

The top banner features the Fudan University logo on the left, which includes the university's name in English ('FUDAN UNIVERSITY') and Chinese ('復旦大學') around a central emblem. To the right of the logo are several blue gears of varying sizes. Some gears contain white icons: a building, a graduation cap, a medical cross, a person silhouette, and an atomic symbol. The background of the banner is light blue with a dark blue curved shape on the right side.

Big Data Analytics & Applications

Bin Li

School of Computer Science

Fudan University

About the Course

- Introduce real-world data analytics applications
- Introduce typical machine learning solutions
- Syllabus of the course
 - Part 1: Introduction to Machine Learning
 - Part 2: **Text** Data Analytics
 - Part 3: **Network** Data Analytics
 - Part 4: **Image** Data Analytics
- Assessment methods
 - Interactions and presentations (30%)
 - Course project with final report (70%)

About the Students

■ Prerequisites of this course

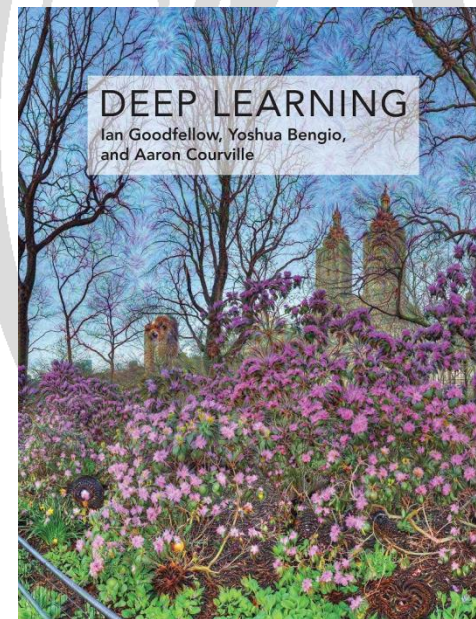
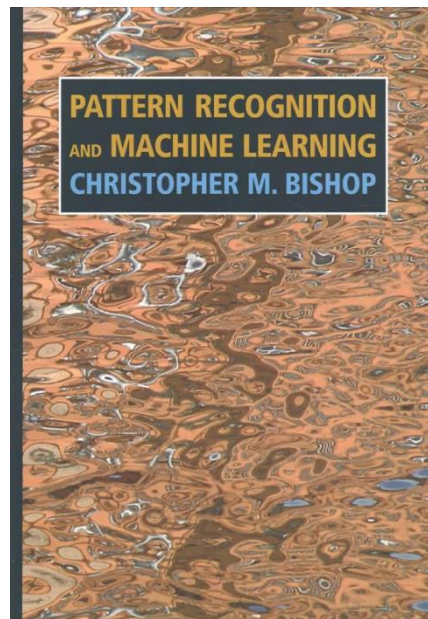
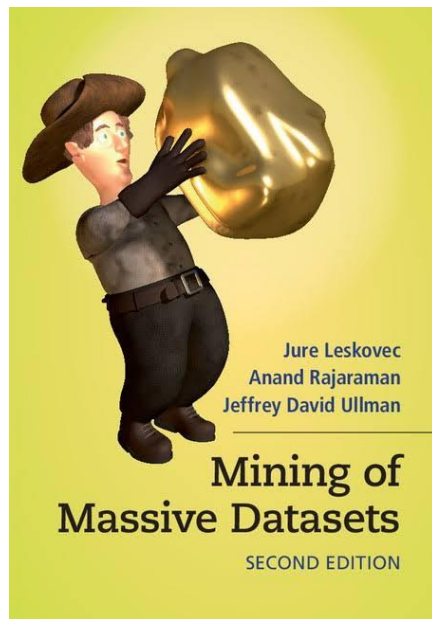
- ❑ Calculus
- ❑ Linear Algebra
- ❑ Probability Theory and Statistics
- ❑ Programming
- ❑ English Speaking & Writing

■ Course project requirements

- ❑ Select one data analytics problem introduced in the course
(Collect real-world new dataset by yourself – **extra point**)
- ❑ Use toolbox (e.g., scikit, TensorFlow) to solve the problem
(Develop new method by yourself – **extra point**)
- ❑ Write project report in English

Reference Books

- No textbook for this course
- Reference books and some papers mentioned in the course



Data Explosion

- Can you tell some Chinese counterparts of these Apps?



Big Data (why called "Big")

- Big data^[1] can be described by the "3V" characteristics
 - ▣ **Volume**: The quantity of generated and stored data
 - ▣ **Variety**: The type and nature of the data
 - ▣ **Velocity**: Big data is often available in real-time



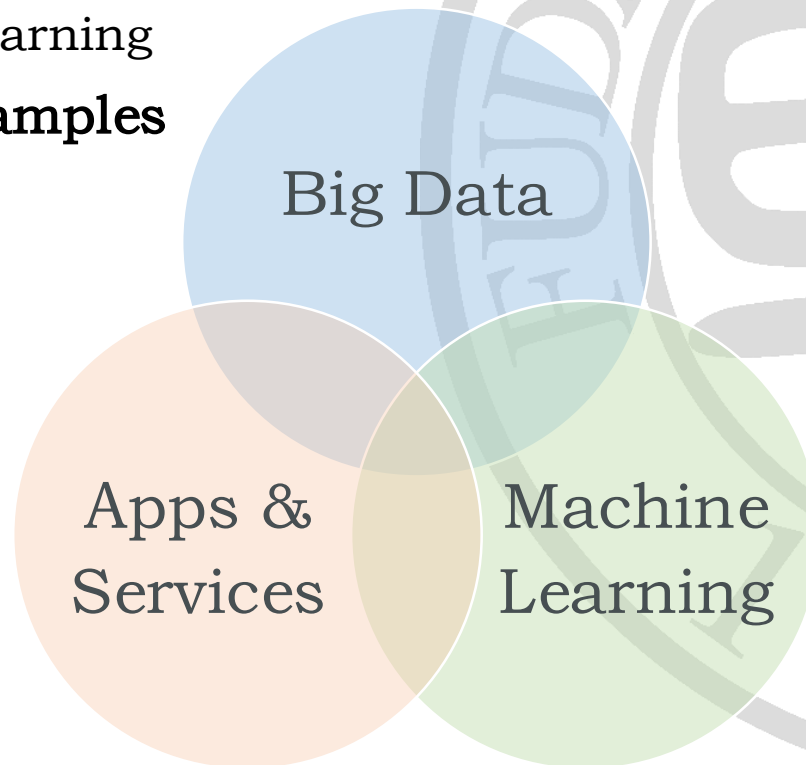
[1] https://en.wikipedia.org/wiki/Big_data

Big Data Analytics (BDA)

- Three elements for big data analytics

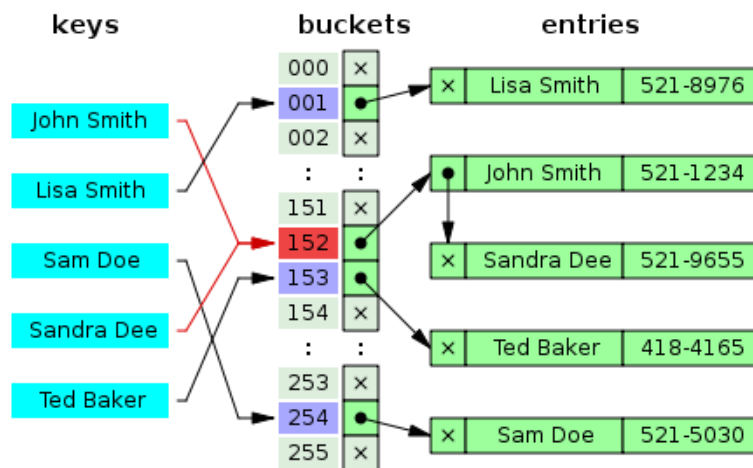
- **Source**: Big Data
- **Demand**: Applications & Services
- **Tools**: Machine Learning

- Try to find some examples



BDA Example: LSH

- **Locality-sensitive hashing** (LSH) reduces the dimensionality of high-dimensional data such that similar items map to the same buckets with high probability.
 - ❑ Similarity search
 - ❑ Duplicate detection
 - ❑ Dimension reduction
 - ❑ Preprocessing for machine learning tasks



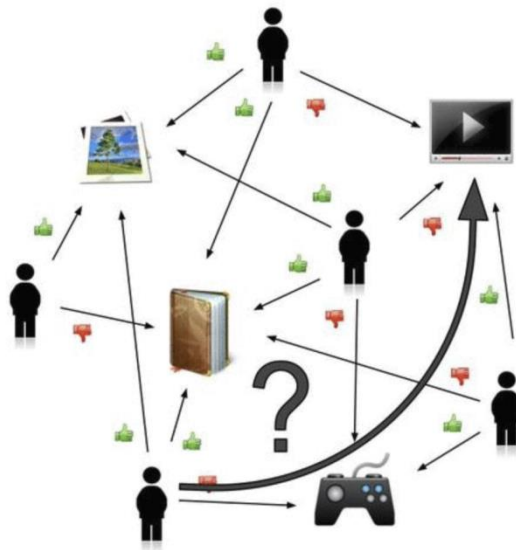
[illegible]


























- Document (e.g., news) automatic categorization
- Any bag-of-words objects grouping tasks



BDA Example: Collaborative Filtering

- Collaborative filtering makes predictions (**filtering**) about the interests of a user by collecting preferences or taste information from many users (**collaborating**)
 - Recommender systems
 - User behavior prediction
 - Marketing



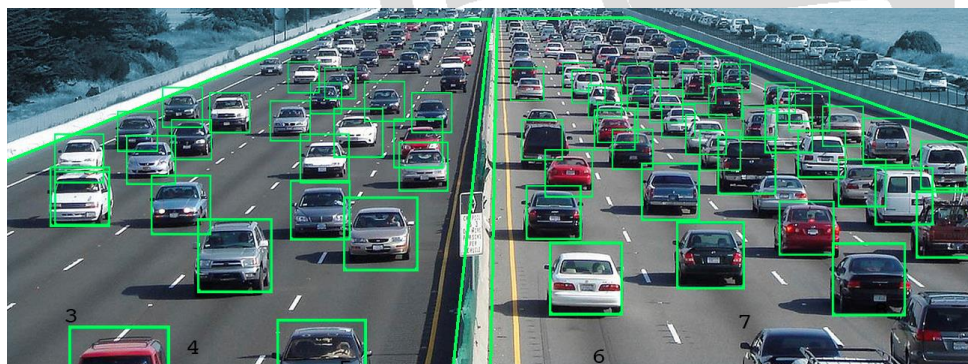
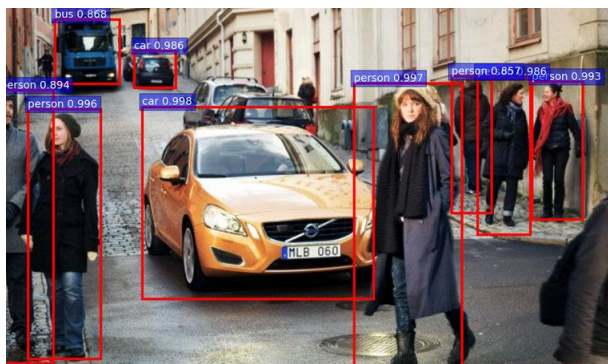
BDA Example: Social Network Analysis

- Social network analysis (SNA) is the process of investigating **social structures** where nodes denote users and links denote relationships.
 - ❑ Community detection
 - ❑ Link prediction
 - ❑ Social capital detection
 - ❑ etc.



BDA Example: Object Detection

- **Detect certain objects** from massive images or video streams using bounding boxes
 - ❑ Security surveillance
 - ❑ Intelligent transport
 - ❑ Image annotation and retrieval
 - ❑ Internet censorship
 - ❑ etc.



BDA Example: OCR

- **Optical character recognition** (OCR) is a technology that extracts and recognizes texts from images.

- ❑ Document scanning
- ❑ Text recognition in the wild
- ❑ Commercial analysis
- ❑ Internet censorship
- ❑ etc.



复仇者联盟

> Prediction: 复仇者联盟

专方高格

> Prediction: 专方高格

DEF

> Prediction: OEF

香酥芝麻味

> Prediction: 香酥芝麻味

冷加工糕点

> Prediction: 冷加工糕点



> Prediction: 七

Three Elements: Examples

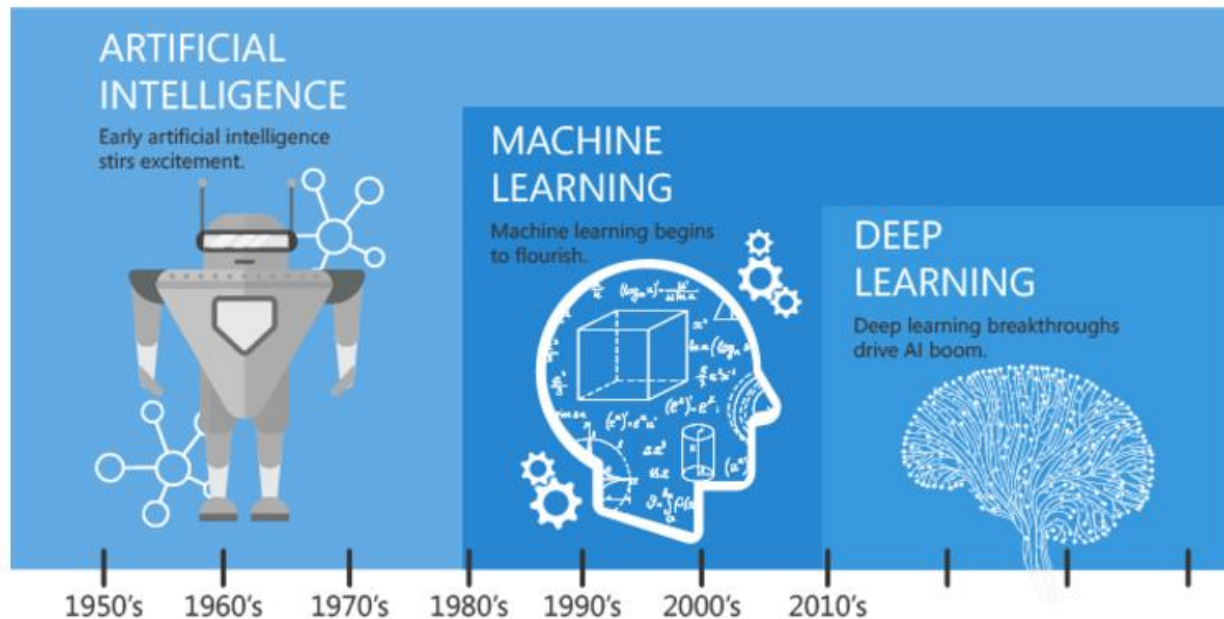
■ Three elements for big data analytics

- ❑ **Source**: Big Data
- ❑ **Demand**: Applications & Services
- ❑ **Tools**: Machine Learning

Demand	Source	Tools
Similarity Search	Webpages	Locality-Sensitive Hashing (LSH)
Topic Modeling	Documents	Latent Dirichlet Allocation (LDA)
Object Detection	Surveillance Video	Covolutional Neural Networks
OCR	Commercials	Covolutional Neural Networks
Recommendation	User Ratings	Matrix Factorization
Community Detection	Social networks	Infinite Relational Models

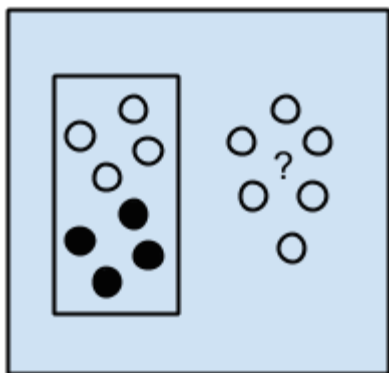
Machine Learning

- Machine learning often uses **statistical techniques** to give computers the ability to "learn" with data.
- Within the field of data analytics, machine learning is a method used to devise **models** and **algorithms** that lend themselves to prediction.



Machine Learning Problems

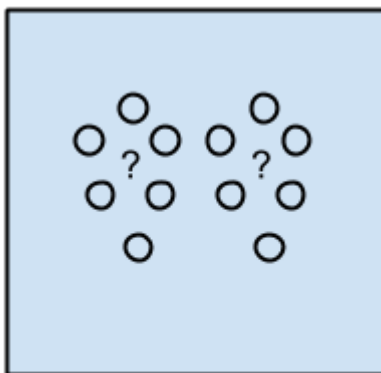
■ Problem settings in machine learning



Supervised Learning
Algorithms

Classification
Regression

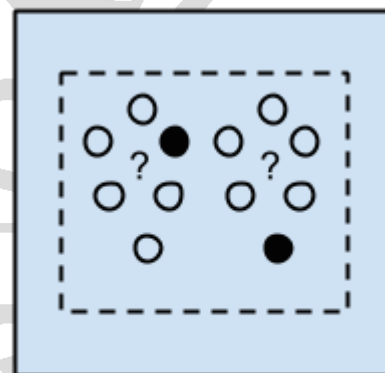
Examples: Logistic
Regression, Neural
Network, etc.



Unsupervised Learning
Algorithms

Clustering
Dimensionality Reduction

Examples: K-Means, Principal
Component Analysis, etc.

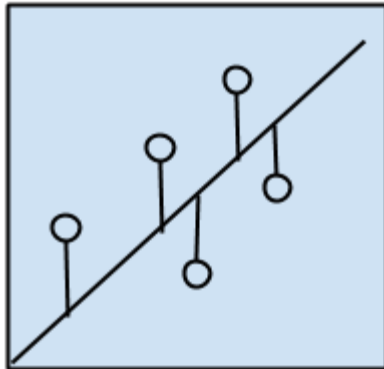


Semi-supervised
Learning Algorithms

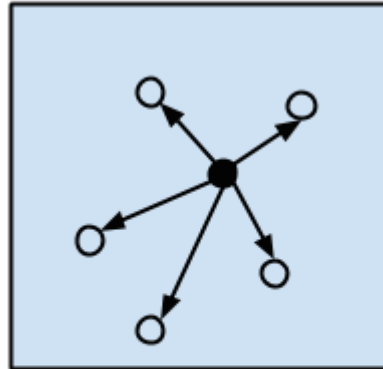
Classification
Regression

Examples: Transductive
SVM, etc.

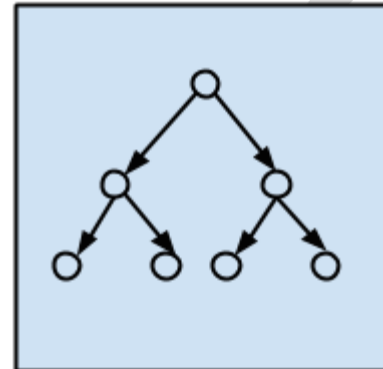
Typical ML Algorithms



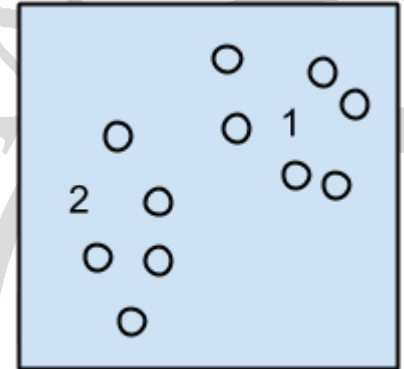
Regression Algorithms



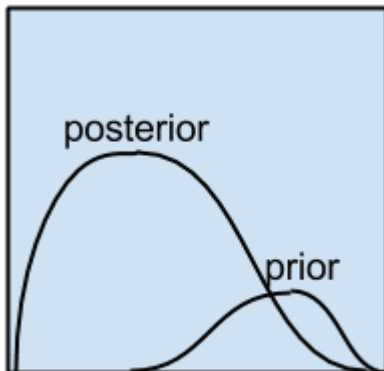
Instance-based Algorithms



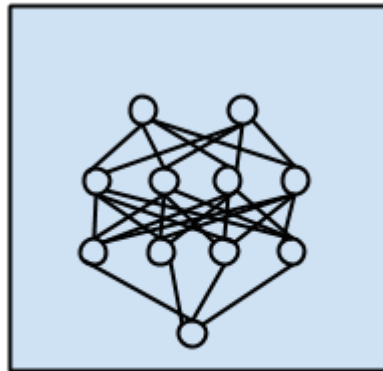
Decision Tree Algorithms



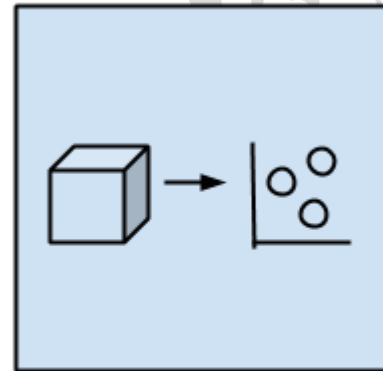
Clustering Algorithms



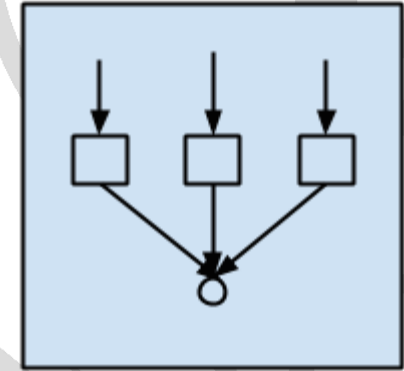
Bayesian Algorithms



Deep Learning Algorithms



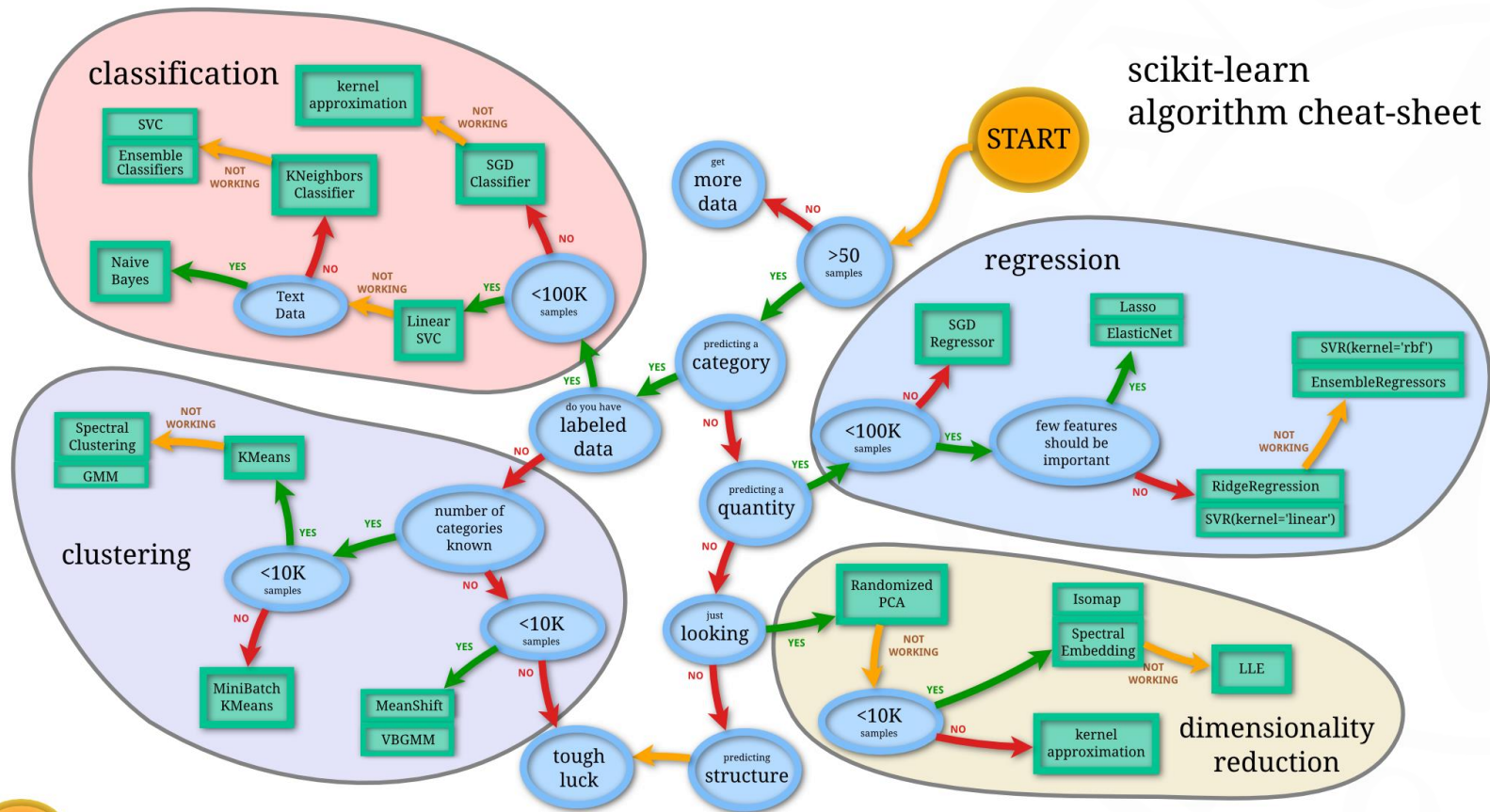
Dimensional Reduction Algorithms



Ensemble Algorithms

A Map of ML Tool Selection

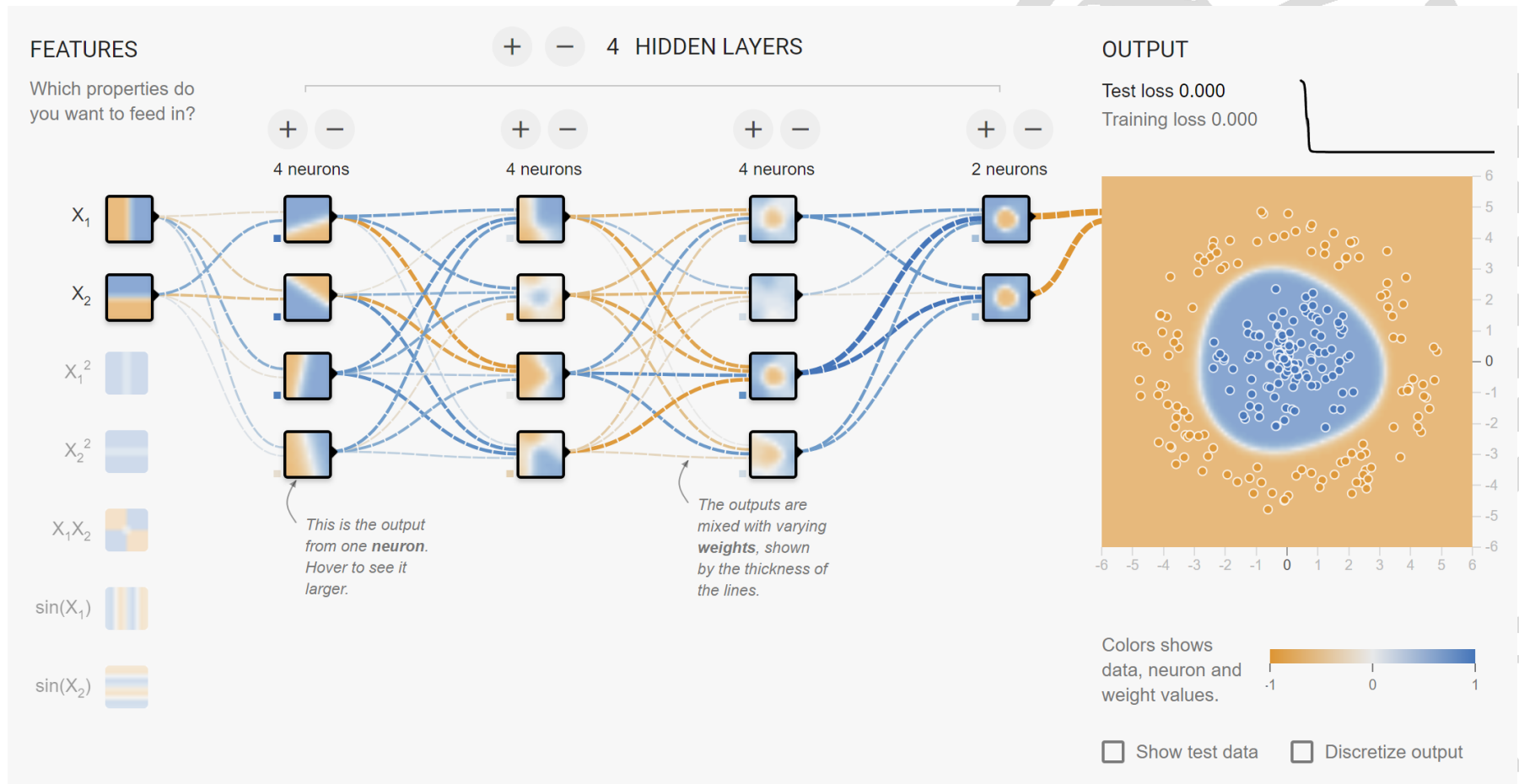
scikit-learn
algorithm cheat-sheet



Toolboxes for DBA

- For big data analytics with **traditional machine learning** tasks: <http://scikit-learn.org>
- For big data analytics with **deep learning** tasks: <https://www.tensorflow.org> or <https://pytorch.org/>

Tensorflow Playground

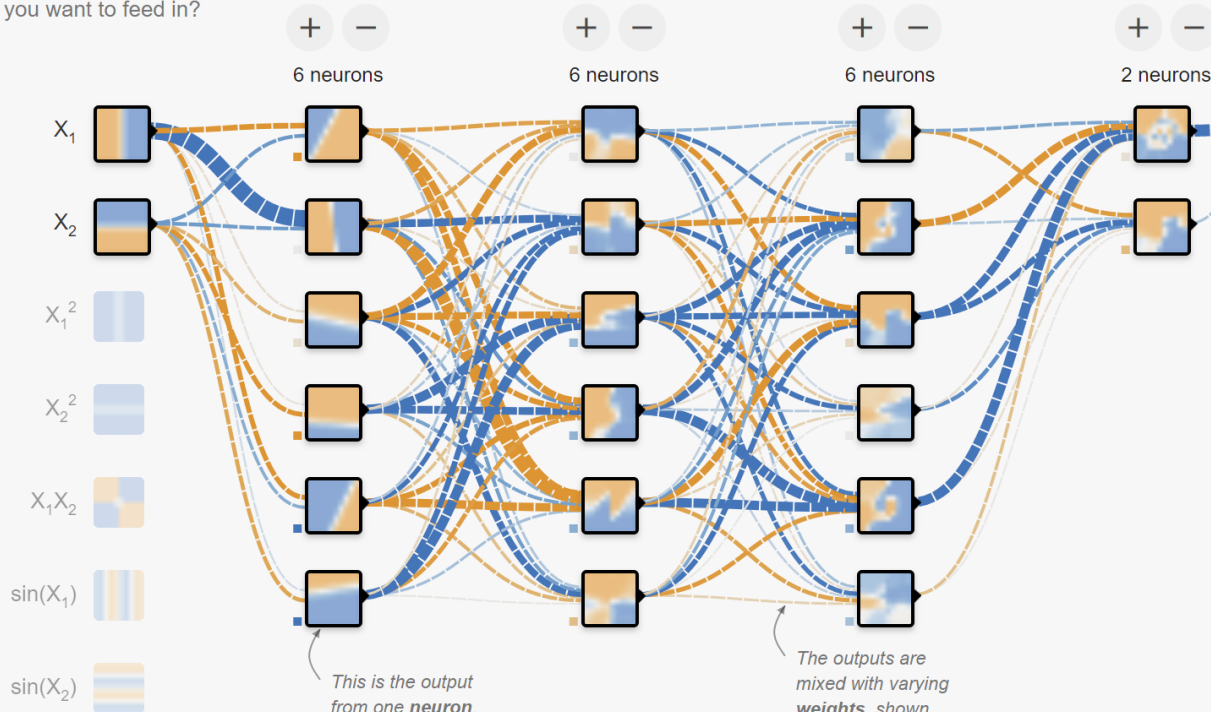


<https://playground.tensorflow.org>

Tensorflow Playground

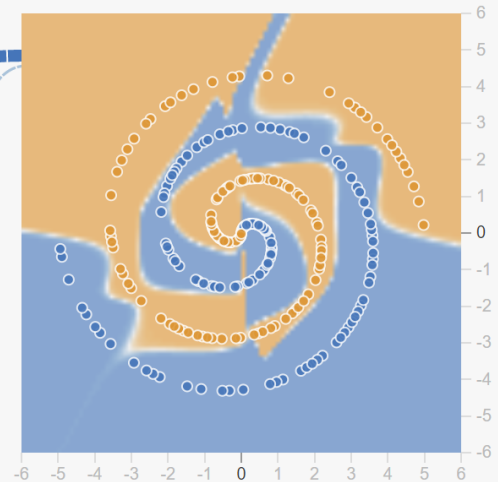
FEATURES

Which properties do you want to feed in?



OUTPUT

Test loss 0.044
Training loss 0.008



Colors shows data, neuron and weight values.

-1 0 1

☐ Show test data ☐ Discretize output

<https://playground.tensorflow.org>



Thanks

Email: libin@fudan.edu.cn