# Outlier Detection

Ninghao Liu

University of Georgia

April 10, 2024

# Overview

1. Background and Problem Definition

2. Outlier Detection (OD) Strategies
   - Supervised, Unsupervised, Semi-supervised

3. Outlier Detection Methods
   - {Statistical, Proximity, Clustering, Classification} based
   - k-means for Outlier Detection
   - Isolation Forest

4. Dealing with High-Dimensional Data in OD
   - Autoencoders for OD

# 1. A Real-World Application

Imagine you are a transaction auditor in a credit card company:

- Protect customers from credit card fraud.
- Pay special attention to card usages that are different from typical cases.
- Not-that-typical cases: the purchase occurs *far* from the owner's resident city, the purchase amount is much *bigger* than usual, very *frequent* purchases in a short time, ...
- Detect such transactions as soon as they occur and contact the card owner for verification.

# 1. Problem Definition

An **outlier** is a data object that deviates significantly from the rest of the objects, as if it were generated by a different mechanism.
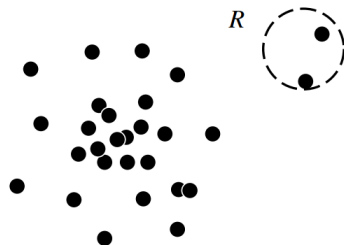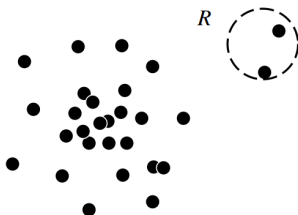


Figure 1: The objects in region R are outliers.

# 1. Problem Definition

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

Please note that:

- Outliers are different from *noisy* data.
- It may be important to justify *why* the outliers are detected.
- Outlier detection is also related to *novelty detection* (but they are not the same).

# 2. Outlier Detection Strategies

Outlier detector training strategies:

- Supervised Learning
- Unsupervised Learning
- Semi-supervised Learning
- ...

# 2. Outlier Detection Strategies

## 2.1 Supervised Learning for OD

Given a dataset $\{(x_i, y_i)\}$, $1 \leq i \leq N$, train a classifier $f : \mathcal{X} \to \{0, 1\}$. For a new instance $x$, it is classified as an outlier if $f(x) = 1$, or a normal instance if $f(x) = 0$.

# 2. Outlier Detection Strategies

## 2.1 Supervised Learning for OD

Given a dataset $\{(x_i, y_i)\}$, $1 \leq i \leq N$, train a classifier $f : \mathcal{X} \to \{0, 1\}$. For a new instance $x$, it is classified as an outlier if $f(x) = 1$, or a normal instance if $f(x) = 0$.

Challenges:

- Data imbalance: Methods are needed for handling data imbalance.
- Evaluation gap: Classification maximizes prediction accuracy, which may not be of interest in applications.
- Labeling difficulty: The sampled data for labeling may not sufficiently represent the outlier distribution.

# 2. Outlier Detection Strategies

## 2.2 Unsupervised Learning for OD

Given a dataset $\{x_i\}$, $1 \le i \le N$, train a detector $f : \mathcal{X} \to \{0, 1\}$. For a new instance $x$, it is regarded as an outlier if $f(x) = 1$, or a normal instance if $f(x) = 0$.

# 2. Outlier Detection Strategies

## 2.2 Unsupervised Learning for OD

Given a dataset $\{x_i\}$, $1 \leq i \leq N$, train a detector $f : \mathcal{X} \rightarrow \{0, 1\}$. For a new instance $x$, it is regarded as an outlier if $f(x) = 1$, or a normal instance if $f(x) = 0$.

- Implicit assumption: The normal objects are "clustered".
- This assumption may not be true in some cases.
  - Outliers may also be "clustered".
  - Normal objects are diversely distributed.
- Computational challenge: We have to process a large population of non-target data entries (i.e., the normal objects) before one can touch the real meat (i.e., the outliers).
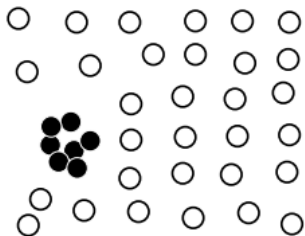
# 2. Outlier Detection Strategies



Figure 2: The black objects form outliers collectively.

# 2. Outlier Detection Strategies

## 2.3 Semi-supervised Learning for OD

Given a dataset $\{x_i\}$ and $\{y_j\}$, $1 \leq i \leq N$, $j \in \mathcal{J}$, $|\mathcal{J}| \ll N$, train a detector $f : \mathcal{X} \rightarrow \{0, 1\}$.

- Most of the instances are unlabeled, so effective unsupervised methods serve as the important basis.
- Both labeled normal/outlier instances are possible.

# 3. Outlier Detection Methods

- Statistical Methods
- Clustering-Based Methods
- Proximity-Based Methods
- ...

# 3.1 Statistical Methods

**The general idea**: Learn a probability density function $p(x)$, denoting the probability that $x$ is generated by the normal mechanism.

- Parametric methods: Assume $p$ follows certain priori distribution (e.g., gaussian) with parameters $\theta$ to be estimated.
- Non-parametric methods: Do not assume a priori distribution.

# 3.1 Statistical Methods

**Steps of parametric methods:**
- Find a certain distribution to work with.
- Estimate the parameters $\theta$ of the distribution.
- Calculate $p(\boldsymbol{x}; \theta)$, and decide if the probability is low enough.

Example 1:

**Univariate outlier detection using maximum likelihood.** Suppose a city's average temperature values in July in the last 10 years are, in value-ascending order, 24.0°C, 28.9°C, 28.9°C, 29.0°C, 29.1°C, 29.1°C, 29.2°C, 29.2°C, 29.3°C, and 29.4°C.

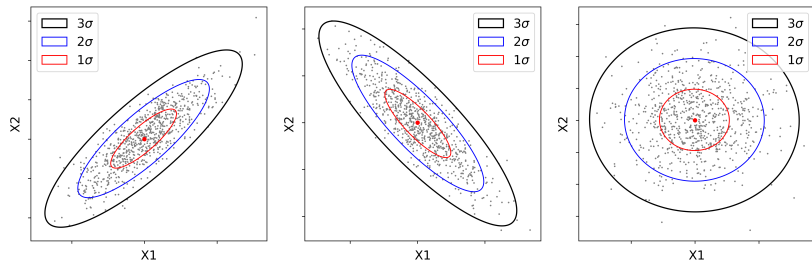$$p(\boldsymbol{x}) = \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp\left(-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}\right)$$

# 3.1 Statistical Methods

**Steps of parametric methods:**

- Find a certain distribution to work with.
- Estimate the parameters $\theta$ of the distribution.
- Calculate $p(\boldsymbol{x}; \theta)$, and decide if the probability is low enough.

Example 2:

$$p(\boldsymbol{x}) = (2\pi)^{-\frac{k}{2}} \cdot \det(\boldsymbol{\Sigma})^{-\frac{1}{2}} \cdot \exp\left(-\tfrac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right)$$
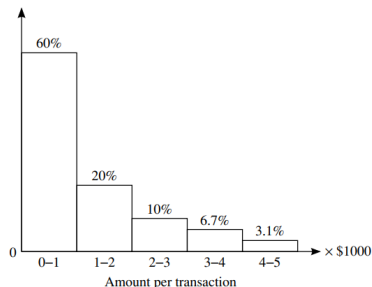
# 3.1 Statistical Methods

**Commonly used non-parametric techniques:**
- Histogram
- Kernel based density estimation

Example:

**Outlier detection using a histogram.** *AllElectronics* records the purchase amount for every customer transaction. Figure 12.5 uses a histogram (refer to Chapters 2 and 3) to graph these amounts as percentages, given all transactions.



Q: Should a transaction amounted $6500 be regarded as an outlier?

# 3.1 Statistical Methods

**Steps of histogram methods:**

- Build the histogram from data.
- Given an instance $x$, check it against the histogram, and assign it with a outlier score of $\frac{1}{hist(x)}$.
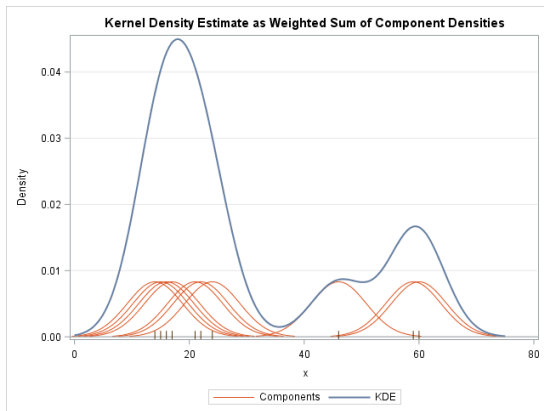
**Challenges**

- How to choose the bin size?

# 3.1 Statistical Methods

**Another Method**: Treat an observed object as an indicator of high probability density in the surrounding region.

# 3.1 Statistical Methods

**Another Method**[1]: Treat an observed object as an indicator of high probability density in the surrounding region.



Kernel Density Estimate as Weighted Sum of Component Densities

---
[1]https://blogs.sas.com/content/iml/2016/07/27/
visualize-kernel-density-estimate.html

# 3.1 Statistical Methods

## Kernel density estimation

Given samples $\{x_1, x_2, ..., x_N\}$ and a test instance $x$,

$$p(x) = \frac{1}{Nh} \sum_{i=1}^{N} K(\frac{x - x_i}{h}) \tag{1}$$

where $K()$ is the kernel function, $h$ is the bandwidth parameter.

A kernel $K()$ is a non-negative real-valued function that satisfies:

- $\int_{-\infty}^{+\infty} K(u) du = 1$,
- $K(u) = K(-u)$.

A frequently used kernel is the RBF kernel: $K(x) = \exp\left(-\gamma \cdot \|x\|_2^2\right)$.

# 3.1 Statistical Methods

Properties of statistical methods:

- Statistically justifiable
- Handling **high-dimensional data** is challenging
- Computational cost: For simple parametric models, it is linear time cost. For kernel density estimation, the model learning cost can be up to quadratic.

# 3.3 Clustering-Based Methods

Possible scenarios:

- Does the instance belong to any cluster?
    - If no, then an outlier.
- Is there a large distance between the instance and the cluster to which it is closest?
    - If yes, then an outlier.
- Is the instance part of a small or sparse cluster?
    - If yes, then the whole cluster is outlying.

Usually **two steps**: (1) clustering, (2) outlier score measuring.

- Use existing outlier detection methods for the second step.

# 3.3 Clustering-Based Methods

**A review of k-mean clustering**:

- Given a set of observations $(x_1, x_2, ..., x_N)$, k-means clustering aims to partition the $N$ observations into $k(\leq n)$ sets $S = \{S_1, S_2, ..., S_k\}$.
- Formally, the objective is to find:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^{k} \sum_{x \in S_i} \|x - \mu_i\|^2, \qquad (2)$$

where $\mu_i$ is the mean of $S_i$.

- This is equivalent to:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^{k} \frac{1}{2|S_i|} \sum_{x,y \in S_i} \|x - y\|^2, \qquad (3)$$

i.e., minimizing the pairwise squared deviations of points in the same cluster.

# 3.3 Clustering-Based Methods

**Equivalence between the two objectives**:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^{k} \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu_i}\|^2$$

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^{k} \frac{1}{2|S_i|} \sum_{\mathbf{x,y} \in S_i} \|\mathbf{x} - \mathbf{y}\|^2$$

# 3.3 Clustering-Based Methods

**A review of k-means clustering**:

- Initiate the $k$ mean vectors $\{\boldsymbol{\mu_1}, \boldsymbol{\mu_2}, ..., \boldsymbol{\mu_k}\}$.
- Alternate between the two steps as below.

  1. **Assignment:** Assign each instance to the cluster with the nearest mean vector:

  $$S_i^{(t)} = \{\boldsymbol{x}_p : \|\boldsymbol{x}_p - \boldsymbol{\mu}_i\| \leq \|\boldsymbol{x}_p - \boldsymbol{\mu}_j\|, \forall j, 1 \leq j \leq k\} \quad (4)$$

  2. **Update:** Recalculate mean vectors

  $$\boldsymbol{u}_i = \frac{1}{|S_i^{(t)}|} \sum_{\boldsymbol{x}_p \in S_i^{(t)}} \boldsymbol{x}_p \quad (5)$$

- The algorithm converges when the assignments no longer change.

# 3.3 Clustering-Based Methods

**Some additional details in k-means**:

- How to initialize?
  - The Forge method.
  - The Random Partition method.
- How to find the best $k$?
  - Through validation.
  - We could use the Gap Statistic [2].
- Limitations of k-mean?

---

[2]"Estimating the number of clusters in a data set via the gap statistic." Journal of the Royal Statistical Society. 2001.
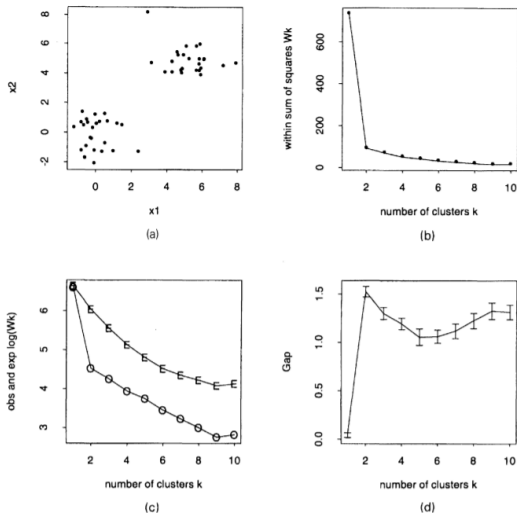
# 3.3 Clustering-Based Methods



**Fig. 1.** Results for the two-cluster example: (a) data; (b) within sum of squares function $W_k$; (c) functions $\log(W_k)$ (O) and $\hat{E}_n^*\{\log(W_k)\}$ (E); (d) gap curve
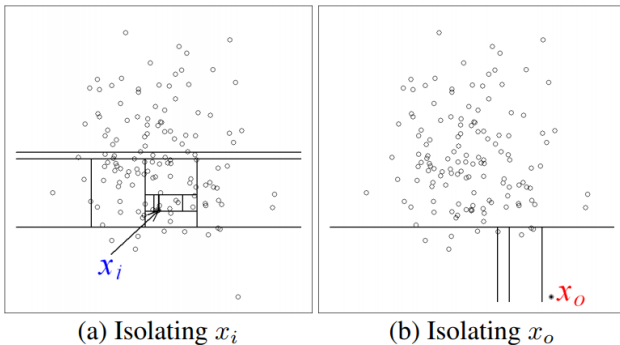
# 3.3 Clustering-Based Methods

**Properties of clustering-based methods**:

- Obtain a summary of data
- Good explainability
- The first step could be time consuming
- The detection performance largely depends on clustering quality
    - Gaussian Mixture Model
    - Hierarchical Clustering
    - Spectral Clustering
    - Deep Model Based Clustering
    - ...

# 3.5 iForest

**Isolation Forest (iForest)** [3]

- **Isolation** means "separating an instance from the rest of the instances".
- Easily being separated $\rightarrow$ More likely to be an outlier.



(a) Isolating $x_i$       (b) Isolating $x_o$
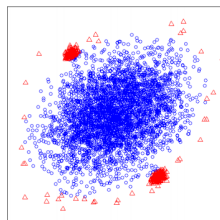
---

[3]"Isolation forest." ICDM. 2008.

# 3.5 iForest

---

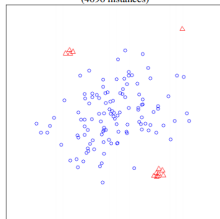**Algorithm 1** : $iForest(X, t, \psi)$

---

**Inputs**: $X$ - input data, $t$ - number of trees, $\psi$ - su
sampling size

**Output**: a set of $t$ *iTrees*

1:  **Initialize** $Forest$
2:  set height limit $l = ceiling(\log_2 \psi)$
3:  **for** $i = 1$ to $t$ **do**
4:      $X' \leftarrow sample(X, \psi)$
5:      $Forest \leftarrow Forest \cup iTree(X', 0, l)$
6:  **end for**
7:  **return** $Forest$

---



(a) Original sample
(4096 instances)



(b) Sub-sample
(128 instances)

# 3.5 iForest

**Algorithm 2** : $iTree(X, e, l)$

**Inputs**: $X$ - input data, $e$ - current tree height, $l$ - height limit

**Output**: an iTree

1: **if** $e \geq l$ or $|X| \leq 1$ **then**
2:      return $exNode\{Size \leftarrow |X|\}$
3: **else**
4:      let $Q$ be a list of attributes in $X$
5:      randomly select an attribute $q \in Q$
6:      randomly select a split point $p$ from $max$ and $min$ values of attribute $q$ in $X$
7:      $X_l \leftarrow filter(X, q < p)$
8:      $X_r \leftarrow filter(X, q \geq p)$
9:      return $inNode\{Left \leftarrow iTree(X_l, e + 1, l),$
10:                     $Right \leftarrow iTree(X_r, e + 1, l),$
11:                     $SplitAtt \leftarrow q,$
12:                     $SplitValue \leftarrow p\}$
13: **end if**

There are two types of nodes: inNode and exNode.

Not always separate each point.

# 3.5 iForest

Complexity $O(nt\log\psi)$

---

**Algorithm 3** : $PathLength(x, T, e)$

---

**Inputs** : $x$ - an instance, $T$ - an iTree, $e$ - current path length; to be initialized to zero when first called

**Output**: path length of $x$

1: **if** $T$ is an external node **then**
2:   return $e + c(T.size)$ {$c(.)$ is defined in Equation 1}
3: **end if**
4: $a \leftarrow T.splitAtt$
5: **if** $x_a < T.splitValue$ **then**
6:   return $PathLength(x, T.left, e + 1)$
7: **else** {$x_a \geq T.splitValue$}
8:   return $PathLength(x, T.right, e + 1)$
9: **end if**

---

# 3.5 iForest

The outlier score is an instance $x$ is:

$$s(x, \psi) = 2^{-\frac{E(h(x))}{c(\psi)}} \tag{6}$$

- $h(x)$: a single path length from root to $x$.
- $c(\psi)$: normalization
- iTrees have an equivalent structure to Binary Search Tree (BST).
- The estimation of $E(h(x))$ for external node terminations is the same as the unsuccessful search in BST.
- $c(\psi) = $ the average path length of unsuccessful search in BST.
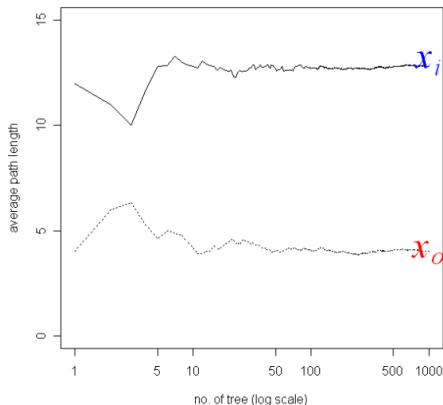
# 3.5 iForest

The outlier score is an instance $x$ is:

$$s(x, \psi) = 2^{-\frac{E(h(x))}{c(\psi)}} \tag{7}$$

- If instances return $s \approx 1$, then they are outliers.
- If instances have $s \ll 0.5$, then they are quite safe to be seen as normal instances.
- If all the instances return $s \approx 0.5$, then the entire sample does not really have any distinct outliers.

# 3.5 iForest

Two variables to choose in iForest:

- number of trees to build $t$ (default: 100)
- sub-sampling size $\psi$ (default: 256)

# 4 Outlier Detection in High-Dimensional Data

**Dimensionality Curse**: Various phenomena that arise when analyzing data in high-dimensional spaces that do not occur in low-dimensional ones [4].

- The *pairwise distances* between different samples in the space converging to the same value as the dimensionality of the data increases.
- As the dimensionality increases, the distance between objects may be heavily dominated by noise.
- The *proximity* or *similarity* between samples may not be qualitatively relevant in higher dimensions.

---

[4]https://en.wikipedia.org/wiki/Curse_of_dimensionality

**Dimensionality Reduction**: compress data

- Classical approaches
    - E.g., Principle Component Analysis (PCA)
- Deep learning approaches.
    - E.g., Robust Deep Autoencoder (RDA) [5]

---

[5]Zhou, Chong, and Randy C. Paffenroth. "Anomaly detection with robust deep autoencoders." KDD. 2017.

# 4.1 PCA

- A feature/dimensionality reduction technique.
- A relevant (but different) technique is feature selection.
- The goal of PCA is to transform the high-dimensional features to a lower-dimensional space.
- PCA is unsupervised.

**What is a good data transformation?**

# 4.1 PCA

**What is a good data compression (transformation)?**

- High compression rate (smaller data size).
- Key characteristics are maintained (keep more information).

In the deep learning era, the transformed/compressed data is also called as:

- the "representation" of data
- "latent representation"
- "embeddings"

# 4.1 PCA

- Let us start from the simple case where PCA projects the features to a **1-dimensional space**.
- Suppose we have $n$ $p$-dimensional ($p > 1$) instances, $x_1$, $x_2$, ..., $x_n \in \mathbb{R}^p$.
- We want to learn a mapper $a \in \mathbb{R}^p$ that $a^\mathsf{T} x_i = z_i$. Here $z_i$ is where $x_i$ locates in the 1-dimensional space.
- PCA tries to solve:

$$a^* = argmax_{\|a\|=1} \ \underbrace{\frac{1}{n}(z_i - \overline{z})^2}_{\text{variance}}, \tag{8}$$

where PCA tries to find the projection with the maximum variance in reduced data.

# 4.1 PCA

- PCA tries to solve:

$$a^* = argmax_{\|a\|=1} \underbrace{\frac{1}{n}(z_i - \overline{z})^2}_{\text{variance}}. \qquad (9)$$

- Q1: Why $\|a\| = 1$ ?
- Q2: What does the objective imply?
    - Trying to make the data points distinguishable in the new space.
    - Why?

# 4.1 PCA

Since

$$\bar{z} = \frac{1}{n}\sum_{i=1}^{n} z_i = \frac{1}{n}\sum_{i=1}^{n} a^T x_i = a^T\left(\frac{1}{n}\sum_{i=1}^{n} x_i\right) = a^T \bar{x},$$

the problem can be written as

$$
\begin{aligned}
a^* &= \underset{||a||=1}{\arg\max} \; \frac{1}{n}\sum_{i=1}^{n}(z_i - \bar{z})^2 \\
&= \underset{||a||=1}{\arg\max} \; \frac{1}{n}\sum_{i=1}^{n}(a^T x_i - \bar{z})^2 \\
&= \underset{||a||=1}{\arg\max} \; \frac{1}{n}\sum_{i=1}^{n}(a^T x_i - a^T \bar{x})^2 \\
&= \underset{||a||=1}{\arg\max} \; \frac{1}{n}\sum_{i=1}^{n} a^T(x_i - \bar{x})(x_i - \bar{x})^T a \\
&= \underset{||a||=1}{\arg\max} \; a^T \underbrace{\left(\frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})(x_i - \bar{x})^T\right)}_{p \times p \text{ covariance matrix}} a \\
&= \underset{||a||=1}{\arg\max} \; a^T C a,
\end{aligned}
$$

where $C = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})(x_i - \bar{x})^T$, denotes the covariance matrix.

# 4.1 PCA

The problem is

$$\max_{a} \; a^\top C a$$
$$s.t., \; \|a\|^2 = 1. \tag{10}$$

How to solve this?

# 4.1 PCA

**Eigen-Decomposition**

A square $n \times n$ matrix $\boldsymbol{C}$ with $n$ linearly independent eigenvectors can be factorized as:

$$\boldsymbol{C} = \boldsymbol{U} \boldsymbol{\Lambda} \boldsymbol{U}^{-1} \tag{11}$$

- $\boldsymbol{U}$ is a $n \times n$ matrix whose columns are eigenvectors of $\boldsymbol{C}$,
- $\boldsymbol{\Lambda}$ is a diagonal matrix whose diagonal elements are the corresponding eigenvalues.

If $\boldsymbol{C}$ is symmetric, then

$$\boldsymbol{C} = \boldsymbol{U} \boldsymbol{\Lambda} \boldsymbol{U}^{\mathsf{T}} = \sum_{i=1}^{n} \lambda_i \boldsymbol{u}_i \boldsymbol{u}_i^{\mathsf{T}}. \tag{12}$$

# 4.1 PCA

Q: What about mapping data into a $k$-dimensional space?

A: The PCA problem is then defined as:

$$A^* = \underset{A \in \mathbb{R}^{p \times k} : A^\mathsf{T} A = I_k}{\operatorname{argmax}} \operatorname{trace}(A^\mathsf{T} C A) \tag{13}$$

where $A = [a_1, a_2, ..., a_k] \in \mathbb{R}^{p \times k}$. PCA requires different projection vectors to be orthogonal.

$$\operatorname{trace}(A^\mathsf{T} C A) = \sum_{i=1}^{k} (a_i^\mathsf{T} C a_i), \tag{14}$$

which is a natural extension of 1-D PCA.

Q: How to solve this?

# 4.1 PCA

When $k = 2$, we want to solve

$$\max_{a_2} \ a_2^\mathsf{T} C a_2$$
$$s.t., \ \|a_i\|^2 = 1 \tag{15}$$
$$a_2^\mathsf{T} a_1 = 0.$$

How to solve this?

# 4.1 PCA

We then have the following theorem.

**Theorem.** (*Ky Fan*) Let $H \in \mathbb{R}^{n \times n}$ be a symmetric matrix with eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n,$$

and the corresponding eigenvectors $U = [u_1, \ldots, u_n]$. Then

$$\lambda_1 + \cdots \lambda_k = \max_{A \in \mathbb{R}^{n \times k} : A^T A = I_k} \operatorname{trace}\left(A^T H A\right).$$

And the optimal $A^*$ is given by $A^* = [u_1, \ldots, u_k]Q$ with $Q$ an arbitrary orthogonal matrix. ■

# 4.1 PCA

**❶** Note that in Eqn. (2), the covariance matrix $C$ is a symmetric matrix. Given the above theorem, we directly obtain

$$\lambda_1 + \cdots \lambda_k = \underset{A \in \mathbb{R}^{n \times k} : A^T A = I_k}{\arg\max} \ \text{trace}\left(A^T C A\right),$$

$$A^* = [u_1, \ldots, u_k] Q,$$

where $\lambda_1, \ldots, \lambda_k$ are the $k$ largest eigenvalues of the covariance matrix $C$, and the solution $A^*$ is the matrix whose columns are corresponding eigenvectors.

**❷** It also follows from the above theorem that solutions to PCA are not unique, and they differ by an orthogonal matrix. We used the special case where $Q = I$, *i.e.*, $A^* = [u_1, \ldots, u_k]$.

# 4.1 PCA for Outlier Detection

Using PCA, we obtain the mappers $A = \{u_i\}$, $1 \leq i \leq k$. An original data instance $x$ will be mapped to $z = [z_1, z_2, ..., z_k]$. After mapping all the data, outlier detection could be done in the $z$-space.

- Apply the detection methods we already discussed.
- Choose the mappers corresponding to higher or lower eigenvalues?
- How many $k$'s to choose?

# 4.1 PCA for Outlier Detection

There is another way to understand PCA from the **data reconstruction** perspective.

- Assume $X$ is centered.
- Data mapping: $Z = A^\mathsf{T}X$.
- Reconstructed data: $\tilde{X} = AZ = AA^\mathsf{T}X$.
- Reconstruction error: $\|X - \tilde{X}\|_F = \|X - AA^\mathsf{T}X\|_F$.
- It has also been shown that:

$$AA^\mathsf{T}X = B^* = \underset{B:\mathrm{rank}(B)\leq k}{\mathrm{argmin}} \|X - B\|_F. \qquad (16)$$

PCA projects the high-dimensional features to a lower-dimensional space with minimum reconstruction error.

  - Keep most of the information = Be able to recover.

- Instances with large reconstruction errors $\|x - AA^\mathsf{T}x\|_2$ are outliers.

# 4.1 PCA for Outlier Detection

In previous parts, we need to compute $C$ first before doing PCA.
This step could be costly if $p$ or $n$ is large.
Can this step be avoided?

**Singular Value Decomposition:**
The singular value decomposition (SVD) of an $m \times n$ real matrix $R$
(we assume $m \geq n$) can be written as:

$$R = U\Sigma V^{\mathsf{T}}, \tag{17}$$

where $U \in \mathbb{R}^{m \times m}$, $\Sigma \in \mathbb{R}^{m \times n}$, and $V \in \mathbb{R}^{n \times n}$, $UU^{\mathsf{T}} = U^{\mathsf{T}}U = I_{m \times m}$,
and $VV^{\mathsf{T}} = V^{\mathsf{T}}V = I_{n \times n}$.

# 4.1 PCA for Outlier Detection

1. Pre-process data.

$$\begin{aligned}
\boldsymbol{X} &= [\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_n] \in \mathbb{R}^{p \times n}, \\
\bar{\boldsymbol{X}} &= \frac{1}{n} \boldsymbol{X} \mathbf{1}_n \in \mathbb{R}^{p \times 1}, \\
\tilde{\boldsymbol{X}} &= (\boldsymbol{X} - \bar{\boldsymbol{X}} \mathbf{1}_n^T) \in \mathbb{R}^{p \times n},
\end{aligned}$$

2. We know $\boldsymbol{C} = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{x}_i - \overline{\boldsymbol{x}})(\boldsymbol{x}_i - \overline{\boldsymbol{x}})^{\mathsf{T}} = \frac{1}{n} \tilde{\boldsymbol{X}} \tilde{\boldsymbol{X}}^{\mathsf{T}}$.

3. If the SVD of $\tilde{\boldsymbol{X}}$ is $\tilde{\boldsymbol{X}} = \boldsymbol{U} \tilde{\Sigma} \boldsymbol{V}$, the columns of $\boldsymbol{U}$ (left-singular vectors) are eigenvectors of $\tilde{\boldsymbol{X}} \tilde{\boldsymbol{X}}$.

4. We do not perform eigen-value decomposition of $\boldsymbol{C}$, but do singular-value decomposition of $\tilde{\boldsymbol{X}}$, to obtain the mappers $\boldsymbol{U}$.

# 4.1 PCA for Outlier Detection

A more efficient way of computing PCA can be described as:

$$\underbrace{\boldsymbol{X}}_{p \times n} \rightarrow \underbrace{\tilde{\boldsymbol{X}}}_{p \times n} \rightarrow \tilde{\boldsymbol{X}} = \boldsymbol{U}\tilde{\boldsymbol{\Sigma}}\boldsymbol{V}^T \rightarrow \boldsymbol{G} = \boldsymbol{U}_k \in \mathbb{R}^{p \times k}$$

$$\boldsymbol{G}^T \tilde{\boldsymbol{X}} = \boldsymbol{Z} \in \mathbb{R}^{k \times n}$$

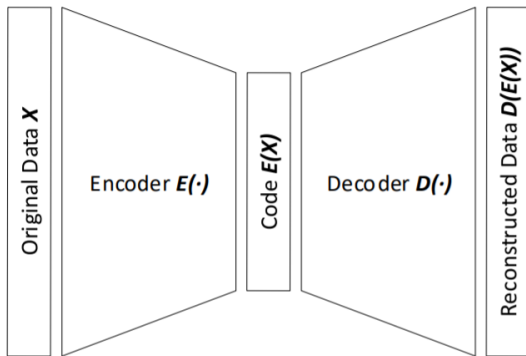# 4.1 Connection between PCA and Autoencoders



Figure 3: An autoendoer.

- What does each module in PCA correspond to?
- How to go deeper?

# 4.2 Outlier Detection with Deep Autoencoders

**Deep Autoencoders**

- A feed-forward multi-layer neural network.
- The desired output is the input itself.
  - Looks trivial?
  - Making hidden representation to be low-dimensional.
- The loss function of training an autoencoder:

$$\min_{D,E} \|\boldsymbol{X} - D(E(\boldsymbol{X}))\| \qquad (18)$$

# 4.2 Outlier Detection with Deep Autoencoders

**Robust PCA (RPCA)**

- Try to reduce the sensitivity of PCA to outliers. Intuition?
- Split the data matrix into two parts

$$X = L + S, \qquad (19)$$

where $L$ is a low-rank matrix, and $S$ is a sparse matrix.

# 4.2 Outlier Detection with Deep Autoencoders

**Robust PCA (RPCA)**

- The problem can be formulated as:

$$\min_{L,S} \|L\|_* + \lambda\|S\|_1$$
$$s.t. \|X - L - S\|_F^2 = 0, \tag{20}$$

  where $\|\cdot\|_*$ is the nuclear norm, $\|\cdot\|_1$ is the L1 norm.

- If we let $D$ and $E$ to be linear mappers, then PCA is the optimal solution for recovering $L$.

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

- The problem can be formulated as:

$$\min_{\boldsymbol{L},\boldsymbol{S}} \|\boldsymbol{L} - D_\theta(E_\theta(\boldsymbol{L}))\|_2 + \lambda\|\boldsymbol{S}\|_1$$
$$s.t.\ \boldsymbol{X} - \boldsymbol{L} - \boldsymbol{S} = 0, \tag{21}$$

  where both the encoder and the decoder are parameterized by $\theta$.

- The first term encourages $\boldsymbol{L}$ to contain the "main features" in data.

- $\lambda$ plays an essential role:
  - a smaller $\lambda \rightarrow$ encourage more data into $\boldsymbol{S}$
  - a larger $\lambda \rightarrow$ include less data in $\boldsymbol{S}$

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

- How to choose $D_\theta$ and $E_\theta$?
    - Any commonly used deep autoencoder architectures could be used.
    - The paper choose a simple architecture:

$$
\begin{aligned}
E_\theta(\mathbf{x}) &= g(\mathbf{W}\mathbf{x} + b_E) \\
D_\theta(\mathbf{x}) &= g(\mathbf{W}^\mathsf{T} E_\theta(\mathbf{x}) + b_D)
\end{aligned}
\tag{22}
$$

- $g(x) = \frac{1}{1+\exp(-x)}$, a nonlinear transformation.

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

- Is $\|S\|_1$ always a good choice?
- What if each anomaly point is not a matrix entry, but the whole row or column?
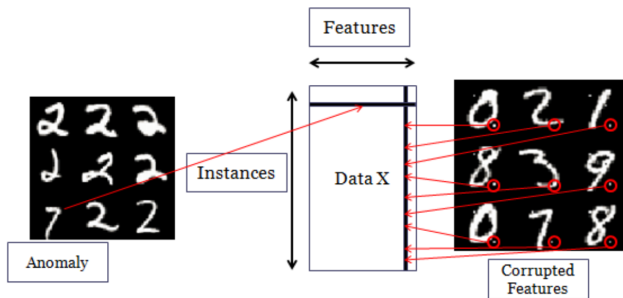


Figure 4: Two examples of group anomalies.

# 4.2 Outlier Detection with Deep Autoencoders

**Regularization**

- $\min \|x\|_1$ vs $\min \|x\|_2$
- Which one encourages a sparse $x$?
- An example of taxes.
- Can we combine them?

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

- $l_{2,1}$ normalization encourages structured/group sparsity.
- $\|\boldsymbol{X}\|_{2,1} = \sum_{j=1}^{n} \|\boldsymbol{x}_j\|_2 = \sum_{j=1}^{n} (\sum_{i=1}^{m} |x_{i,j}|^2)^{1/2}$.
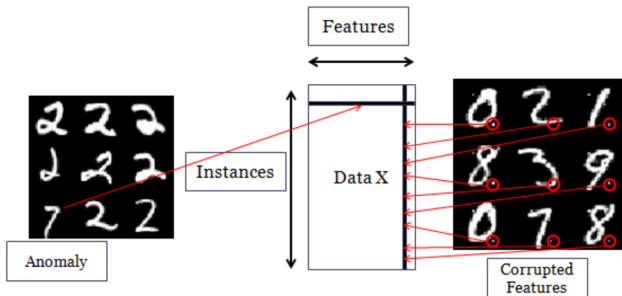- $\|\boldsymbol{X}^\intercal\|_{2,1}$ ?



Figure 5: Two examples of group anomalies.

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

- The problem can be formulated as:

$$\min_{L,S} \|L - D_\theta(E_\theta(L))\|_2 + \lambda(\|S\|_1 \text{ or } \|S\|_{2,1} \text{ or } \|S^\mathsf{T}\|_{2,1})$$

$$s.t. \, X - L - S = 0,$$

- Just denoising: choose $\|S\|_1$.
- To find abnormal data instances: choose $\|S^\mathsf{T}\|_{2,1}$.

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

- How to solve the problem?

$$\min_{L,S} \|L - D_\theta(E_\theta(L))\|_2 + \lambda(\|S\|_1 \text{ or } \|S\|_{2,1} \text{ or } \|S^\mathsf{T}\|_{2,1})$$

$$s.t. \ X - L - S = 0,$$

- Alternating Direction Method of Multipliers (ADMM):
  - Divide the objective into several pieces.
  - Optimizes one of the pieces while keeping the others fixed.
- How to update $L$ with respect to $\min_L \|L - D_\theta(E_\theta(L))\|_2$ ?
- How to update $S$ with respect to $\min_S \|S\|_1$ ?

# 4.2 Outlier Detection with Deep Autoencoders

**Robust Deep Autoencoders (RDA)**

The algorithm of solving

$$\min_{\boldsymbol{L}, \boldsymbol{S}} \|\boldsymbol{L} - D_\theta(E_\theta(\boldsymbol{L}))\|_2 + \lambda(\|\boldsymbol{S}\|_1 \text{ or } \|\boldsymbol{S}\|_{2,1} \text{ or } \|\boldsymbol{S}^\mathsf{T}\|_{2,1})$$

$$s.t. \, \boldsymbol{X} - \boldsymbol{L} - \boldsymbol{S} = 0.$$

1. Initialize $\theta$, $\boldsymbol{S}$
2. while not converge:
   1. $\boldsymbol{L} = \boldsymbol{X} - \boldsymbol{S}$
   2. Minimize $\|\boldsymbol{L} - D_\theta(E_\theta(\boldsymbol{L}))\|_2$ with gradient descent
   3. $\boldsymbol{L} \leftarrow D_\theta(E_\theta(\boldsymbol{L}))$
   4. $\boldsymbol{S} = \boldsymbol{X} - \boldsymbol{L}$
   5. Optimize $\boldsymbol{S}$ using a proximal operator
   6. Check convergence.

# 4.2 Outlier Detection with Deep Autoencoders
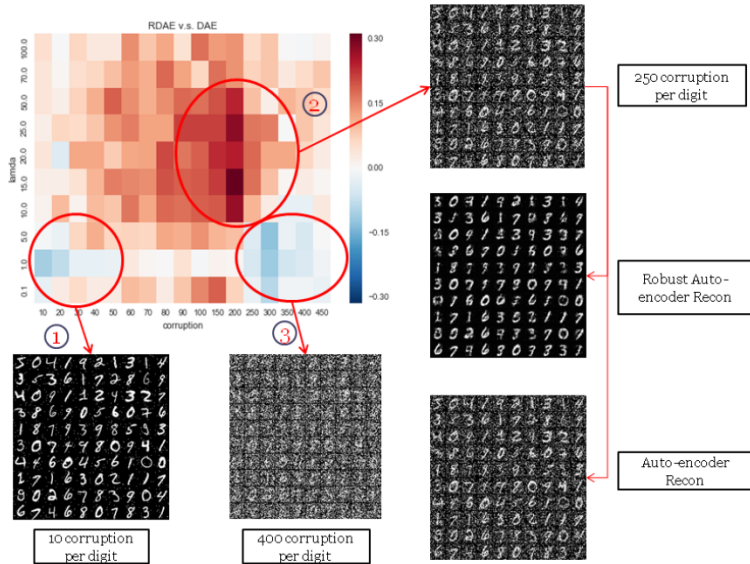
**Experiment Design**

- Datasets
- Research Problems
    - Baseline methods
    - Metrics
    - Visualization
- Effect of hyperparameters or modules

# 4.2 Outlier Detection with Deep Autoencoders

**Experiment Design in RDA**

- Datasets: MNIST
- Research Problem 1: Whether excluding from data an independent sparse matrix using $\|S\|_1$ could help the autoencoder in capturing the essential patterns in data.
    - Baseline methods: Autoencoder
    - Metrics: Classification error
- Experiment design:
    - Extract hidden features from $E_\theta(L)$.
    - Use the hidden features to build a classifier to classify digits.
    - Lower error $\rightarrow$ Better model.

RDAE v.s. DAE

250 corruption per digit

Robust Auto-encoder Recon

Auto-encoder Recon
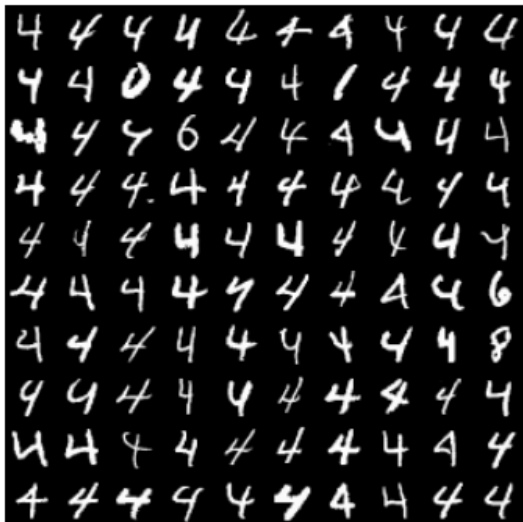
10 corruption per digit

400 corruption per digit

# 4.2 Outlier Detection with Deep Autoencoders

**Experiment Design in RDA**

- Datasets: MNIST
- Research Problem 2: Is $\|S\|_{2,1}$ regularized RDA effective in detecting outliers?
  - Baseline method: Isolation Forest (IS)
  - Evaluation metric: F1-score
- Experiment design:
  - Compare detected outliers with the ground truths.
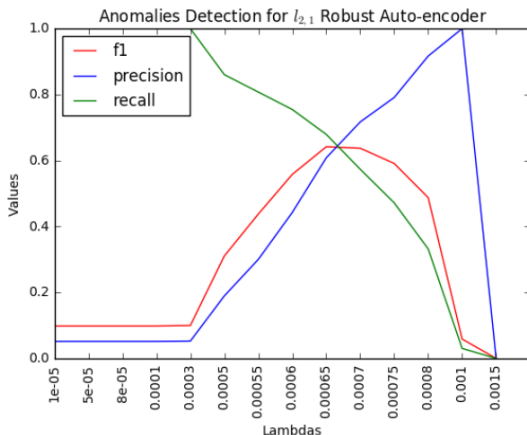  - 0.64 (RDA) vs 0.37 (IS)
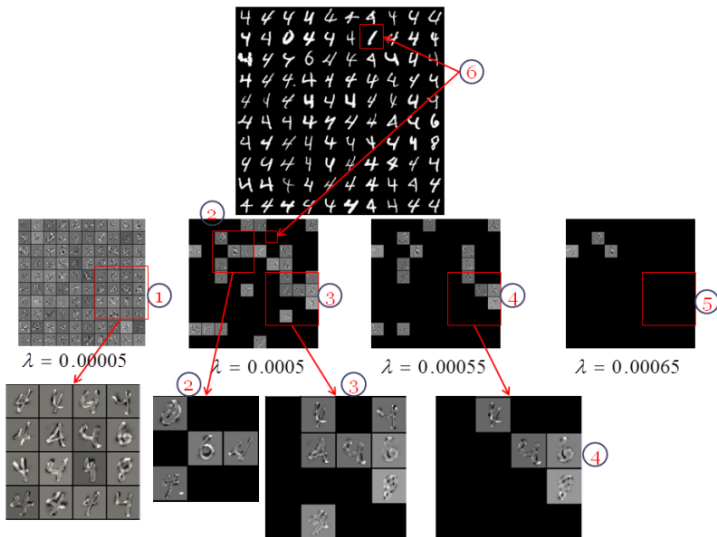
# 4.2 Outlier Detection with Deep Autoencoders

# 4.2 Outlier Detection with Deep Autoencoders

**Experiment Design in RDA**

- Datasets: MNIST
- Research Problem 3: Effect of $\lambda$?



Anomalies Detection for $l_{2,1}$ Robust Auto-encoder

$\lambda = 0.00005$

$\lambda = 0.0005$

$\lambda = 0.00055$

$\lambda = 0.00065$

# 4.2 Outlier Detection with Deep Autoencoders