

# Scientific computation and anomaly recognition

Chen Zhang<sup>1,2</sup>

<sup>1</sup>RTL CRI\_SZ UIH @ Shenzhen, <sup>2</sup>PMJL RUJC-CRI UIH @ Shanghai

Git URL: [CubicZebra](#)



## 1 Scientific computation

- Motivation: full support of algorithm
- Significance: unveil the black box of AI
- Supremum: reduce nonsense efforts

## 2 Anomaly recognition

- Introduction: integration of scientific computation
- Comprehension: concept and applicable scope
- Applications: principles and implementations

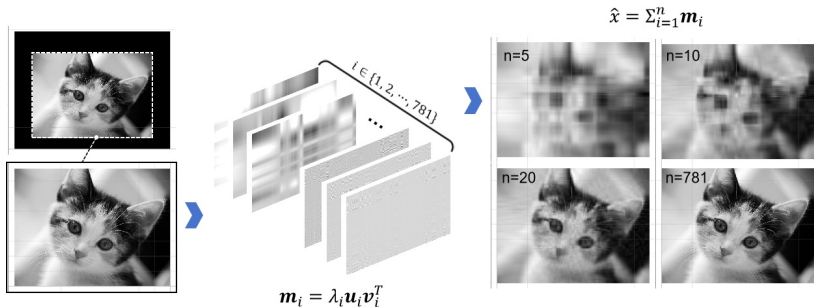
## 3 Summary

# Essence: utilizing the valuable pattern from data

- Program in CS:  
data (digital) + algorithms (business)
- Data mining: (purification, **critical**)  
data (conceptual) + algorithms (simplication)
- AI application: (utilization)  
data (representation) + algorithms (criterion)

# Data mining

- Pattern: minimal representation of data<sup>1</sup>



# AI application

- Mathematical statistics:  
hypotest, ANOVA, Bayesian stats, statistical learning, ...
- Machine learning:
  - analytical: PCA, SVM, conjugate gradient descent, ...
  - randomness: RF, ensemble, stochastic gradient descent, ...
- Deep learning:
  - Conv+Pool (+randomness & variation, data aug & mining)
  - MLP ( $p = f(\mathbf{x}) \rightarrow f$ , set criterion of regressor)

# Scientific computation: base implementation of algorithms

- Signal decomposition
  - eigen and siglar value decomp.,
  - tensor decomp. & synthesis (CP, tucker, train, ring);
- Optimization
  - linear or nonlinear OLS,
  - gradient descent: (quasi-)newton iter, conjugate, stochastic,
  - stochastic process;
- Statistics
  - inference: MLE, Bayes (posteriori & prediction dis),
  - measurement: univariate/multivariate hypothesis tests,
  - sampling: MH, MCMC, stat simulation;
- Linear algebra, ODR, PDE, ...

# For rational AI approaches

- Comprehension on algorithm frame
  - What: base components of scientific computation methods,
  - Why: reason of frame architecture design,
  - How: concrete computation steps of each component;
- Modification on algorithm frame
  - Parameter: inherent adaption support of frame,
  - Component: change features of frame for certain purpose;
- Creating customized frame
  - Analysis: data  $\Rightarrow$  concerned info  $\Rightarrow$  sc methods,
  - Architect: sc implements  $\Rightarrow$  atomic ops  $\Rightarrow$  pipe;



# For rational AI approaches

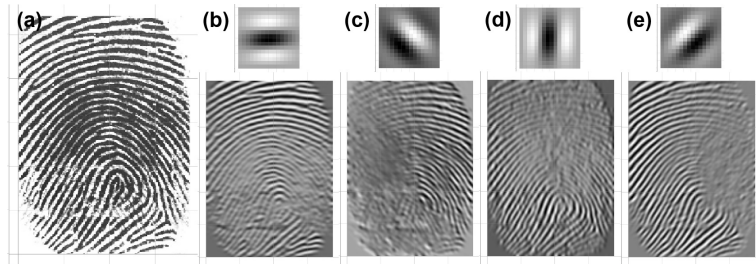
Two suggestions in practice:

1. **Deep into principles:** despite autonomous vehicles, a driver inside should have license at least.
2. **Don't against the current:** modification/design based on *concerned info* and *task objective*.



# The objective-oriented methodology

- Case 1: fingerprint recognition<sup>2</sup>
  - interested info: pattern of texture
  - sc methods: spatial gabor filtering (2D)



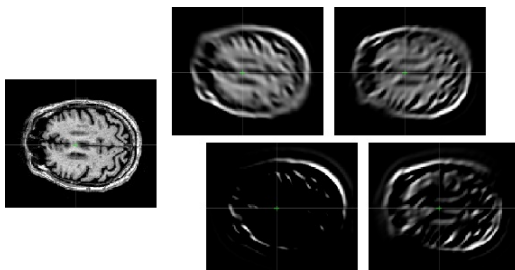
# The objective-oriented methodology

- Case 2: vessel measurement in pancreatic carcinoma study
  - interested info: space angle
  - sc methods: morphological ops + linalg computation



# The objective-oriented methodology

- Case 3: brain-related study (pipe reusability in Case 1)
  - interested info: texture pattern of brain
  - sc methods: spatial gabor filtering (3D)



# Primer concepts in statistics

- Why we need statistics
  - data: listing all acquired observations
  - statistics: description on data via sth. (e.g. dis)
- Samples & population
  - estimation: study on samples  $\Rightarrow$  conclusion on population
  - bias: diff between samples & population<sup>3</sup>

---

<sup>3</sup>Unbias estimation

# Primer concepts in statistics

- Parametric or non-parametric<sup>4</sup>
  - parametric: at least info of a dis., strong assumption,
  - -diff, the statistical illusion (e.g. median/mean)

$$A = \{1, 2, 3, 4, 5\} \text{ and } B = \{1, 2, 3, 4, 5000\}$$

Data does not lie. People do.

—Lee Baker, *Truth, Lies & Statistics: How to Lie with Statistics*

---

<sup>4</sup>Parametric and non-parametric

# Primer concepts in statistics

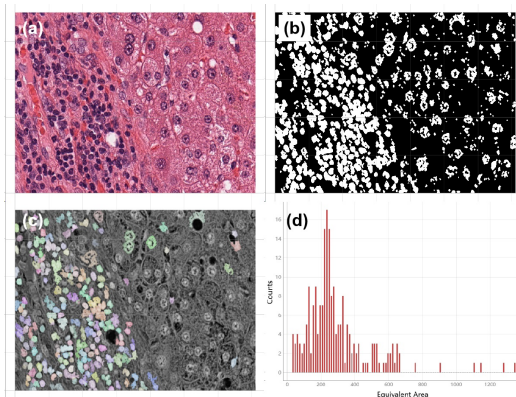
- Consideration of sufficiency
  - def:  $f(\mathbf{x}|\theta) = g(T(\mathbf{X}|\theta))h(\mathbf{x})$
  - design intention of StratifiedKFold in scikit-learn
  - interpretation of bootstrap (CS) in statistics (convergency)

# Proposals in managing data practice

- Diagnosis on samples (assume  $\mathbf{A} \subset \mathbf{B} \subset \mathbf{C}$ )
  - $s(\mathbf{A})$  diff from  $s(\mathbf{B}) \Rightarrow$  increase nums. | test variants,
  - $\mathbf{A} \subset \mathbf{B} \rightarrow f$  be ineffective from  $\mathbf{B}$  to  $\mathbf{C} \Rightarrow$  data aquisition;
- Applicability of algorithm
  - evaluate applicability from frame to task-objective,
  - modification | design  $\Rightarrow$  parameterization;
- Further works
  - cleaning, preprocessing, training, validation, etc.

# Determination of algorithm design

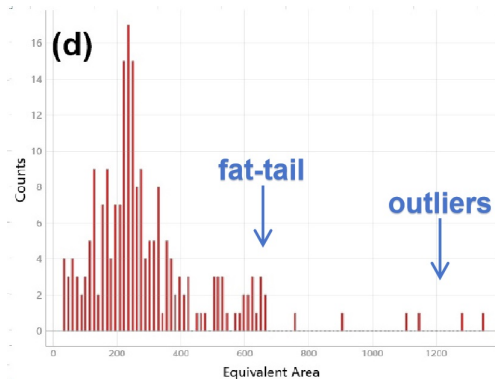
- Case 4: statistics on nuclei of cancer cells
  - objective:  $\mathcal{N}(\mathbf{x}, \Sigma)$  dis of nuclei, statistical approaches





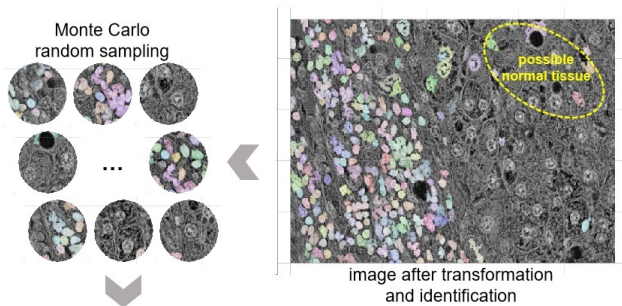
# Determination of algorithm design

- Case 4: statistics on nuclei of cancer cells
  - diagnosis on dis  $\Rightarrow$  fat-tailed & outliers



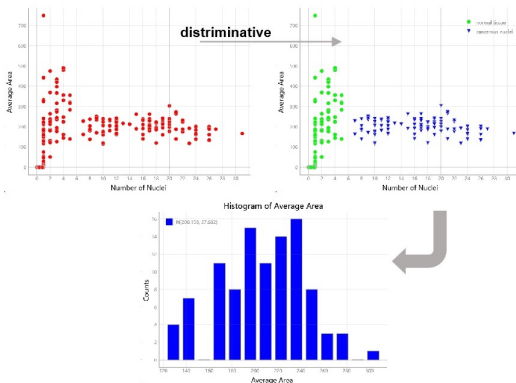
# Determination of algorithm design

- Case 4: statistics on nuclei of cancer cells
  - +MCMC aug => params converge in prob



# Determination of algorithm design

- Case 4: statistics on nuclei of cancer cells
  - +MCMC aug => params converge in prob



# Comb. & App. of multi sc meta implementations

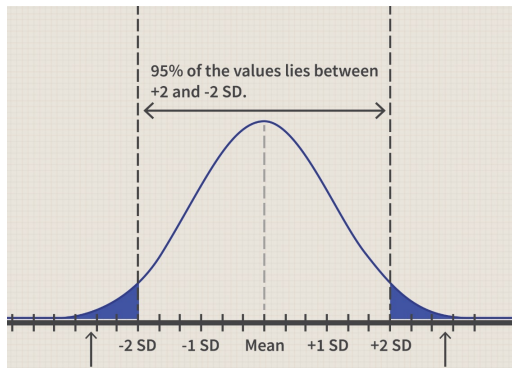
- Dis-related methods (e.g. T2, Naive Bayes)
  - mathematical statistics, hypothesis tests
  - linear algebra (supporting multivariate)
  - Bayes statistics and computation
- Data-related methods (e.g. neighbors)
  - data structure (CS), KDTree for query
  - optimizations, for numeric solution
  - signal decomposition
  - linear algebra as well
- Other methods. . .

# Comb. & App. of multi sc meta implementations

- Security:
  - risky transaction, anti-fraud, ...
- Quality ensurance:
  - quality testing, mechanical fault monitoring, ...
- Networks:
  - spam recog., attack monitoring, ...
- Medical:
  - **less reported**

# Essentials: conventional hypothesis testing

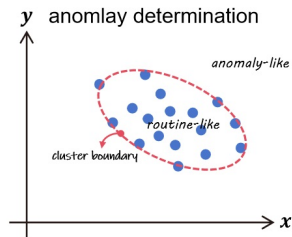
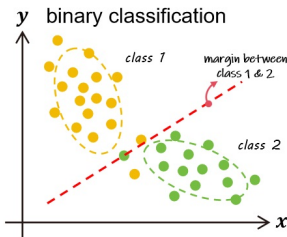
- One-side modeling: stats on one-class  $\rightarrow$  normal/anomaly



accept/rejection regions<sup>5</sup>

# Essentials: concept and applicable scope

- Anomaly detection<sup>6</sup>:
  - stats on one-class  $\rightarrow$  accept/rejection  $\rightarrow$  routine/anomaly-like



# Essentials: concept and applicable scope

*Why not binary classification?*

- **binary classification**

1. narrow vs. narrow
2. large data required
3. low generalization

- **anomaly detection**

1. narrow vs. generic
2. low data required
3. high generalization

**Source: limited capability for unknown pattern**



# Hotelling T2:<sup>7</sup> multivariate student-T

- **Hypothesis:**

- null: case is derived from the identical population,
- alternative: case is not derived from the identical population;

- **T2 statistic:**

$$T^2 = \frac{N-M}{(N+1)M} (\mathbf{x}' - \hat{\boldsymbol{\mu}})^\top \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x}' - \hat{\boldsymbol{\mu}}) \sim F(M, N - M)$$

- **Criterion:**

$$(\mathbf{x}' - \hat{\boldsymbol{\mu}})^\top \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x}' - \hat{\boldsymbol{\mu}}) \sim \chi^2(x|M, 1)$$

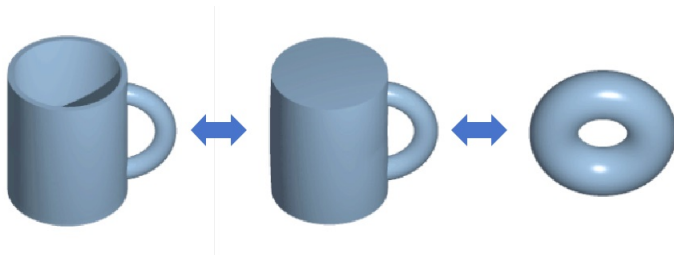
calculation on basis of  $\chi^2(x|M, 1)$

---

<sup>7</sup>Hotelling T-squared

# Large margin nearest neighbors (LMNN)

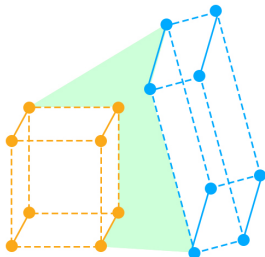
- Essential concept: *homeomorphism*<sup>8</sup>



# LMNN: role of Riemannian space

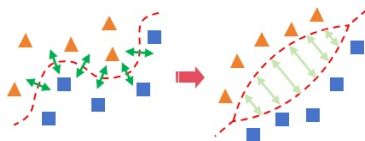
Euclidean (global ops)

$$\mathbf{x}' = f(\mathbf{x}) = \mathbf{A}\mathbf{x}$$



Riemannian (local ops)

$$\mathbf{x}' = f(\mathbf{x}) = \mathbf{B}\mathbf{x}$$



- The determination of  $\mathbf{B}$ ?

# LMNN: problem description<sup>9</sup>

- **Optimization objective:**

$$\Psi(\mathbf{R}) = \frac{1}{N} \sum_{c=1}^s \sum_{n=1}^N \left[ w_c \cdot \psi_1^{(n)}(\mathbf{R}) + \sum_{m \in \{c\}^C} w_m \cdot \psi_2^{(n)}(\mathbf{R}) \right]$$

- **Constraints:**

$$\text{s.t. } \mathbf{R} \succeq 0$$

- **Supports:**

linalg, matrix decomp., data structure, ...



---

<sup>9</sup>Empirical distribution and neighbors

# LMNN: optimization & Riemannian space

$$\arg \min_R \Psi(R) \rightarrow R^*$$

- matrix decomposition: LDLt

$$R^* = L_m \Lambda_m L_m^\top = L \Lambda \Lambda^\top L^\top$$
$$\therefore L' = L \Lambda$$

- Cartesian to Riemannian space:

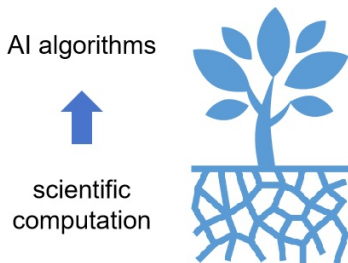
$$f(x) = L'x$$

# Anomaly detection: decoupling the AI algorithms

- As for general case:
  - normal/anomaly determination;
- As for anomaly cases:
  - expertise model for subtypes (most researchers mainly do);
- As for interested subtypes:
  - further investigation for principles/mechanisms;
- Well recognized → full utilization:
  - model to export (suggested) decision, etc.

# AI algorithm in practice

- The relations between AI algorithm and scientific computation



1. **SC**  $\rightarrow$  AI  $\rightarrow$  App.
  - technical methods set
  - rational combination
2. objective  $\rightarrow$  AI frame
  - high generalization
  - expertise for problem
  - interpretability

# Mottos:

- There is neither elixir for all diseases in this world, nor generic solution for all questions.<sup>10</sup>
- Invocation without comprehension is just like tree without root, stream without source.
- Scientific computation and rationality: the motivation and reason of finding the information underlying the data.



Thanks.