

“达观杯” 文本智能处理挑战赛

TNT_000_

目录

CONTENT

- ① 团队介绍
- ② 赛题分析
- ③ 总体框架
- ④ 特征构建
- ⑤ 模型构建
- ⑥ 模型融合
- ⑦ 总结展望

团队介绍

成员介绍



杨亚涛-中山大学
InplusLab



罗志鹏-北京大学
Microsoft



肖小粤-中山大学
InplusLab



何嘉伟
湖南大学



丁晓菲
湖南大学

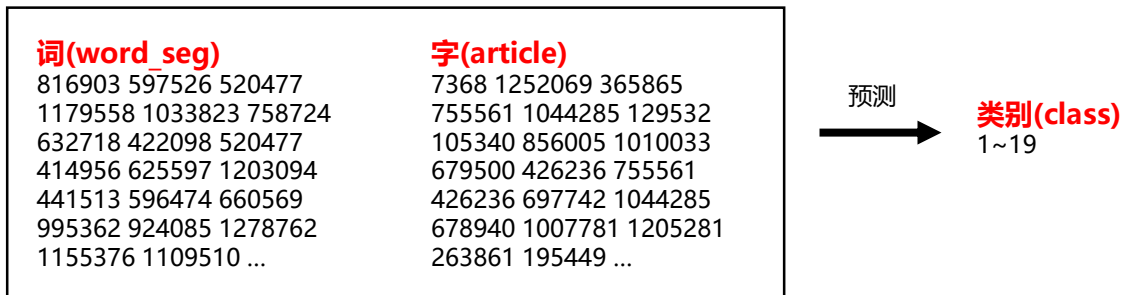
团队成绩

● 天池大数据	CIKM CUP 2018	1 st /1027
● KDD CUP 2018	Last 10-Day Prediction	1 st /4000+
● KDD CUP 2018	Second 24-Hour Prediction	1 st /4000+
● DataCastle	微博传播热度预测大赛	1 st /668
● Kesci	BOT2016人工智能聊天机器人大赛	1 st /78
● KDD CUP 2018	Main Track	2 nd /4000+
● DataFountain	2017CCF 让AI当法官	2 nd /415
● DataFountain	唯品会用户购买行为预测	2 nd /542
● Kesci	拍拍贷魔镜算法大赛	2 nd /485
● DataCastle	微额贷款人品预测大赛	3 rd /1666
● DataFountain	2016CCF 搜狗用户画像	3 rd /894
● DataFountain	2016CCF 客户用电异常行为分析	4 th /888
● DataFountain	2017CCF 360人机大战	5 th /888
● KDD CUP 2017	Tasks1	6 th /3582
● 滴滴研究院	第一届滴滴算法大赛	8 th /7000+
● 天池大数据	阿里聚安全算法挑战赛	9 th /1124

赛题分析

赛题任务

任务：通过**长文本**的**字和词**的序列数据，预测文本类别

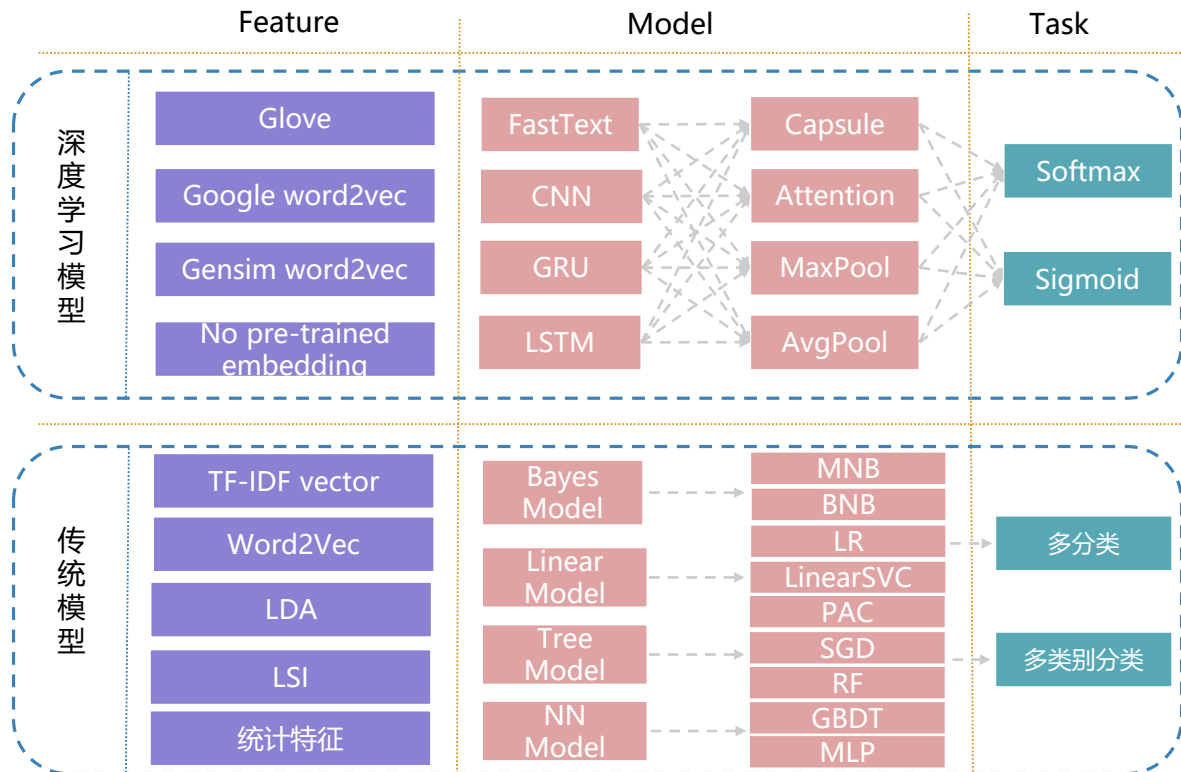


评价指标

Macro-F1:
$$\langle F1 \rangle = \frac{1}{n} \sum_i^n F1_i = \frac{1}{n} \sum_i^n \frac{2 \cdot P_i \cdot R_i}{P_i + R_i}$$

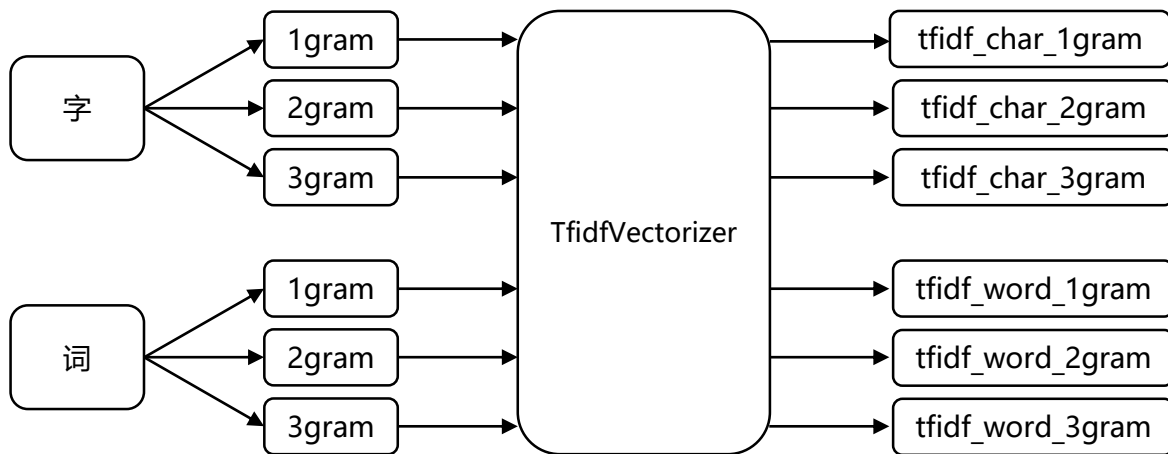
总体框架

模型架构

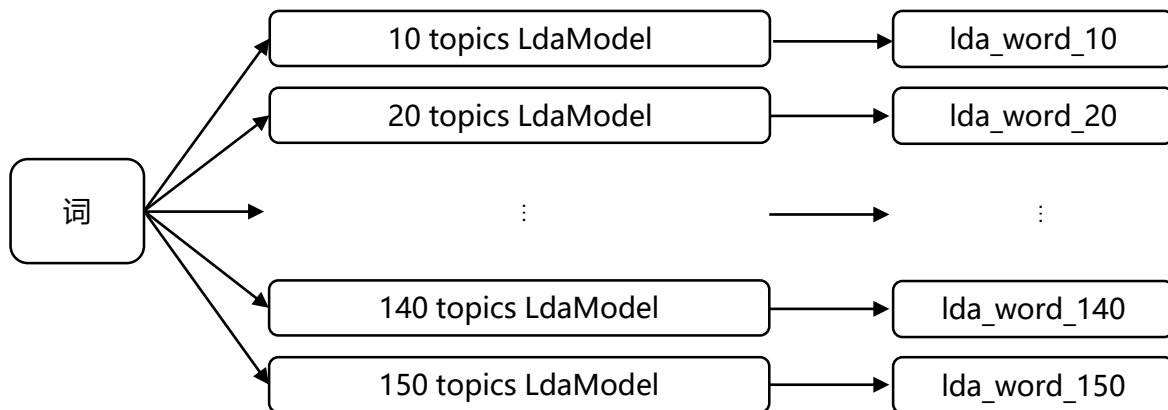


特征构建

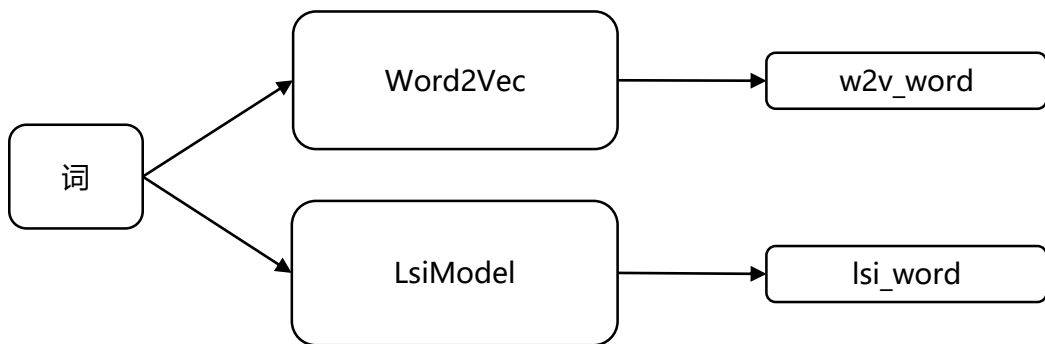
TF-IDF特征



LDA特征



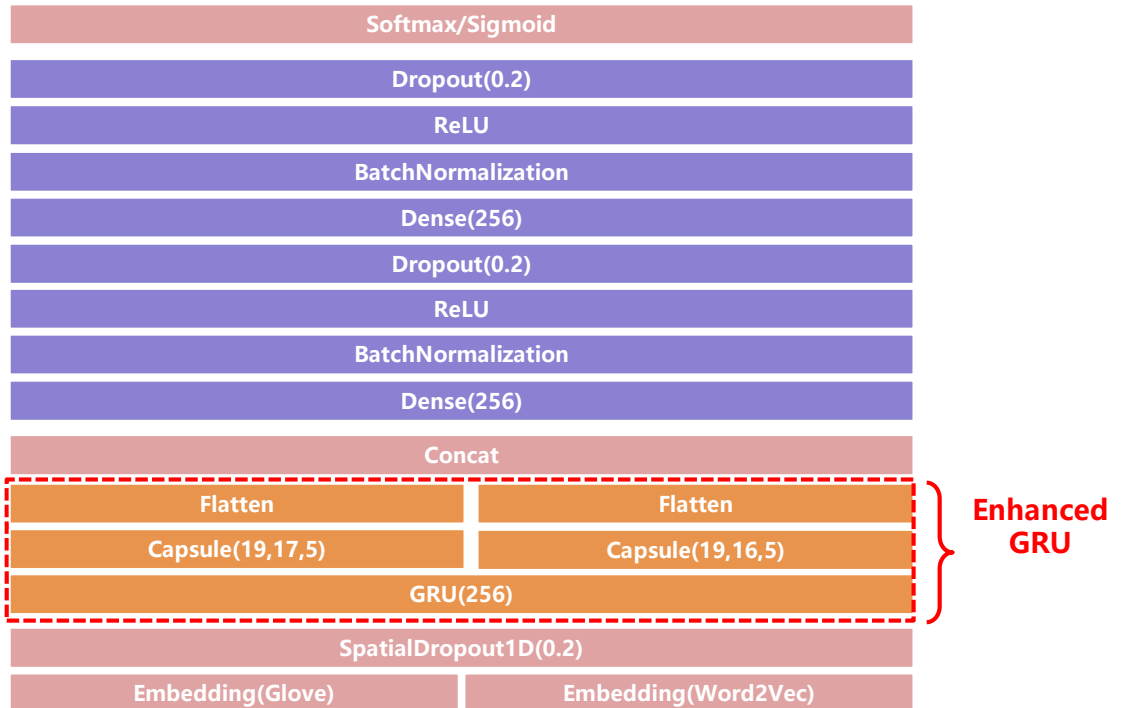
LSI特征+Word2Vec特征



模型构建

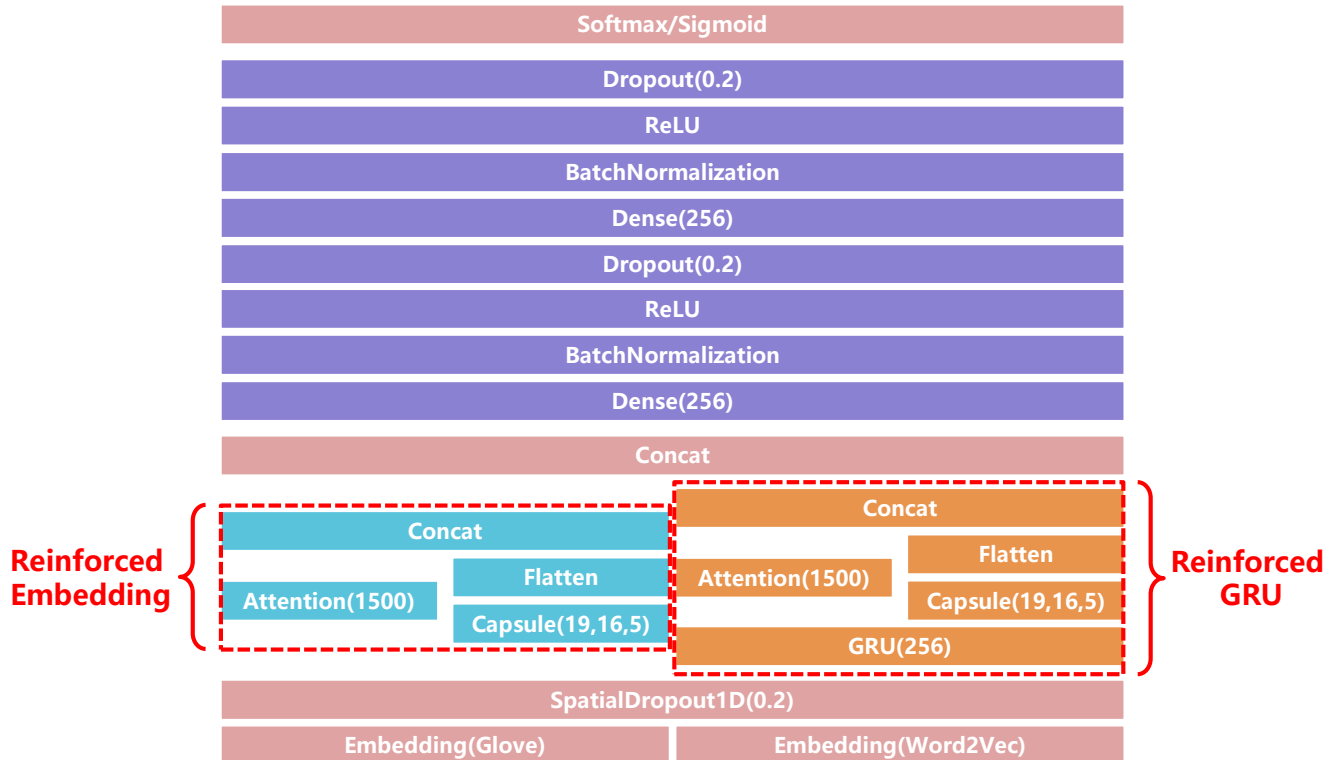
Hybrid NN-1

0.8013



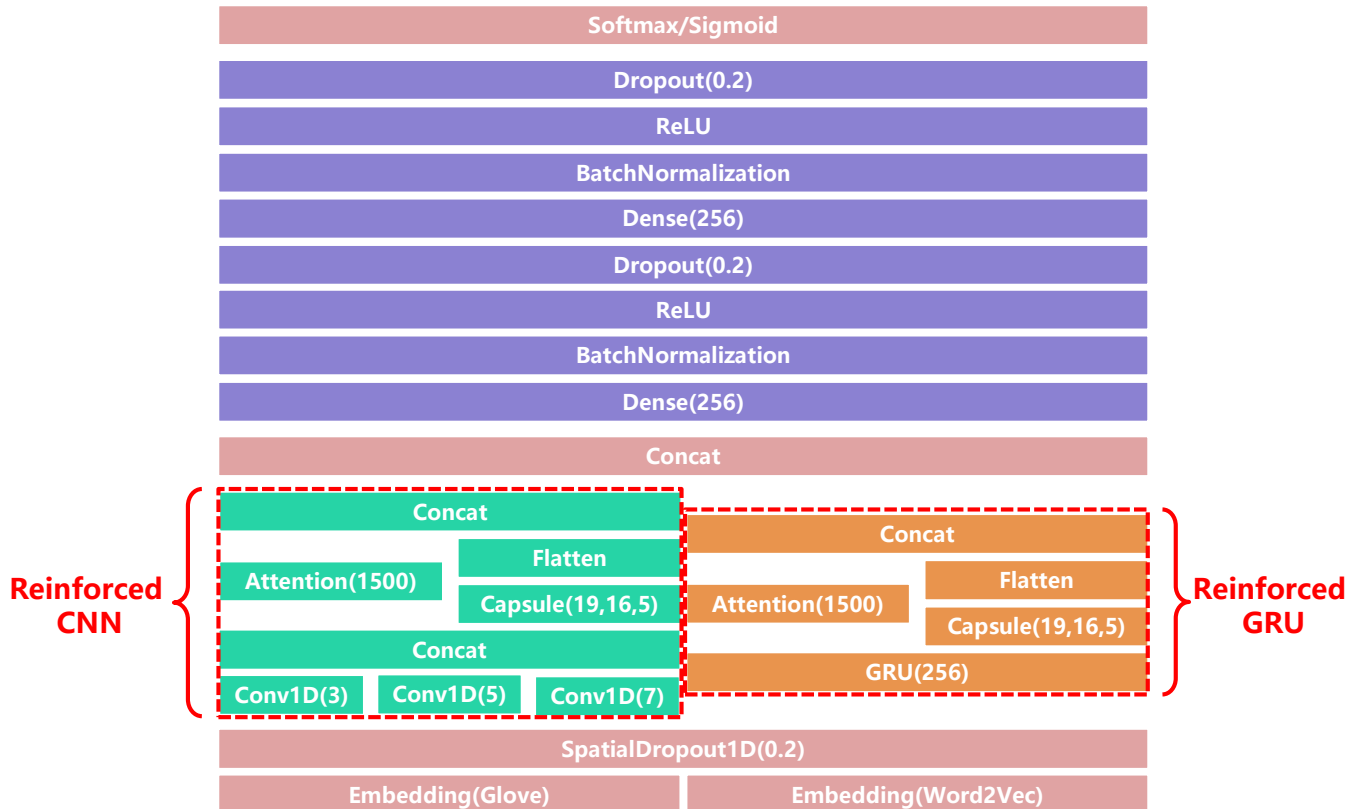
Hybrid NN-2

0.8017



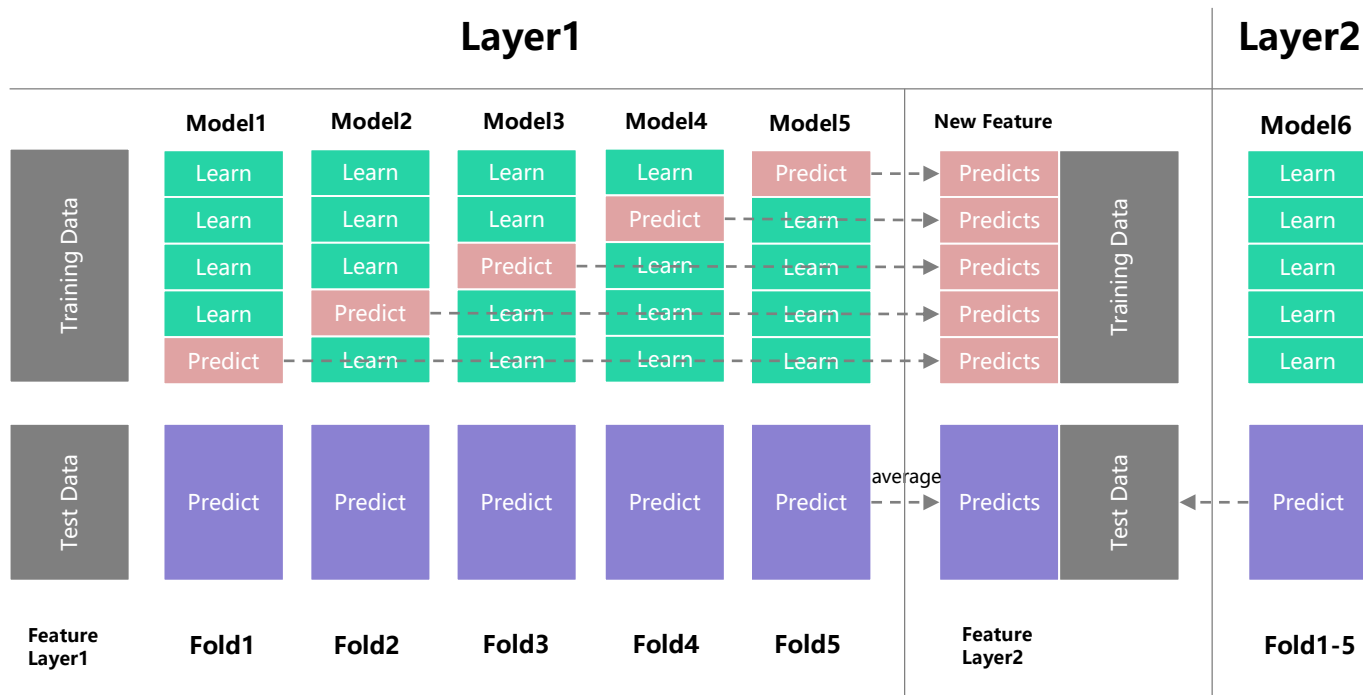
Hybrid NN-3

0.8019



模型融合

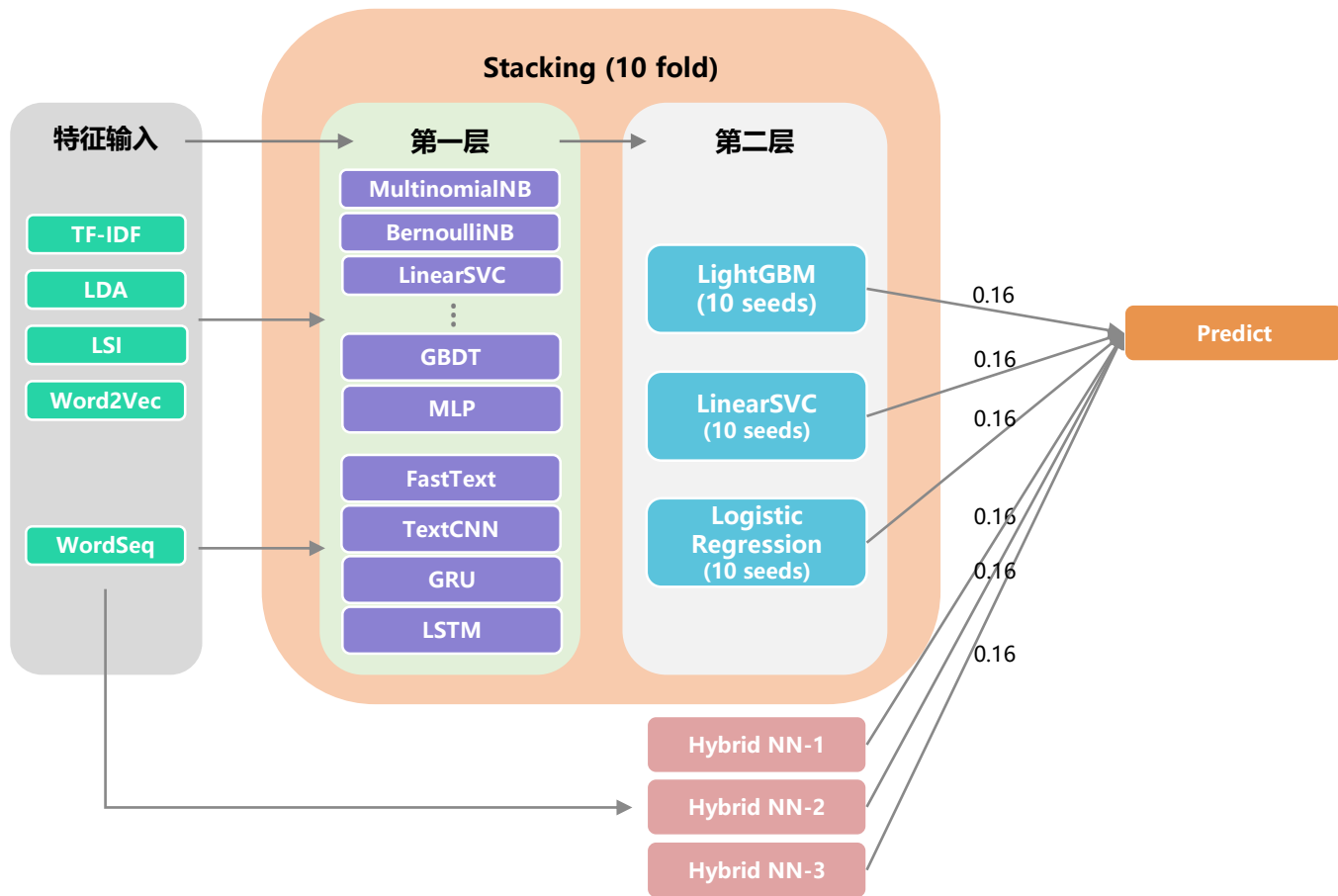
Stacking框架



10 folds

10 seeds avg

融合策略



总结展望

总结

- 网络结构的创新和改进对本赛题的效果是明显的
- 预训练的Embedding能加快网络的训练，并且效果俱佳
- 传统模型对于融合的提升是巨大的

展望

- 尝试其他的网络结构，例如：DPCNN，HAN等
- 混合网络的参数调优
- 融合系数的优化

Thanks!

TNT_000_