} + p_{i|j}) / (2n)$  ( $n$  is the total number of data points) - guaranteed to be between 0 and 1. Data-specific variances  $\sigma_i^2$  will be specified later...

2. Assume that their projections onto a 2D plane are samples from a **Student t-distribution** with one degree of freedom (= the **Laplace distribution**)

$$q_{ij} = \frac{\left(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2\right)^{-1}}{\sum_{k \neq i} \left(1 + \|\mathbf{y}_k - \mathbf{y}_j\|^2\right)^{-1}}$$

A fat-tailed distribution

Large distances are suppressed as a power law

# Probabilistic dimension reduction with t-SNE

3. Fix the data-dependent variances  $\sigma_i^2$  from the requirement that distributions of points  $\mathbf{x}_j$  centered around  $\mathbf{x}_i$  have a fixed perplexity (specified by the user!). Perplexity controls the effective number of neighbors for each point.)

$$Perp(\mathbf{P}_i) = 2^{H(\mathbf{P}_i)}, \quad H(\mathbf{P}_i) = -\sum_j p_{j|i} \log_2 p_{j|i}$$

4. Use the Kullback-Leibler (KL) divergence to quantify the collective dissimilarity between the set of points in  $D$  dimensions with their projections onto a plane:

$$C = KL(P \parallel Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

5. Optimize  $C$  with respect to the choice of points  $\{\mathbf{y}_i\}$  using a gradient descent method, with the following value of the gradient:

$$\frac{\partial C}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j), \quad i = 1, \dots, n$$

# The basic t-SNE algorithm

Laurens van der Maaten and Geoffrey Hinton, “Visualizing Data using t-SNE”, Journal of Machine Learning Research 9 (2008), 2579-2605.

---

**Algorithm 1:** Simple version of t-Distributed Stochastic Neighbor Embedding.

---

**Data:** data set  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ ,

cost function parameters: perplexity  $Perp$ ,

optimization parameters: number of iterations  $T$ , learning rate  $\eta$ , momentum  $\alpha(t)$ .

**Result:** low-dimensional data representation  $\mathcal{Y}^{(T)} = \{y_1, y_2, \dots, y_n\}$ .

**begin**

    compute pairwise affinities  $p_{j|i}$  with perplexity  $Perp$  (using Equation 1)

    set  $p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$

    sample initial solution  $\mathcal{Y}^{(0)} = \{y_1, y_2, \dots, y_n\}$  from  $\mathcal{N}(0, 10^{-4}I)$

**for**  $t=1$  **to**  $T$  **do**

        compute low-dimensional affinities  $q_{ij}$  (using Equation 4)

        compute gradient  $\frac{\delta C}{\delta \mathcal{Y}}$  (using Equation 5)

        set  $\mathcal{Y}^{(t)} = \mathcal{Y}^{(t-1)} + \eta \frac{\delta C}{\delta \mathcal{Y}} + \alpha(t) (\mathcal{Y}^{(t-1)} - \mathcal{Y}^{(t-2)})$

**end**

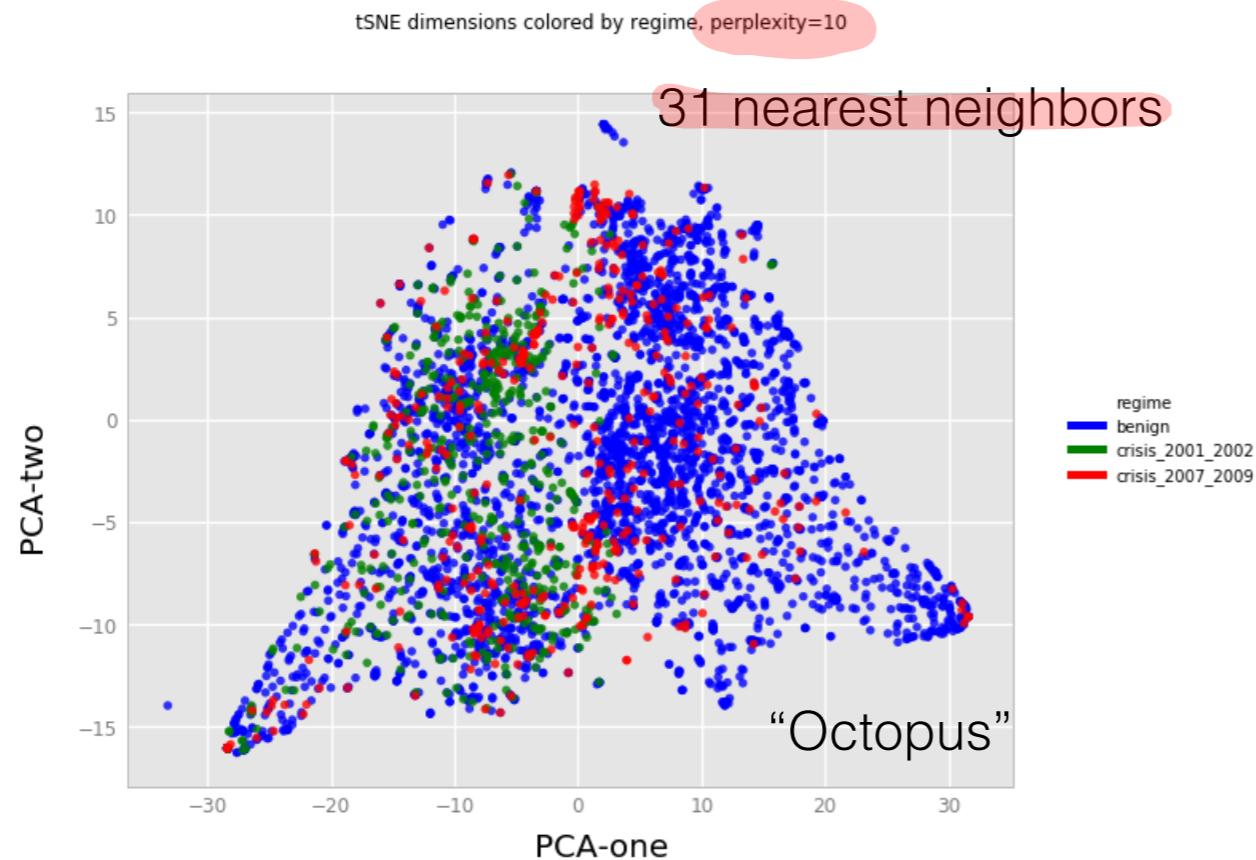
**end**



Momentum term

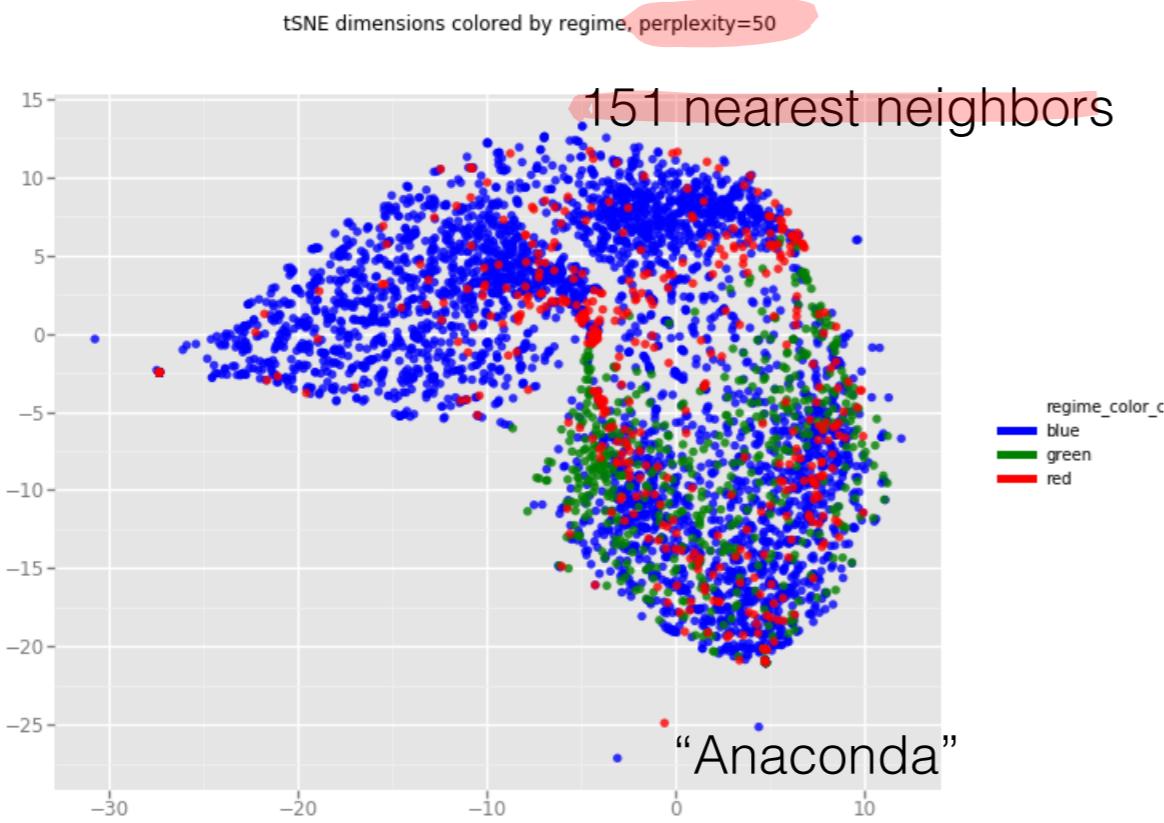
# t-SNE for stock returns

Apply t-SNE to the DJI stock return data, for different values of perplexity



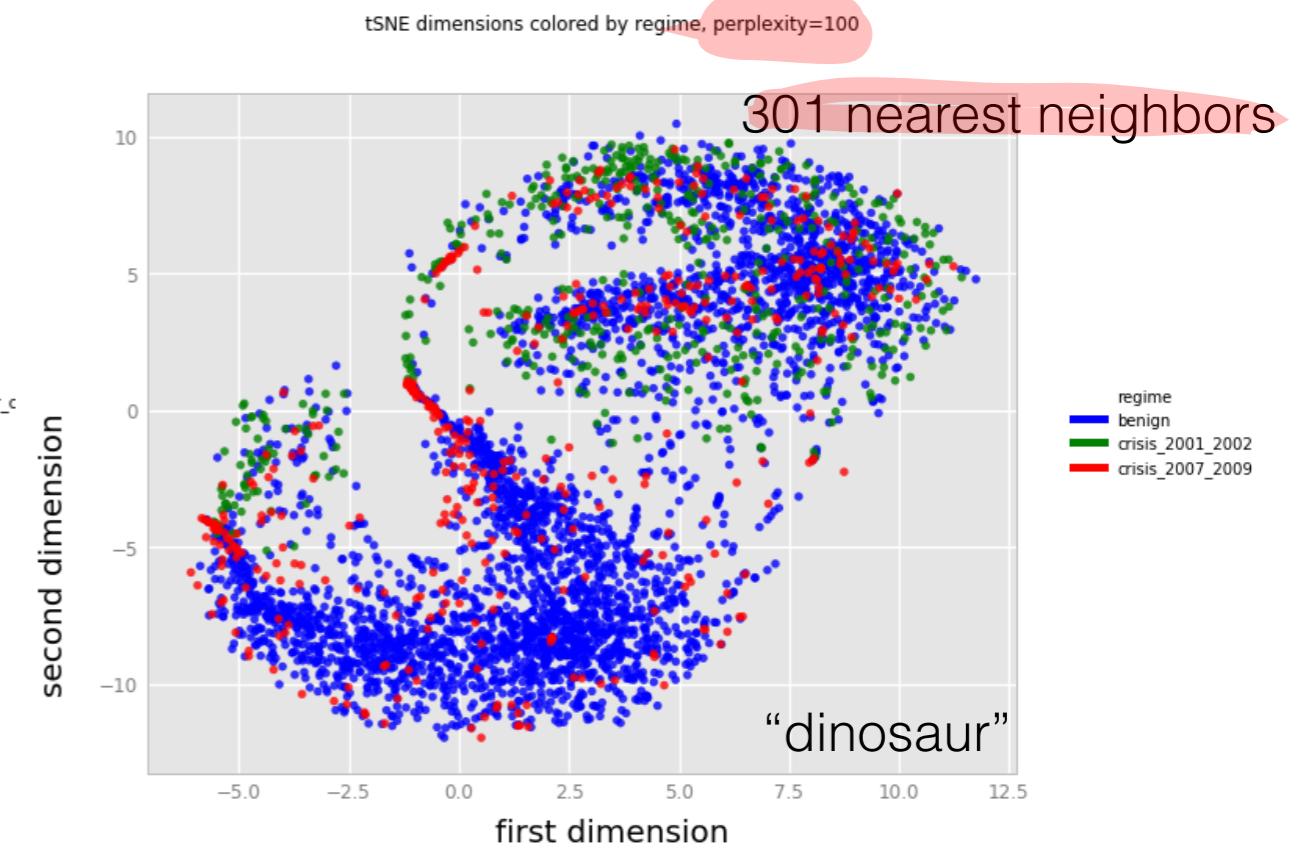
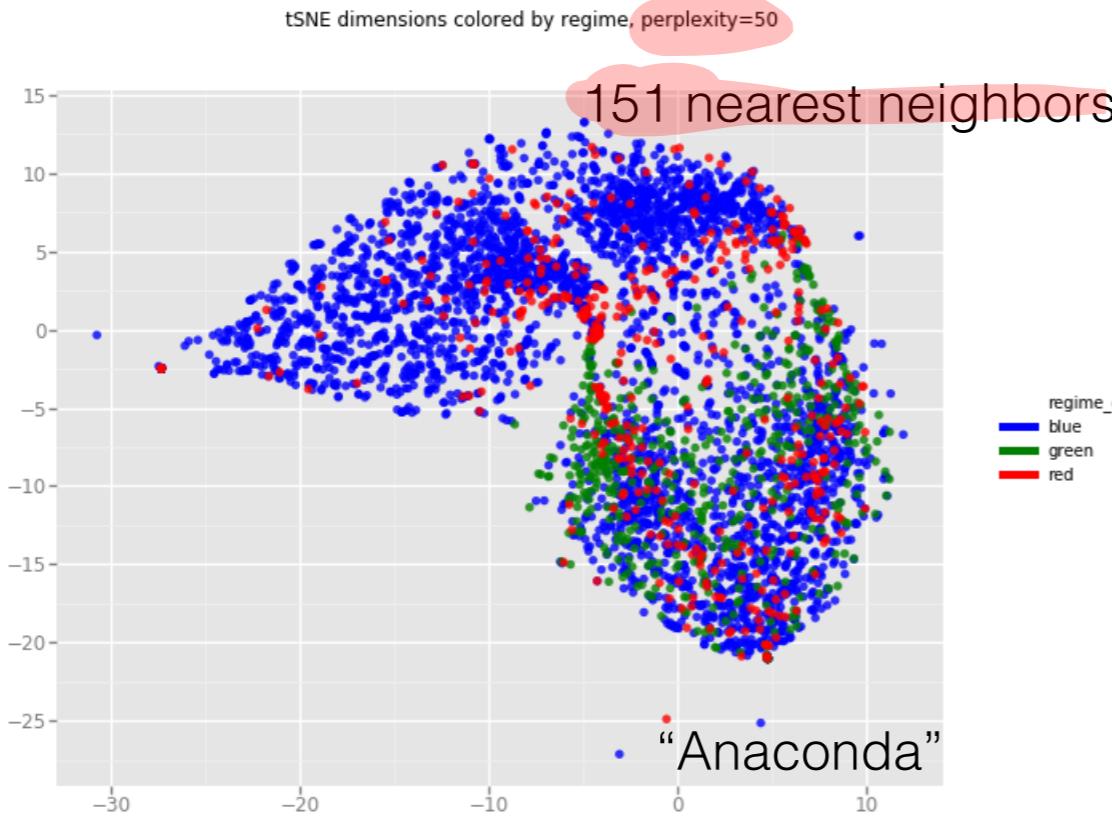
# t-SNE for stock returns

Apply t-SNE to the DJI stock return data, for different values of perplexity



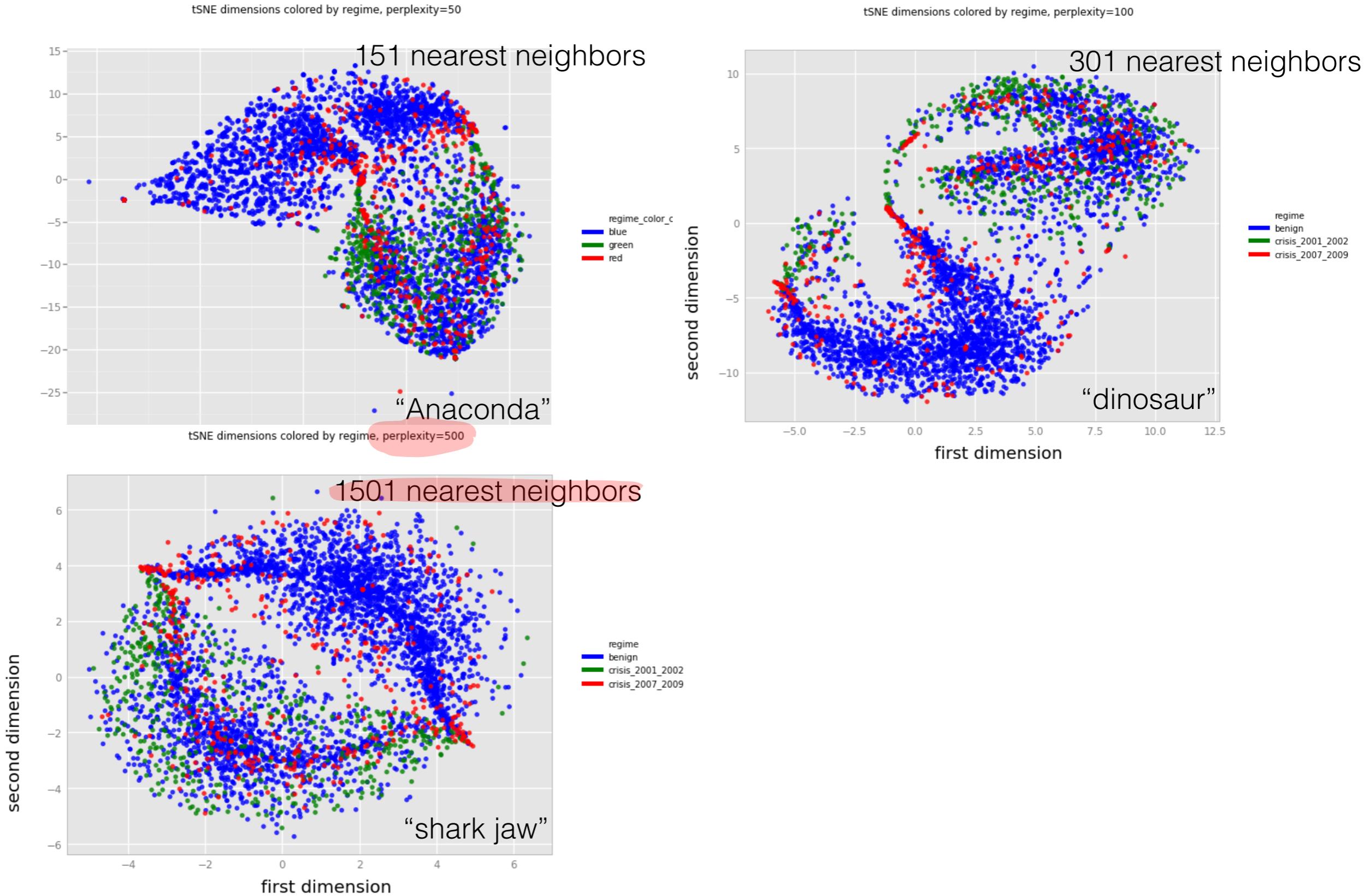
# t-SNE for stock returns

Apply t-SNE to the DJI stock return data, for different values of perplexity



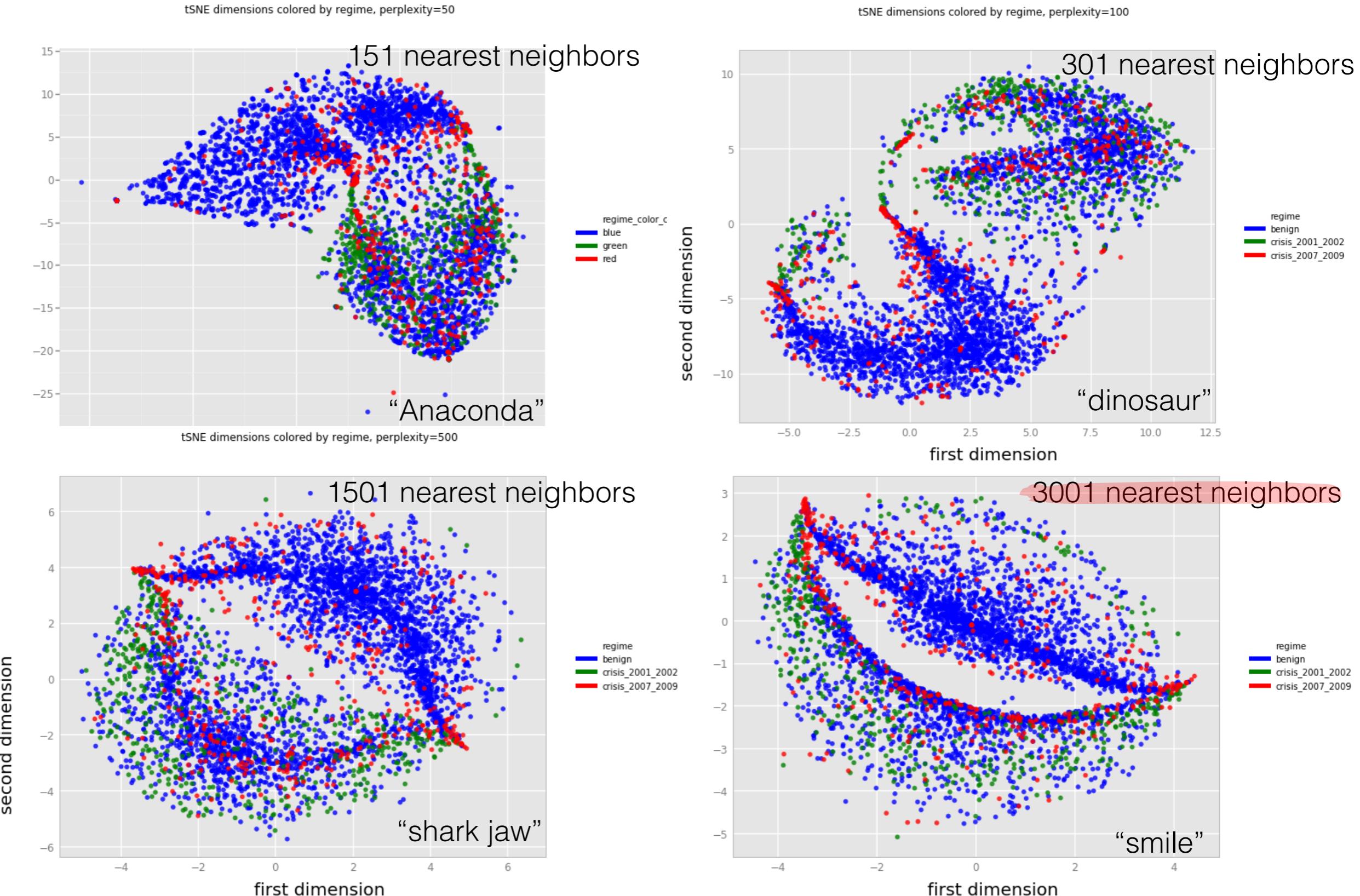
# t-SNE for stock returns

Apply t-SNE to the DJI stock return data, for different values of perplexity



# t-SNE for stock returns

Apply t-SNE to the DJI stock return data, for different values of perplexity



# t-SNE algorithm: further notes

- Visualization results depend on perplexity. What perplexity would be optimal?
-  In its basic form, t-SNE provides visualization of a given dataset but not a dimensional reduction recipe/algorithm for a new, unseen data. New data points cannot be embedded into an existing low dimensional representation.
- t-SNE preserves nearest neighbors but not distance. Running distance or density based clustering algorithms on outputs of the t-SNE can be problematic!
- The basic algorithm is not well scalable with the number of data points (has complexity (scales as  $N^2$ )). Additional tricks are required to apply it to large datasets, e.g. a tree-based algorithm of van der Maaten (2014) has  $O(N \log N)$  complexity

# Guided Tour of Machine Learning in Finance

## Week 3: Unsupervised Learning

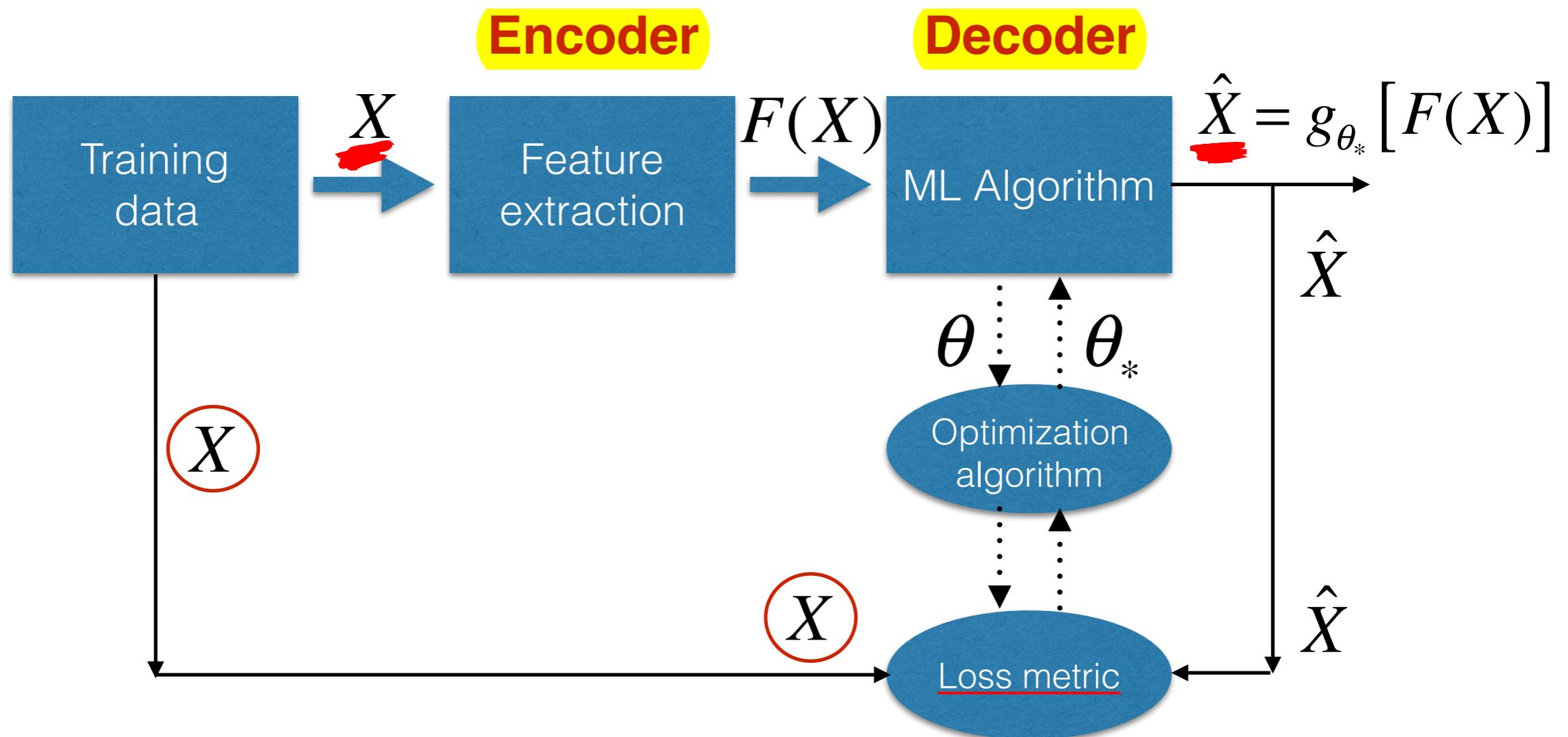
### 3-1-5-Autoencoder

Igor Halperin

NYU Tandon School of Engineering, 2017

# Unsupervised learning diagram: Autoencoder

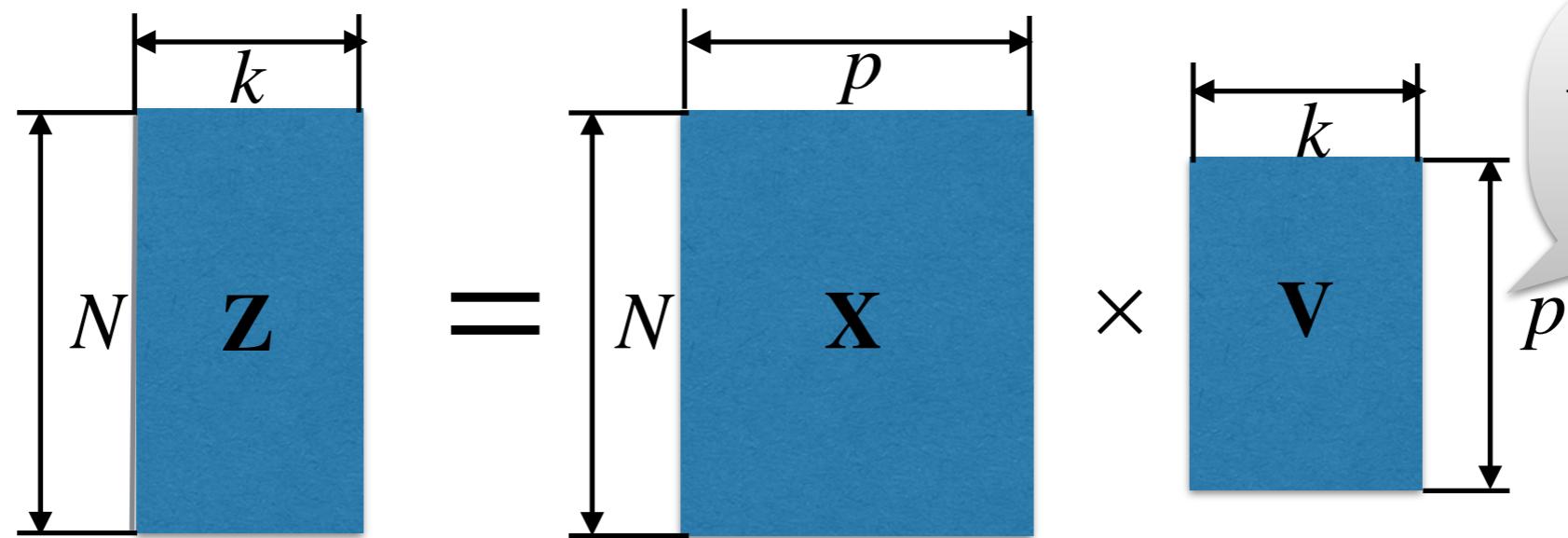
**Autoencoder:** replace  $Y$  with  $X$ !



# PCA as a dimension reduction method

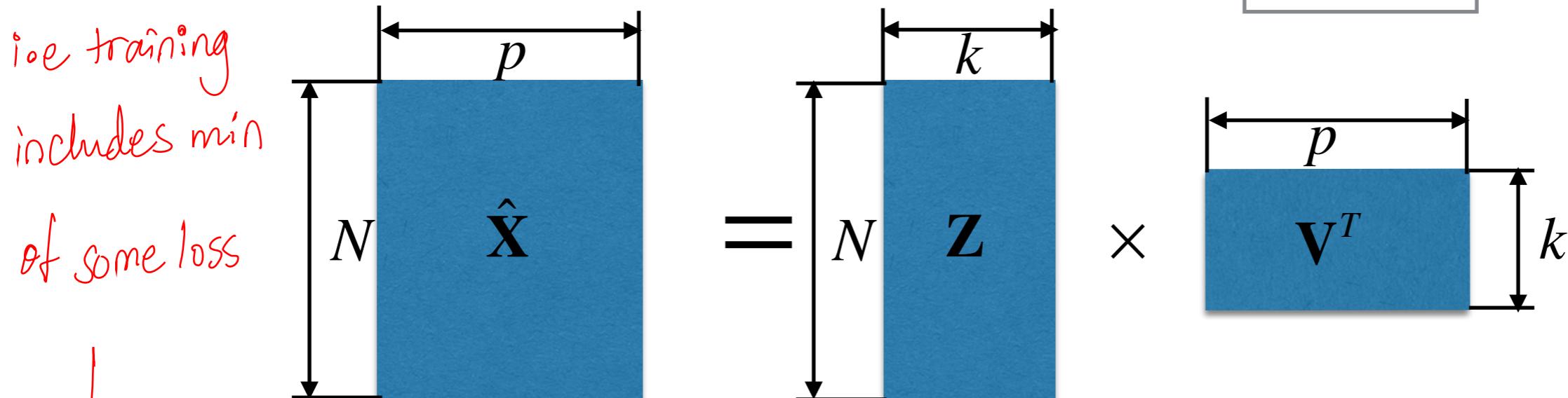
We had a linear transform (linear encoder) of the data  $\mathbf{Z} = \mathbf{X}\mathbf{V}$  parametrized by a  $p \times k$  orthogonal matrix  $\mathbf{V}$  with  $\mathbf{V}\mathbf{V}^T = 1$

Special  
Linear  
Autoencoder



If  $k = p$ ,  $\mathbf{V} = \mathbf{U}$   
this is just a rotation  
of the eigenvalue  
decomposition

A decoded signal (inverse transform) is obtained as  $\hat{\mathbf{X}} = \mathbf{Z}\mathbf{V}^T$

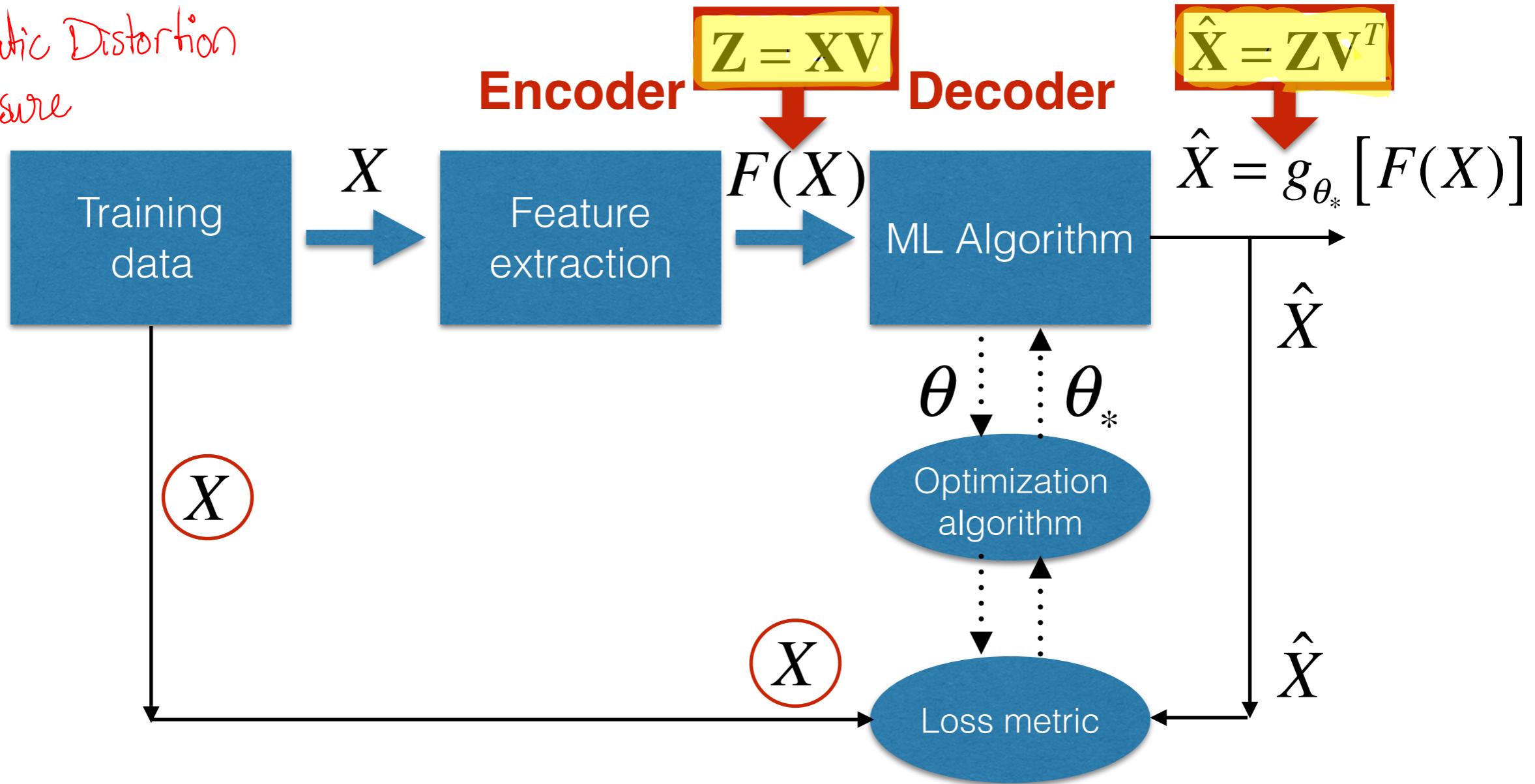


2<sup>nd</sup> interpretation → reconstruction error minimization

## Autoencoder

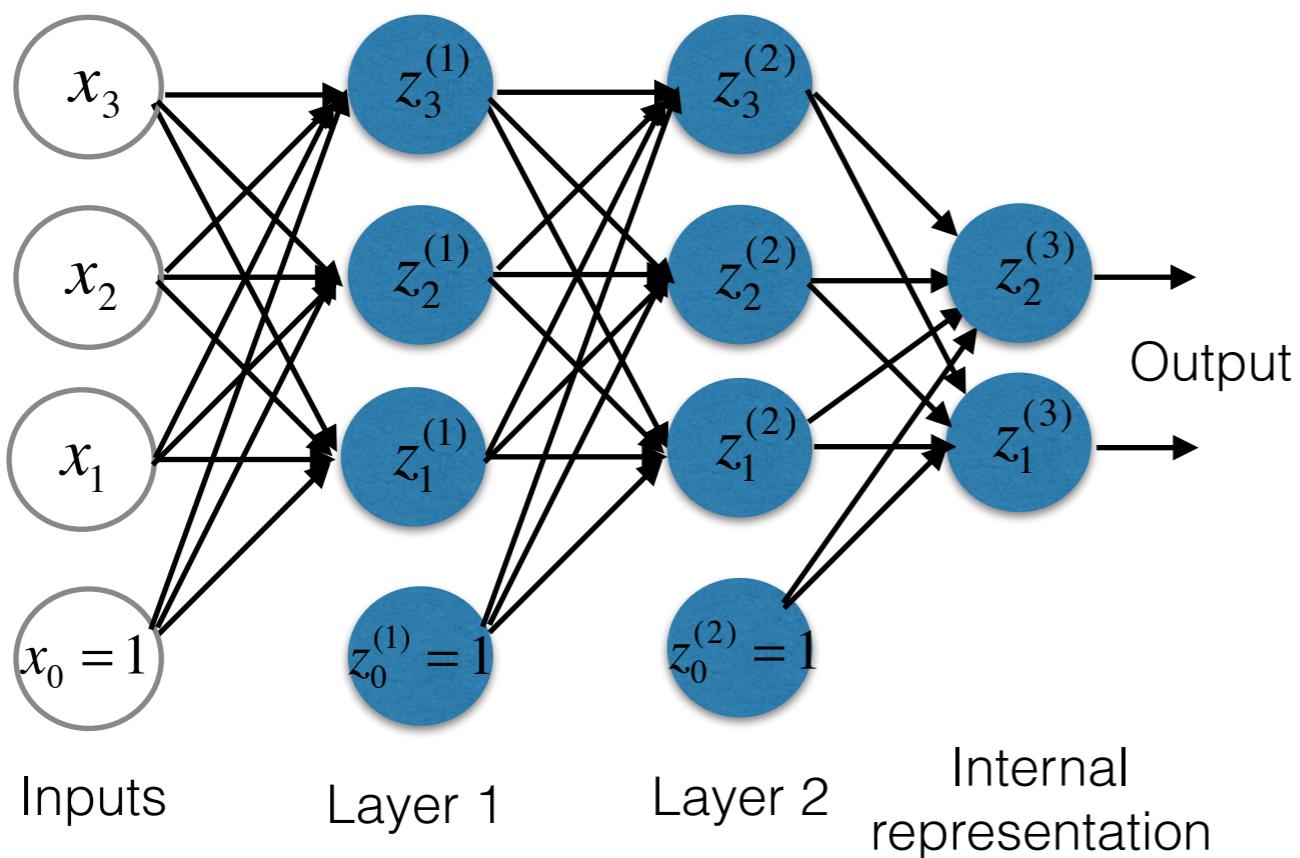
min the MSE error bet  $X$  & reconstructed  $\hat{X}$

Quadratic Distortion Measure



# Neural Encoder

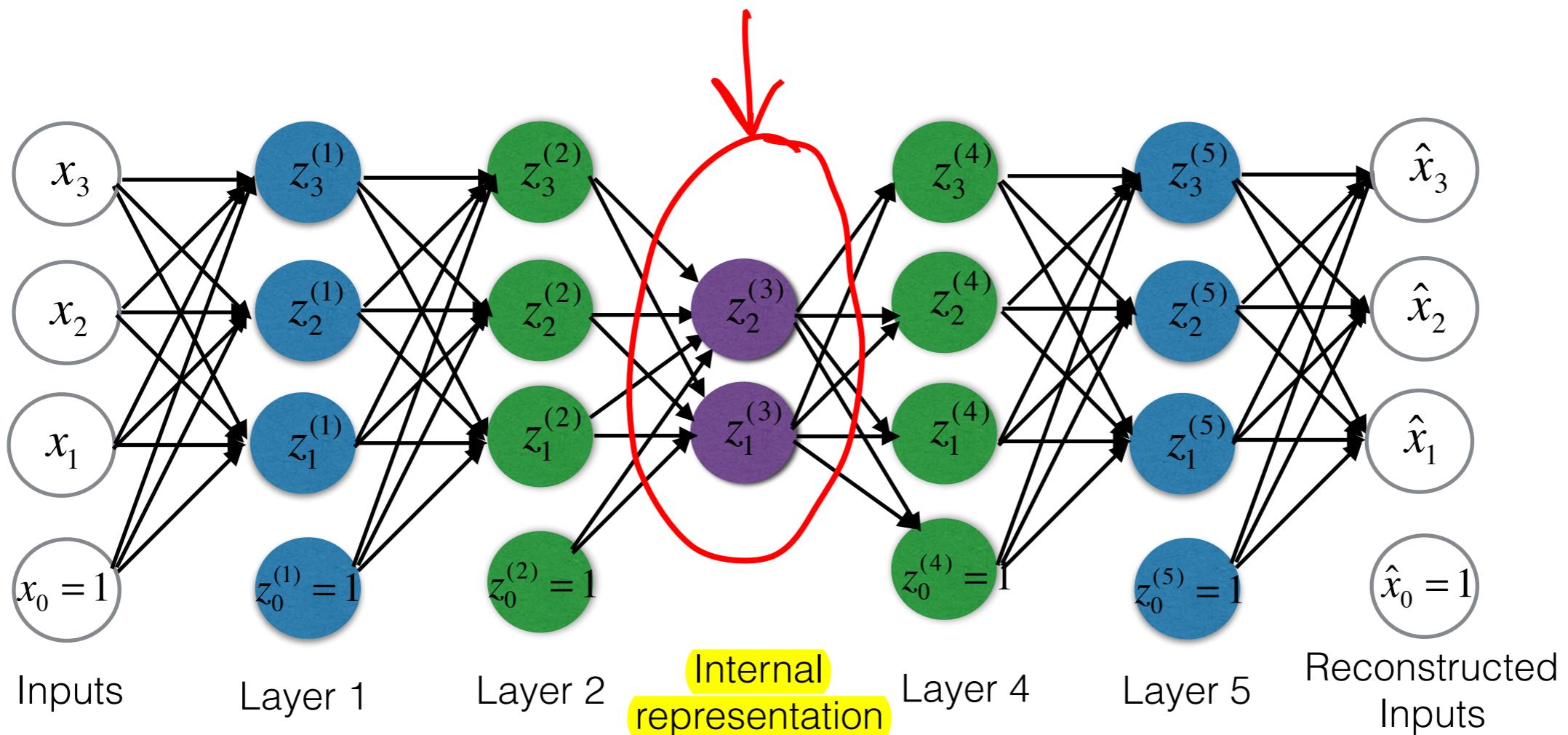
- A Neural Encoder is composed of layers of artificial neurons (perceptrons), with a final layer providing a low-dimensional representation of data



- For linear activation functions in all layers, this is equivalent to a single linear layer, which is equivalent to the Linear Encoder of the PCA
- If we use non-linear activations, then the depth of the network matters. For complex data, many layers might be needed to build a good internal representation.

# Neural Encoder-Decoder

- A Neural Encoder/Decoder pair implements an Autoencoder:

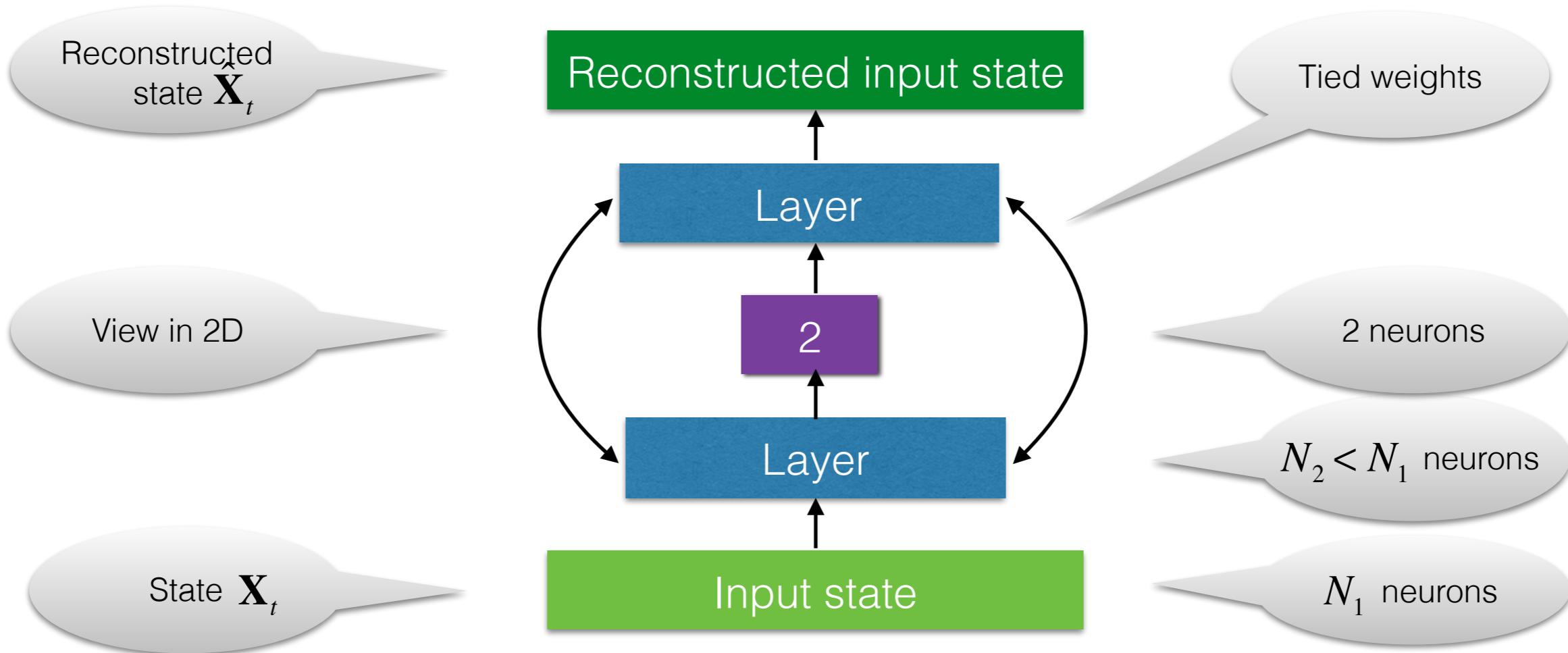


- If the inner-most layer has only two neurons, their activation can be plotted for visualization of data with the autoencoder.
- A Stacked (multi-level) autoencoder may have layers with gradually decreasing number of neurons per layer
- Tied weights in the Encoder and Decoder are often used in practice

# Visualizing data with autoencoder



Autoencoder provides a way to visualize a high-dimensional input vector by a non-linear projection onto a 2D space of neuron activations for the inner-most layer with 2 neurons:



Autoencoder is trained by minimizing the **reconstruction error** between input state  $\mathbf{X}_t$  and reconstructed state  $\hat{\mathbf{X}}_t$

# Control question

Select all correct answers

- \* 1. The PCA implements a Linear Encoder that converts correlated inputs  $\mathbf{X}$  into uncorrelated features  $\mathbf{Z}$  by a linear transform  $\mathbf{Z} = \mathbf{X}\mathbf{V}$
2. If the orthogonal matrix  $\mathbf{V}$  has dimension  $p \times p$ , i.e. it keeps all eigenvectors of correlation matrix of  $\mathbf{X}$ , then  $\mathbf{Z}$  preserves the total variation of  $\mathbf{X}$ .
3. When the data is noisy or non-stationary,  $\mathbf{Z}$  can have a higher variance than  $\mathbf{X}$ .
4. The PCA is a probabilistic method that enables simulating from data.

**Correct answer:** 1, 2