



# Data Mining



## Chapter 8: Outlier Detection

**Yunming Ye, Baoquan Zhang**

**School of Computer Science**

**Harbin Institute of Technology, Shenzhen**

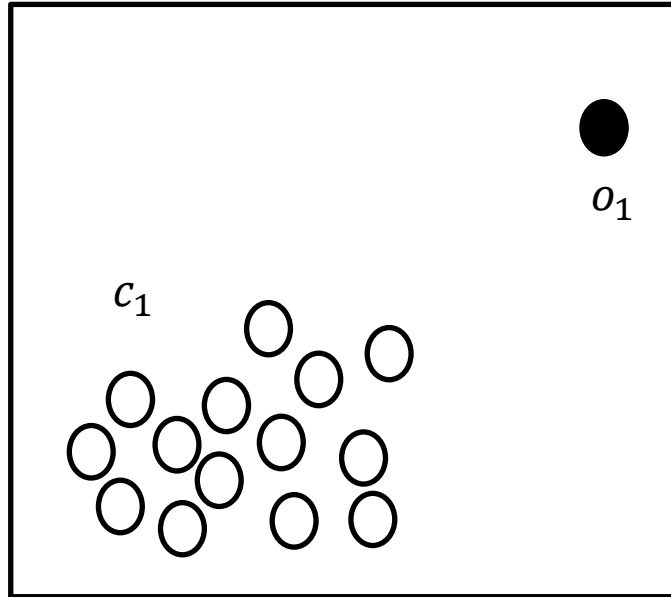
# Agenda

- Introduction to Outlier Detection
- Maximum Likelihood Method
- One-class SVM and Isolation Forest
- Reconstruction Methods

## **8.1 Introduction to Outlier Detection**

# What Is Outlier Detection?

- An outlier is a data object that *deviates significantly* from the rest of the objects



- **Outlier** detection (**anomaly detection**) is the process of finding data objects with behaviors that are very different from expectation

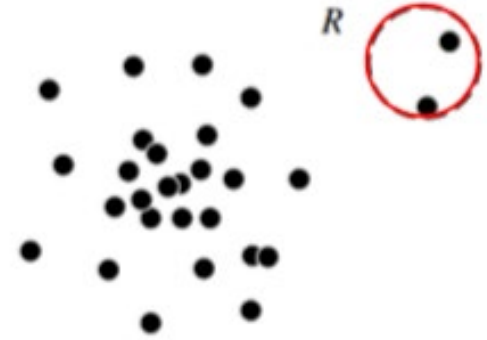
# Applications of Outlier Detection

- Network intrusion detection
- Insurance / Credit card fraud detection
- Healthcare Informatics / Medical diagnostics
- Industrial Damage Detection
- Image Processing / Video surveillance
- Novel Topic Detection in Text Mining
- ...

# Types of outliers

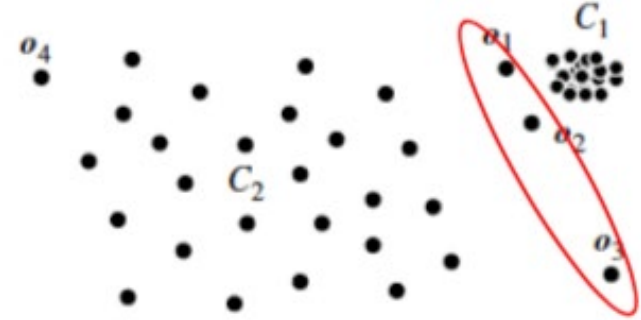
- **Global outliers (point anomalies)**

- A global outlier deviates significantly from the rest of the data set



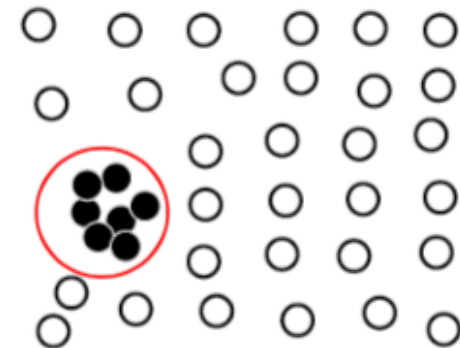
- **Contextual outliers**

- A contextual outlier deviates significantly w.r.t a specific context of the object



- **Collective outliers**

- A collective outlier refer to objects as a whole deviate significantly from the entire data set



# Challenges of Outlier Detection

- **Modeling normal objects and outliers effectively**
- **Application-specific outlier detection**
- **Handling noise in outlier detection**
- **Understandability**

# Methods for Outlier Detection

- Supervised methods
  - train classifier with “normal” and “abnormal” data
  - Challenge: unbalanced data
- Unsupervised methods
  - Proximity-based: an outlier’s nearest neighbors should be far away



# Nearest-Neighbor Based Approach

- Compute the distance between every pair of data points
- There are various ways to define outliers:
  - Data points for which there are fewer than  $p$  neighboring points within a distance  $D$
  - The top  $n$  data points whose distance to the  $k$ th nearest neighbor is greatest
  - The top  $n$  data points whose average distance to the  $k$  nearest neighbors is greatest

# Methods for Outlier Detection

- Supervised methods

- train classifier with “normal” and “abnormal” data
- Challenge: unbalanced data

- Unsupervised methods

- Proximity-based: an outlier’s nearest neighbors should be far away
- Clustering-based: normal data belonging to large and dense clusters
- One-class Method
  - ✓ Statistical method : data normality from some statistical model
  - ✓ Other one-class methods
- Reconstruction method
- .....

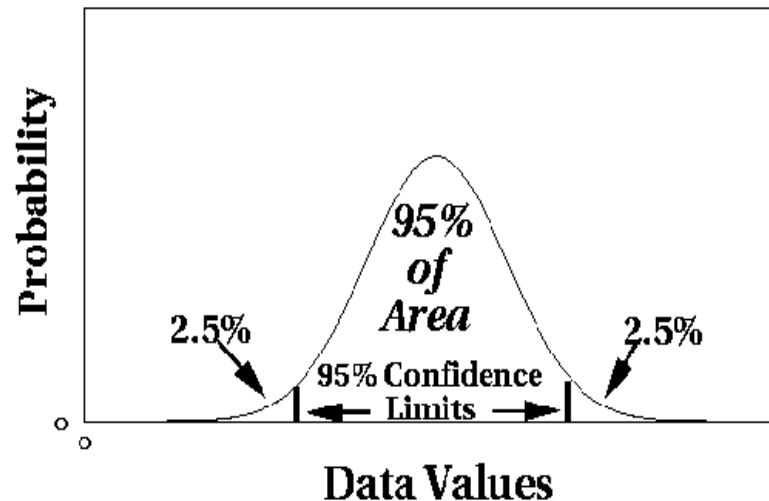
# Outlier detection using one-class model

- **Only modeling the normal class (with large amount of objects).**
  - Centroid-based method
  - Statistical method
  - One-class SVM
  - Isolation forest
  - .....

## **8.2 Maximum Likelihood Method**

# Statistical Approaches

- Assume a parametric model describing the distribution of the data (e.g., normal distribution)
- Apply a statistical test that depends on
  - Data distribution
  - Parameter of distribution (e.g., mean, variance)
  - Number of expected outliers (confidence limit)



# Maximum Likelihood Method: problem definition

- Given a data set  $\mathbb{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$
- Assume the probability density function of  $\mathbb{X}$  is known to be  $f_{\theta}(x)$ 
  - $\theta$  is the parameters, and to be learned from data

- Likelihood of  $\mathbb{X}$ :
$$L(\theta) = \prod_{i=1}^m f_{\theta}(\mathbf{x}^{(i)}) = f_{\theta}(\mathbf{x}^{(1)})f_{\theta}(\mathbf{x}^{(2)})\dots f_{\theta}(\mathbf{x}^{(m)})$$

- Maximize:
$$\theta^*: \operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} \left( \prod_{i=1}^m f_{\theta}(\mathbf{x}^{(i)}) \right)$$

- Given a data set  $\mathbb{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$   $\mathbf{x}^{(i)} \in \mathbf{R}^d$
- Assume the distribution of  $\mathbb{X}$  be Gaussian:

$$\theta: \boldsymbol{\mu}, \boldsymbol{\sigma} \quad f_{\boldsymbol{\mu}, \boldsymbol{\sigma}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\boldsymbol{\sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- Maximize:

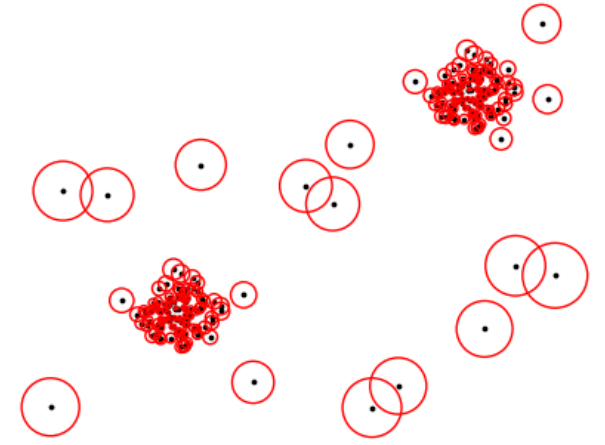
$$\theta^*: \operatorname{argmax}_{\theta: \boldsymbol{\mu}, \boldsymbol{\sigma}} L(\theta) = \operatorname{argmax}_{\boldsymbol{\mu}, \boldsymbol{\sigma}} \left( \prod_{i=1}^m f_{\boldsymbol{\mu}, \boldsymbol{\sigma}}(\mathbf{x}^{(i)}) \right)$$

$$\boldsymbol{\mu}^* = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}$$

$$\boldsymbol{\sigma}^* = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu}^*)(\mathbf{x}^{(i)} - \boldsymbol{\mu}^*)^T$$

# Detection Phase

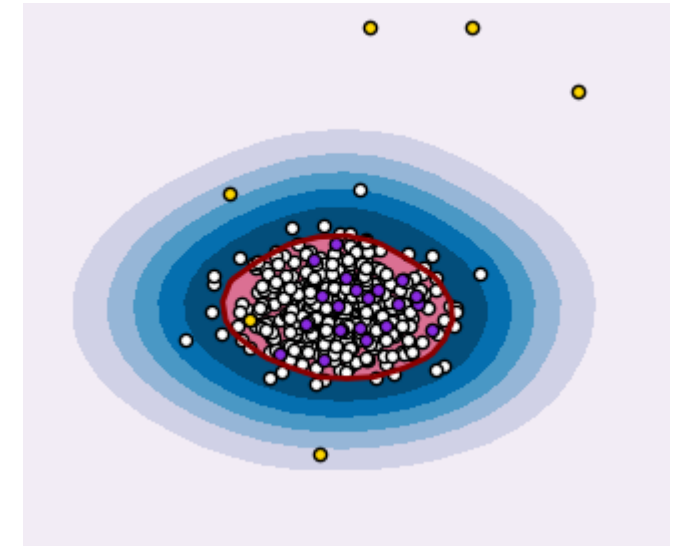
- Given a threshold value  $\delta$
- Decide whether a given data object  $x$  is an outlier



$$f_{\mu^*, \sigma^*}(x) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\sigma^*|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^*)^T \sigma^{*-1} (x - \mu^*) \right\}$$

$$\left\{ \begin{array}{ll} f_{\mu^*, \sigma^*}(x) \geq \delta & \text{Normal data object} \\ f_{\mu^*, \sigma^*}(x) < \delta & \text{Outlier} \end{array} \right.$$

$\delta$  determine the results

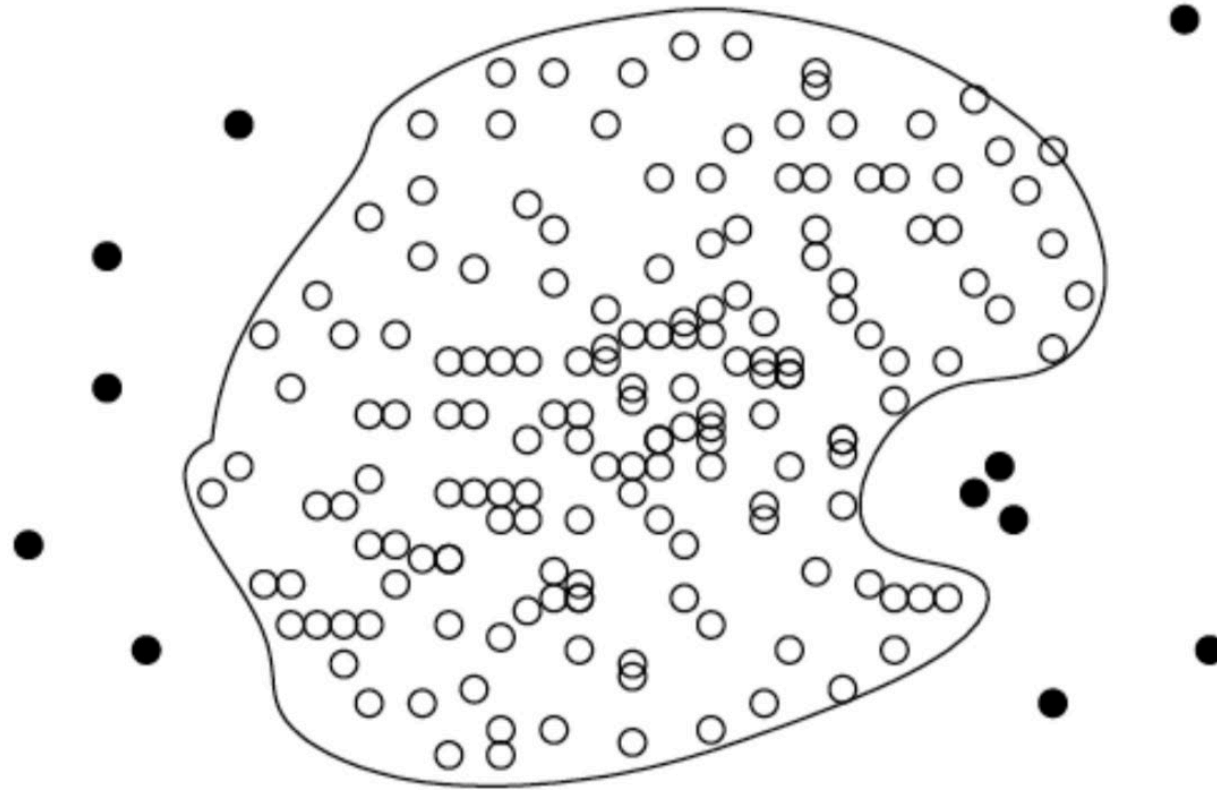




## 8.3 One-class SVM and Isolation Forest

# One-class SVM: Basic Idea

- **Learning the Boundary of the input data (i.e. normal class with large amount of objects).**

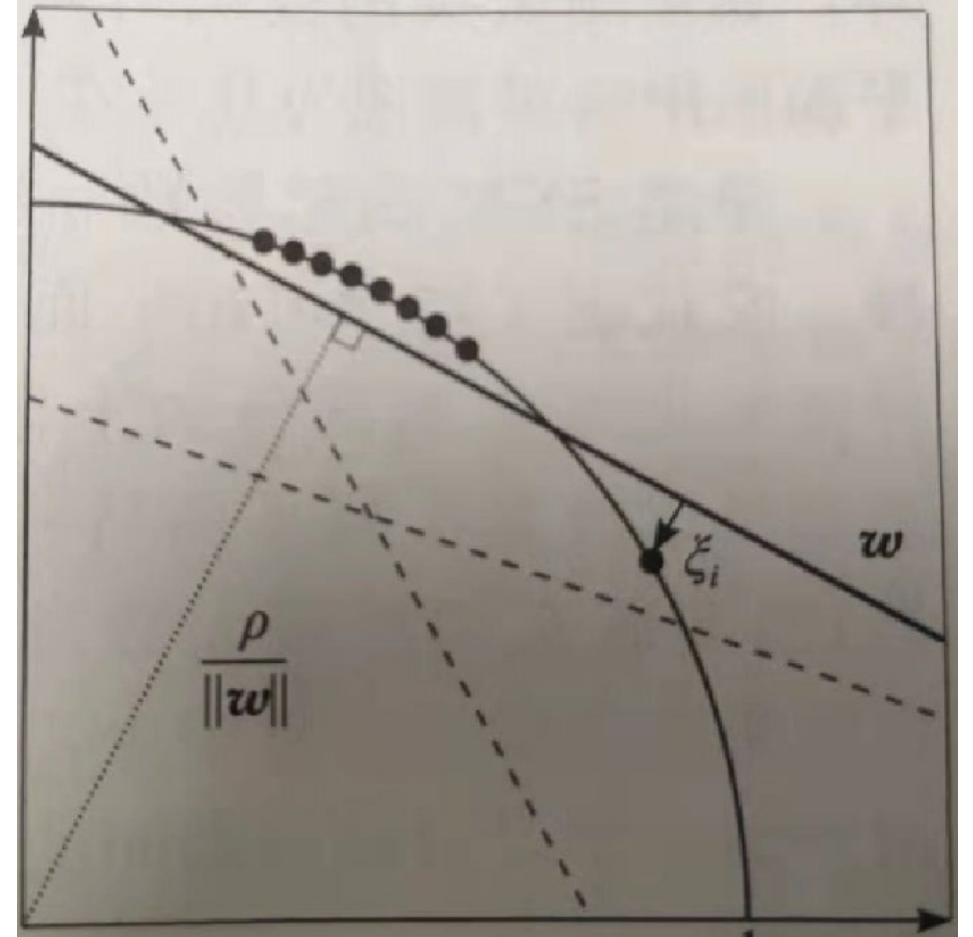


# One-class SVM: $\nu$ -SVM

- How to define the second class?
  - $\phi(x)$  : Projection to high-dimensional feature space
  - With Gaussian kernel:

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right)$$

$$K(\mathbf{x}, \mathbf{x}) = \langle \phi(\mathbf{x}), \phi(\mathbf{x}) \rangle = \|\phi(\mathbf{x})\|^2 = 1$$



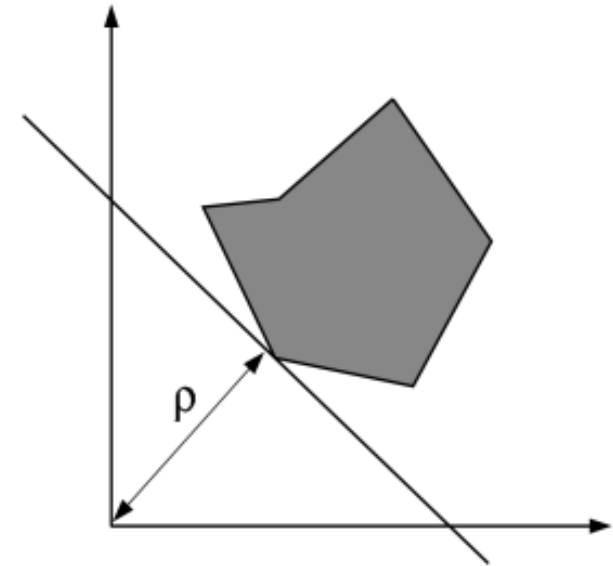
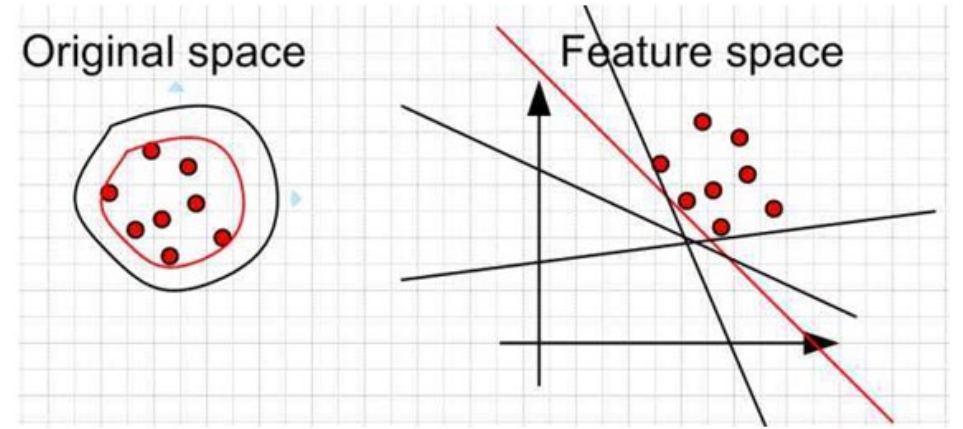
- the origin of the feature space is the second class!

- Define the separating hyperplane:

$$w \cdot \phi(x) = \rho$$

- Outlier detection function:

$$f(x) = \text{sgn}(\langle w, \phi(x) \rangle - \rho)$$



- Learning task:

- Learned model:  $w \cdot \phi(x) = \rho$

- Objective function:

$$\left. \begin{aligned} \min_{w, \xi, \rho} \quad & \frac{1}{2} \|w\|^2 + \frac{1}{vm} \sum_{i=1}^m \xi_i - \rho \\ \text{subject to: } & \langle w, \phi(x_i) \rangle \geq \rho - \xi_i, \xi_i \geq 0 \\ & 0 < v \leq 1 \end{aligned} \right\}$$

$$K(x, y) = \exp \left( -\frac{\|x - y\|^2}{2\sigma^2} \right)$$

$$w^* = \sum_{i=1}^m \lambda_i^* \phi(x_i)$$

$$\begin{aligned} \rho^* &= \sum_{i=1}^m \lambda_i^* \phi(x_i) \phi(x_s) \\ &= \sum_{i=1}^m \lambda_i^* K(x_i, x_s) \end{aligned}$$

- Outlier detection:

$$f(x) = \text{sgn}(\langle w, \phi(x) \rangle - \rho)$$

$$= \text{sgn} \left( \sum_{i=1}^m \lambda_i^* \phi(x_i) \phi(x) - \rho \right) = \text{sgn} \left( \sum_{i=1}^m \lambda_i^* K(x_i, x) - \rho \right)$$

# One-class SVM: *SVDD*

- Support Vector Domain Description:

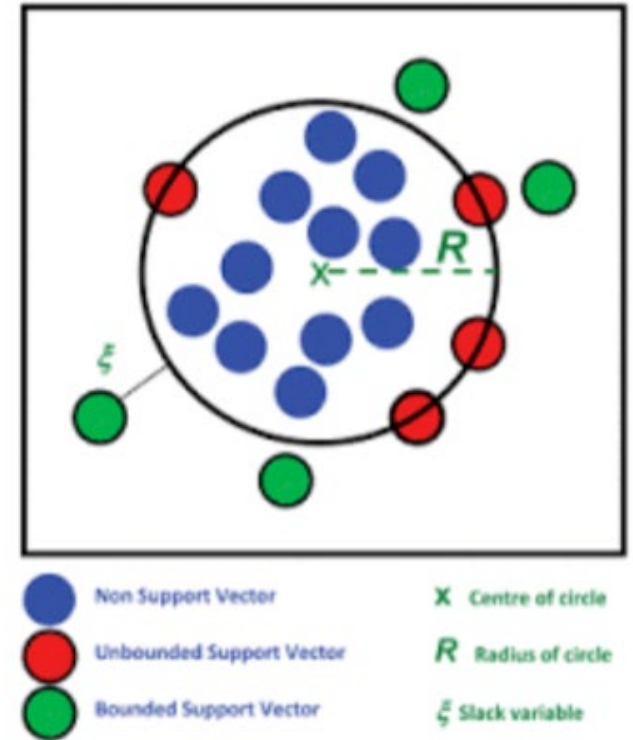
Constraining "normal" data in a ball with relative small radius

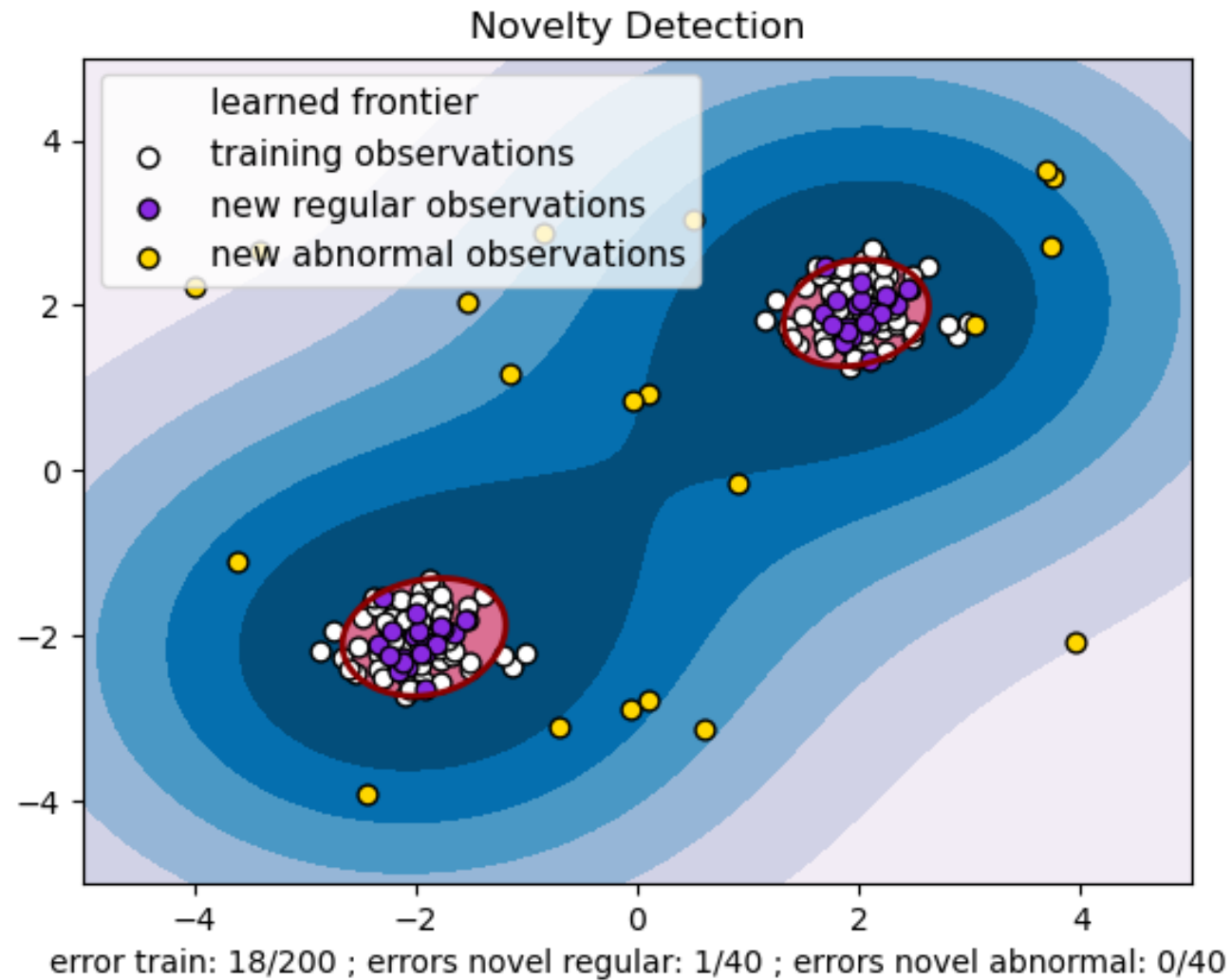
$$\min_{R, \xi, \mathbf{c}} \mathbf{R}^2 + \frac{1}{vm} \sum_{i=1}^m \xi_i \quad R \in \mathbb{R}, \xi \in \mathbb{R}^m, \mathbf{c} \in \mathcal{H}$$

subject to:

$$\|\phi(x_i) - \mathbf{c}\|^2 \leq R^2 + \xi_i, \quad \xi_i \geq 0$$

$$0 < v \leq 1$$

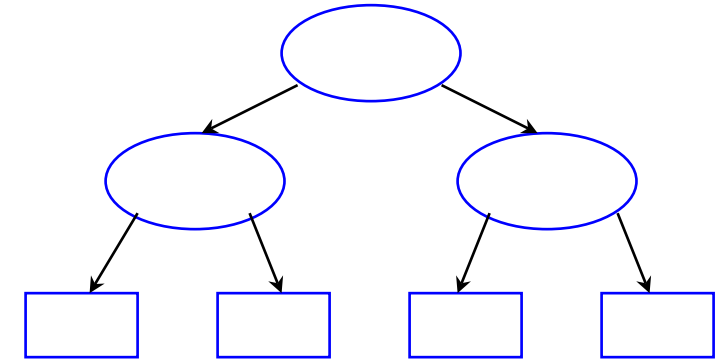
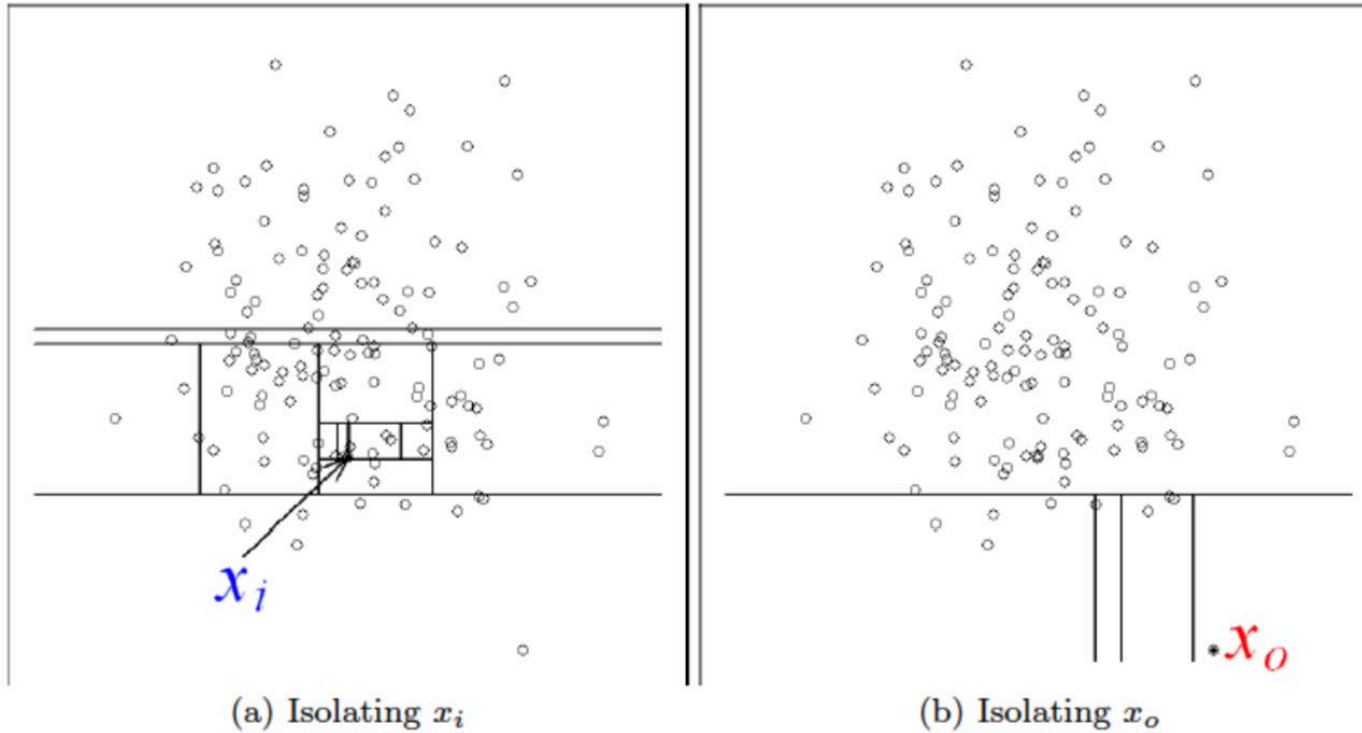




[https://scikit-learn.org/stable/auto\\_examples/svm/plot\\_oneclass.html](https://scikit-learn.org/stable/auto_examples/svm/plot_oneclass.html)

# Isolation-based outlier detection

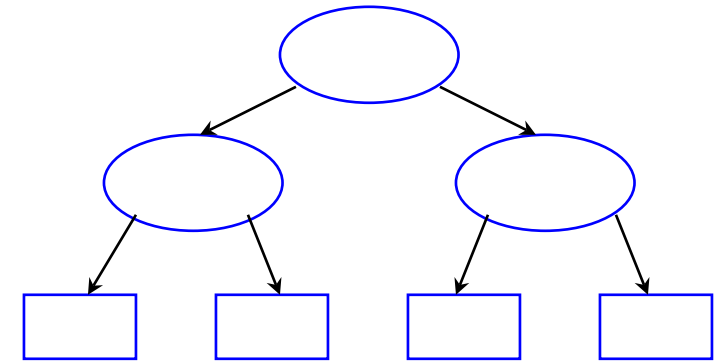
- Outliers are few and different
- when randomly split the space into small region, an outlier is **more likely to be *ISOLATED***





- F. T. Liu, K. M. Ting and Z. H. Zhou, ***Isolation-based Anomaly Detection. ACM TKDD***, 2011

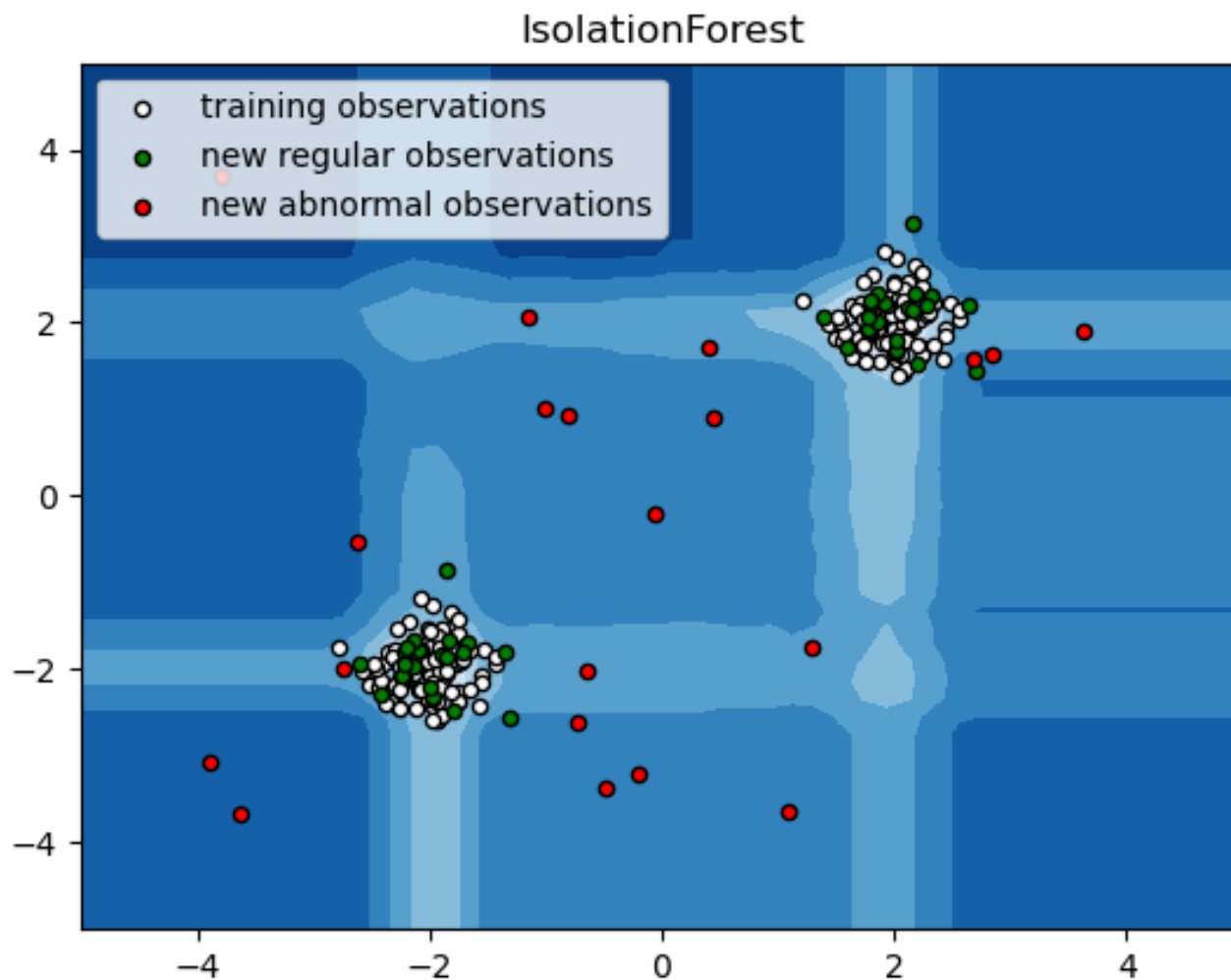
**Key idea:** Modeling “Isolation” using tree structure, and characterize the outlier suspiciousness with the path length from the root to the isolated object



- General steps:
  - Randomly subsample the original data
  - On each subsampled data, grow an *iTree* by
    - ✓ Randomly pick an attribute and a split value between min and max
    - ✓ Split the data into two subtrees
    - ✓ This process iterates until “isolation” is reached (no more points to be split or instances share the same value)
  - Compute outlier score by consulting the average path length from root to the isolated objects

# Output of outlier detection

- Label
  - Each test instance is given a normal or anomaly label
- Score
  - Each test instance is assigned an anomaly score
    - ✓ Allows the output to be ranked
    - ✓ Requires an additional threshold parameter



[https://scikit-learn.org/stable/auto\\_examples/ensemble/plot\\_isolation\\_forest.html#sphx-glr-auto-examples-ensemble-plot-isolation-forest-py](https://scikit-learn.org/stable/auto_examples/ensemble/plot_isolation_forest.html#sphx-glr-auto-examples-ensemble-plot-isolation-forest-py)

## 8.4 Reconstruction Method

# Reconstruction Method: basis idea

$$x_i \longrightarrow z_i \longrightarrow \hat{x}_i \quad x_i, \hat{x}_i \in R^n \quad z_i \in R^{d'}$$

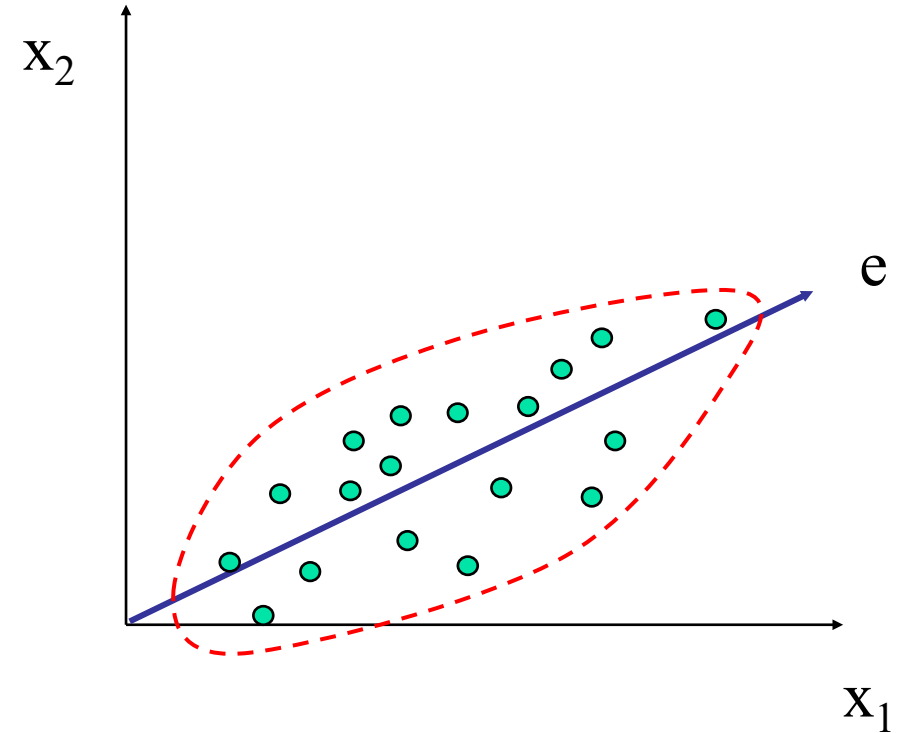
$$d' \ll n$$

Minimize:

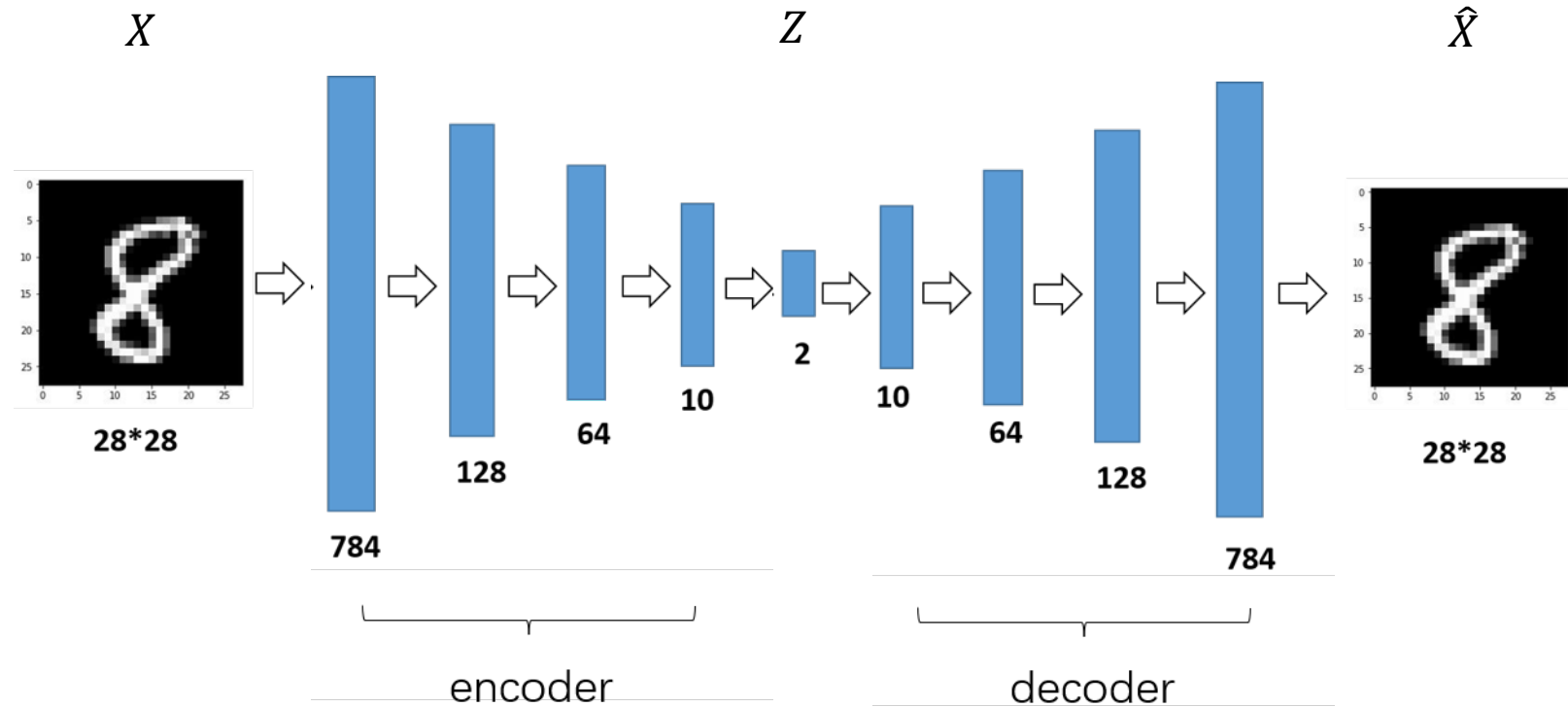
$$e = \|x_i - \hat{x}_i\|^2$$

# Reconstruction with PCA

- Find a projection that captures the largest amount of variation in data.
- Project and reconstruct  $x_i$  with PCA
- Limitation of PCA method
  - **Can only model linear combination of original features**



# Reconstruction with Autoencoder



# ● Training objective

➤ encoder: map  $x$  to low dimensional  $z$ .

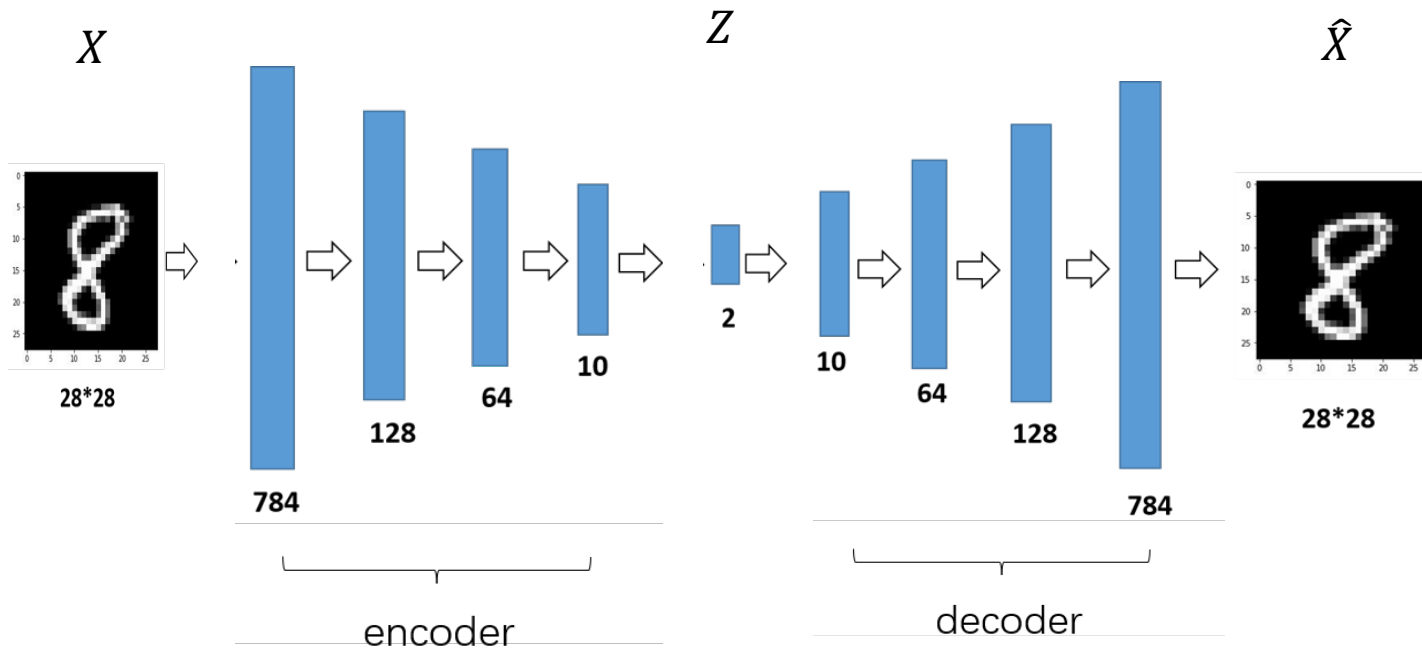
$$z = f_{\phi}(x)$$

➤ decoder: reconstruct  $x$  based on  $z$ .

$$\hat{x} = g_{\theta}(z)$$

➤ Objective function: minimize reconstruction error.

$$\begin{aligned} L &= \sum_{i=1}^N \|x_i - \tilde{x}_i\|_2^2 \\ &= \sum_{i=1}^N \|x_i - g_{\theta}(f_{\phi}(x_i))\|_2^2 \end{aligned}$$



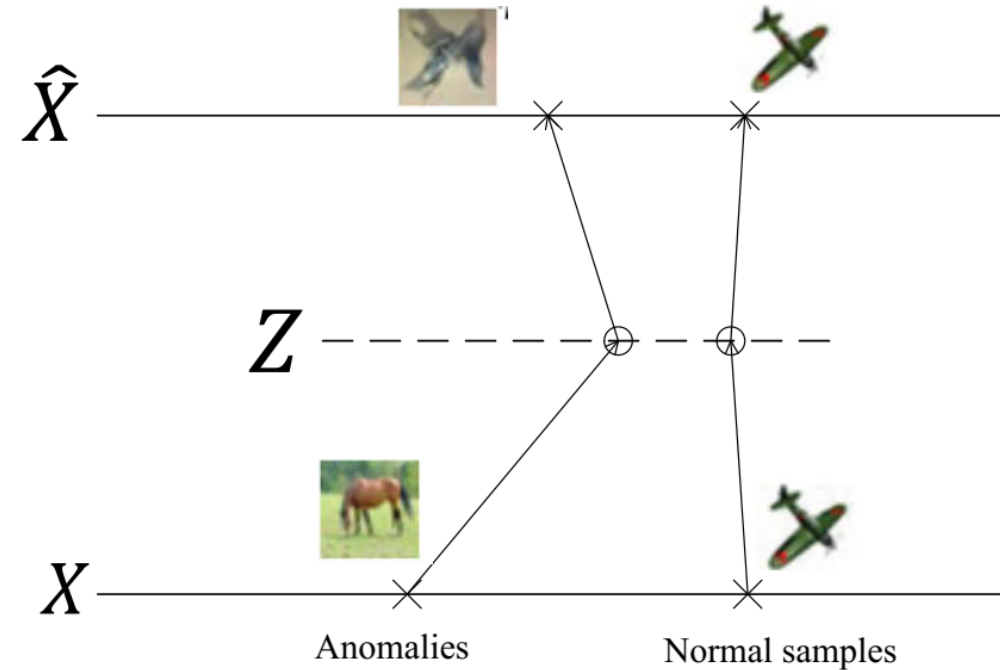


## ● Outlier detection

- Core idea: compare reconstruction errors with given threshold to detect anomalies.

error > threshold **anomaly**

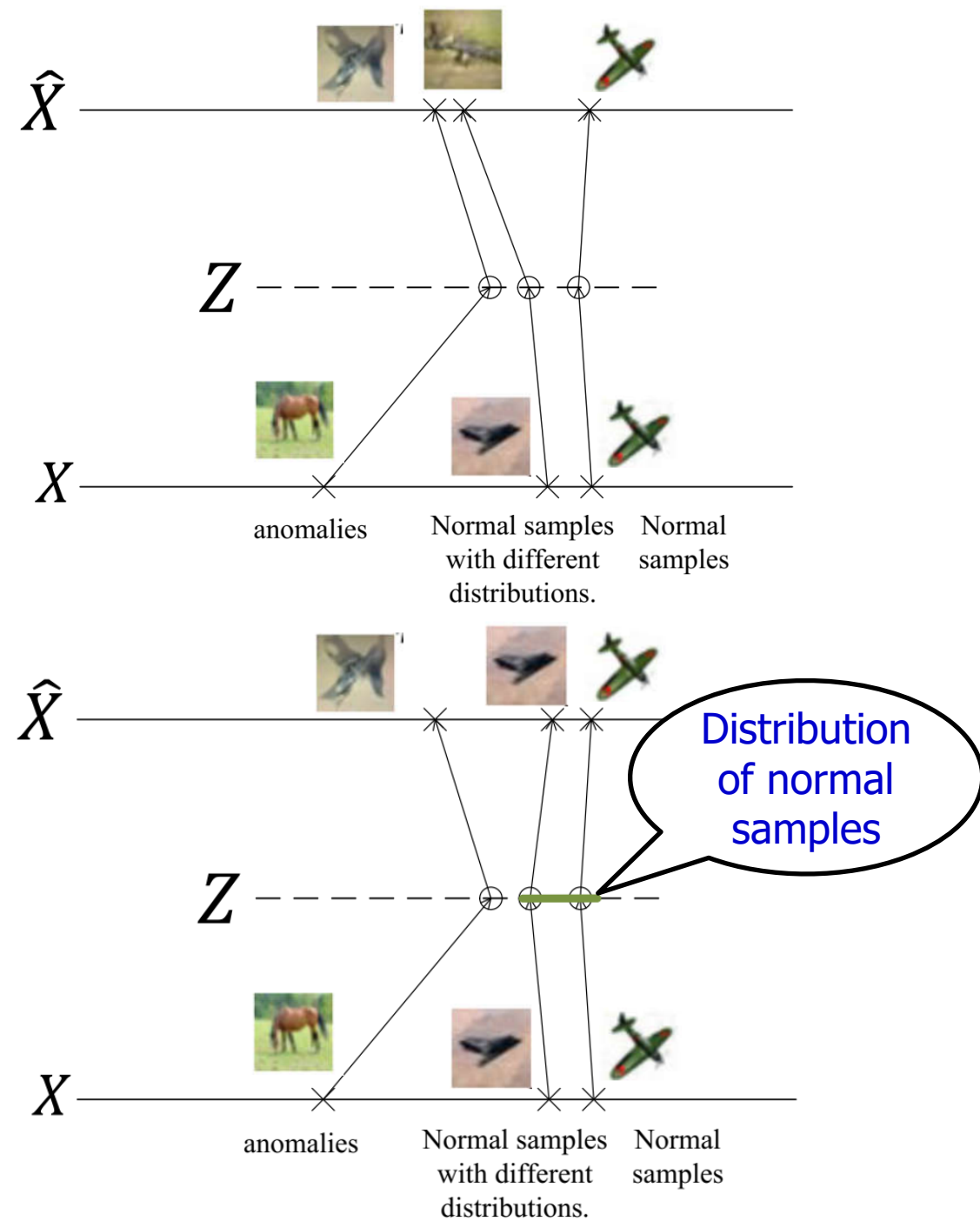
error < threshold **normal**



## ● Shortages of AE

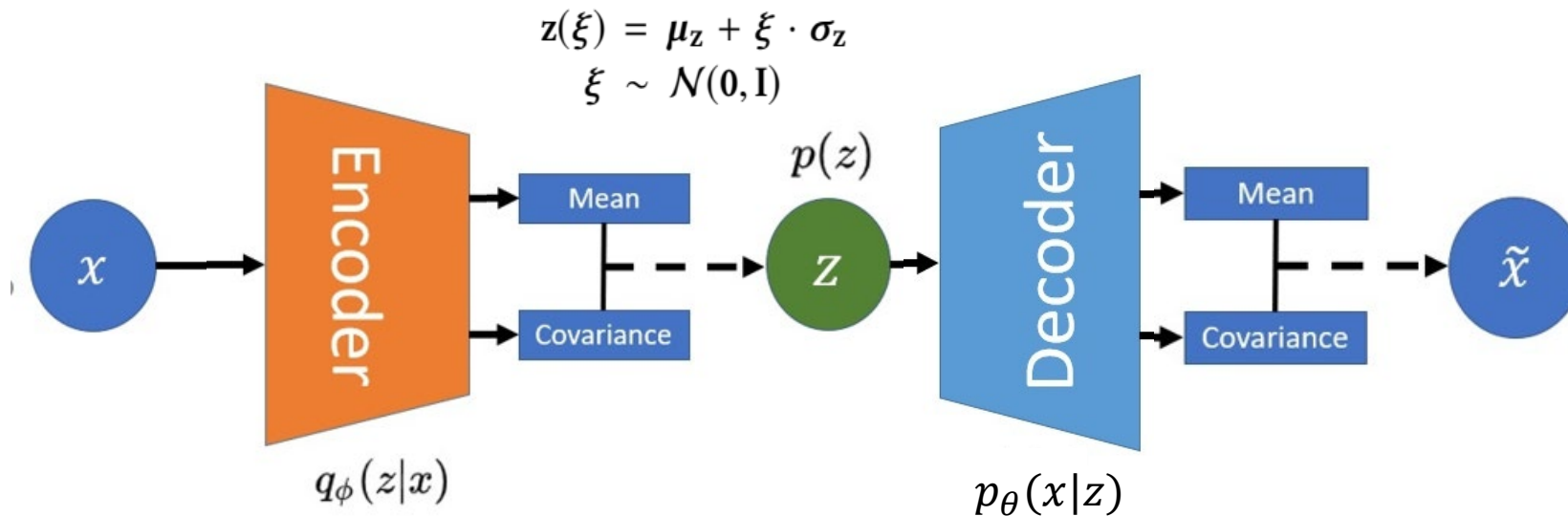
$$z = f_{\varphi}(x)$$
$$\hat{x} = g_{\theta}(z)$$

- Learn one-on-one mapping between  $x$  and  $z$ .
- Can not handle variance in normal samples. Low generalization capability.
- Normal samples may also be falsely judged as anomalies.



# Reconstruction with Variational Autoencoder

- Core idea:  
model the **parameters of distributions** of  $z$  and  $\tilde{x}$  rather than their values.  
———main difference with autoencoders.



- Encoder: learn distributional parameters of  $z$  based on  $x$ .

- Decoder: learn distributional parameters of  $\tilde{x}$  based on  $z$ .

# Reading list

- An J, Cho S. Variational autoencoder based anomaly detection using reconstruction probability[J]. Special Lecture on IE, 2015, 2(1): 1-18.
- B. Zhou, S. Liu, B. Hooi, X. Cheng, J. Ye, Beatgan: Anomalous rhythm detection using adversarially generated time series, in: S. Kraus (Ed.), IJCAI 2019.
- Alexander Geiger, Dongyu Liu, Sarah Alnegheimish, Alfredo Cuesta-Infante, Kalyan Veeramachaneni. TadGAN: Time Series Anomaly Detection Using Generative Adversarial Networks. In proceedings of IEEE International Conference on Big Data, 2020

# Acknowledgements

- Some text, figures and formulations are from WWW. Thanks for their sharing. If you have copyright claim please contact with me at [yym@hit.edu.cn](mailto:yym@hit.edu.cn).
- This lecture is distributed for nonprofit purpose.

# **Thank You for Your Attention**

Contact me at: [yym@hit.edu.cn](mailto:yym@hit.edu.cn)

Tel: 26033008, 13760196623

Address: Rm.1402, H# Building