



全球

World Of Tech 2017

2017年12月1日-2日 • 深圳中洲万豪酒店

软件开发技术峰会

DEVELOPMENT



深度学习Deep CTR 套件

胡南炜

微博 架构师

 目录

- 01 概述
- 02 训练与分析
- 03 导出与预测
- 04 总结

01

概述

- ◆ 2016-06 Google, Play 应用推荐引入 Wide & Deep , 3.9%
- ◆ 2016-09 Google, YouTube 视频推荐引入 DNN 模型
- ◆ 2017-05 Twitter, 时间序引入深度学习 , 互动率和时长显著 增长
- ◆ 2017-07 美团 , 搜索推荐引入深度学习 , 3%AUC

背景

- ◆ 标准化 更省心
- ◆ 工具化 更简单
- ◆ 服务化 更快速

◆ 功能

- ◆ 输入：csv格式，数据清洗等
- ◆ 特征：连续、离散、文本、标签...
- ◆ 模型：DNN、Wide&Deep
- ◆ 预测：开箱即用
- ◆ 分布式计算

02 训练与分析

配置

训练

分析

02 训练与分析

配置

训练

分析

配置

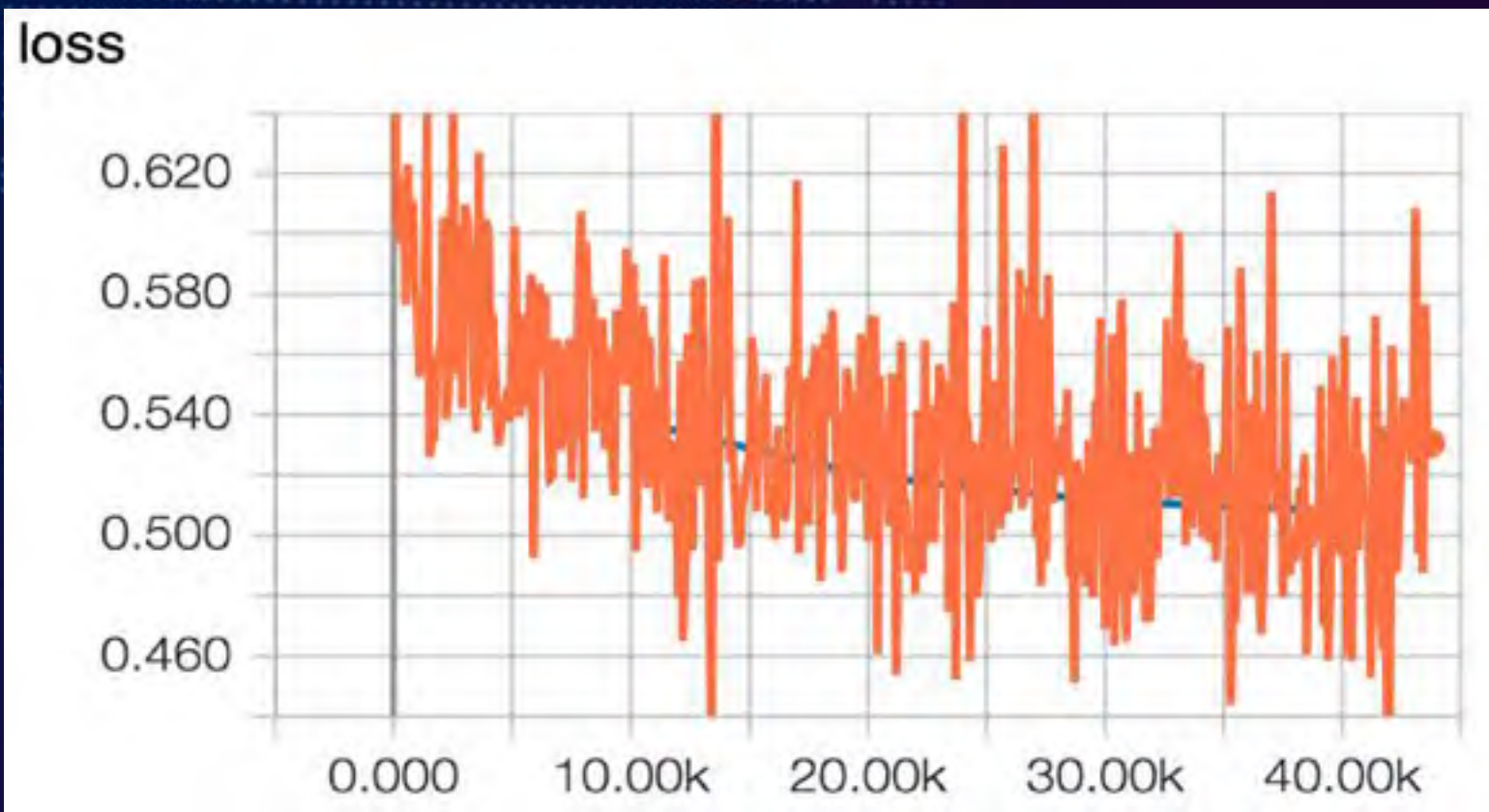
```
1 <config version="0.0.3" model="prometheus">
2   <inputs>
3     <train_data>./data/train_ori.csv</train_data>
4     <test_data>./data/test_ori.csv</test_data>
5     <batch_size>256</batch_size>
6     <header>True</header>
7     <label_index>3</label_index>
8   </inputs>
9   <columns>
10    <column name="index" id="0" cls="continuous" dtype="float" unused="True"></column>
11    <column name="uid" id="1" cls="category" dtype="int">
12      <num_buckets>6841</num_buckets>
13      <dimension>32</dimension>
14    </column>
15    <column name="mid" id="2" cls="category" dtype="int">
16      <num_buckets>3953</num_buckets>
17      <dimension>32</dimension>
18    </column>
19    <column name="score" id="3" cls="continuous" dtype="int"></column>
20    <column name="time" id="4" cls="continuous" dtype="int" unused="True"></column>
21    <column name="gender" id="5" cls="category" dtype="str">
22      <vocab_list>M, F</vocab_list>
23      <dimension>16</dimension>
24    </column>
25    <column name="age" id="6" cls="category" dtype="int">
26      <num_buckets>57</num_buckets>
27      <dimension>16</dimension>
28    </column>
29    <column name="occupation" id="7" cls="category" dtype="int">
30      <num_buckets>22</num_buckets>
31      <dimension>16</dimension>
32    </column>
33  </columns>
34 </config>
```

NORMAL TST/main_feed] config/config.xml utf-8[unix] 1K : 1/76 : 1

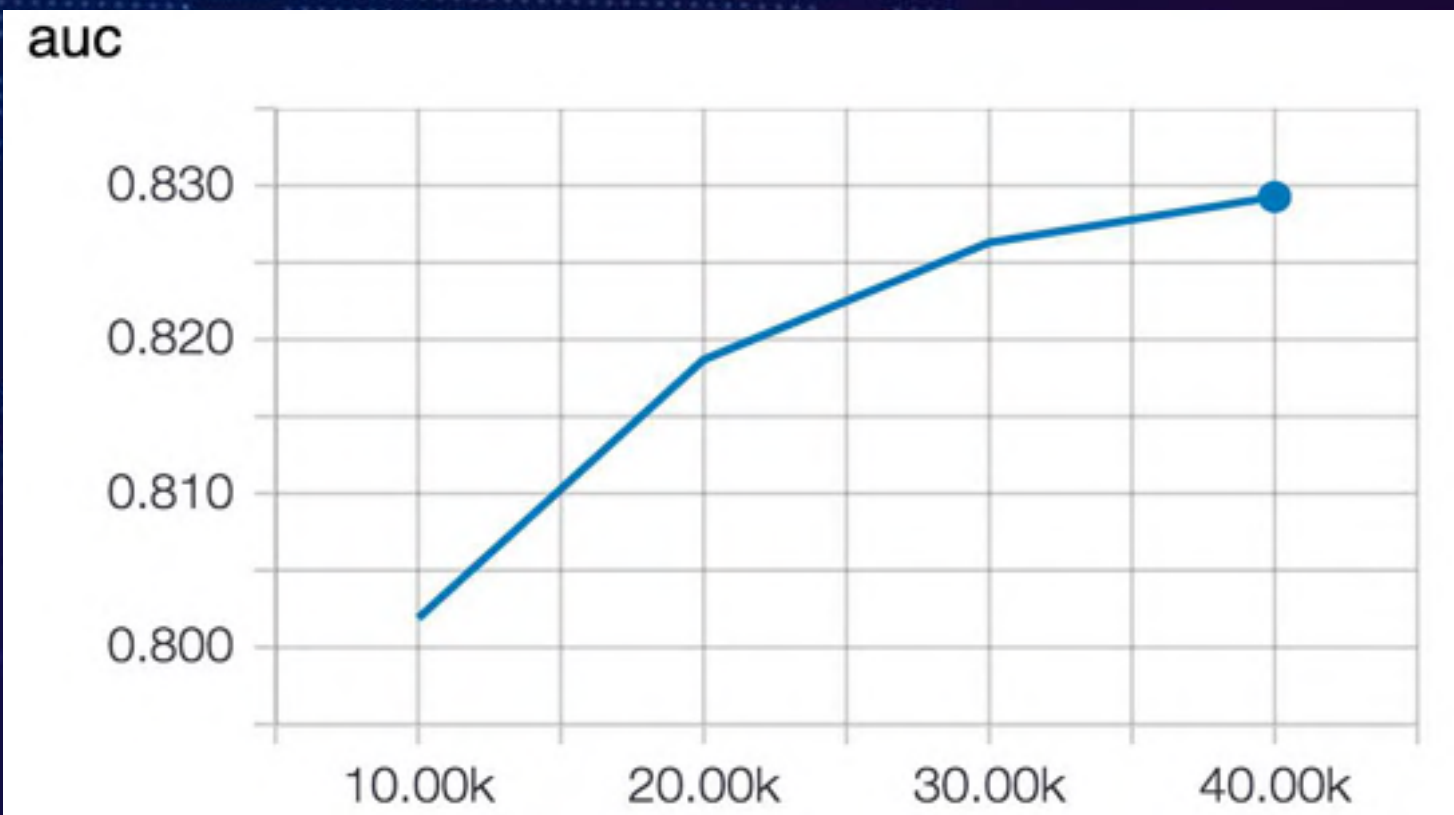
[prom] 0:/acai@77-112-153-hadoop-test:/first_features_96- 1:vim 2:vim "bj-m-203058a,local" 18:46 22-Sep-17



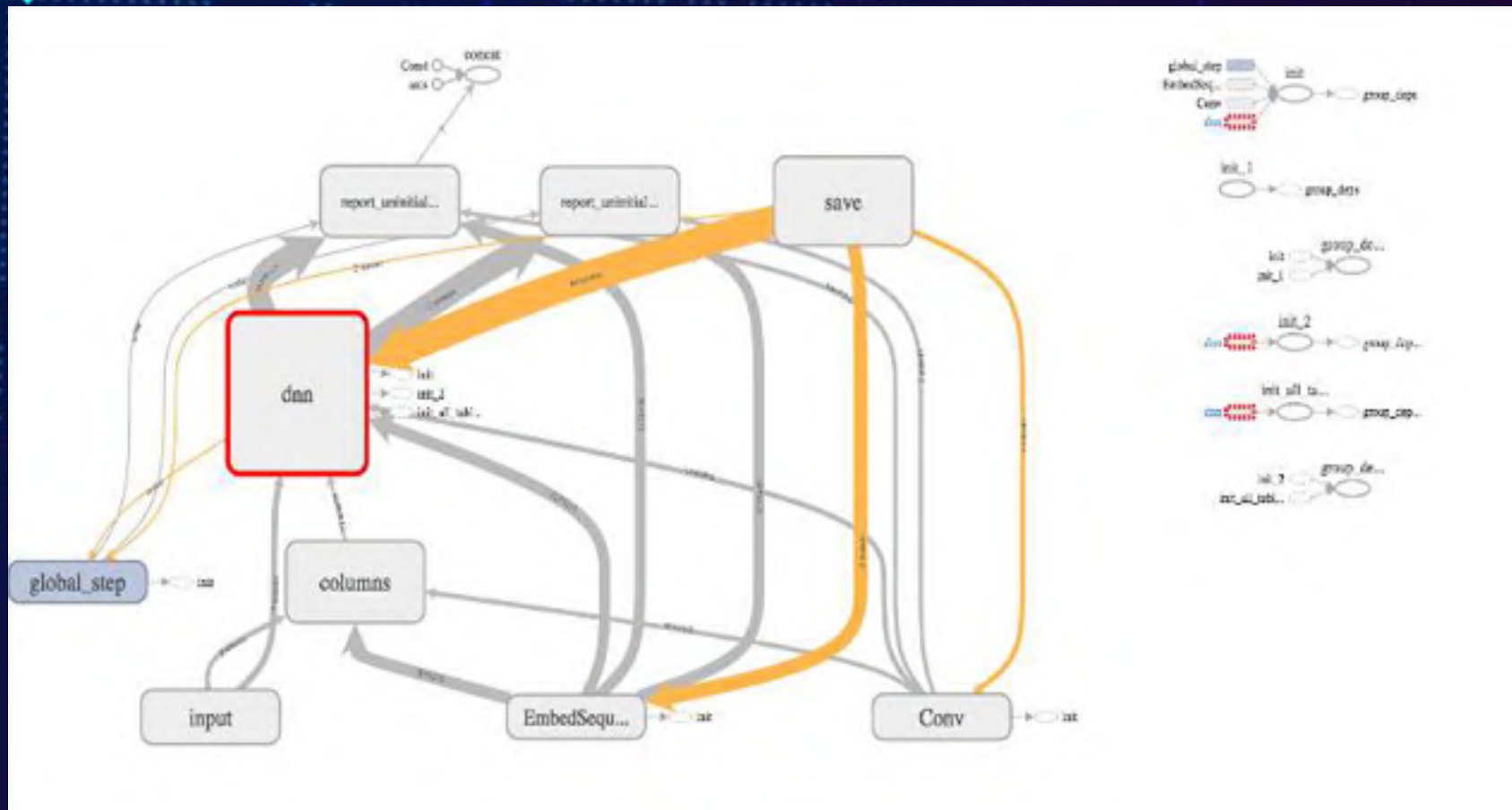
训练



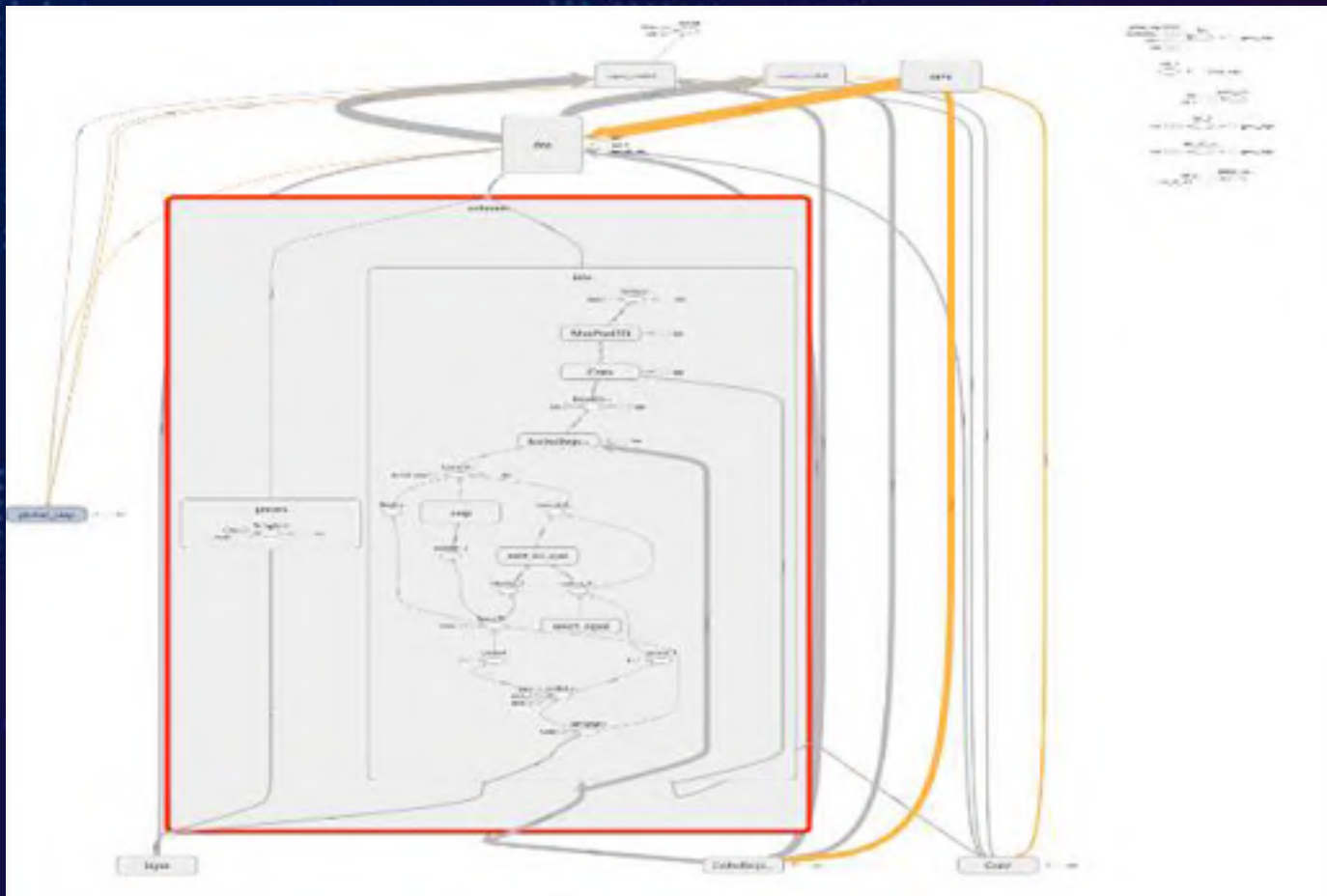
训练



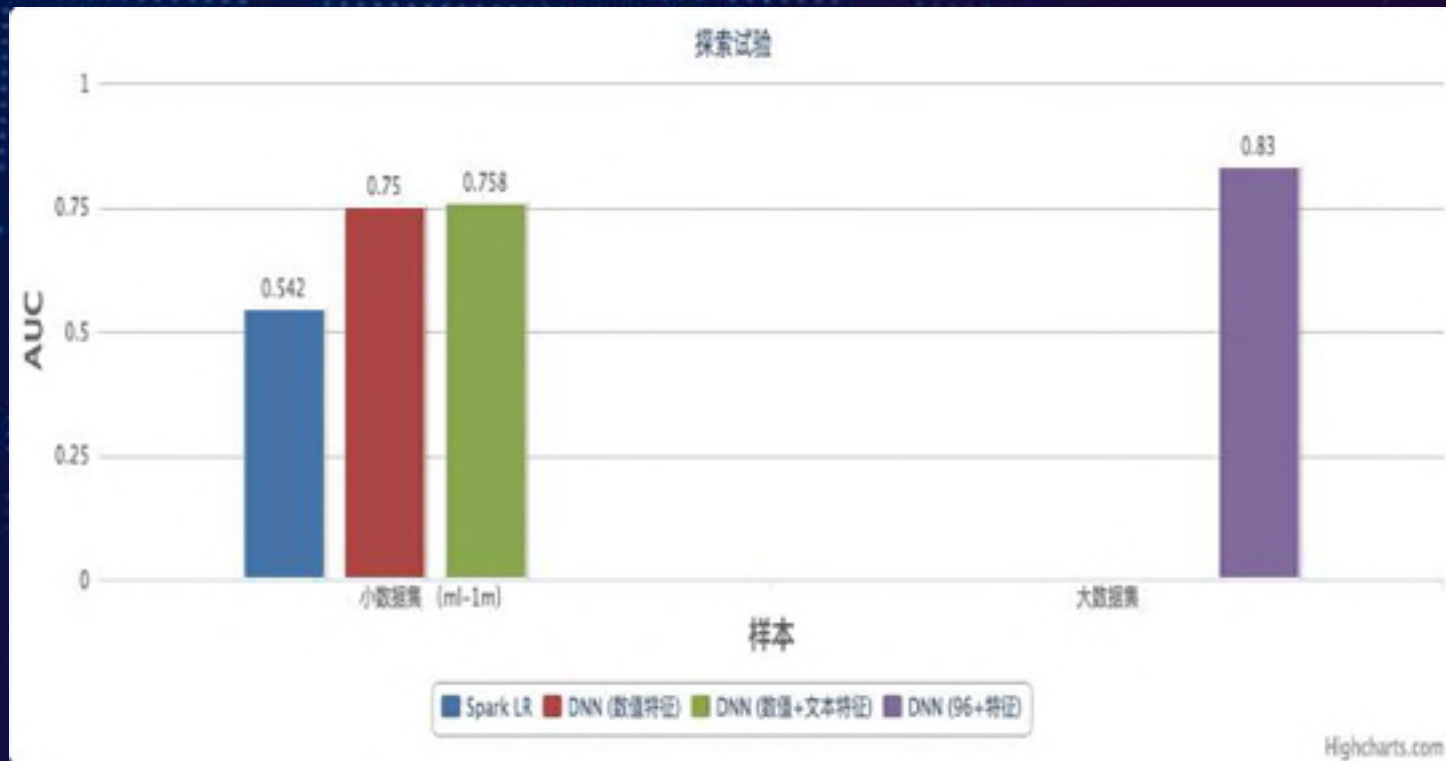
◆ 分析



◆ 分析



分析



03

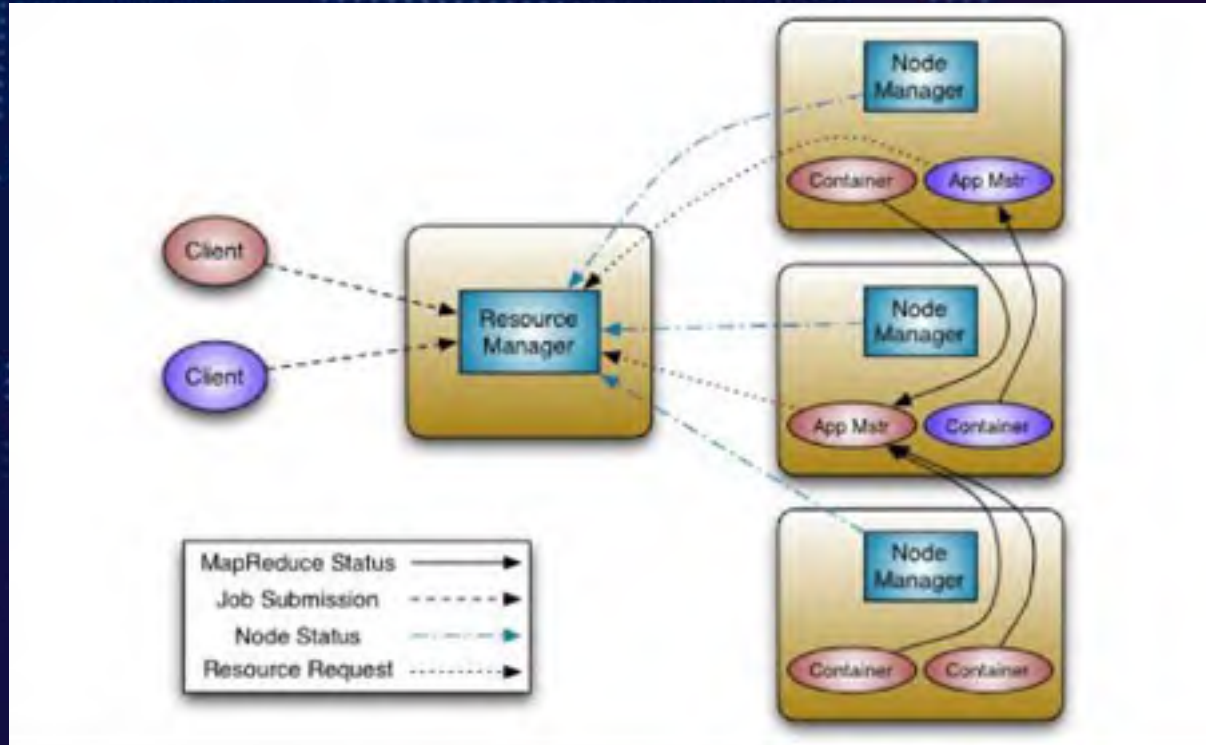
导出与预测

```
<config version="0.0.3" model="prometheus">
  <output>
    <export_dir_base>file:///tmp/facai/prometheus/export</export_dir_base>
  </output>
</config>
```

```
client = dnn_tf-serving_client.DnnClient(FLAGS.server)
inputs_dict = {
    "uid": _int64_feature(value=1),
    "mid": _int64_feature(value=1193),
    "time": _float_feature(value=0),
    "gender": _bytes_feature(value='M'.encode()),
    "age": _int64_feature(value=25),
    "job": _int64_feature(value=0),
    "zip_code": _float_feature(value=0),
    "title": _bytes_feature(value='Shine'.encode()),
    "genres": _bytes_feature(value='Drama'.encode())}
result = client.do_predict(FLAGS.model, inputs_dict, FLAGS.timeout)
```

04

分布式计算-TensorFlow on yarn



04 为什么TensorFlow on yarn

- ◆ 集群资源的统一管理和分配
- ◆ 作业统一管理，运行状态实时跟踪，在线的log查看
- ◆ 作业进程的资源隔离
- ◆ 利用微博现有的Hadoop集群

04 TensorFlow on yarn-资源分配

- ◆ AppMaster自动分配PS hosts和Worker Hosts , 并自动分配端口号
- ◆ 自动分配TensorBoard可视化所需的host和端口号
- ◆ 对训练数据自动进行数据分片
- ◆ 启动PS Task和Worker Task , 启动TensorBoard

04 TensorFlow on yarn-模型训练

- ◆ TensorFlow根据PS hosts和Worker hosts组装成ClusterSpec
- ◆ Worker task从HDFS读取其分片的训练数据，进行训练
- ◆ 训练过程中定期将checkpoint保存到HDFS中
- ◆ 通过TensorBoard实时查看计算情况

04 TensorFlow on yarn-资源回收

- ◆ 训练结束后模型保存到HDFS
- ◆ AppMaster强制结束PS进程
- ◆ AppMaster强制结束TensorBoard进程

04 TensorFlow on yarn-示例启动脚

```
# Start Script for launching the Tensorflow job in cluster.
# =====
hadoop jar agent.jar com.weibo.flyhorse.yarn.Client hdfs:///user/weibo/agent.jar TF " \
--command_file train.py \
--command_path hdfs:///user/weibo/train.py \
--input hdfs://77-16-121-hadoop:9000/user/weibo/feed_data_v0/train/20170703 \
--output hdfs://77-16-121-hadoop:9000/user/weibo/tensorflow/train/context/output \
--num_of_worker 10 \
--num_of_ps 5 \
--memory_of_worker 4096 \
--memory_of_ps 1024 \
--num_epochs 1 \
--board_enable true"
```

◆ 总结

- ◆ 标准化：支持全部微博特征
- ◆ 工具化：简化了深度学习引入排序模型
- ◆ 服务化：支持Tf-Serving预测服务

Thank you!