

# 大数据机器学习

中国科学院自动化研究所

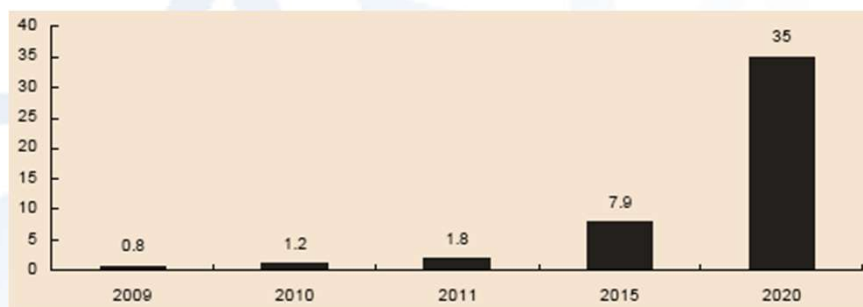
中科院大学人工智能学院

2019年12月19日

# 课程提纲

- 一、何谓大数据时代
- 二、大数据机器学习
- 三、典型应用举例
- 四、开源软件平台

## 数据量增加

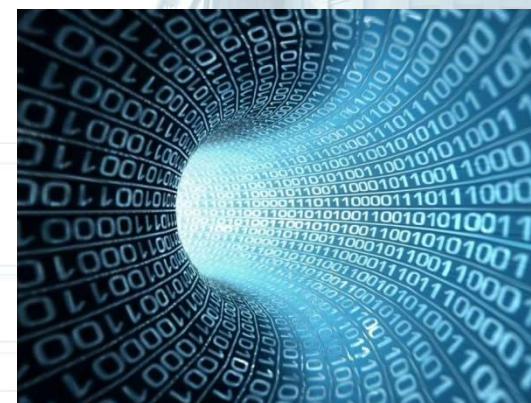


根据IDC 监测，人类产生的数据量正在呈指数级增长，大约每两年翻一番，这个速度在2020 年之前会继续保持下去。这意味着人类在最近两年产生的数据量相当于之前产生的全部数据量

TB  $\Rightarrow$  PB  $\Rightarrow$  EB  $\Rightarrow$  ZB

## 数据结构日趋复杂

大量新数据源的出现则导致了非结构化、半结构化数据爆发式的增长

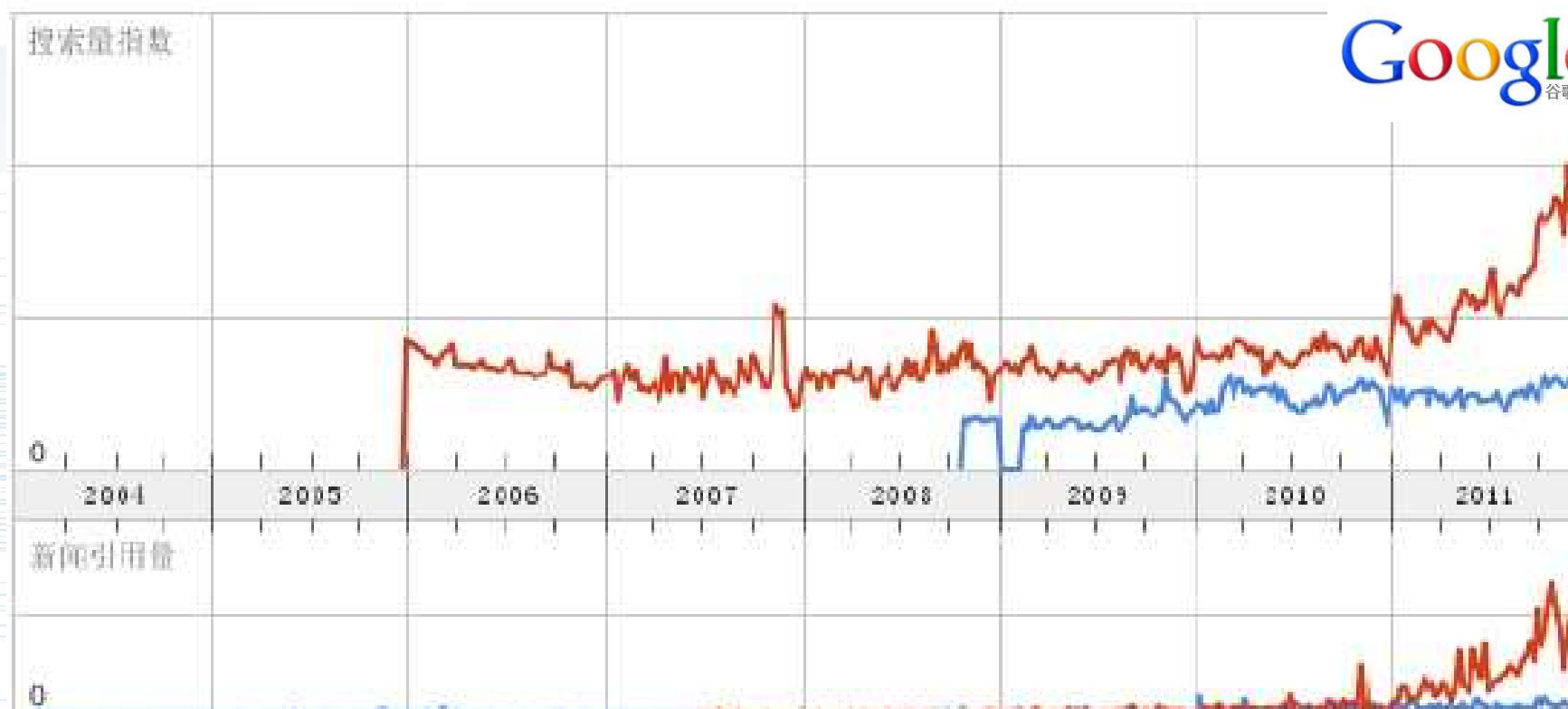


- 这些由我们创造的信息背后产生的这些数据早已经远远超越了目前人力所能处理的范畴
- 大数据时代正在来临…

EMC<sup>2</sup>

2011年5月，在“云计算相遇大数据”为主题的  
EMC World 2011 会议中，EMC 抛出了Big Data概念

● internet of things ● big data





# “大数据”成为学术界关注焦点



**nature**

Vol 455 | Issue no. 7209 | 4 September 2008

How to cope with the flood of data-- one of the **most daunting challenges** facing modern science.

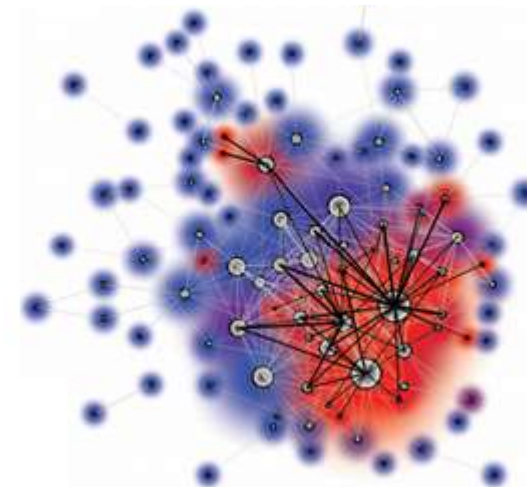


Detecting Novel Associations in Large Data Set  
David N. Reshef, *et al.*  
*Science* 334, 1518 (2011);  
DOI: 10.1126/science.1205438

## Detecting Novel Associations in Large Data Sets

David N. Reshef,<sup>1,2,3\*</sup> Yakir A. Reshef,<sup>2,4\*</sup> Hilary K. Finucane,<sup>5</sup> Sharon R. Grossman,<sup>2,6</sup> Gillean McVean,<sup>3,7</sup> Peter J. Turnbaugh,<sup>6</sup> Eric S. Lander,<sup>2,8,9</sup> Michael Mitzenmacher,<sup>10</sup> Pardis C. Sabeti<sup>2,4</sup>‡

Identifying interesting relationships between pairs of variables in large data sets is increasingly important. Here, we present a measure of dependence for two-variable relationships: the maximal information coefficient (MIC). MIC captures a wide range of associations both functional and not, and for functional relationships provides a score that roughly equals the coefficient of determination ( $R^2$ ) of the data relative to the regression function. MIC belongs to a larger class of maximal information-based nonparametric exploration (MINE) statistics for identifying and classifying relationships. We apply MIC and MINE to data sets in global health, gene expression, major-league baseball, and the human gut microbiota and identify known and novel relationships.



# 大数据开启了一次重大的时代转型

- 望远镜让我们能够感受宇宙
- 显微镜让我们能够观测微生物
- 大数据正在改变我们的生活和理解世界的方式，  
成为新发明和新服务的源泉
- 更多的改变正蓄势待发！！！！

# 课程提纲

一、何谓大数据时代

二、大数据机器学习

三、典型应用举例

四、开源软件平台

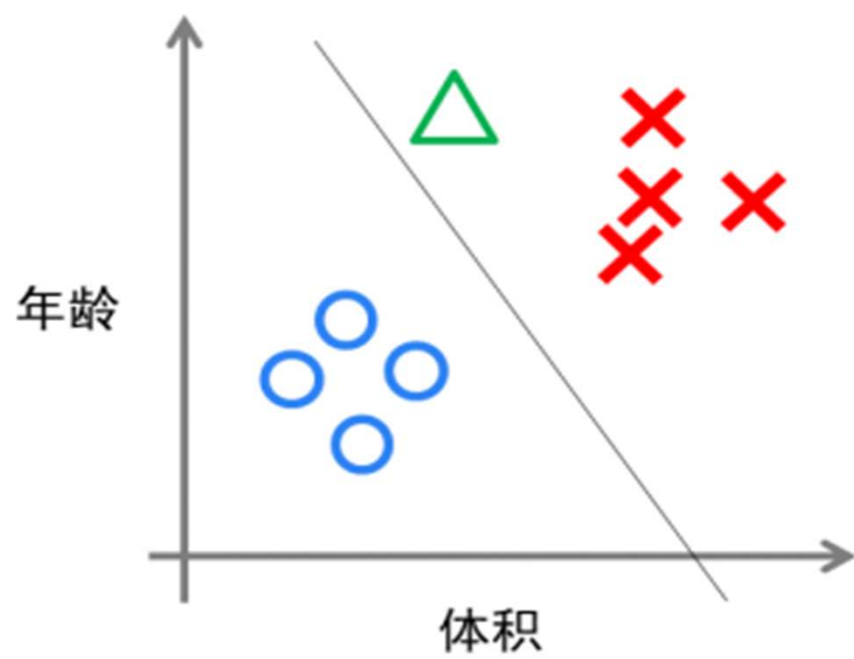


2014年5月16日，吴恩达加入百度，担任百度公司首席科学家，负责百度研究院、Baidu Brain计划

Geoffrey Hinton:Deep Learning祖师爷，多伦多大学教授，负责Google大脑

Facebook人工智能实验室主任、NYU数据科学中心创始人、深度学习界的泰斗Yann LeCun





目标：预测肿瘤的性质

输入：肿瘤的体积，  
患者的年龄

输出：良性或恶性



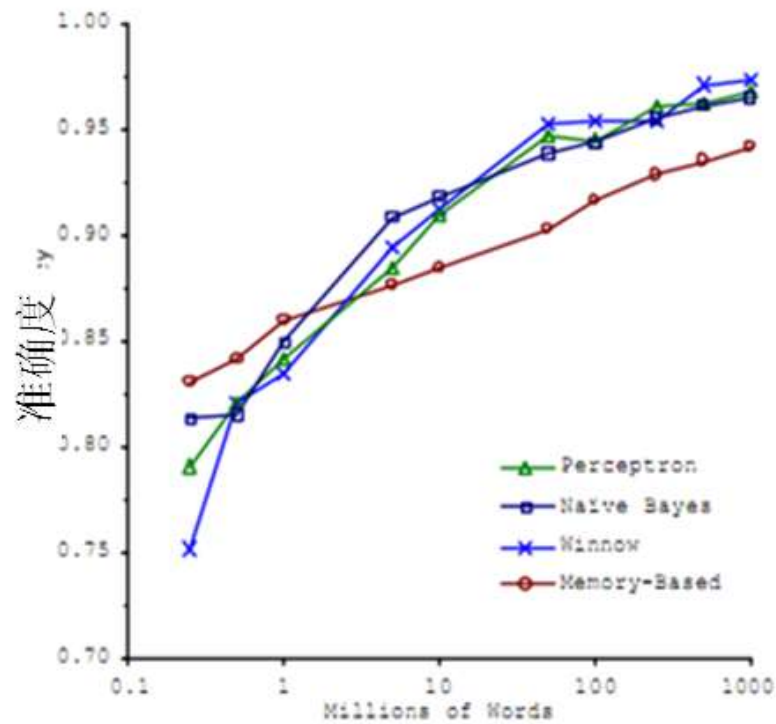
Google

## 淘汰赛赛事预测

更新时间：07月11日

🔵 赛前预测    🟡 准确预测





训练数据大小 (百万)

It's not who has the best algorithm that wins .

It's who has the most data.

# 一种流行的大数据机器学习方法

深度学习介绍【下次课详细讲】



2006年，Geoffrey Hinton在科学杂志《Science》上发表了一篇文章，论证了两个观点：

1. 多隐层的神经网络具有优异的特征学习能力，学习得到的特征对数据有更本质的刻画，从而有利于可视化或分类
2. 深度神经网络在训练上的难度，可以通过“逐层初始化”来有效克服

2006年，Geoffrey Hinton在科学杂志《Science》上发表了一篇文章，论证了两个观点：



2012年6月,《纽约时报》披露了Google Brain项目,这个项目是由Andrew Ng和Map-Reduce发明人Jeff Dean共同主导,用16000个CPU Core的并行计算平台训练一种称为“深层神经网络”的机器学习模型,在语音识别和图像识别等领域获得了巨大的成功,Andrew Ng就是文章开始所介绍的机器学习的大牛

2013年4月,《麻省理工学院技术评论》杂志将深度学习列为2013年十大突破性技术(Breakthrough Technology)之首

2015年10月阿尔法围棋以5：0完胜欧洲围棋冠军、职业二段选手樊麾。2016年3月挑战世界围棋冠军、职业九段选手李世石。根据日程安排，5盘棋将分别于3月9日、10日、12日、13日和15日举行，即使一方率先取得3胜，也会下满5盘。最后以4:1结束了这场“战争”。



4月5日

00:03

“弯弯\_2016”在优酷上传视频20160403北京望京798和颐酒店女生遇袭]

00:06

@弯弯\_2016 在微博上发布第一条信息分享该视频

00:12

@弯弯\_2016 建立微博话题 #和颐酒店女生遇袭#

15:32

接连发布10条微博，持续质疑涉事的酒店、携程、警方等方面

20:10

发布[整理了我被劫持的经过和事态发展到现在的结果]

被@所长别开枪是我 @牛文文 @我的厕所读物 @八哥专用 等大V转发，迅速引爆舆论

21:27

@弯弯\_2016 称酒店终于来电，欲给钱删微博

@平安北京 以回复留言形式表示已通过相关单位开展核实

舆论反响不大，微博仅有少量转发

大V纷纷转发，舆论爆发

微博截图



# 课程提纲

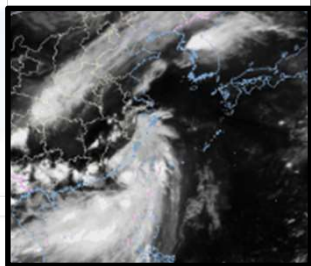
- 一、何谓大数据时代
- 二、大数据机器学习
- 三、典型应用举例
- 四、开源软件平台

- 气象数据真实存在的大数据
- 大数据颠覆了降水计算模式
- 精准降水估计高效学习算法

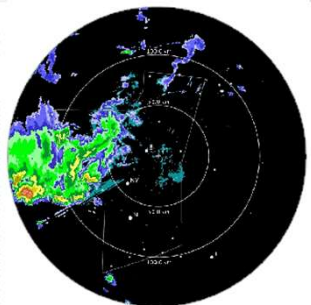
# 多源气象探测信息构成了名副其实的大数据

海量实时、多源异构、复杂多变

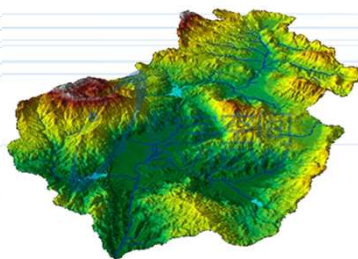
卫星云图



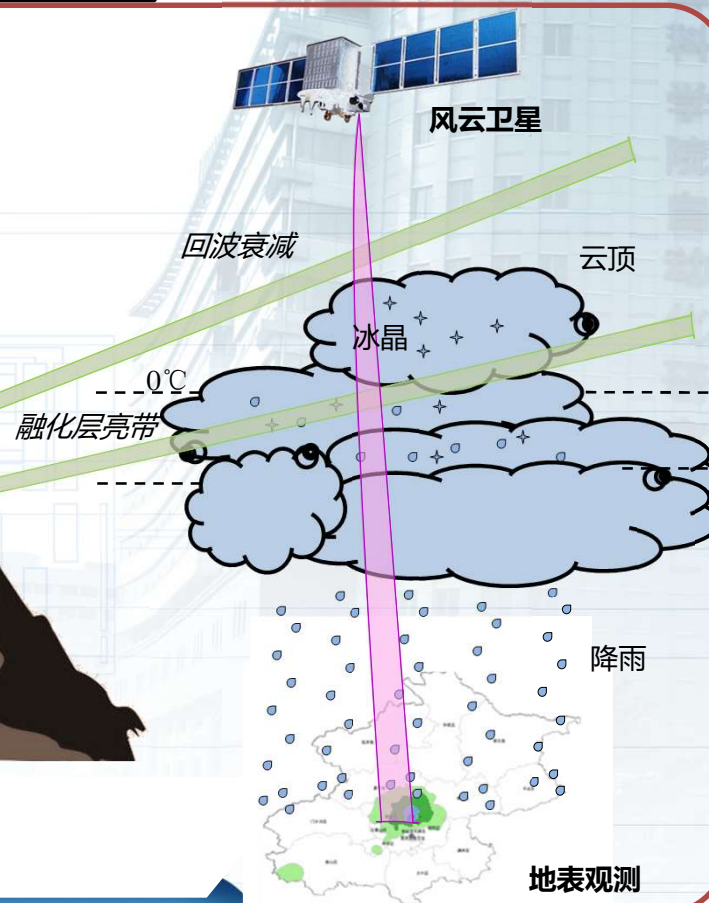
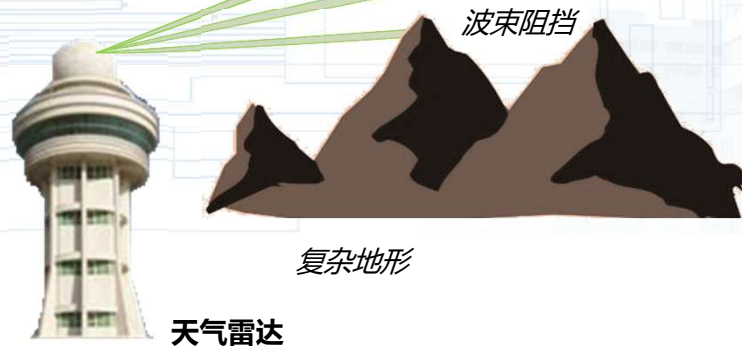
雷达回波



地理信息



- 23颗气象卫星92类卫星产品  
5000万文件3PB卫星数据
- 21类雷达产品超10PB数据  
600万观测/雷达/6分钟
- 6万自动气象站实时数据  
全国30m级DEM数据



- 气象数据真实存在的大数据
- 大数据颠覆了降水计算模式
- 精准降水估计高效学习算法



# 超短临降水预报服务迫切需求



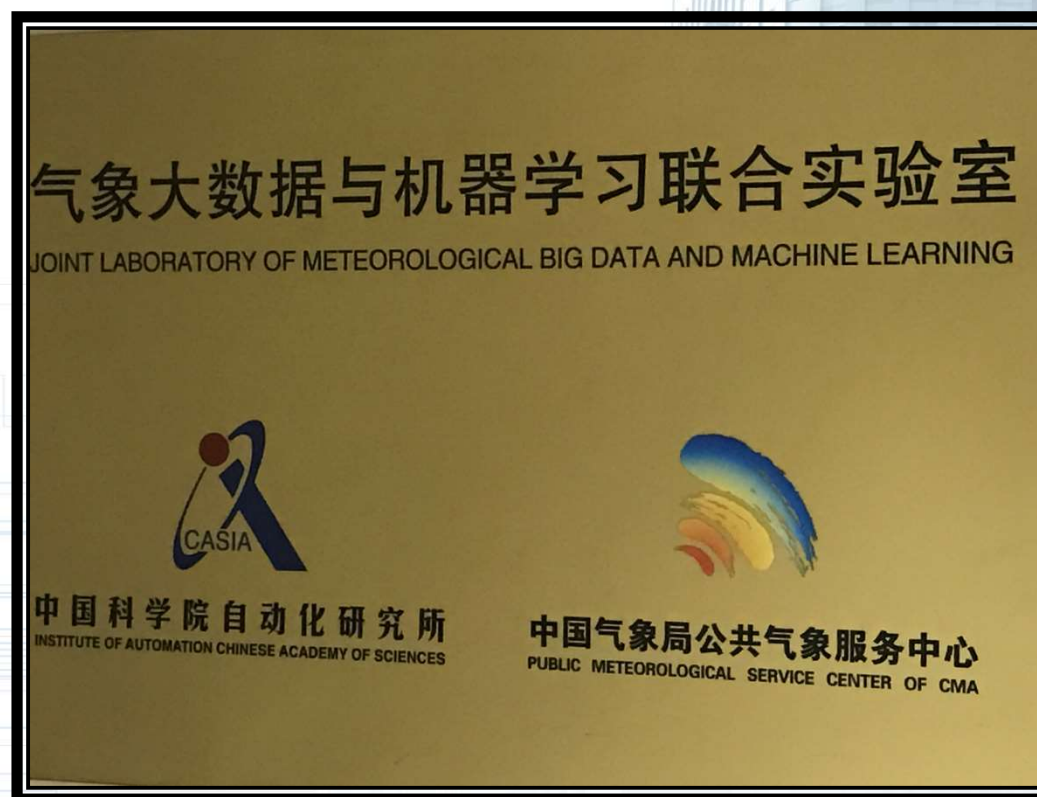
气象观测技术飞速发展，但我国目前还无法做到  
分钟级精确降水预报





# 超短临降水是联合实验室业务需求

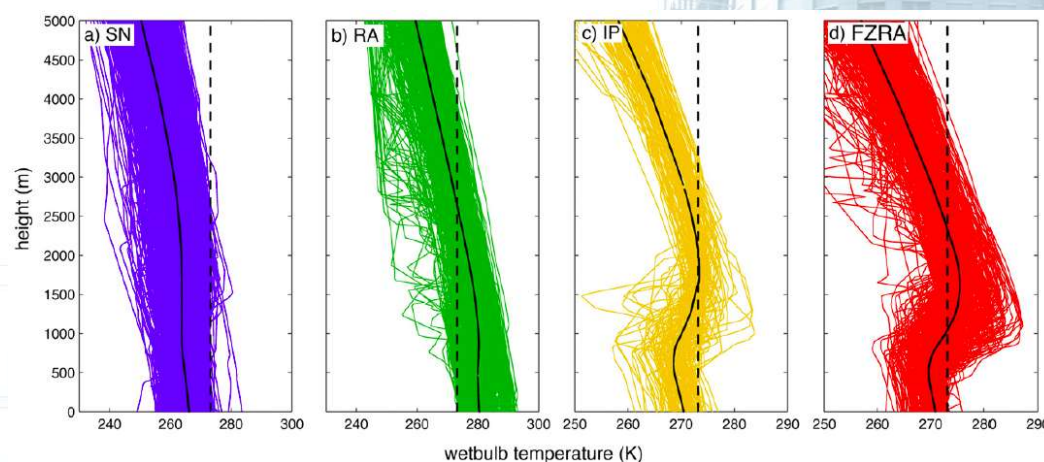
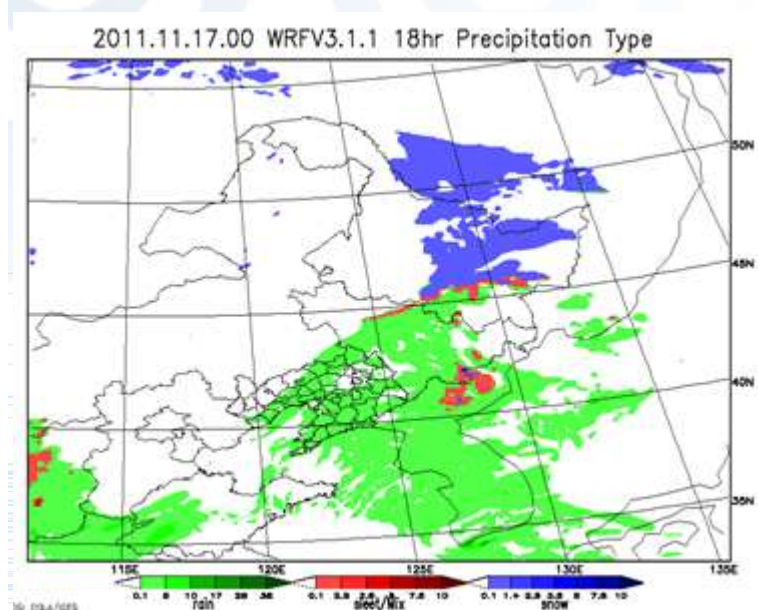
- 2015年6月，国家气象局公共气象服务中心和中国科学院自动化研究所共同组建了“气象大数据与机器学习联合实验室”
- 超短临降水预报(分钟级)是联合实验室业务需求





# 降水相态识别：特别是春秋降水估计

如果将雨误判成雪，可能会将小雨级别的降水误报为大雪



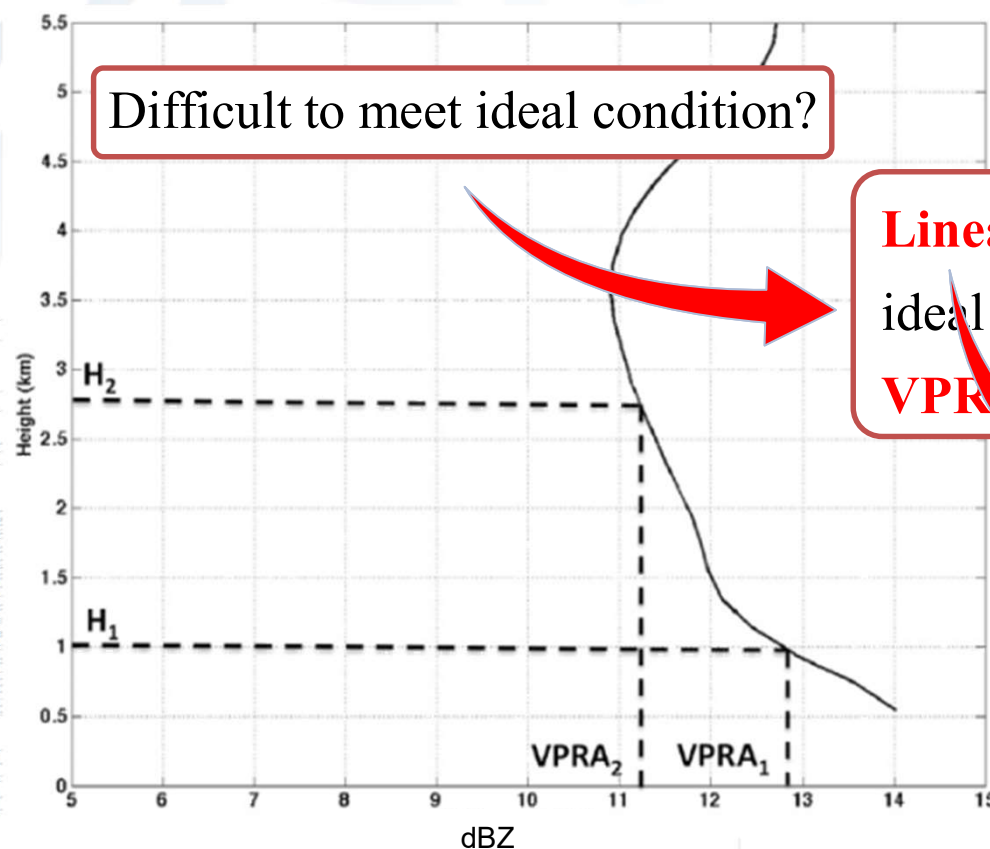
在我国探空气球非常少

降水相态主要与空间上不同高度的温度分布有关

- 气象数据真实存在的大数据
- 大数据颠覆了降水计算模式
- 精准降水估计高效学习算法

# VPR Structure

$$\log(Z) = b \log(R) + \log(a) \quad \rightarrow \quad R = f(Z)$$



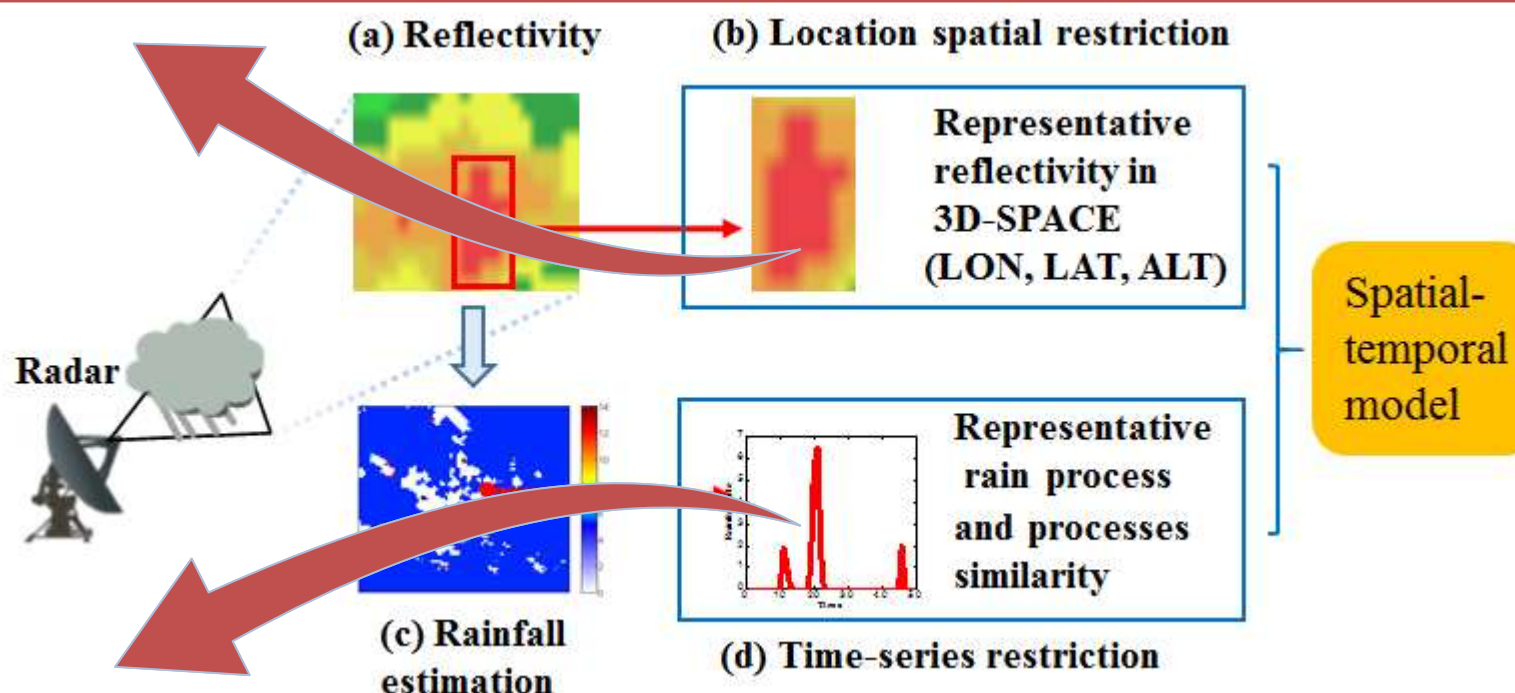
$$\rho_{\log(Z^h), \log(R)} = \frac{\text{cov}(\log(Z^h), \log(R))}{\sigma_{\log(Z^h)} \sigma_{\log(R)}}$$

Which level to select? → Consider VPR structure in statistical model



# Spatial and temporal structure

**spatial structure**  $\min_{\Theta} f(\Theta, x', y) = \sum \left( \sum_{v \in S_j} \left( [y(\Theta) - \bar{y}_j(\Theta)]^2 \right) - \sum_{i \in L, R} \sum_{v \in S_j} \left( [y(\Theta) - \bar{y}_j(\Theta)]^2 \right) \right)$



**temporal structure**  $\min_{\omega, \lambda} g(\omega, \lambda) = \sum_{k=1}^p \left[ \sum_{j=1}^{q(k)} \log(Z_k(x)) - \sum_{j=1}^{q(k)} \sum_{t=1}^T \omega_{kt} P_{kj}(y_{t-1}, y_t, x_t) + \lambda \|\omega_k\|_1 \right]$

Rainfall space-time characteristic?

→ Consider spatial and temporal structure in statistical model

# Result

Comparison between radar rainfall estimation and raingauge measurement for Hangzhou radar coverage

	RMSE <sub>mm</sub>	MAE <sub>mm</sub>	CC
<i>Z-R</i>	3.05	1.64	0.659
<i>SVR</i>	2.69	1.51	0.717
<i>RF</i>	2.80	1.41	0.718
<i>RANMP</i>	2.72	1.36	0.737
<i>LIST</i>	2.51	1.14	0.763
<i>RANLIST</i>	<b>2.15</b>	<b>1.05</b>	<b>0.829</b>

1

Compared to conventional Z-R relationship:

Improvements of **30% in RMSE**, **36% in MAE** and **26% in CC** are obtained

2

Compared to random forest:

Improvements of **23% in RMSE**, **26% in MAE** and **15% in CC** are obtained

# 实验效果

Data set	$P_{min}$ (%)						$P_{avg}$ (%)					
	Base	SSMOTE	GCS	ABNC	OSMOTE	AMDO	Base	SSMOTE	GCS	ABNC	OSMOTE	AMDO
Contraceptive	43.55 <sub>6.04</sub>	38.12 <sub>6.78</sub>	40.24 <sub>5.81</sub>	36.02 <sub>4.51</sub>	<b>47.10</b> <sub>12.20</sub>	45.36 <sub>6.98</sub>	52.07 <sub>2.15</sub>	48.30 <sub>3.56</sub>	49.24 <sub>0.74</sub>	49.21 <sub>2.19</sub>	51.86 <sub>3.04</sub>	<b>52.52</b> <sub>2.09</sub>
Flare	2.22 <sub>4.97</sub>	20.56 <sub>14.65</sub>	<b>30.28</b> <sub>6.39</sub>	20.56 <sub>12.36</sub>	6.67 <sub>9.94</sub>	25.56 <sub>18.26</sub>	59.84 <sub>1.51</sub>	61.37 <sub>3.40</sub>	63.04 <sub>1.93</sub>	58.02 <sub>1.58</sub>	58.57 <sub>1.80</sub>	<b>63.45</b> <sub>1.74</sub>
Thyroid	96.40 <sub>2.49</sub>	97.01 <sub>2.99</sub>	97.61 <sub>2.47</sub>	95.79 <sub>2.68</sub>	<b>98.81</b> <sub>1.64</sub>	93.99 <sub>2.97</sub>	96.70 <sub>1.58</sub>	96.96 <sub>0.49</sub>	97.34 <sub>0.89</sub>	96.55 <sub>0.74</sub>	<b>97.72</b> <sub>0.81</sub>	95.67 <sub>0.98</sub>
Car	63.08 <sub>22.03</sub>	61.54 <sub>19.61</sub>	89.23 <sub>6.88</sub>	67.69 <sub>26.87</sub>	61.54 <sub>26.09</sub>	<b>98.46</b> <sub>3.44</sub>	77.26 <sub>5.59</sub>	76.81 <sub>4.82</sub>	91.72 <sub>4.12</sub>	78.54 <sub>5.63</sub>	80.37 <sub>9.40</sub>	<b>91.97</b> <sub>1.31</sub>
Nursery	67.71 <sub>3.70</sub>	68.64 <sub>6.64</sub>	77.45 <sub>6.54</sub>	66.48 <sub>3.31</sub>	64.36 <sub>6.46</sub>	<b>93.90</b> <sub>2.85</sub>	90.12 <sub>1.19</sub>	90.32 <sub>1.94</sub>	93.45 <sub>1.78</sub>	89.84 <sub>1.05</sub>	89.31 <sub>1.86</sub>	<b>95.75</b> <sub>1.20</sub>

□ **Xuebing Yang**, Qiuming Kuang, Wensheng Zhang, and Guoping Zhang, “AMDO: an Over-Sampling Technique for Multi-Class Imbalanced Problems,” *IEEE Transactions on Knowledge and Data Engineering*, 2017. (SCI, CCF-A)

TABLE 9  
Means and Standard Deviations of AUCm (%) and MAUC (%) Results from Base, SSMOTE, GCS, ABNC, OSMOTE and AMDO for the Nominal/Mixed-type Data Sets

Data set	AUCm (%)						MAUC (%)					
	Base	SSMOTE	GCS	ABNC	OSMOTE	AMDO	Base	SSMOTE	GCS	ABNC	OSMOTE	AMDO
Contraceptive	64.23 <sub>2.14</sub>	60.41 <sub>3.20</sub>	60.78 <sub>1.44</sub>	61.57 <sub>1.32</sub>	64.04 <sub>3.10</sub>	<b>64.63</b> <sub>2.39</sub>	64.05 <sub>1.61</sub>	61.22 <sub>2.67</sub>	61.93 <sub>0.55</sub>	61.91 <sub>1.64</sub>	63.89 <sub>2.28</sub>	<b>64.39</b> <sub>1.56</sub>
Flare	63.18 <sub>0.96</sub>	67.42 <sub>4.40</sub>	<b>70.28</b> <sub>1.76</sub>	66.68 <sub>3.30</sub>	63.93 <sub>2.43</sub>	69.53 <sub>4.81</sub>	75.91 <sub>0.90</sub>	76.82 <sub>2.04</sub>	77.82 <sub>1.16</sub>	74.81 <sub>0.95</sub>	75.14 <sub>1.08</sub>	<b>78.07</b> <sub>1.04</sub>
Thyroid	97.84 <sub>1.11</sub>	98.08 <sub>0.86</sub>	98.37 <sub>0.94</sub>	97.63 <sub>0.87</sub>	<b>98.82</b> <sub>0.66</sub>	96.84 <sub>1.07</sub>	97.53 <sub>1.19</sub>	97.72 <sub>0.37</sub>	98.00 <sub>0.67</sub>	97.41 <sub>0.55</sub>	<b>98.29</b> <sub>0.60</sub>	96.75 <sub>0.73</sub>
Car	81.64 <sub>6.67</sub>	81.08 <sub>5.80</sub>	94.52 <sub>3.03</sub>	82.99 <sub>7.54</sub>	82.94 <sub>9.49</sub>	<b>95.52</b> <sub>1.08</sub>	84.84 <sub>3.73</sub>	84.54 <sub>3.21</sub>	94.48 <sub>2.75</sub>	85.69 <sub>3.75</sub>	86.92 <sub>6.27</sub>	<b>94.64</b> <sub>0.87</sub>
Nursery	88.53 <sub>1.33</sub>	88.82 <sub>2.32</sub>	92.12 <sub>2.23</sub>	88.15 <sub>1.18</sub>	87.42 <sub>2.24</sub>	<b>96.65</b> <sub>1.20</sub>	93.41 <sub>0.79</sub>	93.54 <sub>1.29</sub>	95.63 <sub>1.18</sub>	93.23 <sub>0.70</sub>	92.87 <sub>1.24</sub>	<b>97.16</b> <sub>0.80</sub>
Splice	95.54 <sub>1.27</sub>	95.36 <sub>1.21</sub>	95.36 <sub>1.29</sub>	93.31 <sub>2.51</sub>	95.45 <sub>1.33</sub>	<b>95.99</b> <sub>1.01</sub>	95.67 <sub>0.96</sub>	95.47 <sub>0.73</sub>	95.43 <sub>0.92</sub>	93.37 <sub>2.00</sub>	95.69 <sub>0.85</sub>	<b>96.09</b> <sub>0.55</sub>
PreND	84.98 <sub>0.00</sub>	90.84 <sub>3.25</sub>	95.00 <sub>0.34</sub>	95.12 <sub>0.47</sub>	93.74 <sub>1.99</sub>	<b>95.15</b> <sub>0.55</sub>	89.68 <sub>0.00</sub>	93.52 <sub>2.53</sub>	96.42 <sub>0.23</sub>	96.63 <sub>0.42</sub>	95.43 <sub>0.78</sub>	<b>96.70</b> <sub>0.48</sub>
Average	82.28	83.14	86.63	83.64	83.76	<b>87.76</b>	85.87	86.12	88.53	86.15	86.89	<b>89.11</b>
Mean rank	3.29	3.57	2.43	3.86	3.14	<b>1.71</b>	4.00	4.29	2.86	4.71	3.43	<b>1.71</b>

The best result is in bold face, and the second best result is in italic.

# 课程提纲

- 一、何谓大数据时代
- 二、大数据机器学习
- 三、典型应用举例
- 四、开源软件平台



# AutoML开源软件平台

- 近年来，机器学习在各个领域都取得了重大突破。金融服务、医疗保健、零售、交通等领域一直在以某种方式使用机器学习系统，而且取得了很好的效果
- 传统的 ML 流程仍依赖于人力，但并非所有企业都有资源来投资经验丰富的数据科学团队，AutoML 可能正是这种困境的一个答案

- **自动机器学习**（AutoML）是将机器学习应用于现实问题的端到端流程自动化的过程
- AutoML 对于没有该领域专业知识的人使机器学习成为可能
- AutoML 通过使不同背景的人能够演进机器学习模型来解决复杂的场景，正在从根本上改变基于 ML 的解决方案给人们的印象



- Mercari 是一款在日本很受欢迎的购物应用程序，它一直使用 AutoML Vision（谷歌的 AutoML 解决方案）对图像进行分类
- Mercari 一直在“开发自己的 ML 模型，在照片上传的用户界面上推荐 12 个主要品牌的品牌名称”

- 在 TensorFlow 上训练的模型达到了 75% 的精度，AutoML Vision 的高级模式拥有 50,000 张训练图像，所以精度更是高达 91.3%，提升了竟有 15%
- 基于如此惊人的结果，Mercari 已经将 AutoML 集成到他们的系统中

## AutoML 优点

- 通过自动执行的重复性任务来 **提高工作效率**。这使得数据科学家能够更多地关注问题而不是模型
- 自动化 ML 管道还有助于 **避免** 可能因手动引入的**错误**
- AutoML 是向 **机器学习民主化** 迈出的一步，它使所有人都能使用 ML 的功能

感谢同学们听课  
欢迎讨论与交流