# 机器学习与数据科学

## Machine Learning and Data Science

主讲: 李春光

www.pris.net.cn/teacher/lichunguang

模式识别与智能系统实验室

信息与通信工程学院 网络搜索教研中心

北京邮电大学

# 引子：从身边的学习说起…

- 思考一下我们日常生活中的学习过程
  - 什么是学习
    - 学习: 获取知识、形成技能，获得适应环境、改变环境的能力的过程
  - 学习是怎样一个过程
    - 通过了解、认识、记忆、理解 "问题与答案" 来进行，进而达到能够自由运用与推广
  - 怎样评价学习的效果
    - 测试
      - 测试题目是原封不动所学过的：考察记忆能力
      - 测试题目是学过的但改变了形式的: 考察运用与推广能力
  - 怎样有效地学习

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 机器学习的定义

A computer program is said to learn from experience E

<span style="color:blue">with respect to some tasks T and</span>

<span style="color:red">performance measure P,</span>

<span style="color:blue">if its performance at tasks in T,</span>

<span style="color:red">as measured by P,</span>

improves with experience E.

E: 数据    P: 性能评价指标    T: 特定任务

# 概念辨析

- ## 数据挖掘
  - 源于数据仓库技术
  - 数据库中的知识发现(Knowledge-Discovery in Databases)

- ## 模式识别
  - 源于自动控制工程，跟人的识别能力相对应的各种应用问题及其相关联的一系列技术

- ## 机器学习
  - 人工智能的偏于理论的分支，从具体应用问题抽象出来的具有一般性的模型、算法与理论

- ## 数据科学
  - 与数据的感知、获取、处理和分析相关的模型、算法和理论
    - 对应于机器学习中的无监督部分，又区别于经典信号处理与分析
    - 区别与联系: 信号处理、应用数学、量子物理

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 学习问题的基本形式

- 按照是否有监督信息使用
  - 有监督学习(Supervised learning)
    - 分类: 输出为离散变量
    - 回归: 输出为连续变量
  - 无监督学习(Unsupervised learning)
    - 聚类: 寻找数据的某种划分
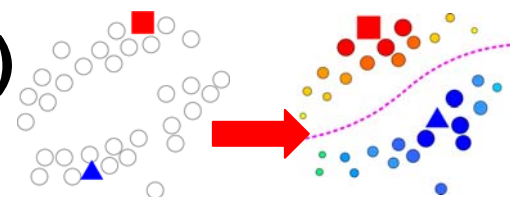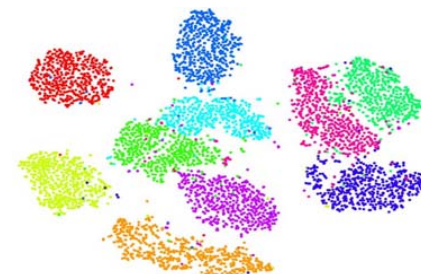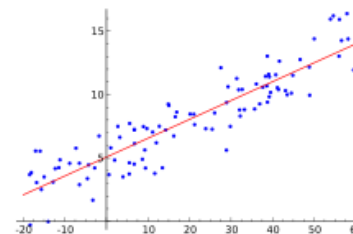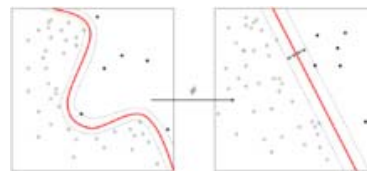    - 降维: 寻找数据的"坐标系"变换
  - 半监督学习(Semi-supervised learning)
- 表示符号说明
  - 训练数据(training data)

$$(X,Y) = \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n) \right\}$$

  - 数据本身

$$X = \left\{ \mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n \right\}$$

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 应用问题举例：I

- 手写数字识别
- 手写汉字识别
- 人脸识别
- 材料、树叶、花卉、食物识别
- 自然场景分类
- 图像物体归类
- 动态场景识别
- 文本分类
- 文本倾向性判断

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 手写数字图像识别

- 



from web

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室
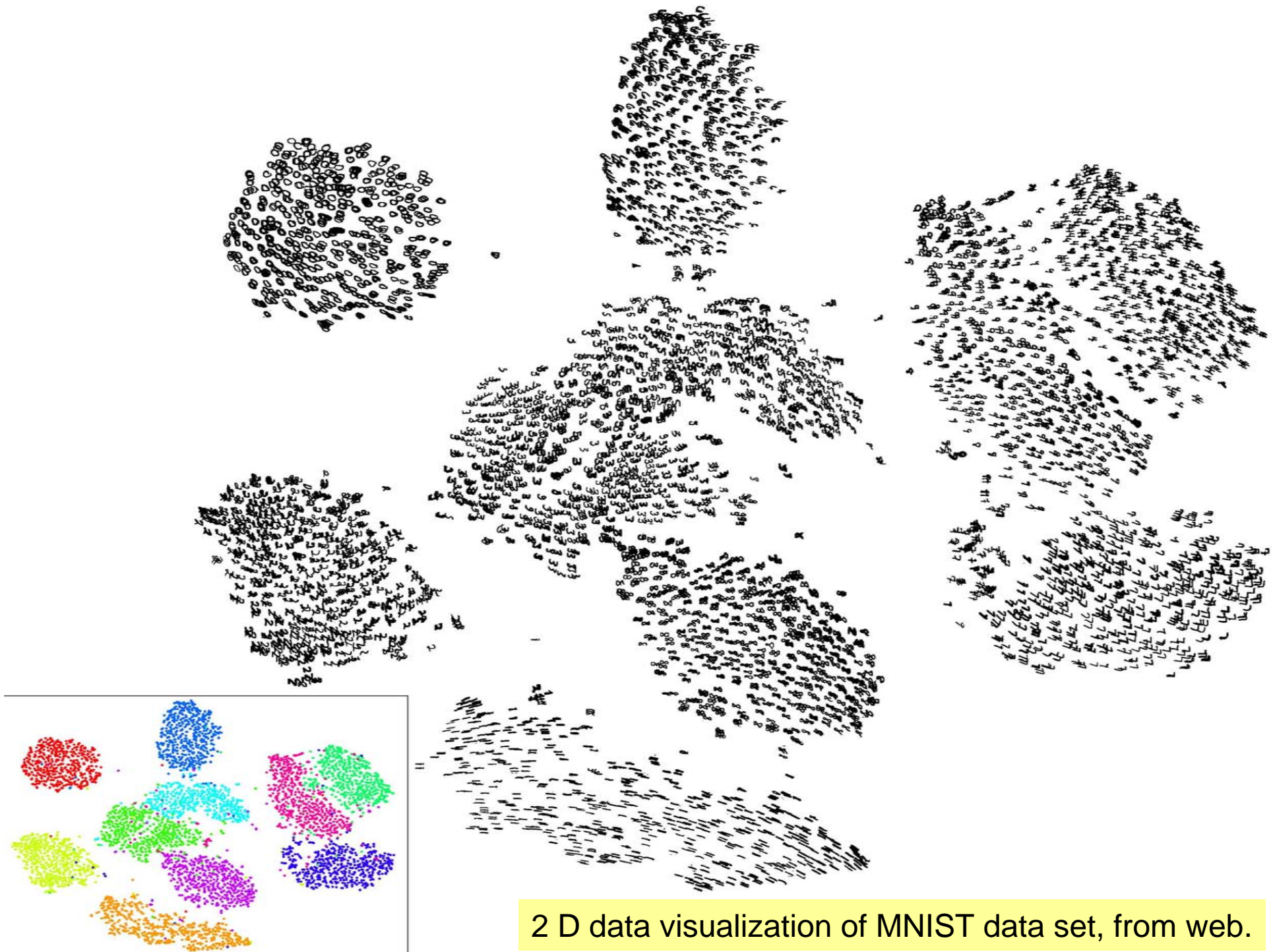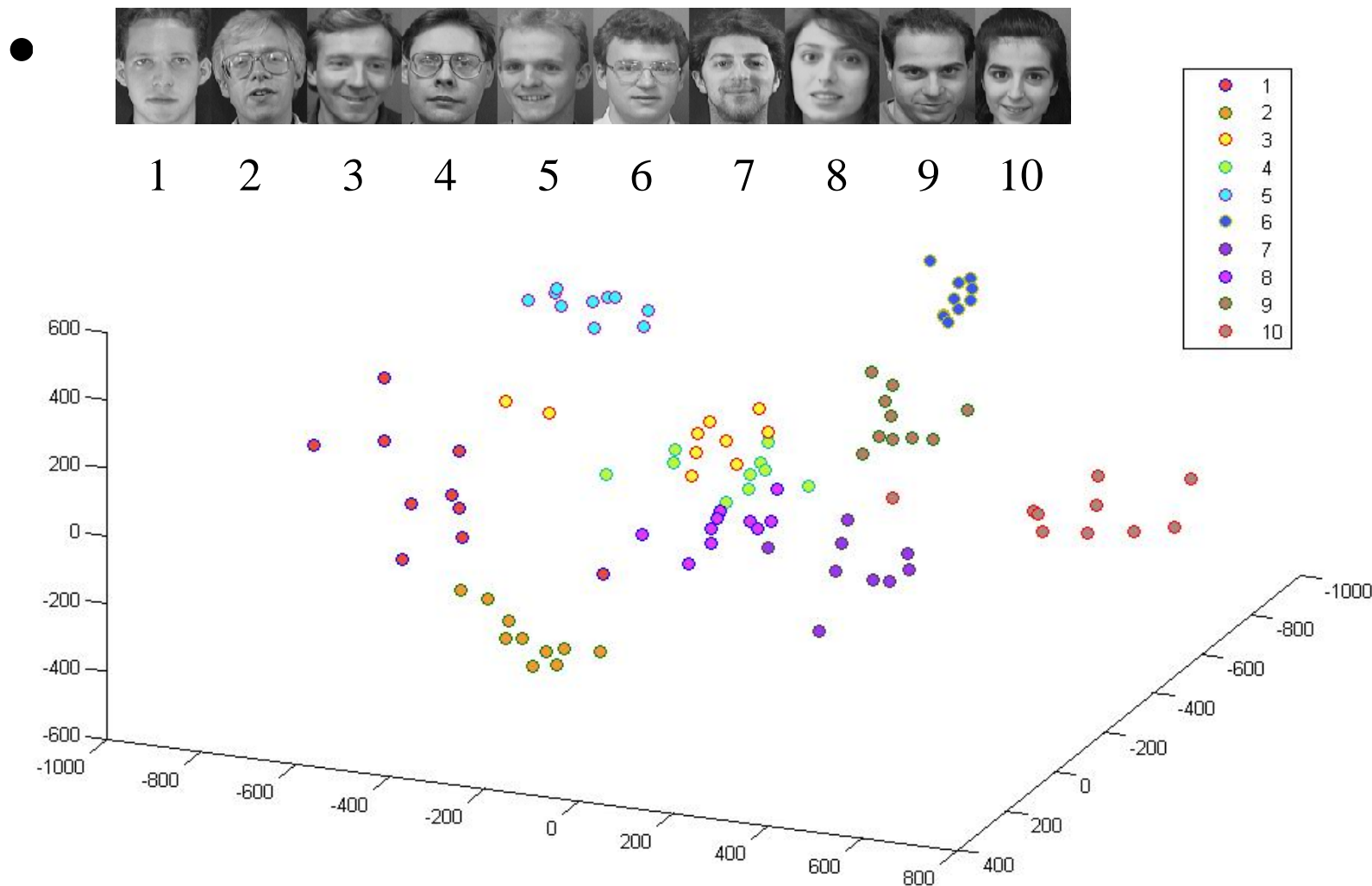
2 D data visualization of MNIST data set, from web.

# 手写汉字图像识别

- 啊阿埃挨哎唉哀皑癌蔼矮艾碍爱
隘鞍氨安俺按暗岸胺案肮昂盎凹
敖熬翱袄傲奥懊澳芭捌扒叭吧笆
八疤巴拔跋靶把耙坝霸罢爸白柏
百摆佰败拜稗斑班搬扳般颁板版
扮拌伴瓣半办绊邦帮梆榜膀绑棒
磅蚌镑傍谤苞胞包褒剥薄雹保堡
饱宝抱报暴豹鲍爆杯碑悲卑北辈

[1] H. Zhang et al., "HCL2000 -- A Large-scale Handwritten Chinese Character Database for Handwritten Character Recognition", ICDAR 2009, pp.286-290

# 人脸识别 (from ORL)



from ORL data set

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# ORL 数据集的可视化(by PCA)



前10个人在前3维特征脸中的投影分布

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 人脸识别 (from UMIST)

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# UMIST 数据集的可视化(by ISOMAP)



UMIST数据库中人脸数目大于30的人的前30张图像在三维嵌入中投影

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 自然场景图像分类



如何获得向量化表达?

From Dataset Scene-15

# 图像物体分类(VOC)



<mark>From Dataset Caltech 256</mark>

如何获得向量化表达?

*机器学习与数据科学 - Machine Learning & Data Science* 模式识别与智能系统实验室

# 材料识别



Glass　　Water　　Leather　　Paper　　Wood

[2] Xianbiao Qi, Rong Xiao, Chun-Guang Li, Yu Qiao, Jun Guo, and Xiaoou Tang, "Pairwise Rotation Invariant Co-occurrence Local Binary Pattern", IEEE TPAMI, Vol. 36, No. 11, Nov. 2014, pp.2199-2213.

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 树叶识别

[2] Xianbiao Qi, Rong Xiao, Chun-Guang Li, Yu Qiao, Jun Guo, and Xiaoou Tang, "Pairwise Rotation Invariant Co-occurrence Local Binary Pattern", IEEE TPAMI, Vol. 36, No. 11, Nov. 2014, pp.2199-2213.

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 花卉识别

●



Fig. 5. Some sample from Oxford Flower 102 data set. (The two images in each column from the same category).

[2] Xianbiao Qi, Rong Xiao, Chun-Guang Li, Yu Qiao, Jun Guo, and Xiaoou Tang, "Pairwise Rotation Invariant Co-occurrence Local Binary Pattern", IEEE TPAMI, Vol. 36, No. 11, Nov. 2014, pp.2199-2213.

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 食物识别

- 

Fig. 7. Sample images from PRID food data set.

[2] Xianbiao Qi, Rong Xiao, Chun-Guang Li, Yu Qiao, Jun Guo, and Xiaoou Tang, "Pairwise Rotation Invariant Co-occurrence Local Binary Pattern", IEEE TPAMI, Vol. 36, No. 11, Nov. 2014, pp.2199-2213.

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 动态纹理 (from YUPENN)



Fig. 6. Sample frames from dynamic scene data set YUPENN. Each image corresponds to a category of video sequence.

[3] Xianbiao Qi, Chun-Guang Li, Guoying Zhao, Xiaopeng Hong, Matti Pietikainen, "Dynamic texture and scene classification by transferring deep image features", Neurocomputing, Vol.171, 2016, pp:1230-1241.
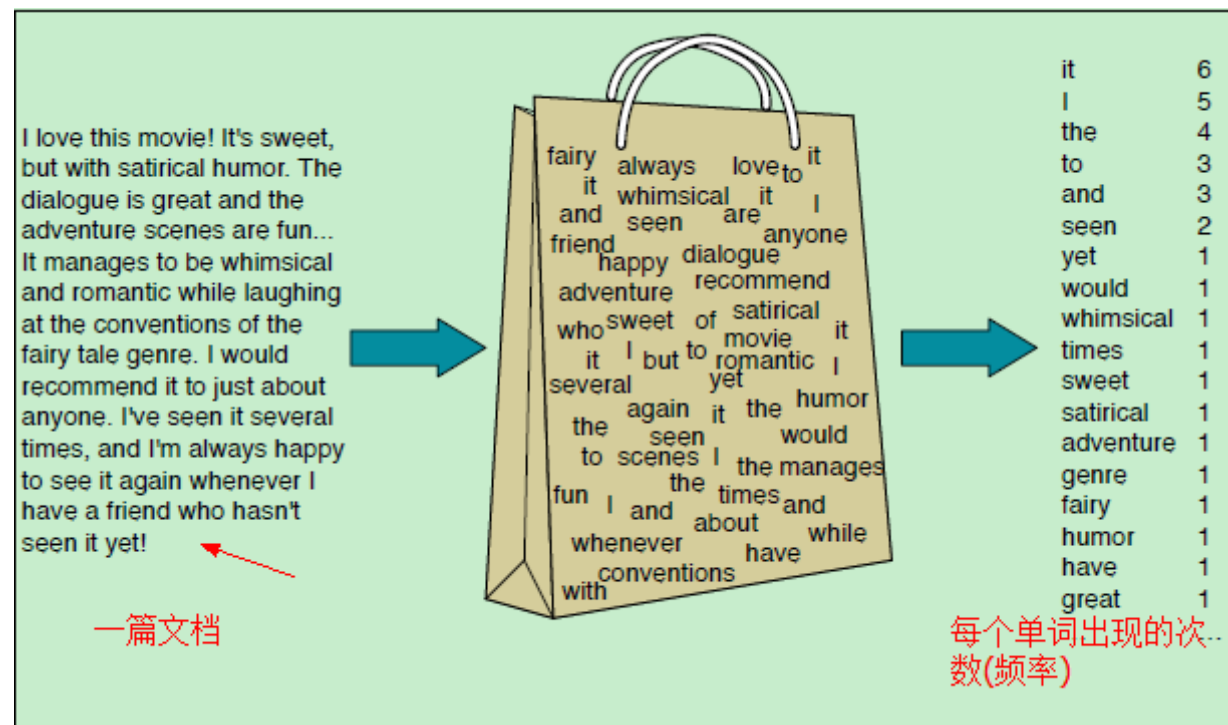
机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 如何构造一个学习问题？

- ## 收集数据
  - 建立数据的表示，有时也叫做特征抽取或学习

- ## 明确特定任务，建立任务的抽象模型
  - 分类 / 回归 / 推荐 / **Ranking** / 聚类 / ...

- ## 定义性能评价指标
  - 误差 / 错误率 ...

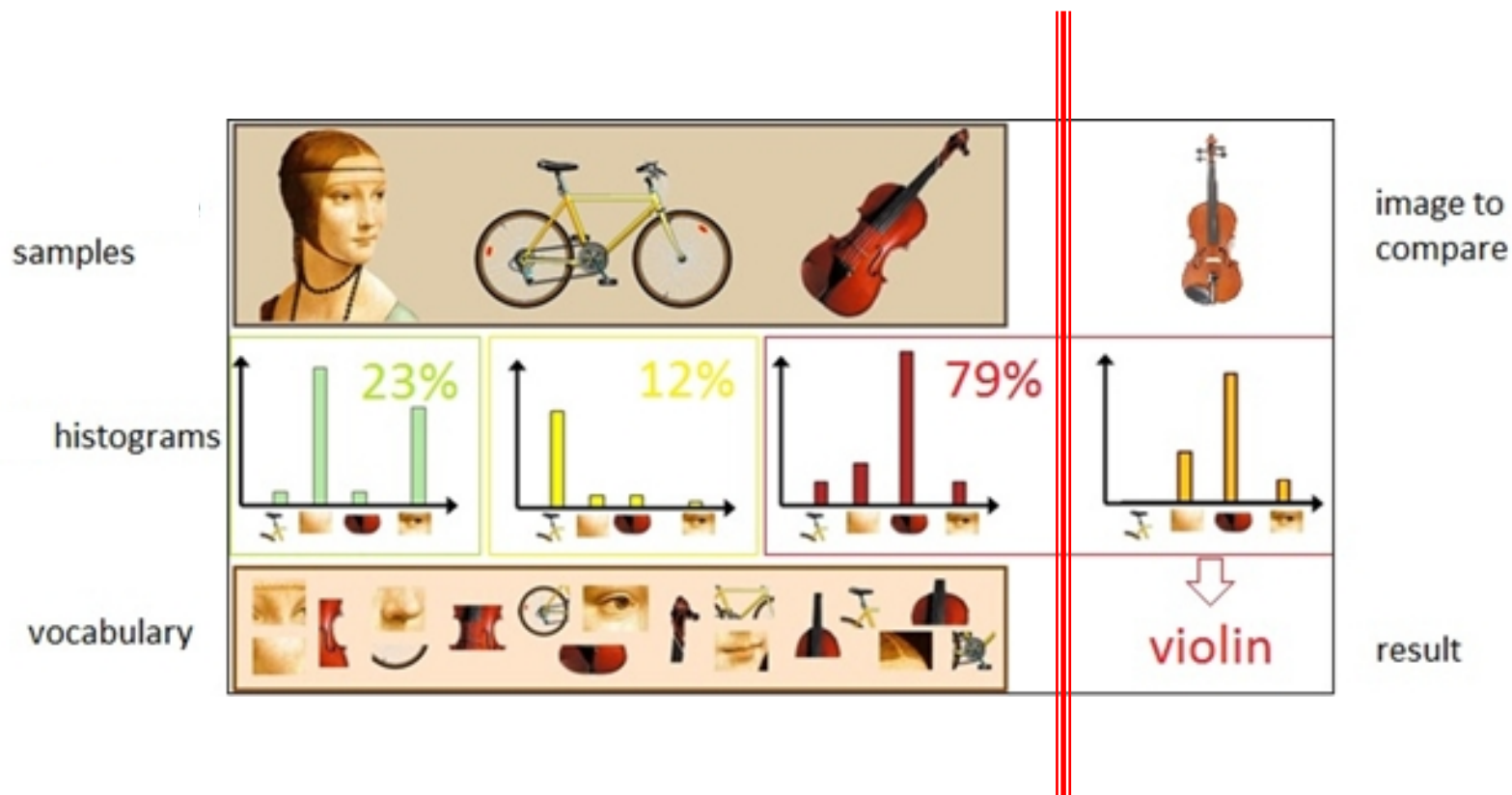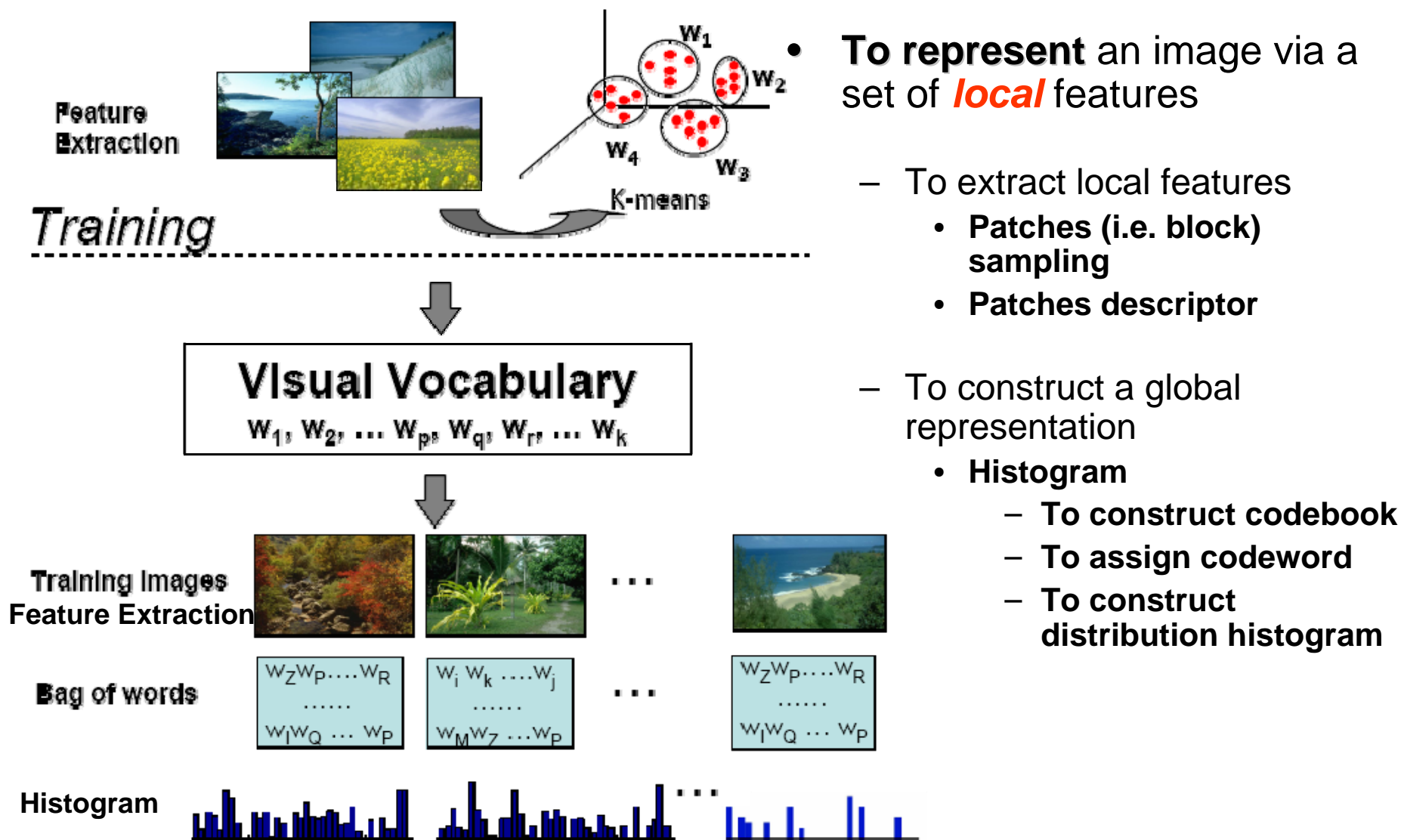# 应用问题举例 1：文本分类

- **基于词袋模型建立直方图特征向量**
  - **Bag of Words**

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 应用问题举例 2：图像检索

- 基于视觉词包模型建立图像的向量表示
  - **Bag of Visual Words**

# 视觉词袋(Bag of Visual Words)



- **To represent** an image via a set of *local* features

  – To extract local features
    - **Patches (i.e. block) sampling**
    - **Patches descriptor**

  – To construct a global representation
    - **Histogram**
      – **To construct codebook**
      – **To assign codeword**
      – **To construct distribution histogram**

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 应用问题举例 3：图像分类

- 一个复杂的应用问题往往需要一系列预处理步骤和特征抽取，也涉及到多种学习问题



- VOC: 图像物体分类(Visual Objects Categorization)
- BoW: 词句(Bag of Words)

机器学习与数据科学 Machine Learning + Data Science 模式识别与智能系统实验室

# 应用问题举例 4：动态场景分类

- **基于迁移的图像特征建立向量化表示**



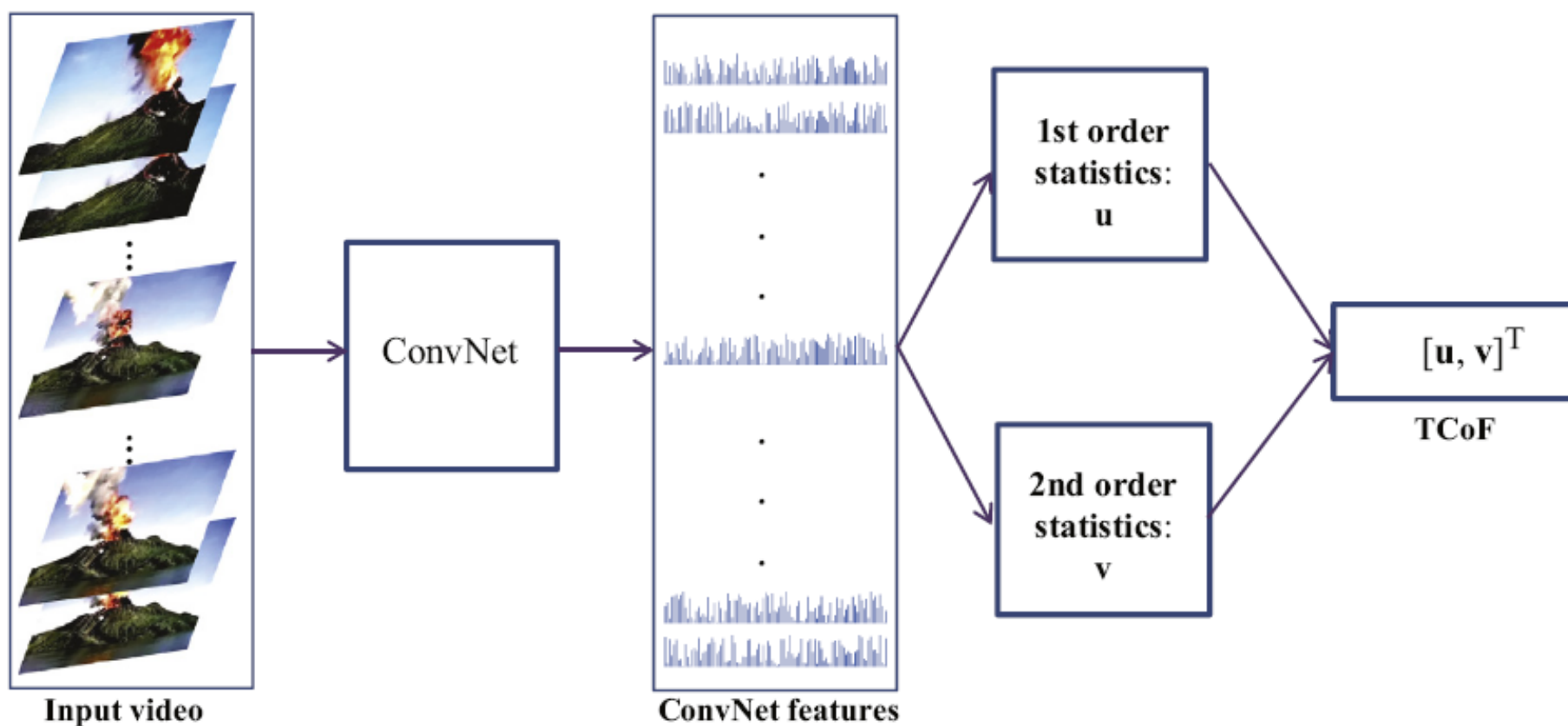**Fig. 4.** An illustration of our TCoF scheme.

# 学习问题的基本形式

- 按照是否有监督信息使用
  - 有监督学习(Supervised learning)
    - 分类: 输出为离散变量
    - 回归: 输出为连续变量
  - 无监督学习(Unsupervised learning)
    - 聚类: 寻找数据的某种划分
    - 降维: 寻找数据的 "坐标系" 变换
  - 半监督学习(Semi-supervised learning)
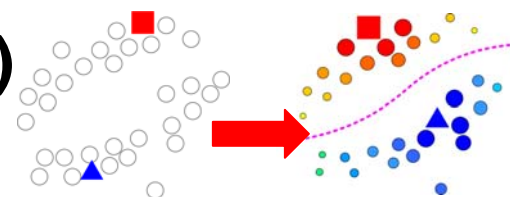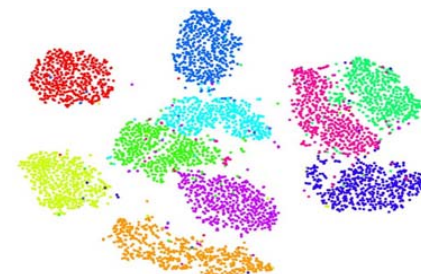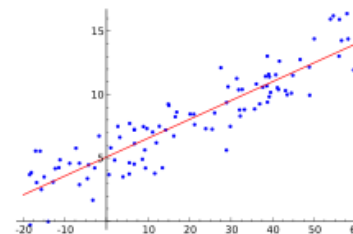- 表示符号说明
  - 训练数据(training data)

$$(X,Y) = \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n) \right\}$$

  - 数据本身

$$X = \left\{ \mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n \right\}$$

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 应用问题举例：II

- 数据降维
- 数据可视化
- 图像修复
- 背景与目标分离
- 歌声与背景音乐分离
- 视频中运动物体分割
- …

# UMIST 数据集的可视化(by ISOMAP)



UMIST数据库中人脸数目大于30的人的前30张图像在二维嵌入中投影

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

- 数据可视化

A

Up-down pose

Lighting direction

Left-right pose

B

[1] J. B. Tenenbaum, V. de Silva, J. C. Langford: "A Global Geometric Framework for Nonlinear Dimensionality Reduction", Science, 2000.

30

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

**Fig. 3.** Images of faces (*11*) mapped into the embedding space described by the first two coordinates of LLE. Representative faces are shown next to circled points in different parts of the space. The bottom images correspond to points along the top-right path (linked by solid line), illustrating one particular mode of variability in pose and expression.

[1] S. Roweis and L. Saul: "Nonlinear dimensionality reduction by locally linear embedding", Science, 2000. 31

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

**A**

master
image
television
paintings
academy
gallery
film
furniture
color
artists
decorative
artist
images
fine
painter
scenes
portrait
artistic
tube
sound
styles
PAINTING
LANDSCAPE
FIGURES
radio
formal
pieces
designs FIGURE
colors
garden
florence
light
inspired
baroque
glass
outstanding
elaborate architect
expression
objects
traditions
subject
design
renaissance
reflected
classical
contemporary
london
paris
medieval
ages ITALIAN
middle ITALY

**B**

LANDSCAPE PAINTING
subjects FIGURES
architectural FIGURE
house
law section
houses courts supreme congress
justice constitution president
architecture federal representatives
office
ITALIAN executive
senate
staff vote
ITALY parties powers
majority election
nuclear weapons party power
commander navy presidential
naval defense political
command american
air russia
military france
force russian
united britain
government forces
front french
battle troops
world allied japan
army british japanese
germany german
war

fought
fighting
captured
killed
defeat
peace
treaty
victory
campaign
invasion
attack

b

[1] J. B. Tenenbaum, V. de Silva, J. C. Langford: "A Global Geometric Framework for Nonlinear Dimensionality Reduction", Science, 2000.
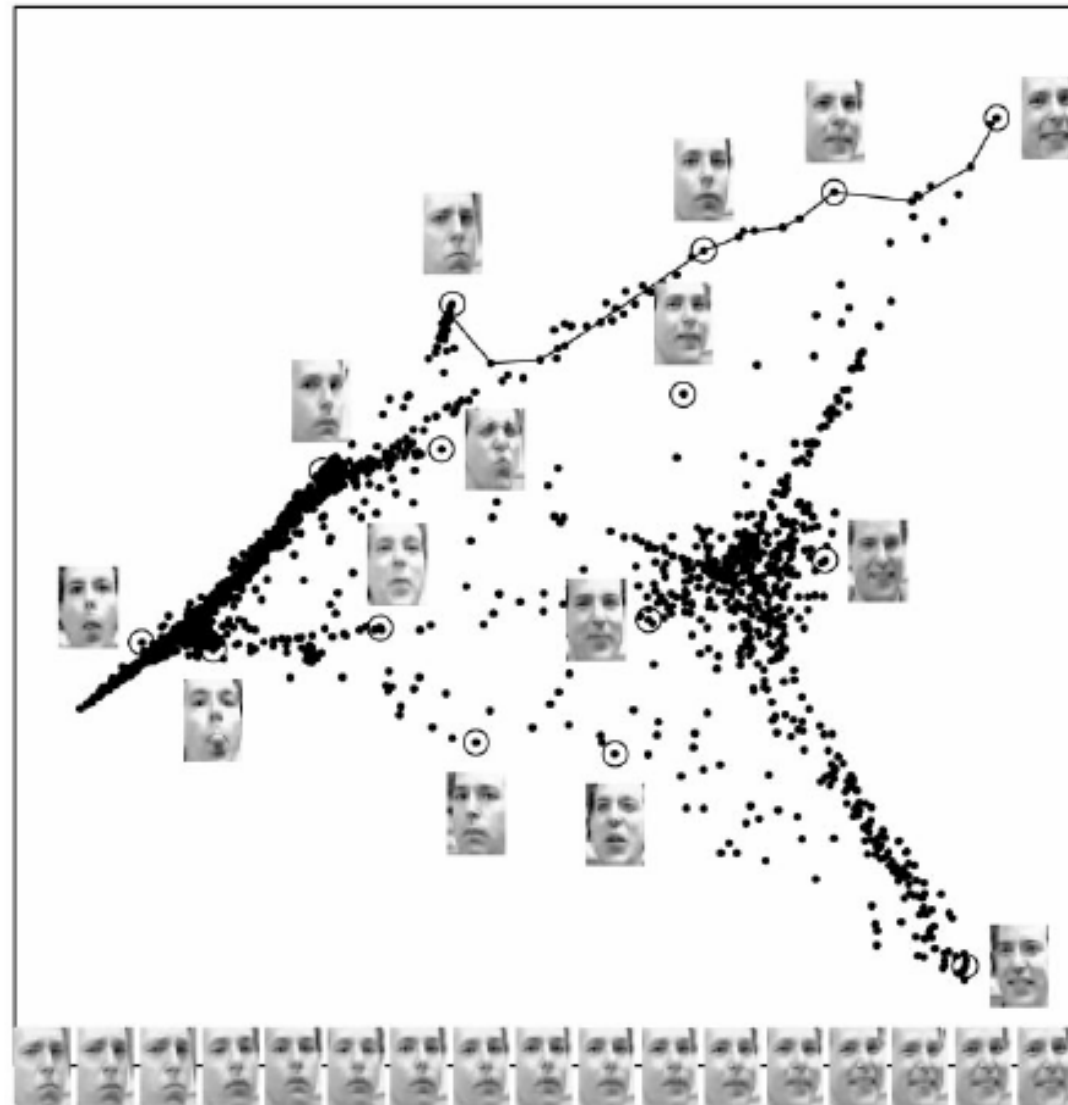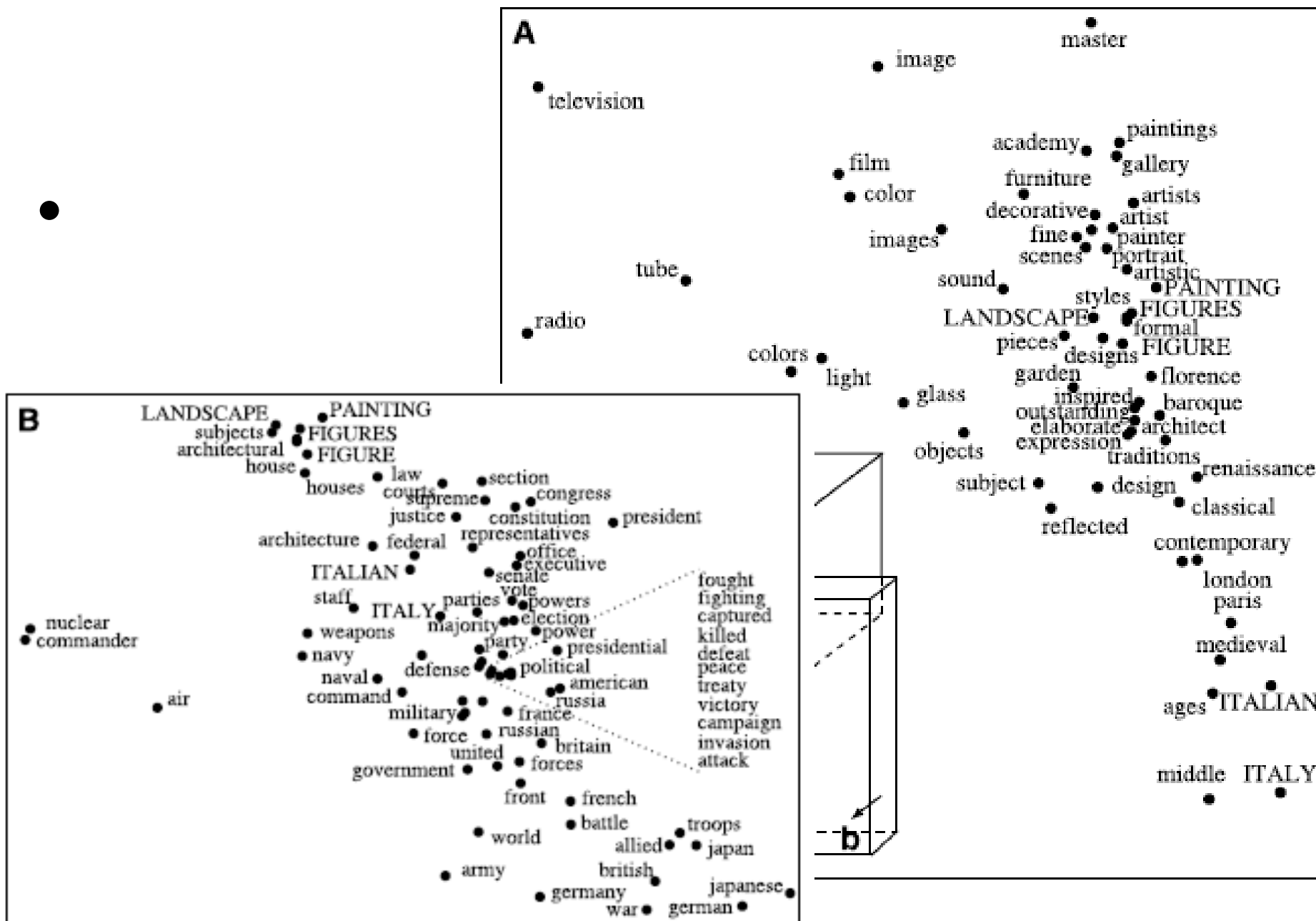
32

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 视频关键帧 (from TRECVID2006)



Sports    Weather    Office    Meeting    Desert

Mountain    Water    Corp. Leader    Police & Security    Military

Animal    Screen    US Flag    Airplane    Car
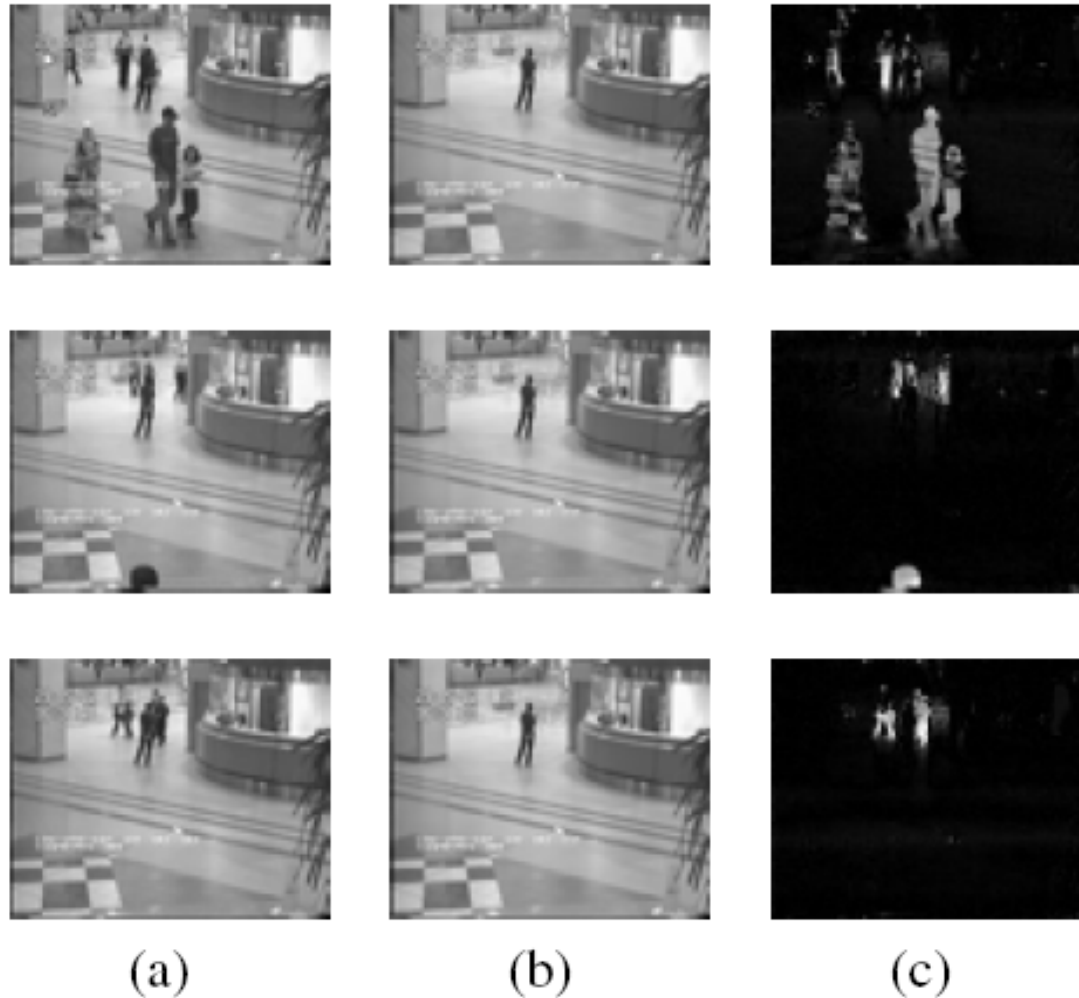
Truck    People Marching    Explosion    Map    Charts
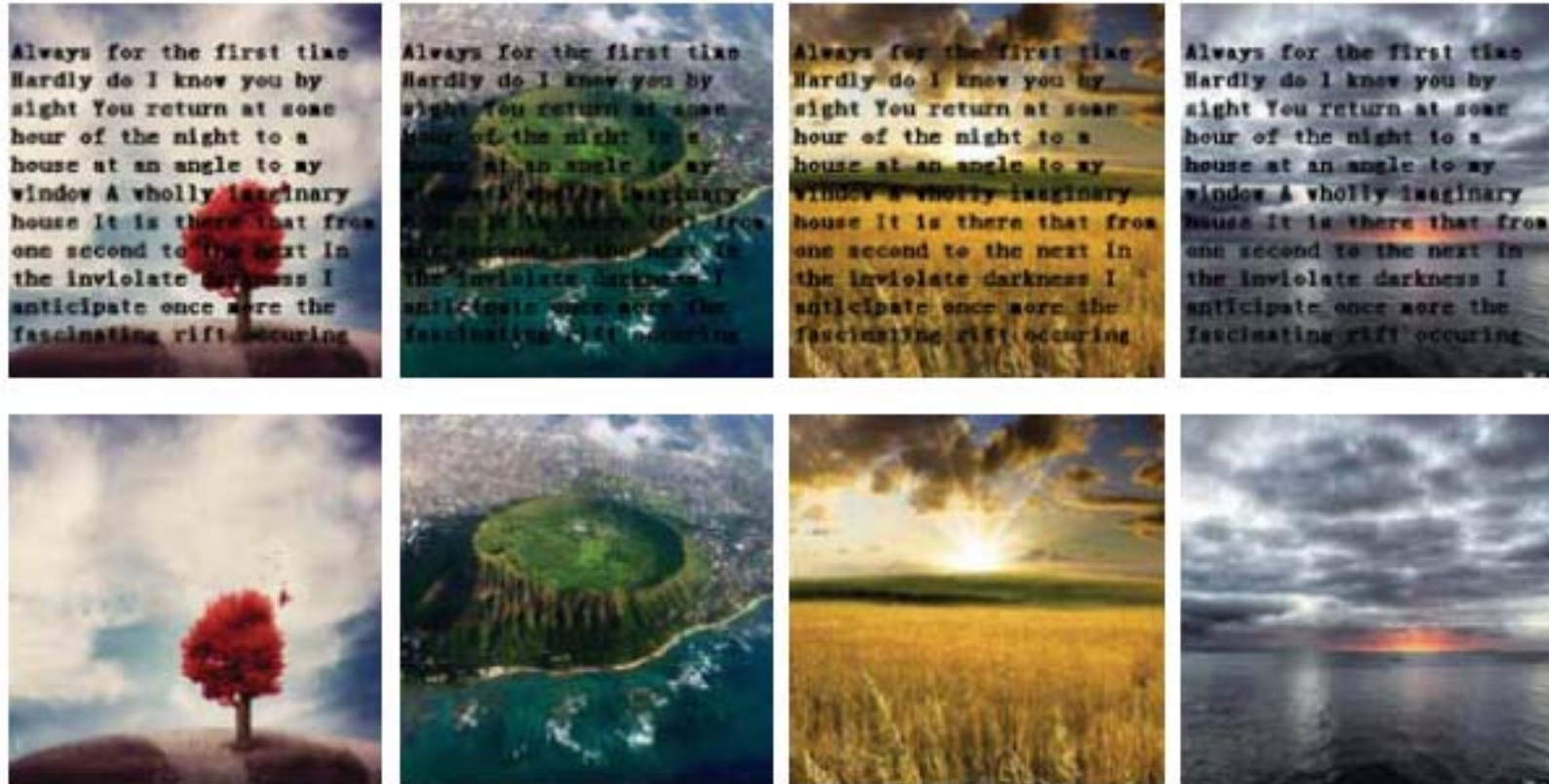
机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 视频中背景与目标分离

$$X \quad = \quad D \quad + \quad E$$

- 



(a)       (b)       (c)

[1] J. Wright, A. Ganesh, S. Rao, Y. Peng, Y. Ma. "Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization", NIPS, pp.2080-2088, 2009.
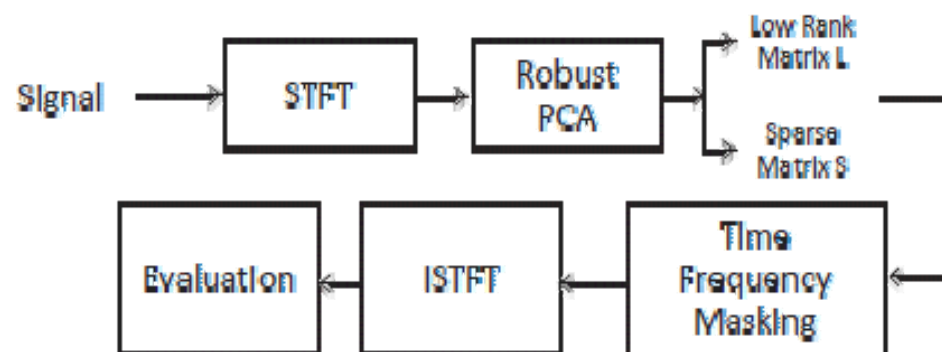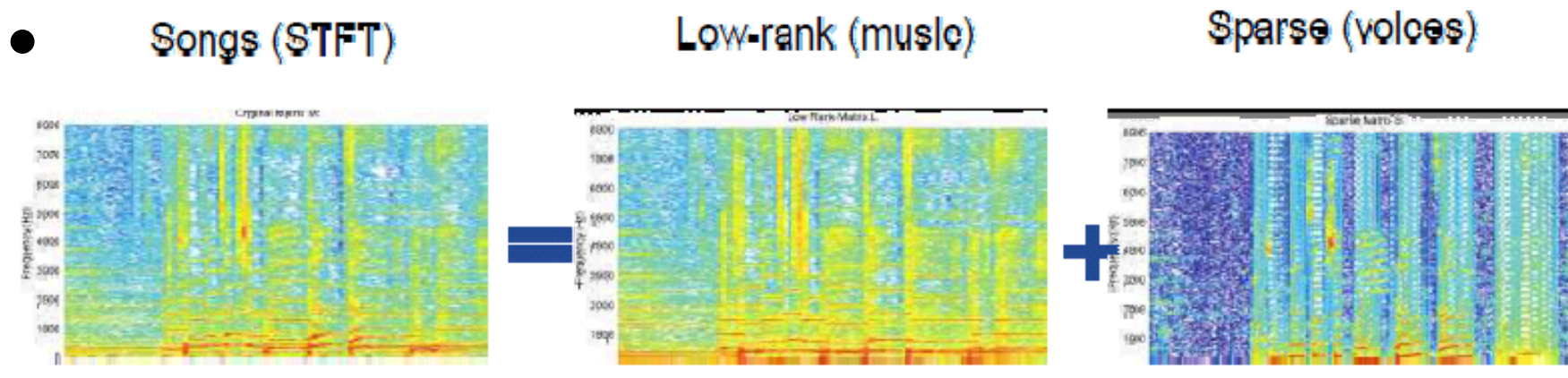
机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室
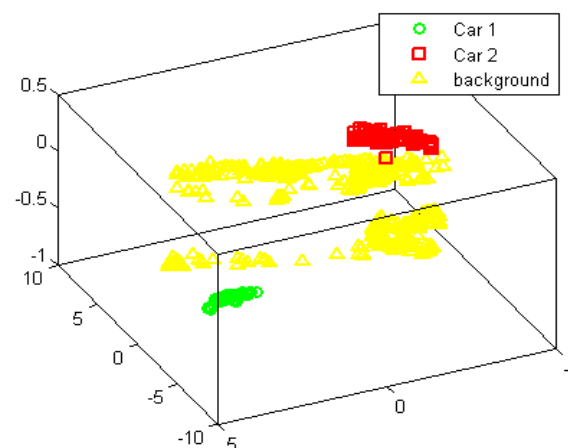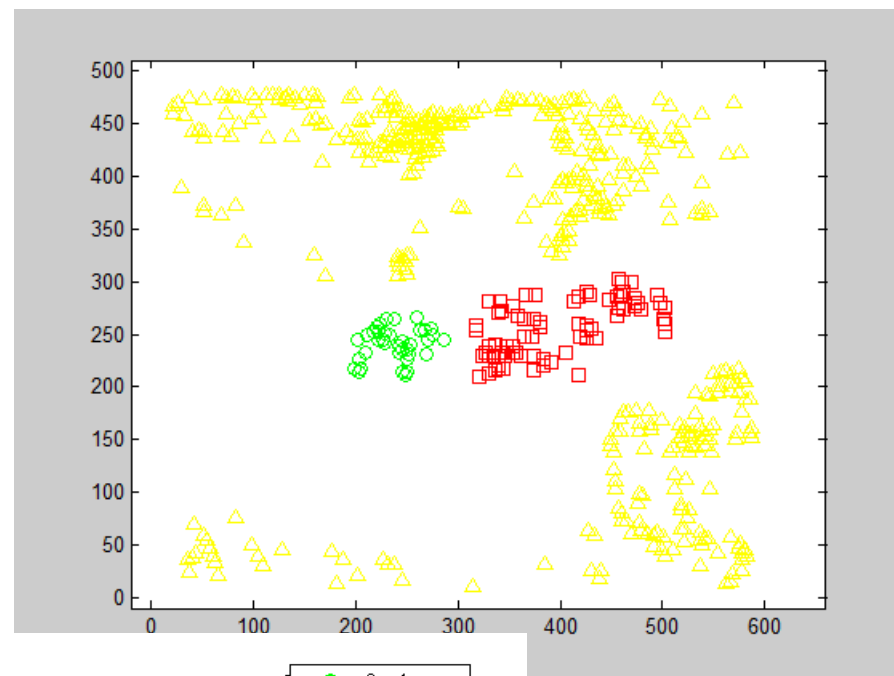
# 图像恢复(Image impainting)

# 歌声与背景音乐的分离

[1] Min, Zhang, Wright, Ma, CIKM 2010 / Sprechmann, Bronstein, Sapiro, ISMIR 2012

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 视频中的运动分割

[1] R. Tron and R. Vidal, "A benchmark for the comparison of 3-D motion segmentation algorithms," in IEEE Conference on Computer Vision and Pattern Recognition, 2007, pp. 1–8.

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 学习问题的其它形式

- 主动学习(Active learning)
  - 被动学习**(Passive learning)**

- 在线学习(Online learning)
  - 成批学习**(Batch learning)**

- 归纳学习(Inductive learning)
  - 传导学习**(Transductive learning)**

- 迁移学习(Transfer learning)
  - 借助一种学习去影响另一种学习

- 强化学习(Reinforcement learning)
  - 决策或控制领域，没有明确的监督信息，只能通过回报(reward)函数对每次决策作出反馈

- 深度学习(Deep learning)

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# Q / A

- Any Question? …

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室

# 专题 一：基于实例的学习

- **内容提要**
  - 引言
  - 最近邻规则 **(Nearest Neighbor Rule)**
    - 非线性回归模型
  - 帕森窗**(Parzen Windows)**
    - 密度估计问题的引出
    - Kernels
    - 瓦森-纳达拉亚估计器(Watson-Nadaraya Estimator)
  - 应用问题举例**:**
    - MNIST数据集 / VOC 与 BoW模型
  - 大数定律？收敛速度**?**

机器学习与数据科学 - *Machine Learning & Data Science* 模式识别与智能系统实验室