

# 在线教育领域的机器学习应用

邓澍军

dengsj@yuantiku.com

2015.04.25



# 提纲

概述

小猿搜题  
之  
拍照搜题

猿题库  
之  
能力预测

猿辅导  
之  
老师推荐

总结



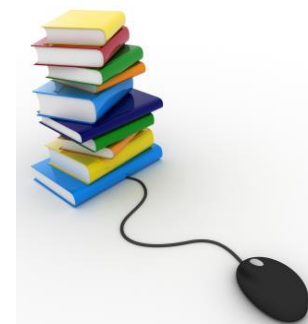
# 机器学习

**10年代**  
Deep Learning

**90-00年代**  
SVM  
Boosting  
随机森林

**70-80年代**  
神经网络

**50-60年代**  
感知机



# 在线教育

**2012-**

移动互联网  
在线教育

**10年代**

互联网公司  
进军在线教育

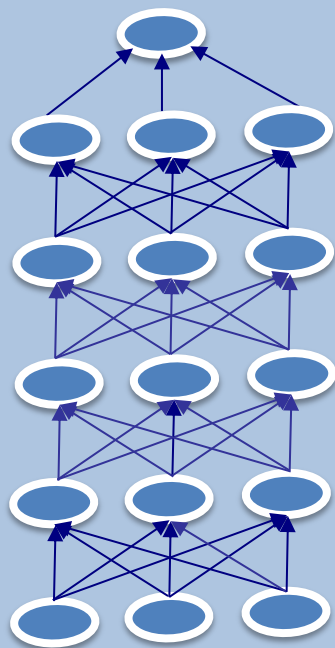
**00年代**

传统教育转  
战线上

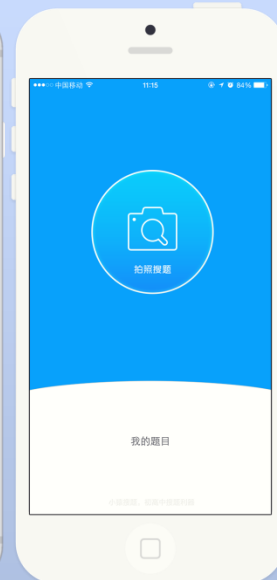
**90年代**

网校：远程  
教育

# 机器学习邂逅在线教育



机器学习



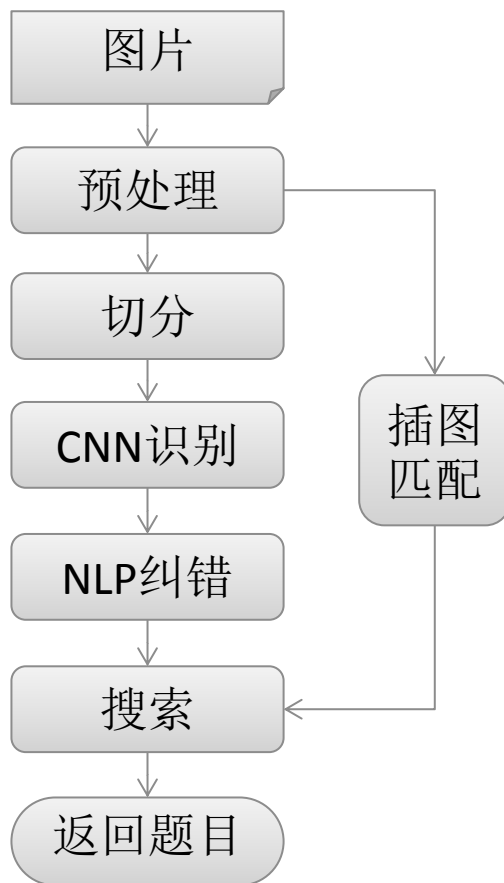
在线教育

应用之一：小猿搜题之拍照搜题

# 小猿搜题之拍照搜题

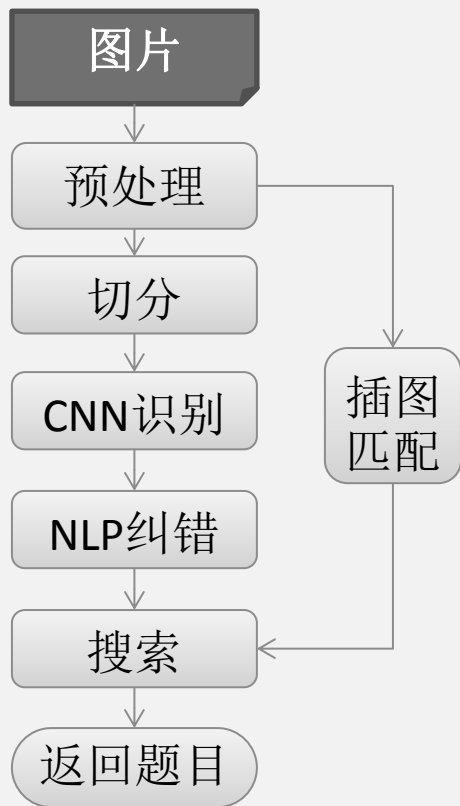


# 小猿搜题技术框架





# 图片类别



## ✓ 图片特征

### ➤ 内容多样

- 语数英等10来个科目

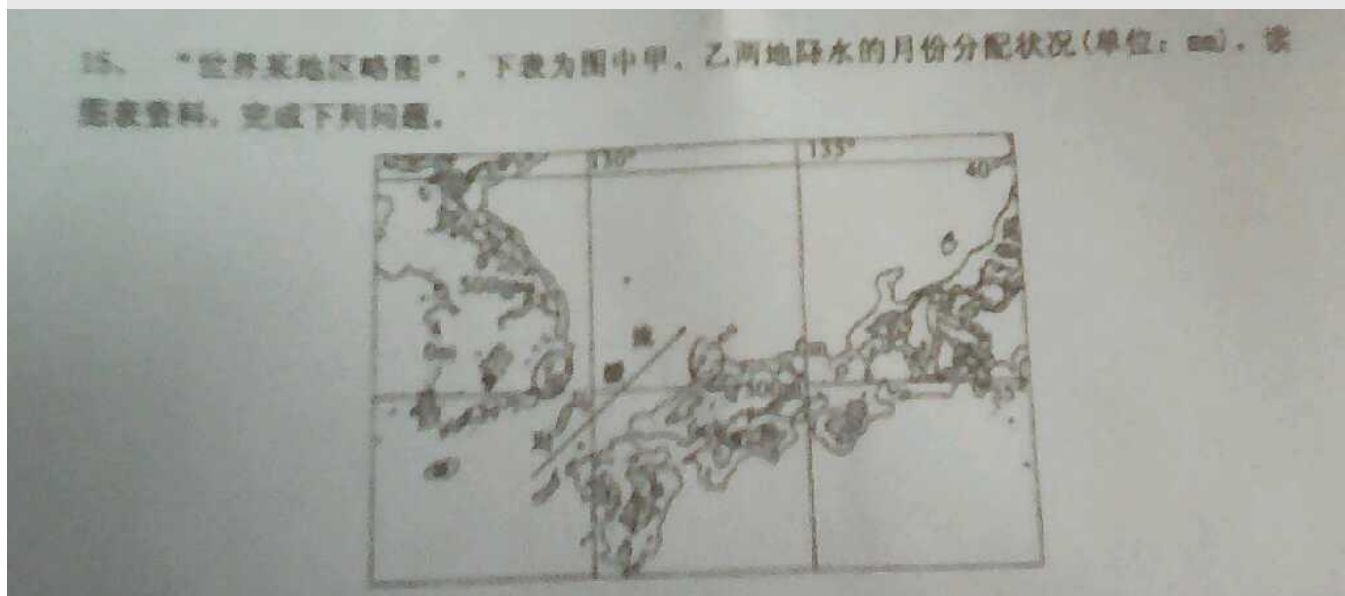
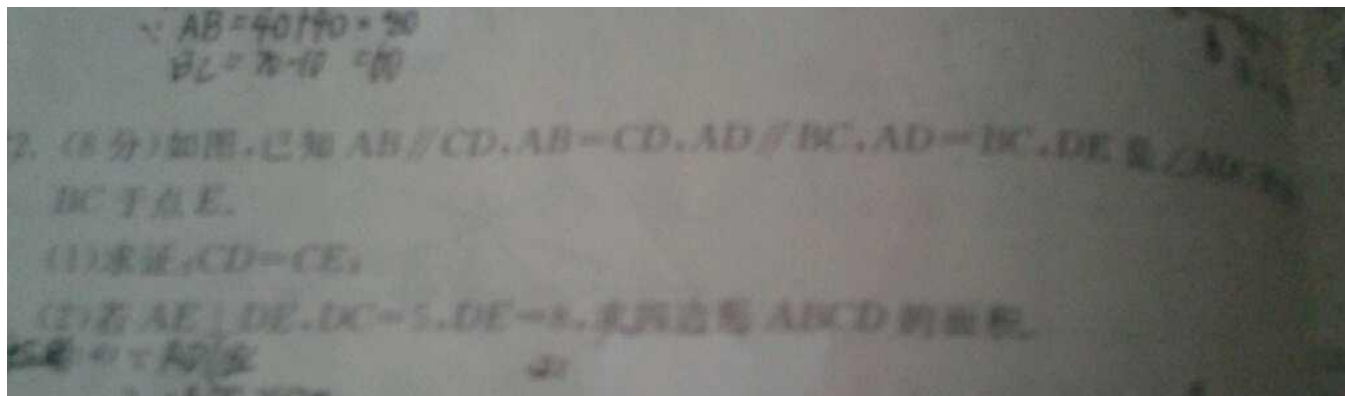
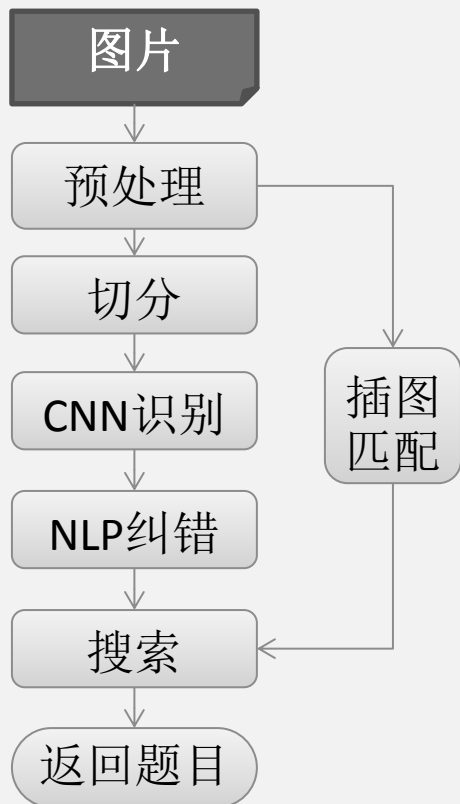
### ➤ 模糊图很多，占**30%+**

- 光照、扭曲、抖动等

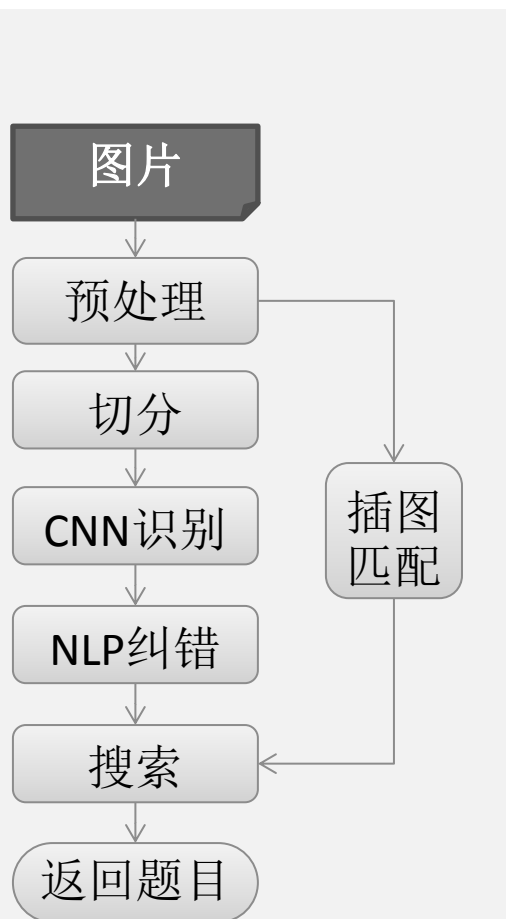
### ➤ 含有大量公式，数学占**50%**

- 上下标、分式、根号等

# 模糊图片



# 公式图片



2. 已知  $|a|$   $|b|$

3. 已知  $x+y=4, xy=1$ , 则  $\frac{x^2 y - xy^2}{x^2 - y^2}$

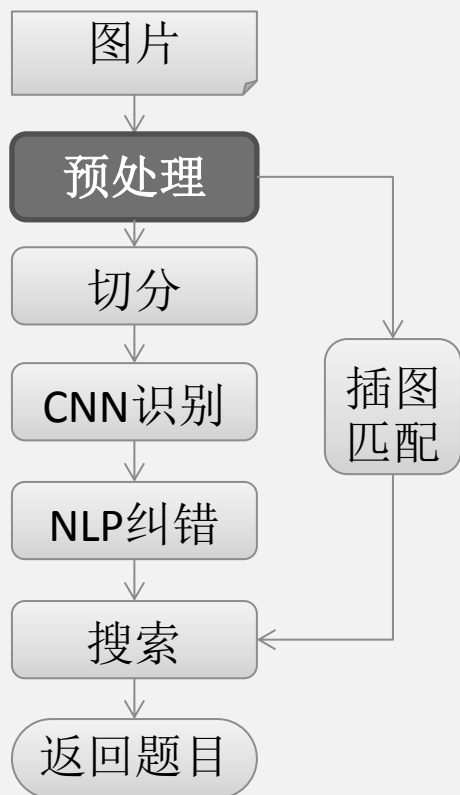
[例 1] 已知  $\alpha$  是三角形的内角, 且  $\sin \alpha + \cos \alpha = \frac{1}{5}$ .

(1) 求  $\tan \alpha$  的值;

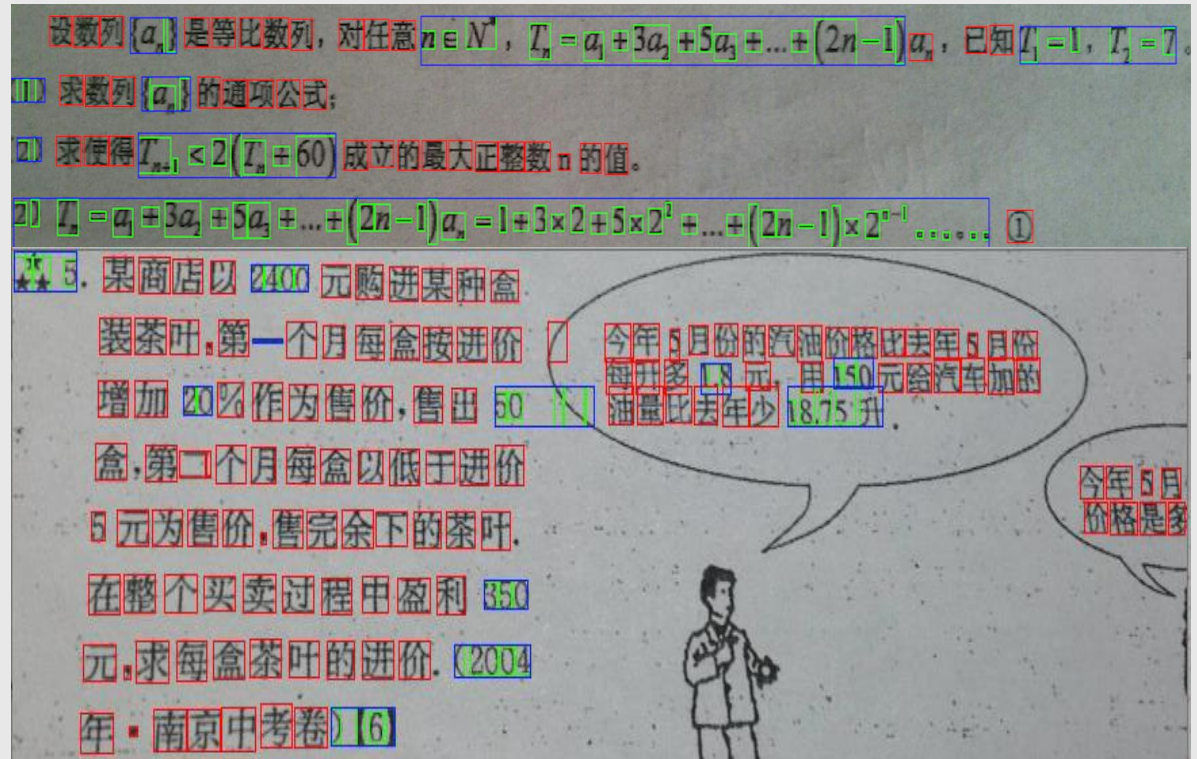
(2) 把  $\frac{1}{\cos^2 \alpha - \sin^2 \alpha}$  用  $\tan \alpha$  表示出来, 并求其值.

(2)  $\frac{2}{\sqrt{3}+1} - (\sqrt{3}-1) + 2\sin 60^\circ - 3\tan 30^\circ$ .

# 图片预处理



# 字符切分



# 复杂公式切分



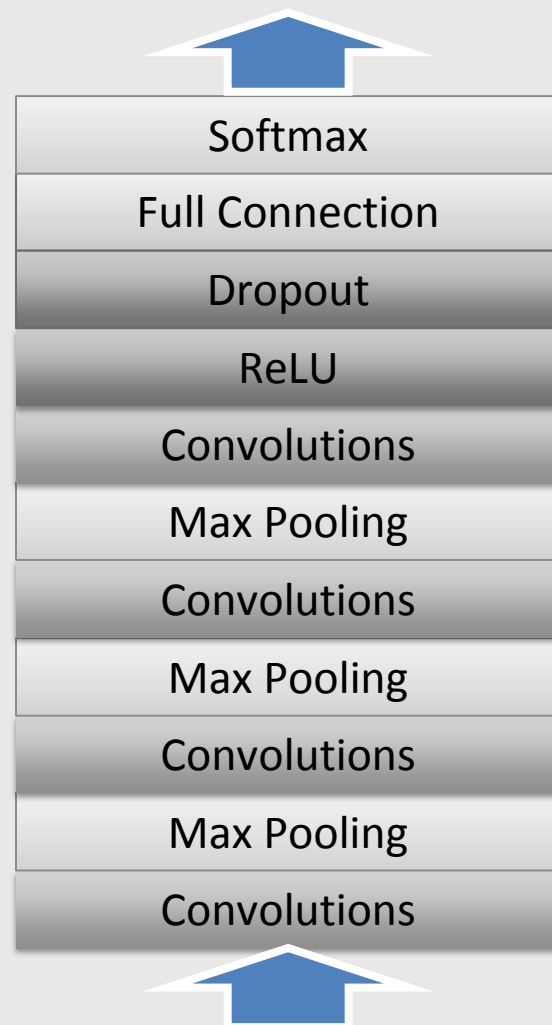
$$(2) \quad y = \sqrt{1 - \left(\frac{1}{2}\right)^x}$$

Diagram illustrating the segmentation of the formula into tokens:  $(2)$ ,  $y$ ,  $=$ ,  $\sqrt{\quad}$ ,  $1$ ,  $-$ ,  $(\frac{1}{2})$ ,  $^x$ .

Diagram illustrating the segmentation of the expression inside the square root:  $1 - \left(\frac{1}{2}\right)^x$ .

Diagram illustrating the final segmentation of the formula into individual characters and symbols:  $1$ ,  $-$ ,  $($ ,  $\frac{1}{2}$ ,  $)$ ,  $^x$ .

# 卷积神经网络模型



# 自动生成标注数据



## ✓ 标注数据自动生成

- 场景融合
- 旋转、拉伸等





# Deep Learning优化



- ✓ Deep Learning的优化算法多种多样，模型最终的效果也不尽相同
  - 小猿搜题中尝试了多种不同优化算法
  - 一般来说，先SGD再采用Gauss-Newton能够在更短时间内收敛

# Deep Learning加速



- ✓ 用GPU K40训练相比CPU模式速度能够提升**5-6**倍
- ✓ GPU K40线上预测速度能够提升**2-4**倍



# 辅助策略——NLP纠错



✓ 利用语言模型进行纠错

➤ 平行回边形 → 平行四边形

➤ 电灯炮 → 电灯泡

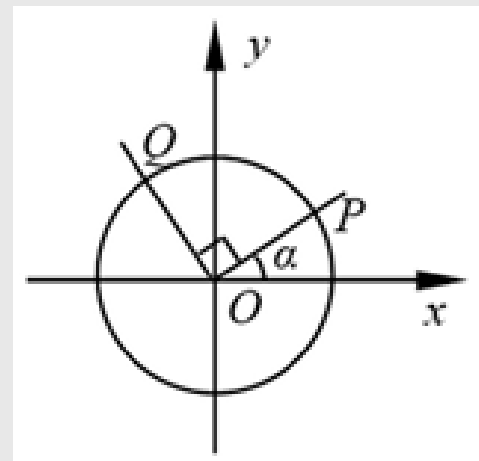
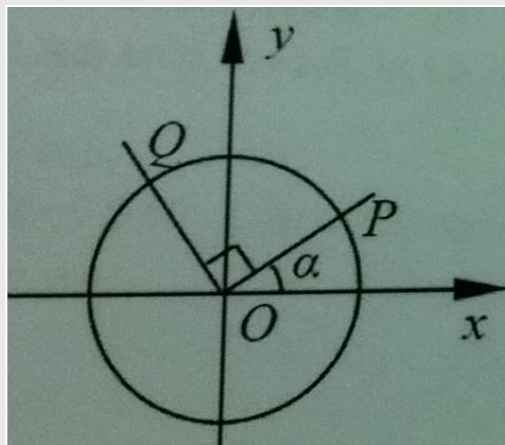
➤ 入 ↔ 人

➤ 1 ↔ 一

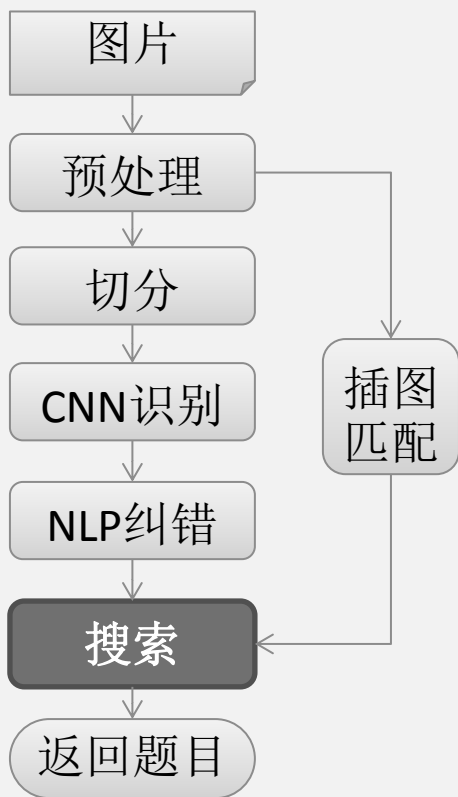
# 辅助策略——插图匹配



✓ 如果题目中的插图匹配(SIFT等特征), 则为加分项



# 搜索



## ✓ 搜索主要模块

➤ 分词

➤ 倒排索引

➤ 排序

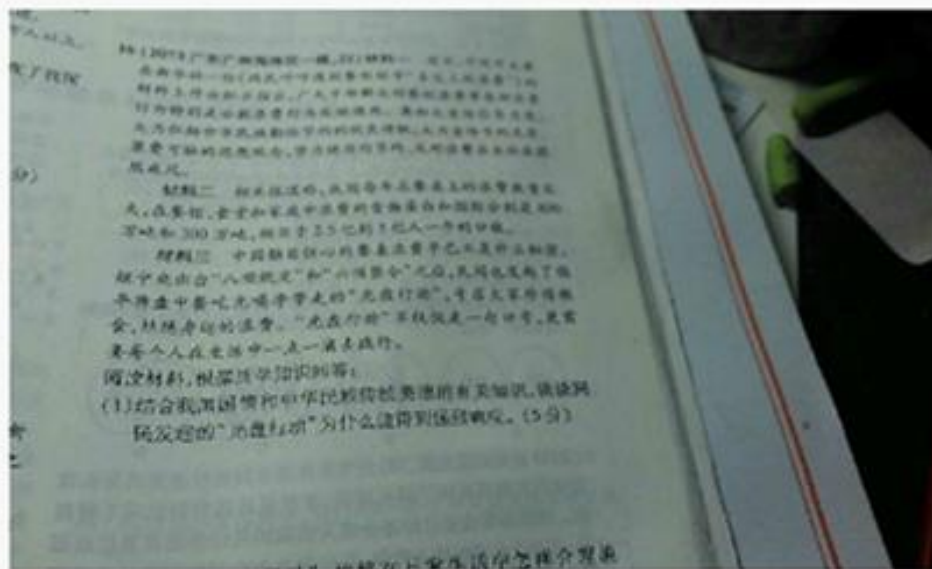
- Learning to Rank: GBRT



# 小猿搜题

初高中拍照搜题利器

## 搜索示例



【题文】材料一：近日，习近平主席在新华社一份《网民呼吁遏制餐饮环节“舌尖上的浪费”》的材料商作出批示指出，广大干部群众对餐饮浪费等各种浪费行为特别是公款浪费行为反映强烈。要加大宣传引导力度，大力弘扬勤俭节约优秀传统文化。

材料二：相关报道称，我国每年在餐桌上的浪费数量巨大。在餐桌上浪费的事务相当于2.5亿到3亿人一年的口粮。

材料三：中国触目惊心的餐桌浪费早已不是秘密。一些网民倡导将盘中餐吃光喝净带走的“光盘行动”。号召大家珍惜粮食，杜绝浪费。“光盘行动”不仅仅是一句口号，更需要去践行。

# 搜索示例

← 题目详情

1. 已知椭圆 $C_1: \frac{x^2}{4} + y^2 = 1$ , 椭圆 $C_2$ 以 $C_1$ 的长轴为短轴, 且与 $C_1$ 有相同的离心率.  
(1) 求椭圆 $C_2$ 的方程;  
(2) 设 $O$ 为坐标原点, 点 $A, B$ 分别在椭圆 $C_1$ 和 $C_2$ 上,  $\vec{OB} = 2\vec{OA}$ , 求直线 $AB$ 的方程.

2014年12月30日 22:01

🔍 搜索结果

题干

已知椭圆 $C_1: \frac{x^2}{4} + y^2 = 1$ , 椭圆 $C_2$ 以 $C_1$ 的长轴为短轴, 且与 $C_1$ 有相同的离心率.

(1) 求椭圆 $C_2$ 的方程;

(2) 设 $O$ 为坐标原点, 点 $A, B$ 分别在椭圆 $C_1$ 和 $C_2$ 上,  $\vec{OB} = 2\vec{OA}$ , 求直线 $AB$ 的方程.

# 小结

## ✓ 小猿搜题之拍照搜题

### ➤ Computer Vision

- 预处理
- 切分
- 训练数据自动生成

### ➤ Deep Learning

- 识别

### ➤ NLP

- 纠错

### ➤ Learning to Rank

- 排序



## 应用之二：猿题库学生能力预测

# 猿题库学生能力预测



# 猿题库学生能力预测(续)



# 传统教育模型

## ✓ 项目反应理论(IRT)

### ➤ 最简单的IRT模型

$$p_{irt}(y = 1|\theta, b) = \frac{1}{1 + e^{-(\theta - b)}}$$

### ➤ 题目难度 $b$

- 标注

### ➤ 学生能力 $\theta$

- 模型参数，优化得到

# 机器学习模型

## ✓ 机器学习模型

$$p(y = 1|w, x) = \frac{1}{1 + e^{-w'x}}$$

### ➤ Offline model

- Logistic Regression

### ➤ Online model

- Follow-the-Regularized-Leader

# 特征

## ✓ 所用特征

### ➤ 用户相关特征

- 学校，地区，目标考试， .....

### ➤ 题目相关特征

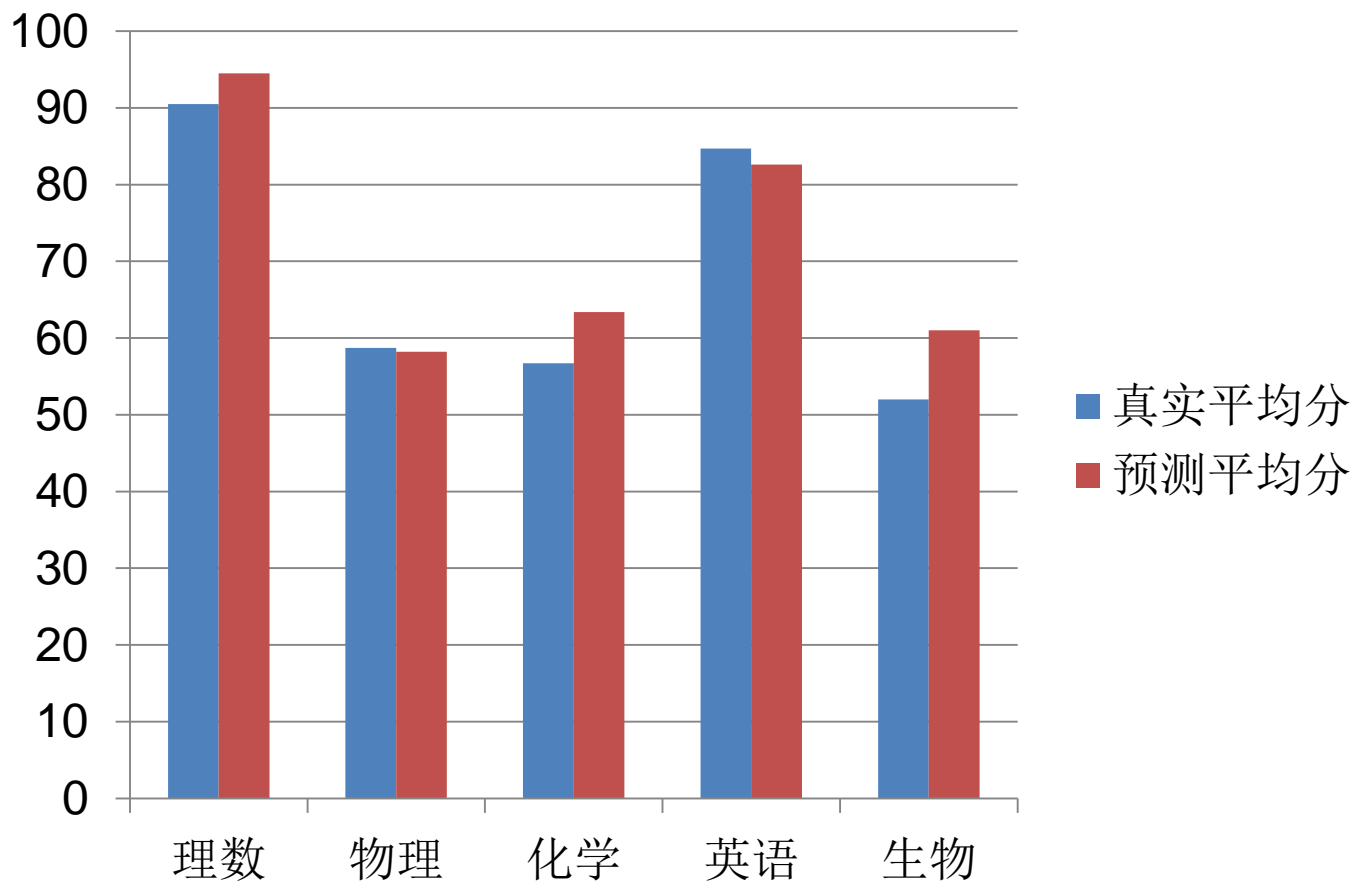
- 知识点，关键词，难度， .....

### ➤ 时序特征

- 距离高考时间， .....

### ➤ 组合特征

# 预测分评估



2014年广东省理科高考

# 小结

- ✓ 传统教育模型

- 项目反应理论 (Item Response Theory)

- ✓ 机器学习模型

- Offline model

- LR (Logistic Regression)

- Online model

- FTRL (Follow-The-Regularized-Leader)



## 应用之三：猿辅导老师推荐

# 猿辅导老师推荐

## 无练习 不辅导



约课



课前练习



在线上课



听课报告

 筛选老师

默认排序 

	<div></div> <p>数学 - 七/八/九年级 清华大学 3年教龄</p>	<p>150元/课时</p> <p>共授7课时 好评率100%</p>
	<div></div> <p>数学 - 七/八/九年级 清华大学 6年教龄 学科带头人</p>	<p>150元/课时</p> <p>共授8课时 好评率100%</p>
	<div></div> <p>数学 - 七/八/九年级 清华大学 6年教龄 竞赛辅导员</p>	<p>150元/课时</p> <p>共授10课时 好评率100%</p>
	<div></div> <p>数学 - 七/八/九年级 7年教龄</p>	<p>150元/课时</p> <p>共授3课时 暂无评价</p>

# 猿辅导老师推荐(续)

## ✓ 推荐系统

- 冷启动: Content-Based
- Item-Based Collaborative Filtering

## ✓ 机器学习

- Logistic Factorization Machine
- Exploitation and Exploration(E&E)

# 总结

- ✓ 小猿搜题之拍照搜题
  - Deep Learning
  - Computer Vision
  - Learning to Rank
- ✓ 猿题库学生能力预测
  - 传统教育领域的项目反应理论（IRT）
  - 计算广告点击率预测模型LR、FTRL
- ✓ 猿辅导老师推荐
  - 推荐系统
  - LFM, E&E

# 未来

- ✓ 教育领域知识图谱
  - 学生的最优能力成长之路
- ✓ 手写识别
  - 手写拍照搜题
  - 解答题
  - 自动判卷
- ✓ 高考机器人
  - 机器自动出题
  - 机器自动做题
- ✓ 智能芯片
- ✓ .....

Q&A?

Thanks!

[dengsj@yuantiku.com](mailto:dengsj@yuantiku.com)

