

OTA场景中的两类对话机器人及其实现方案

2020-05-22

分享人:王创峰

大纲

- OTA场景的两类机器人需求:售前与售后
- 两类场景的性质分析
- 信息引导类建模
- 解决问题类建模
- 总结

两类场景的区别和性质

	售前	售后
类别	信息引导类	解决问题类
诉求	替代翻页,搜索	替代客服
case	我说话,你给我我要的信息	我碰到问题了,你帮我解决,或者告诉我怎么解决.
性质	新的用户界面	解放人力
思考	2C: 界面是命脉,交出去要小心 2B: 参考2C的担忧,是问题也是机会	2C: 人机可替代,协同工作 2B: 壁垒不足,甲方分分钟自己搞

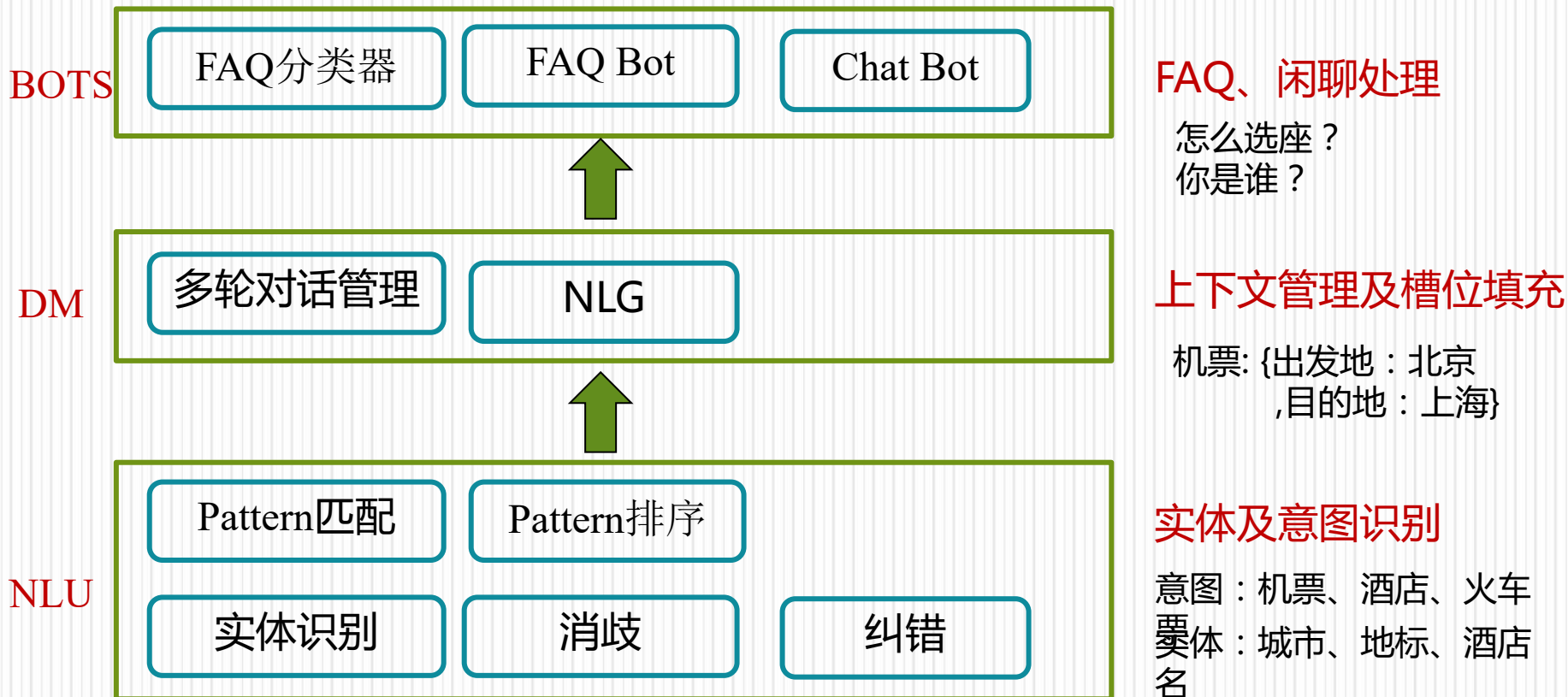
售前,引导类

- 输入: 我要啥啥啥
- 输出: 查询API和参数
- 解决方案: 填槽-状态机
- trick:
 - 纠错
 - 词典
 - 场景间槽位继承
 - 根据业务侧重某些意图
 - 闲聊后置
 - 限制能力
 - ...

QMI遗容



QMI架构



难题

- 危险的闲聊
 - 无法区分被判定为闲聊的问题是否真的涉及业务
 - 先过业务逻辑
 - 闲聊越简单越好
- 长句子
 - 包含多个成分
- 语音杂音录入

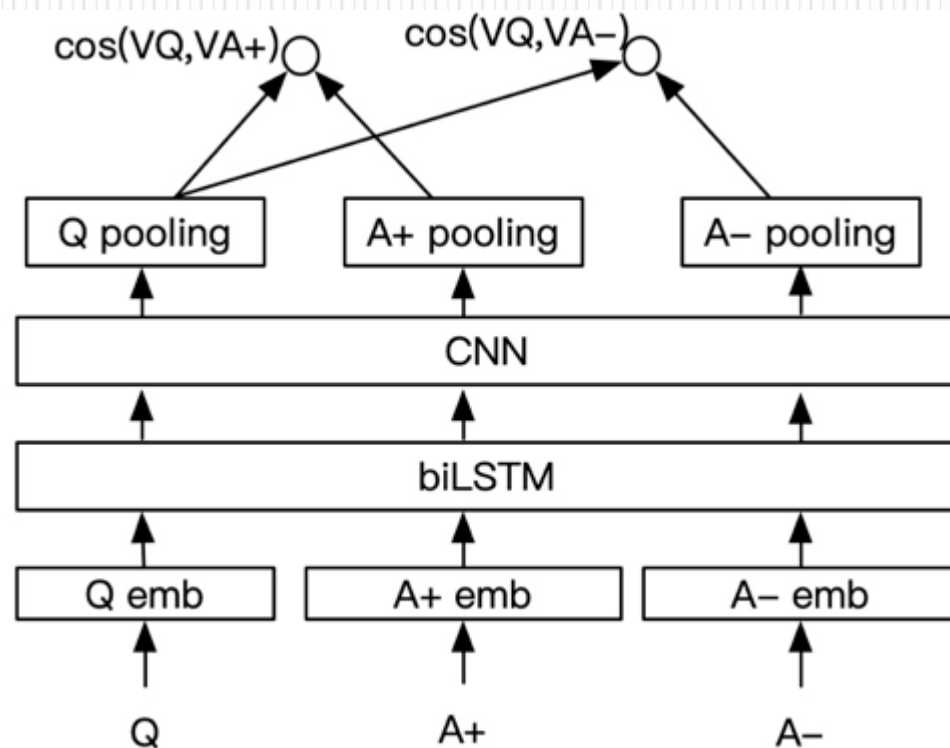
售后,解决问题类

- 特点: 有了标准解决方案的才适合这个场景
- 输入: 我的啥啥怎么回事
- 输出: 标准问题
- 解决方案: 语义相似性匹配
 - DSSM
 - BERT + triplet
- trick:
 - 区分场景,压缩备选空间

售后客服机器人(第四代)

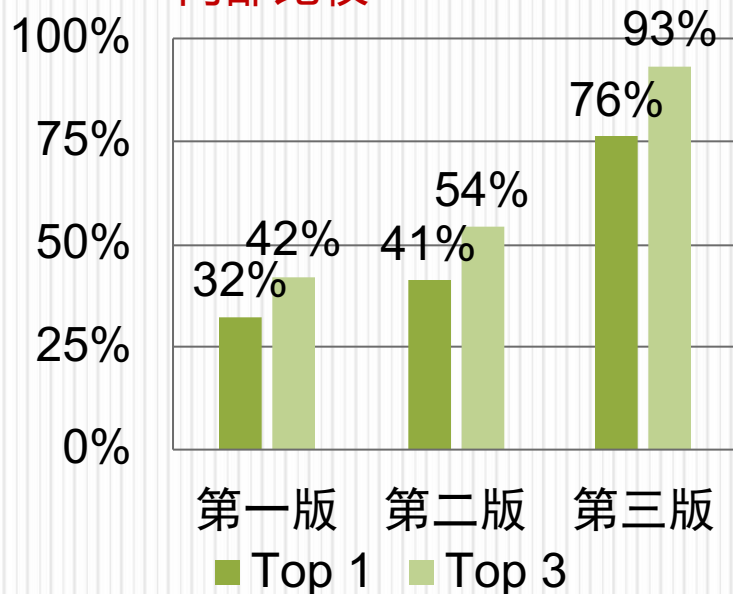


售后机器人模型结构(第三代)

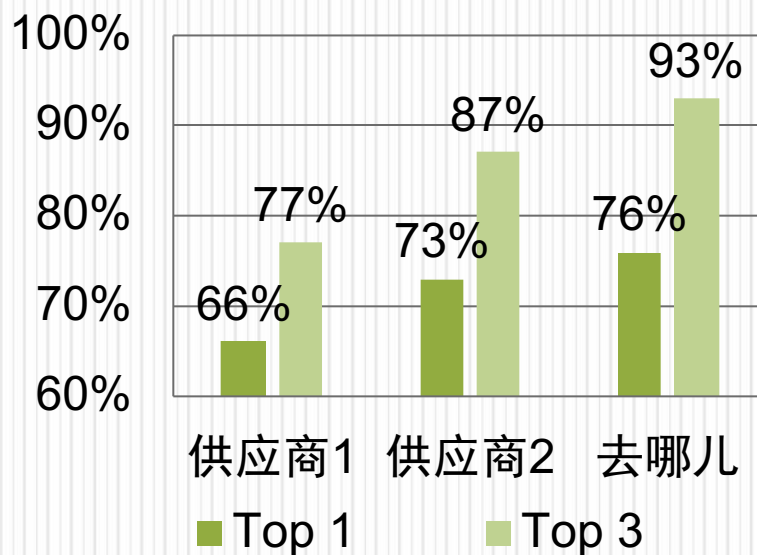


FAQ Bot：售后机器人准确率

内部比较



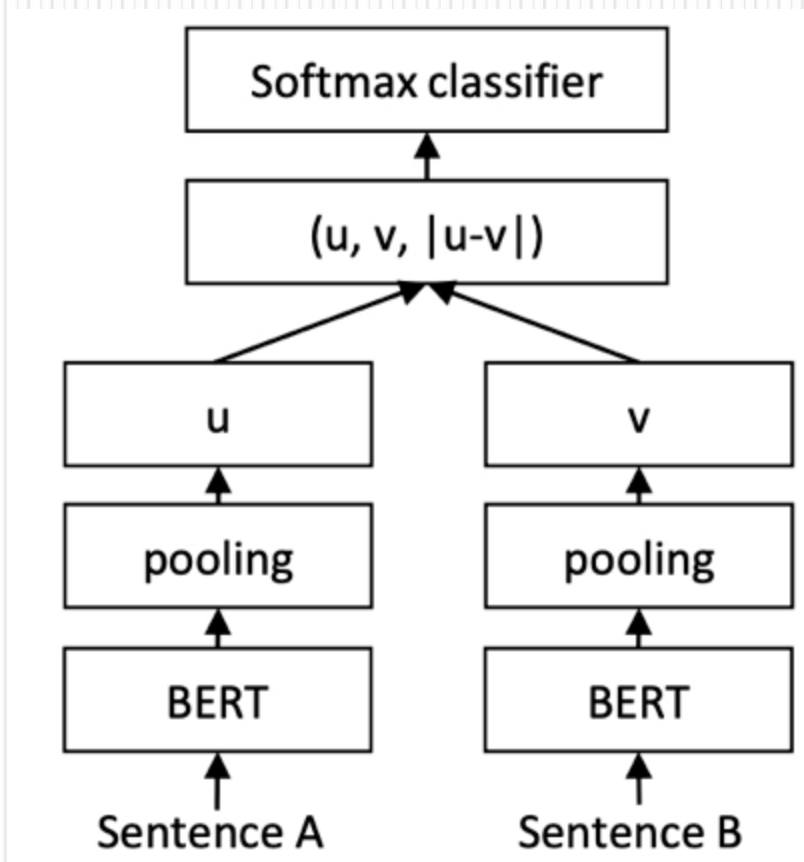
第三方比较



不同版本准确率：第三版比第一版Top 1准确率高出 **44%**

不同供应商准确率：去哪儿自研机器人Top 1准确率比供应商1高**10%** 比供应商2高**3%**

售后机器人模型结构(第四代)



准确率比较

	P@1	P@3	P@5
base	0.6046	0.8207	0.8692
DSSM	0.5632	0.7913	0.8641
ESIM	-	0.3231	0.5095
BERT	0.3619	0.6147	0.7435
base + BERT	0.7626	0.9332	0.9723

推理性能比较

	batch@1	batch@10	all
BERT(8 核 cpu)	0.220 s	1.303 s	
4-Layer+BERT(8 核 cpu)	0.1271 s	0.3528 s	
ALBERT(4 层)	0.030551	0.203293	
ALBERT(12 层)	0.256034	1.600066	
cuBERT(12 层)	0.207686	0.915885	
(base)BiLSTM+CNN(cpu)			0.012
CNN + BERT (gpu)	0.125090	0.166560	
CNN + BERT (cpu)	0.219316	0.854623	
CNN+ALBERT(4,max_seq=64) cpu	0.030551	0.203293	
CNN+ ALBERT (12,max_seq=64) cpu	0.256034	1.600066	
CNN + onnx(cpu)		1.654162	
CNN + cuBERT	0.207686	0.915885	
STS-BERT(gpu)			0.016921
STS-BERT (cpu)			0.049080

总结

- 聊天机器人的两类场景
 - 新用户界面: 替代搜索和翻页
起到信息引导的作用
 - 解放人力: 替代客服
把标准动作自动化
- 信息引导类机器人
 - 复杂繁琐,但有各种trick提升效果
 - 任重道远
- 标准问题的自动问答机器人
 - 能够从各种NLP语义表示的进步中获益
 - 复杂模型的性能是个问题

致谢与部分参考资料

- 致谢
 - 前Leader: 李兆海
 - 前小组一起调参的伙伴们
- 参考:
 - Qunar技术嘉年华(2018)Qdata专场
《Qmi智能机器人》
 - 《深度学习在酒店售后智能问答场景实践》
<https://mp.weixin.qq.com/s/PPUh13SVvk3R0Tn0DgxWlQ>