

# 数据库系统概论新技术篇

## 大数据与机器学习

Part II：大数据机器学习

卢志武

中国人民大学信息学院

2017年4月

# 目录

- ❖ 大数据机器学习的基本概念
- ❖ 大数据机器学习的实现平台
- ❖ 大数据机器学习的总结与反思



# 大数据给机器学习的挑战

- ❖ **数据源多样化**，能否自动地从每种数据源中提取特征？
- ❖ **数据量非常大**，算法的运行效率能否满足实际需求？
- ❖ **数据分布会发生变化**，学习模型的假设是否还成立？



# 大数据机器学习的特点

## ❖ 各种技术的融合

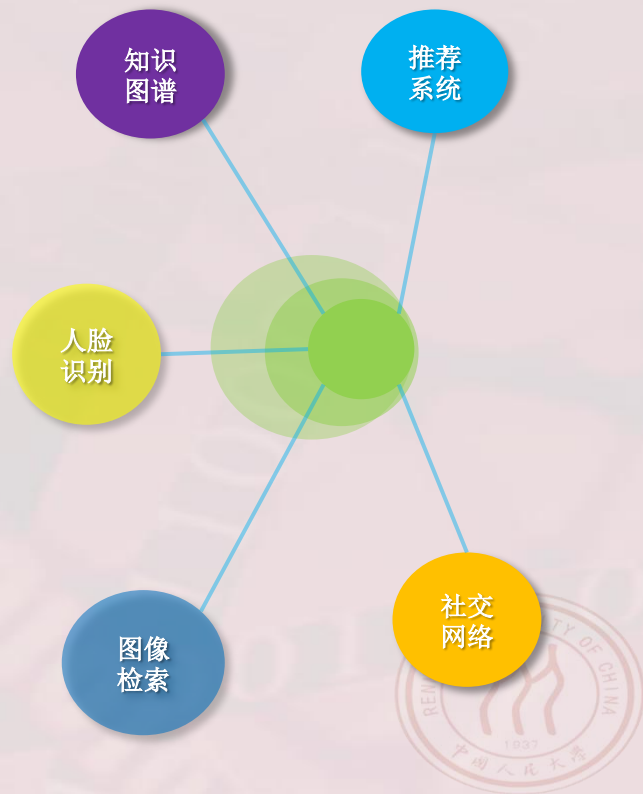
- 单一技术难以满足大数据分析多样性需求。
- **多种技术融合**，有助于提高系统稳定性。

## ❖ 数据理解是难点

- 海量数据高度非结构化，数据的快速准确理解成为关键。
- **深度学习**和**在线学习**在大数据理解领域显得尤为重要。

## ❖ 分类会逐渐弱化

- **快速检索**是有效利用海量数据的前提。
- 大数据时代，分类问题会逐渐弱化，而检索则会变得更重要。



# 大数据机器学习的关键技术

## ❖ 深度学习

- 很好地利用**GPU**等高性能计算设备。
- 解决多源特征的**自动提取**问题。

## ❖ 在线学习

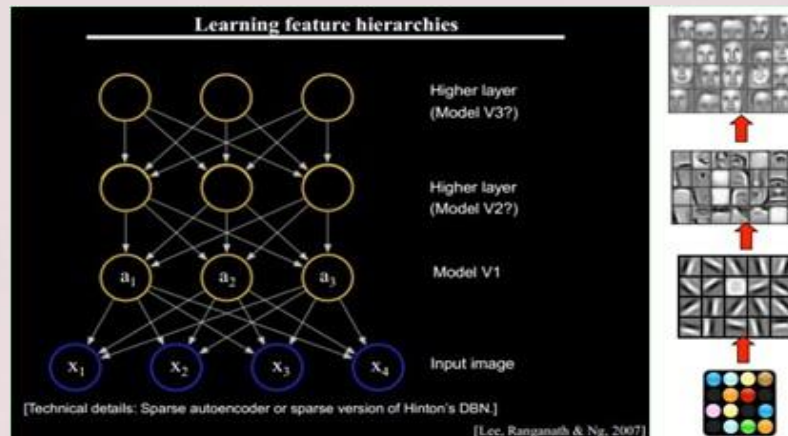
- 解决传统算法**更新模型代价过大**的问题。
- 有效地应对**大数据分布发生变化**的难题。

## ❖ 近似近邻搜索

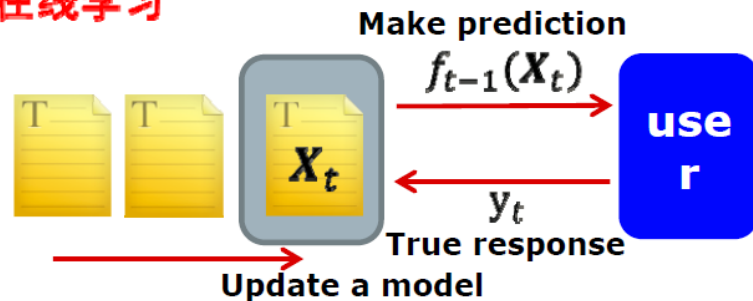
- 使得在**大数据上检索**成为可能。
- 常用方法：**哈希索引**和**基于树的索引**。

## ❖ 并行计算

- **GPU**并行计算：**CUDA**
- **分布式**系统：**Spark**



## 在线学习



# 目录

- ❖ 大数据机器学习的基本概念
- ❖ 大数据机器学习的实现平台
- ❖ 大数据机器学习的总结与反思





# 常用的大数据机器学习平台



# 大数据机器学习平台简介

- ❖ **Hadoop**是一个由**Apache**基金会所开发的分布式并行计算框架。该框架最核心的设计是：**MapReduce**算法和分布式文件系统**HDFS**。
- ❖ **Skytree**是**Skytree**公司开发的机器学习平台，结合先进的机器学习算法，为企业提供大数据高级分析。已用于推荐系统、预测分析、市场细分等。
- ❖ **Spark**是**UC Berkeley AMP lab**所开源的通用并行计算框架，基于**MapReduce**算法实现，拥有**Hadoop MapReduce**所具有的优点，但比**Hadoop**更通用，迭代算法效率更高。
- ❖ **GraphLab**是**CMU**的**Select**实验室提出的一个面向大规模机器学习/图计算的分布式内存计算框架，构建了四种流行的机器学习算法的并行版本。





# 目录

- ❖ 大数据机器学习的基本概念
- ❖ 大数据机器学习的实现平台
- ❖ 大数据机器学习的总结与反思



# 大数据机器学习无处不在



互联网搜索



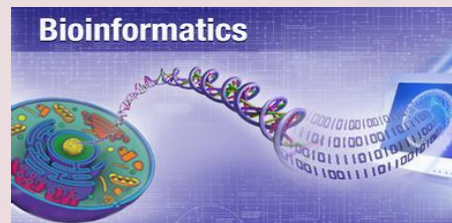
美国大选预测



网络舆情监控



医疗辅助诊断



生物信息学



汽车自动驾驶



# 大数据机器学习的缺陷

- ❖ 大数据可能存在**严重的偏差**，机器学习的结果令人不太满意。
- ❖ 大数据虽很大，但是**总存在例外的情况**。
- ❖ 大数据机器学习的**代价过大**，而人从小样本中就能很好地学习。



# 小数据机器学习

- ❖ 更关注机器学习模型的**可解释性**。
- ❖ 更关注**开放环境**下的机器学习，即机器学习模型的**稳定性**。
- ❖ 更关注**规则与知识**融入到机器学习中。



❖ 主讲人：卢志武

❖ Email: [luzhiwu@ruc.edu.cn](mailto:luzhiwu@ruc.edu.cn)

❖ 电话：010-62515670

