



QCon 全球软件开发大会
INTERNATIONAL SOFTWARE
DEVELOPMENT CONFERENCE

BEIJING 2018

Amazon SageMaker

端到端的托管机器学习平台

演讲者 / AWS 王世帅



基于实践经验总结和提炼的品牌专栏
尽在【极客时间】



重拾极客时间，提升技术认知

通往**年薪百万**的CTO的路上，
如何打造自己的技术**领导力**？

扫描二维码了解详情

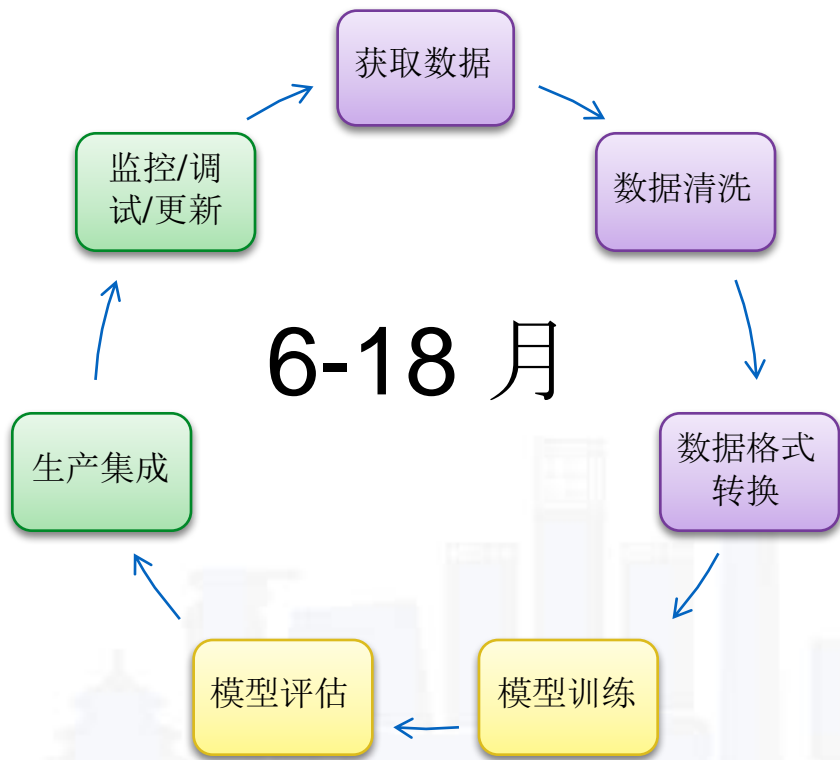


目录

- 为什么我们要推出Amazon SageMaker？
- 什么是Amazon SageMaker？
- 如何使用Amazon SageMaker
- Amazon SageMaker案例分享
- Q&A

为什么我们推出 Amazon SageMaker

应用机器学习艰难而漫长



亚马逊在人工智能领域的大量深度创新



1995

亚马逊自从成立以来一直在人工智能和机器学习领域进行大量投入，并且把我们的知识与能力与客户分享



2017



商品智能推荐



机器人与物流仓储



新产品



供应链管理



智慧呼叫中心

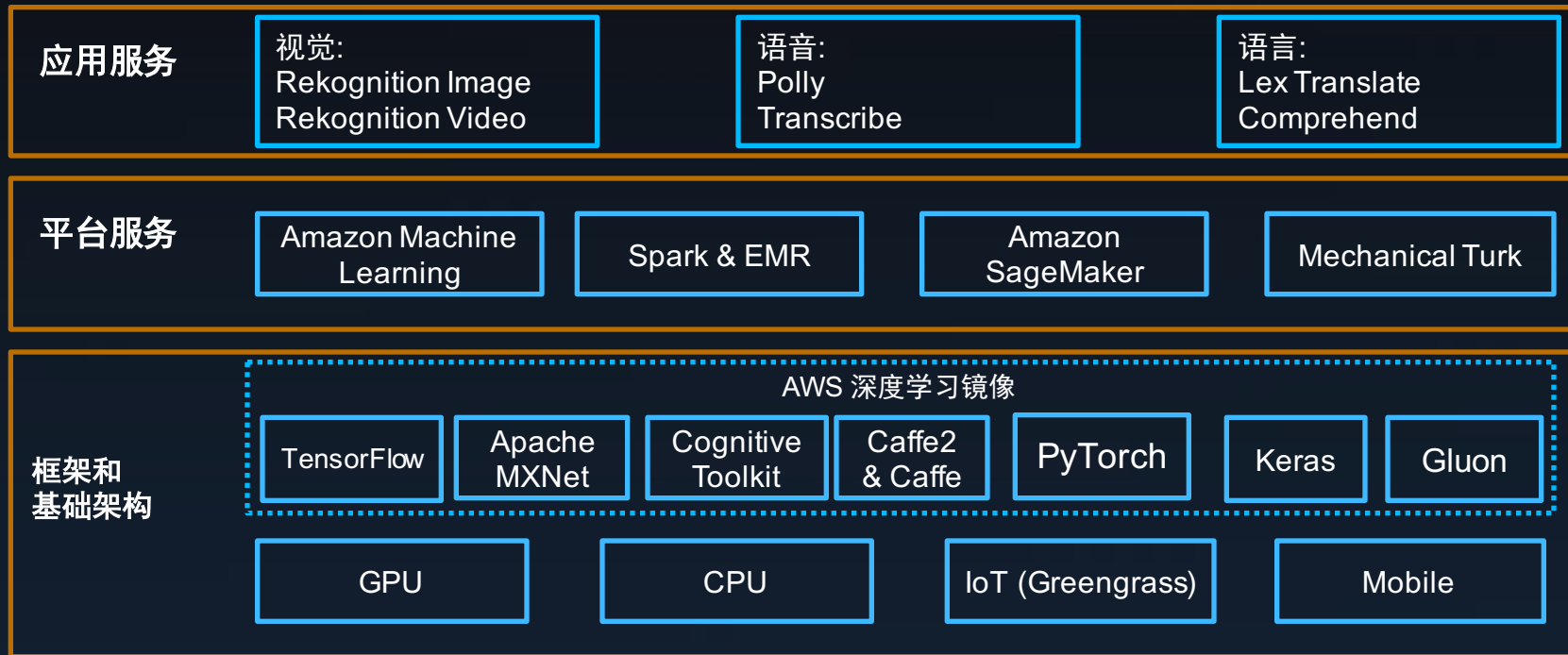


无人值守商店

我们的使命：

为每一个开发者和数据科学家打造机器学习平台

AWS上的机器学习技术堆栈



什么是Amazon Sagemaker ?



Amazon SageMaker

A **fully managed service** that enables **data scientists** and **developers** to quickly and easily **build** machine-learning based models **into production** smart applications.

一个 **全托管服务**，可以帮助 **数据科学家**和 **开发者**快速而轻松地 **构建**基于机器学习的模型的 **生产环境**智能应用

机器学习流程

发现: 业务分析过程

业务问题 -



ML问题转换

- 使用公式来描述正确的问题
 - 需要领域知识

数据调整

数据收集

数据集成

数据准备和清洗

数据可视化和分析

特征工程

模型训练
参数调整

模型校验

符合
业务目标?

否

是

重新训练

监控
调试

- 预测

模型部署

参数调整

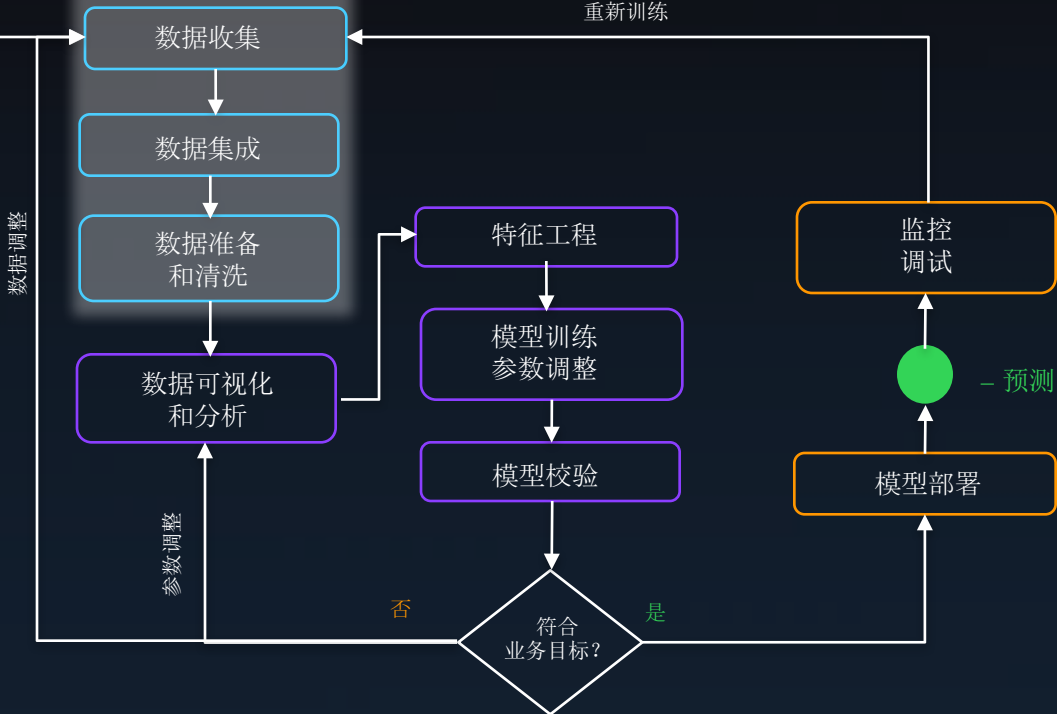
集成: 数据架构过程

业务问题 -

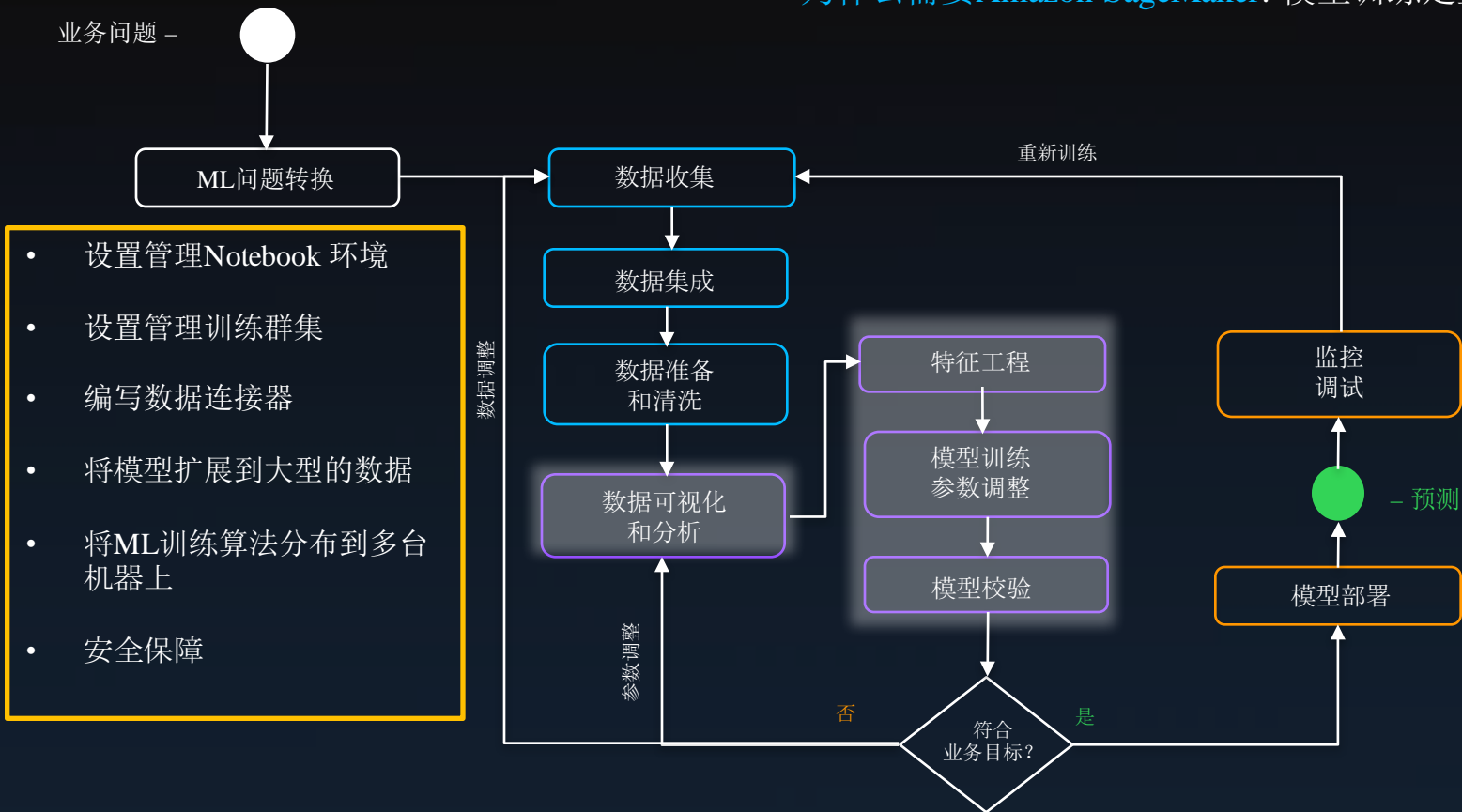


ML问题转换

- 构建数据平台:
 - Amazon S3
 - AWS Glue
 - Amazon Athena
 - Amazon EMR
 - Amazon Redshift Spectrum

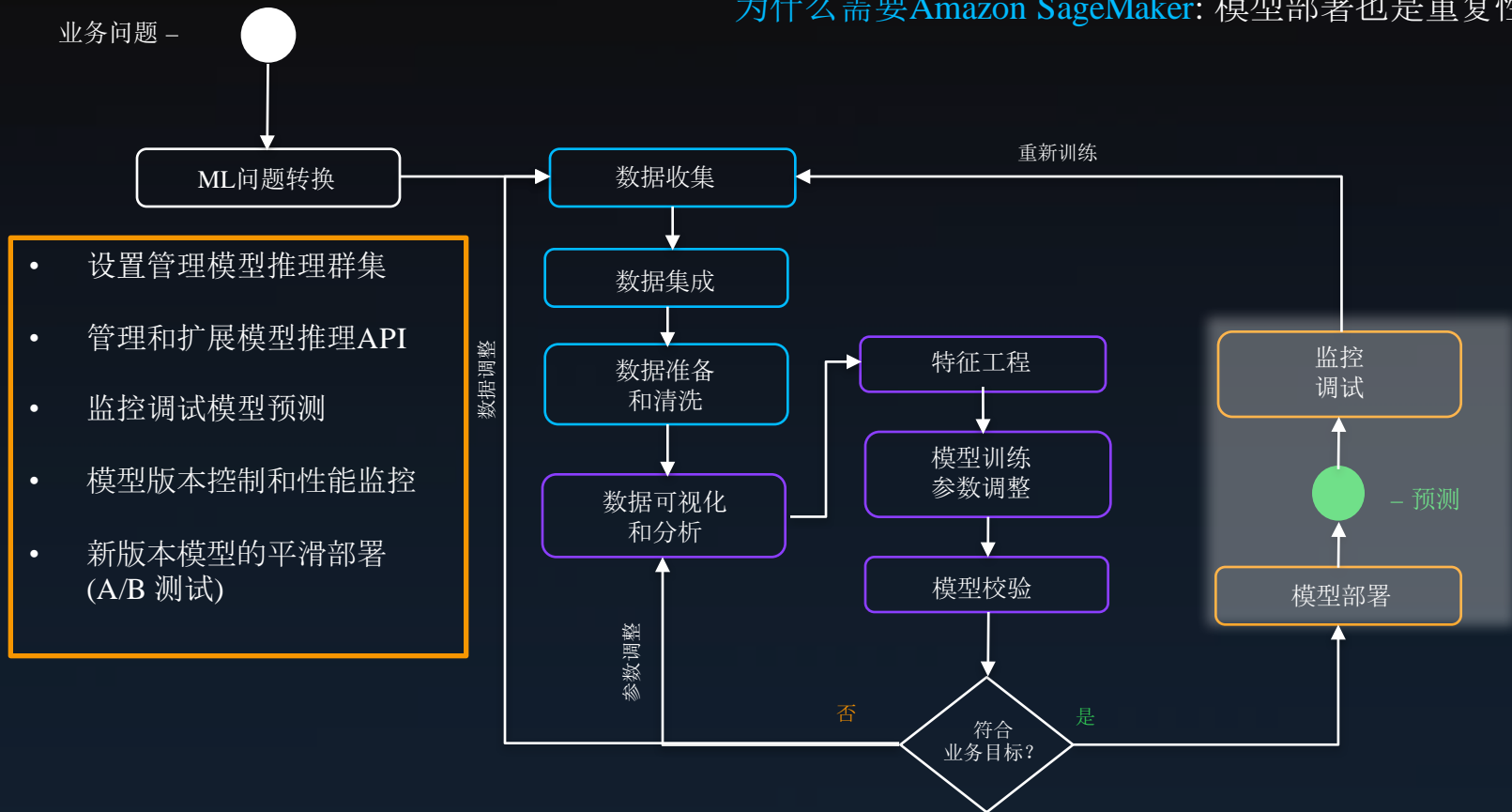


为什么需要Amazon SageMaker: 模型训练是重复性的重体力工作



为什么需要Amazon SageMaker: 模型部署也是重复性的重体力工作

业务问题 -



- 设置管理模型推理群集
- 管理和扩展模型推理API
- 监控调试模型预测
- 模型版本控制和性能监控
- 新版本模型的平滑部署 (A/B 测试)



Amazon SageMaker

1



I

Notebook 实例

2



I

算法

3



I

模型训练服务

4



I

模型部署服务

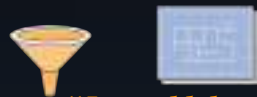
1

零配置 数据探索分析平台



I

Notebook实例



"Just add data"

- 推荐、个性化
- 欺诈检测
- 预测
- 图像分类
- 流失预测
- 市场活动、邮件的定位
- 日志处理, 异常检测
- 语音文字转换
- 其它更多...

2

Amazon SageMaker: 10倍速 算法支持



I
算法



流式数据集
训练费用更低



更快的训练速度

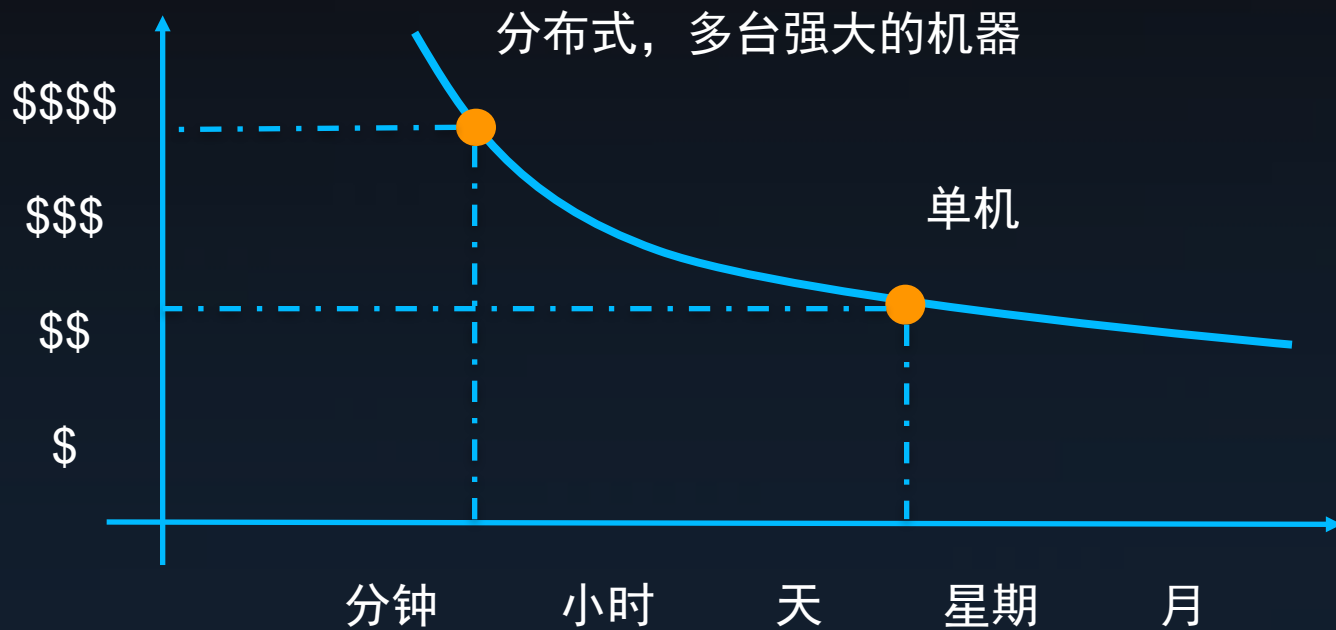


处理超大型数据集时
提供更高的可靠性

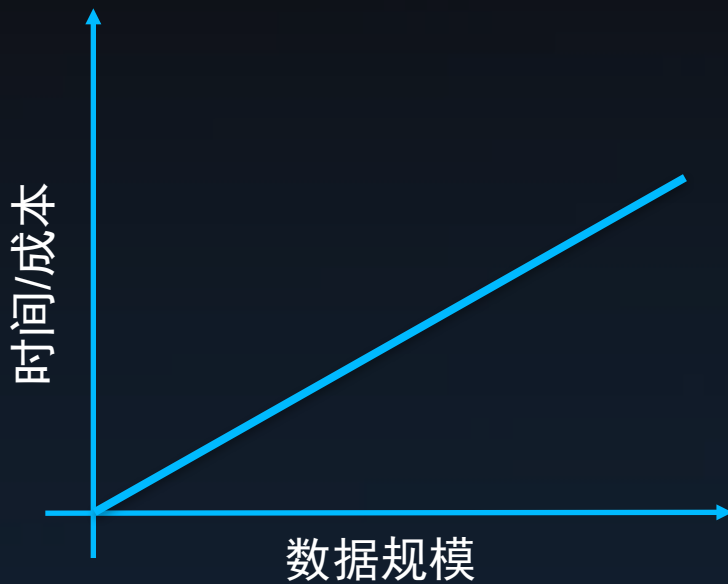
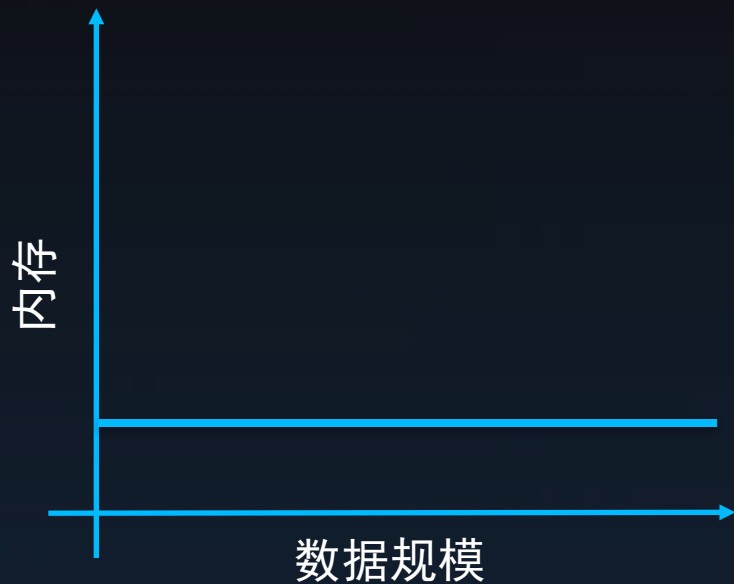


可以选择多种机器学习
算法

费用 vs. 时间



流式数据



费用 vs. 时间



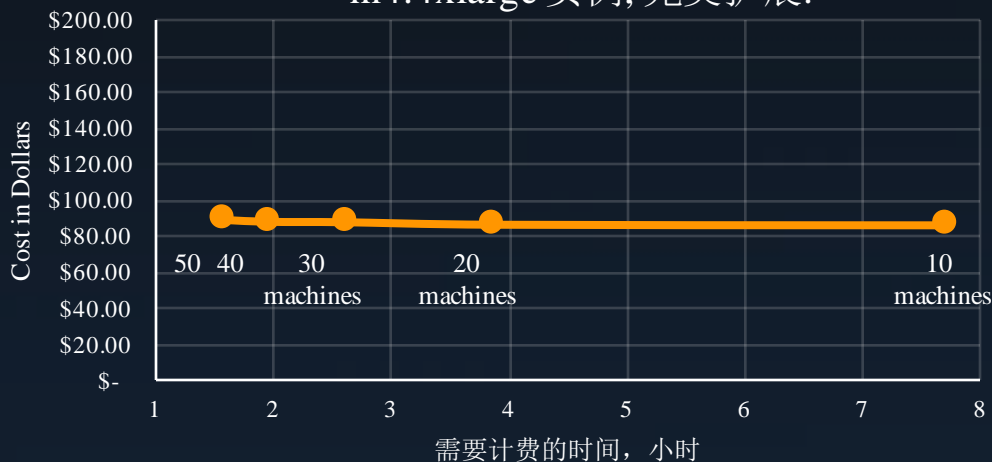
无限扩展的机器学习算法

因子分解机 (Factorization Machines)

$$\tilde{y} = w_0 + \langle w_1, x \rangle + \sum_{i,j>i} x_i x_j \cdot \langle v_i, v_j \rangle$$

| | Log_loss | F1 Score | Seconds |
|------------------|--------------|--------------|------------|
| SageMaker | 0.494 | 0.277 | 820 |
| Other (10 Iter) | 0.516 | 0.190 | 650 |
| Other (20 Iter) | 0.507 | 0.254 | 1300 |
| Other (50 Iter) | 0.481 | 0.313 | 3250 |

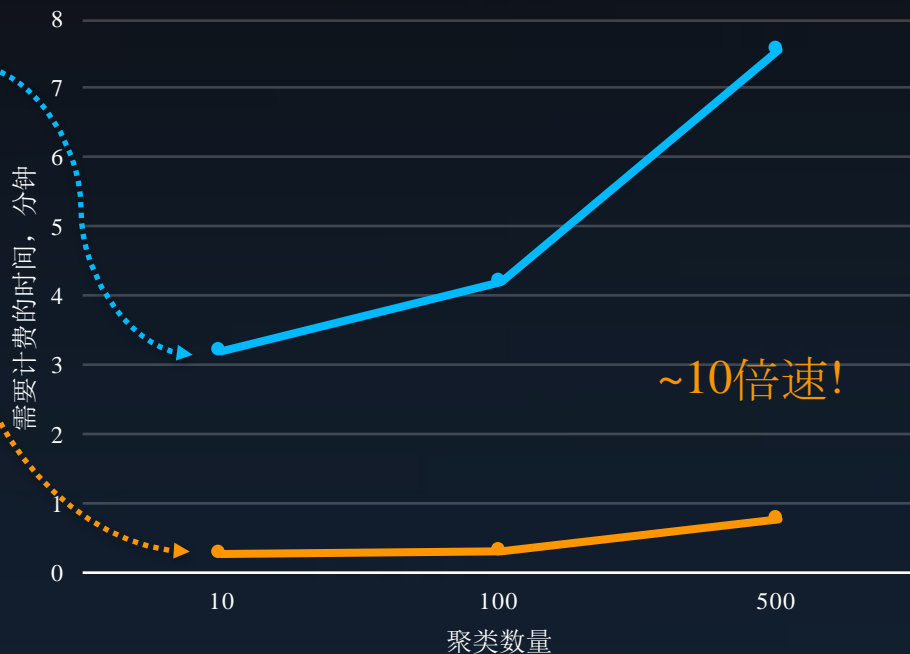
点击预测 1 TB 广告数据集,
m4.4xlarge 实例, 完美扩展.



K-Means 聚类

| | k | SageMaker | Other |
|----------------------|-----|-----------|--------|
| Text 1.2GB | 10 | 1.18E3 | 1.18E3 |
| | 100 | 1.00E3 | 9.77E2 |
| | 500 | 9.18.E2 | 9.03E2 |
| Images 9GB | 10 | 3.29E2 | 3.28E2 |
| | 100 | 2.72E2 | 2.71E2 |
| | 500 | 2.17E2 | Failed |
| Videos 27GB | 10 | 2.19E2 | 2.18E2 |
| | 100 | 2.03E2 | 2.02E2 |
| | 500 | 1.86E2 | 1.85E2 |
| Advertising 127GB | 10 | 1.72E7 | Failed |
| | 100 | 1.30E7 | Failed |
| | 500 | 1.03E7 | Failed |
| Synthetic 1100GB | 10 | 3.81E7 | Failed |
| | 100 | 3.51E7 | Failed |
| | 500 | 2.81E7 | Failed |

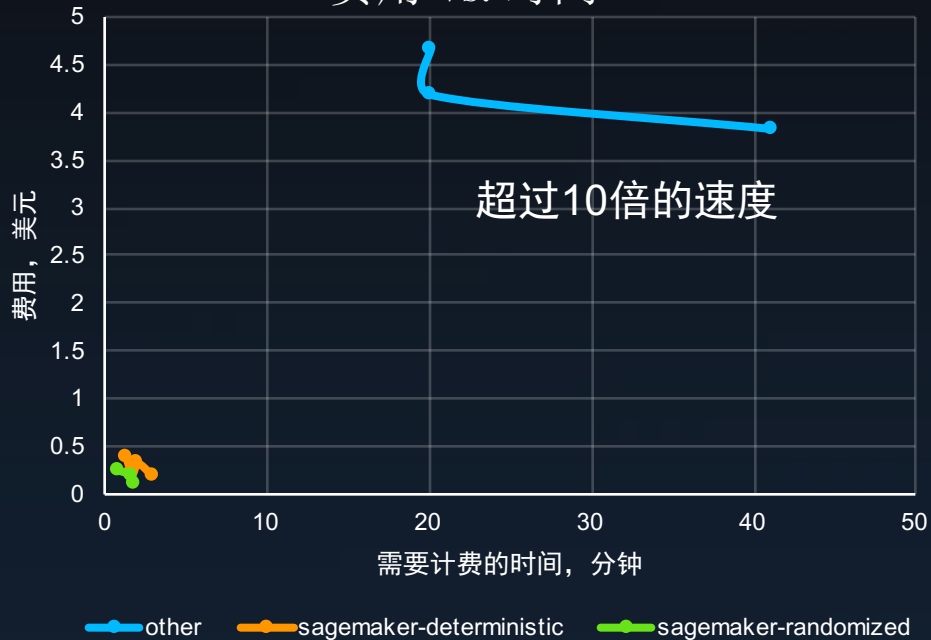
运行时间 vs. 聚类数量



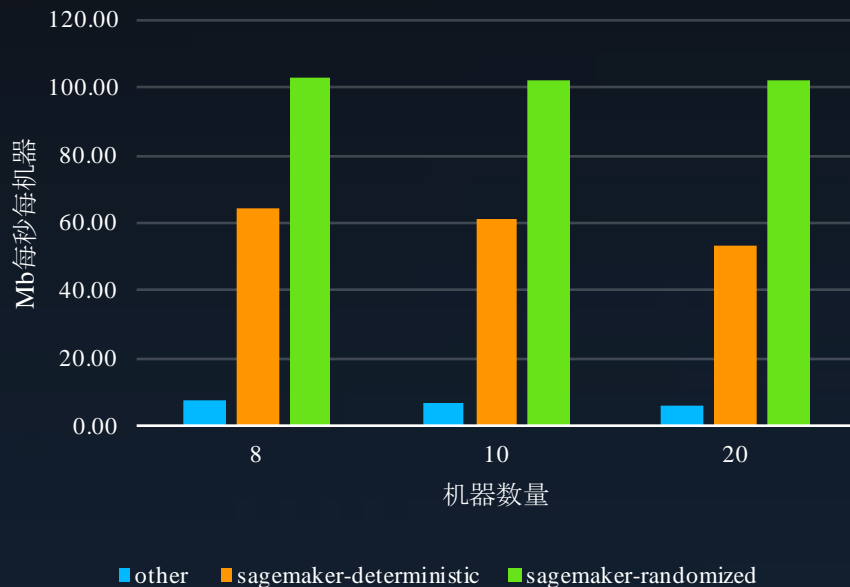
sagemaker other

主成分分析 (PCA)

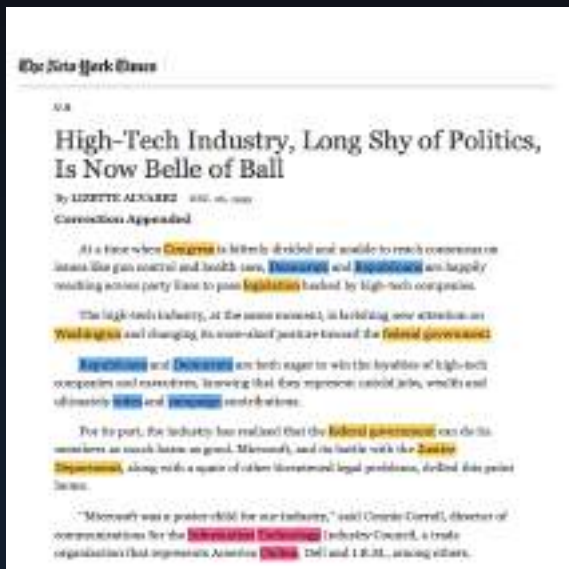
费用 vs. 时间



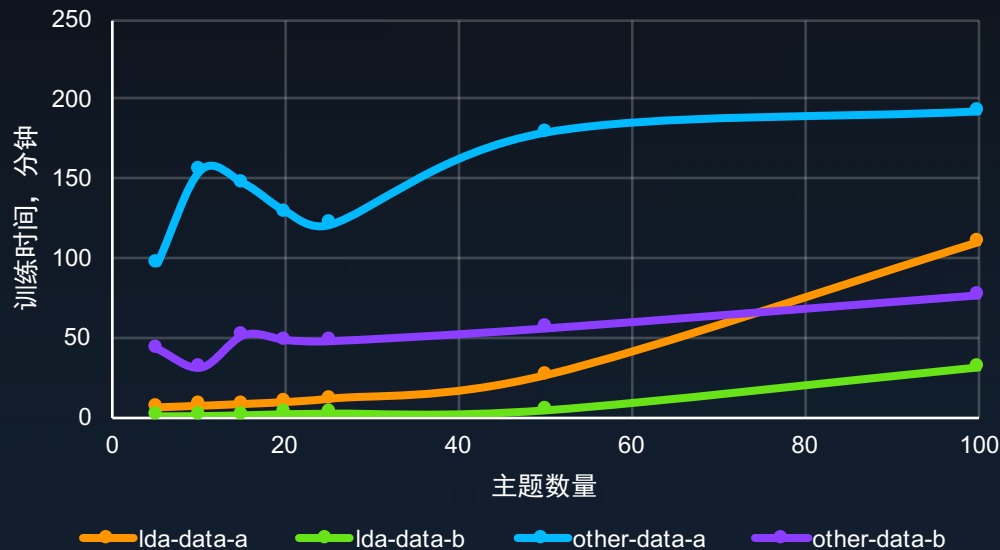
吞吐量 and 扩展性



Spectral LDA



训练时间 vs. 主题数量



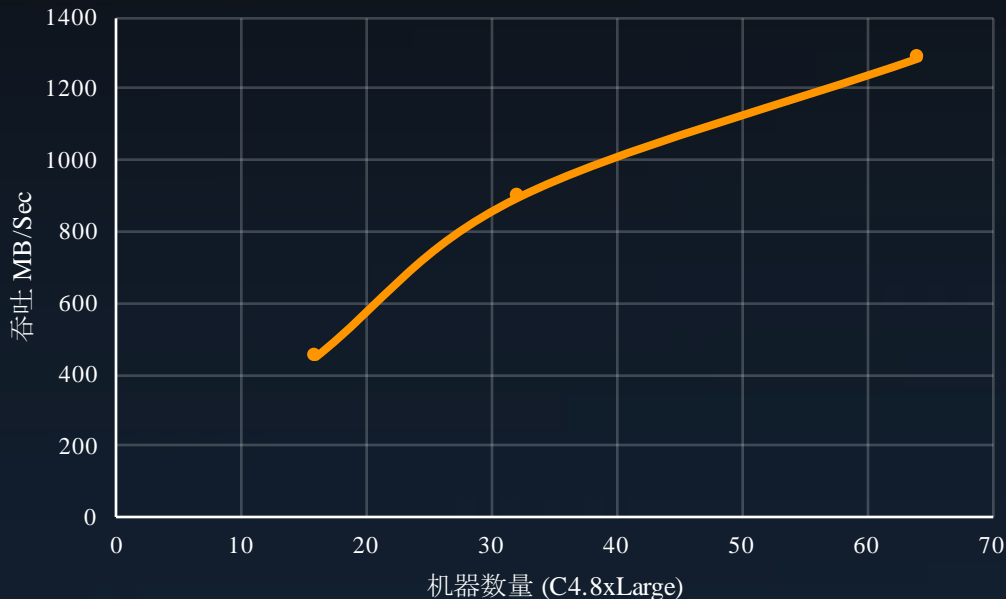
Boosted 决策树

XGBoost 是最常用的boosted决策树实现

现在在 Amazon SageMaker 可以直接使用！



吞吐 vs. 机器数量



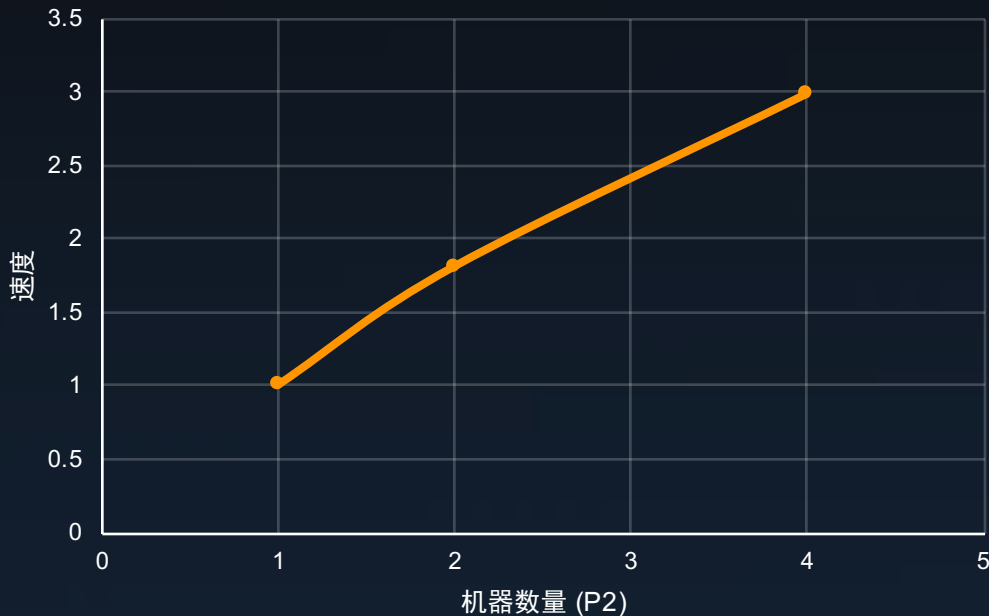
图像分类

通过MxNet 实现的ResNet.

未来会提供更多的网络结构，
如 DenseNet， Inception

迁移学习: 从ImageNet上已经
训练好的模型开始工作

通过水平扩展提升速度



3

全托管的、灵活的 分布式训练



— 安全保障



I

模型训练服务



获取训练数据



训练代码

保存模型



保存推理镜像



Amazon ECR

- Matrix Factorization
- Regression
- Principal Component Analysis
- K-Means Clustering
- Gradient Boosted Trees
- 更多算法！

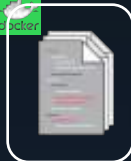
AWS提供的算法



自带脚本 (SM构建容器)



Apache Spark
集成



自带算法 (自己通过容器构建)

CPU

GPU

HPO

全托管



4

轻松将模型部署到Amazon SageMaker



I

ML部署服务

- ✓ 自动扩展的推理API
- ✓ A/B测试
- ✓ 低延迟 & 高吞吐
- ✓ 使用自己的模型
- ✓ Python SDK



Amazon SageMaker

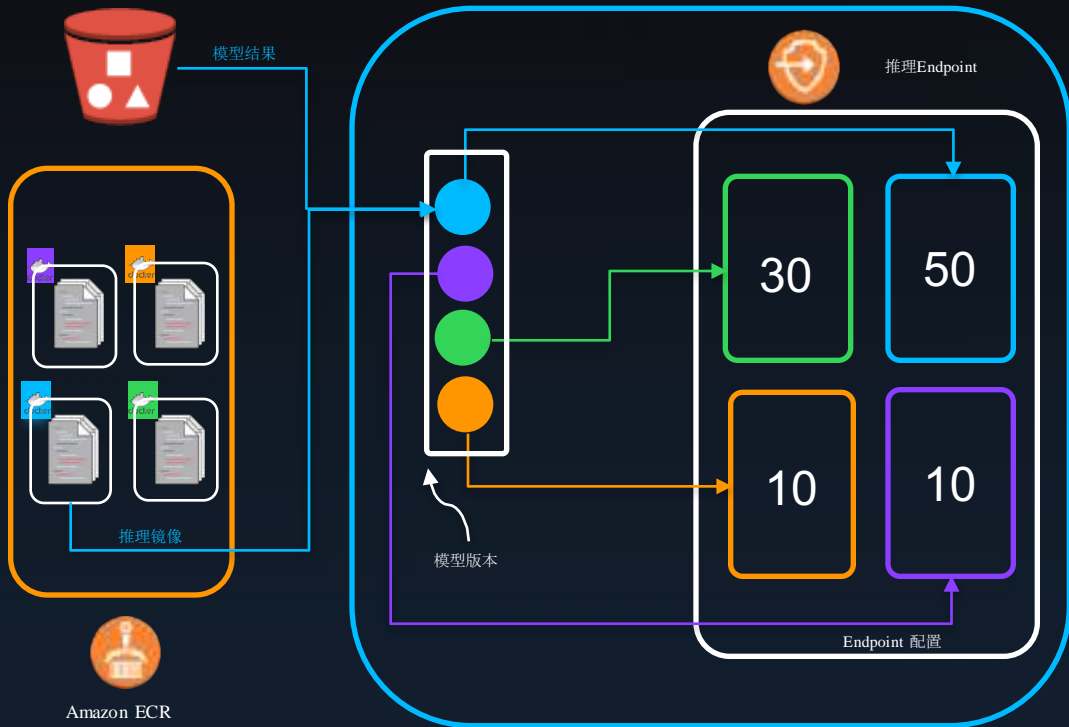
4

轻松将模型部署到Amazon SageMaker



I 模型部署服务

在推理容器中保存多个版本的镜像. Prod 是主要版本, 支持50% 用户流量



实例类型: c3.4xlarge
初始实例数量: 3
模型名称: prod
版本名称: primary
初始版本权重: 50

生产版本

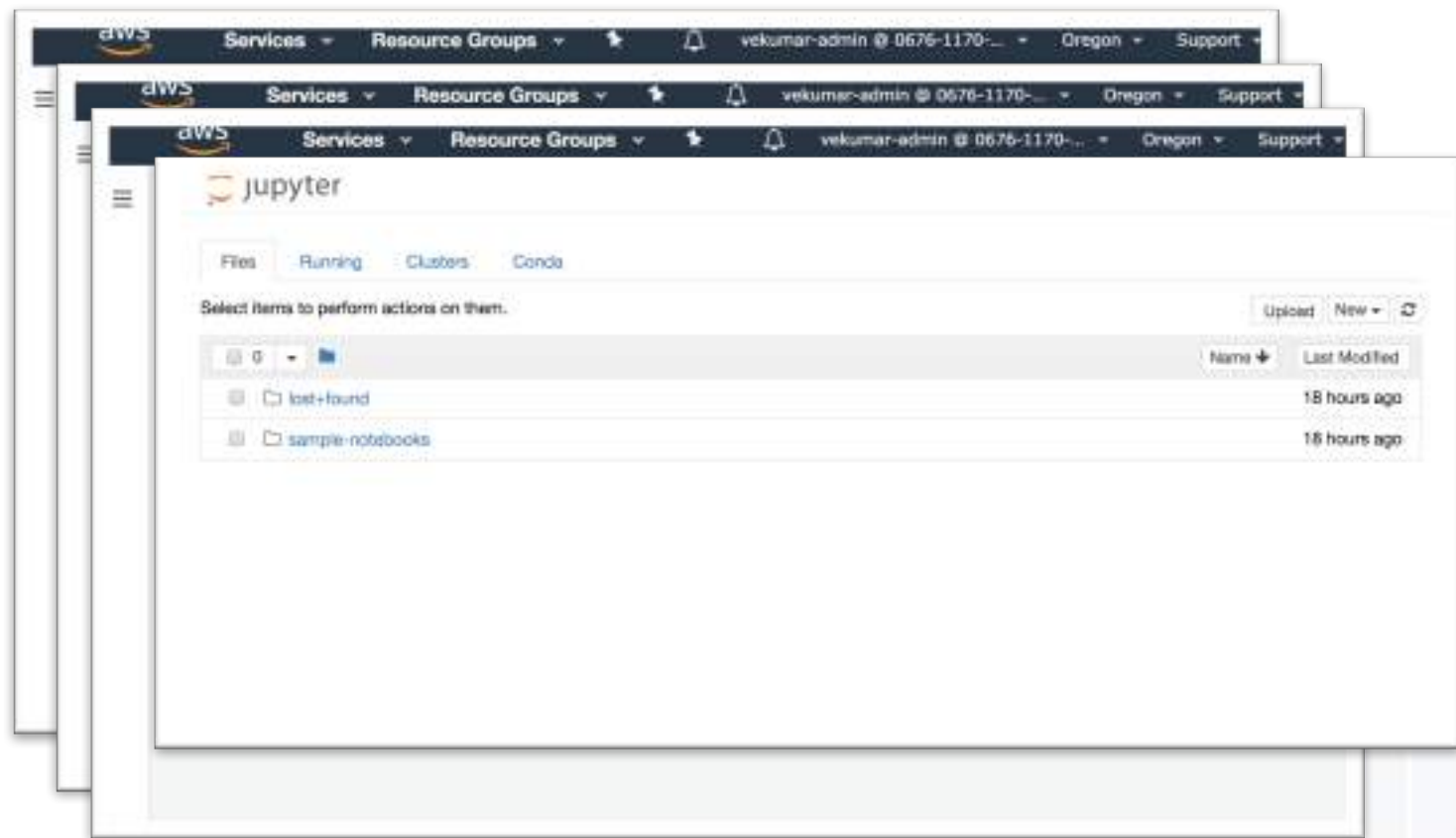
通过Endpoint配置 创建一个
新的 Endpoint



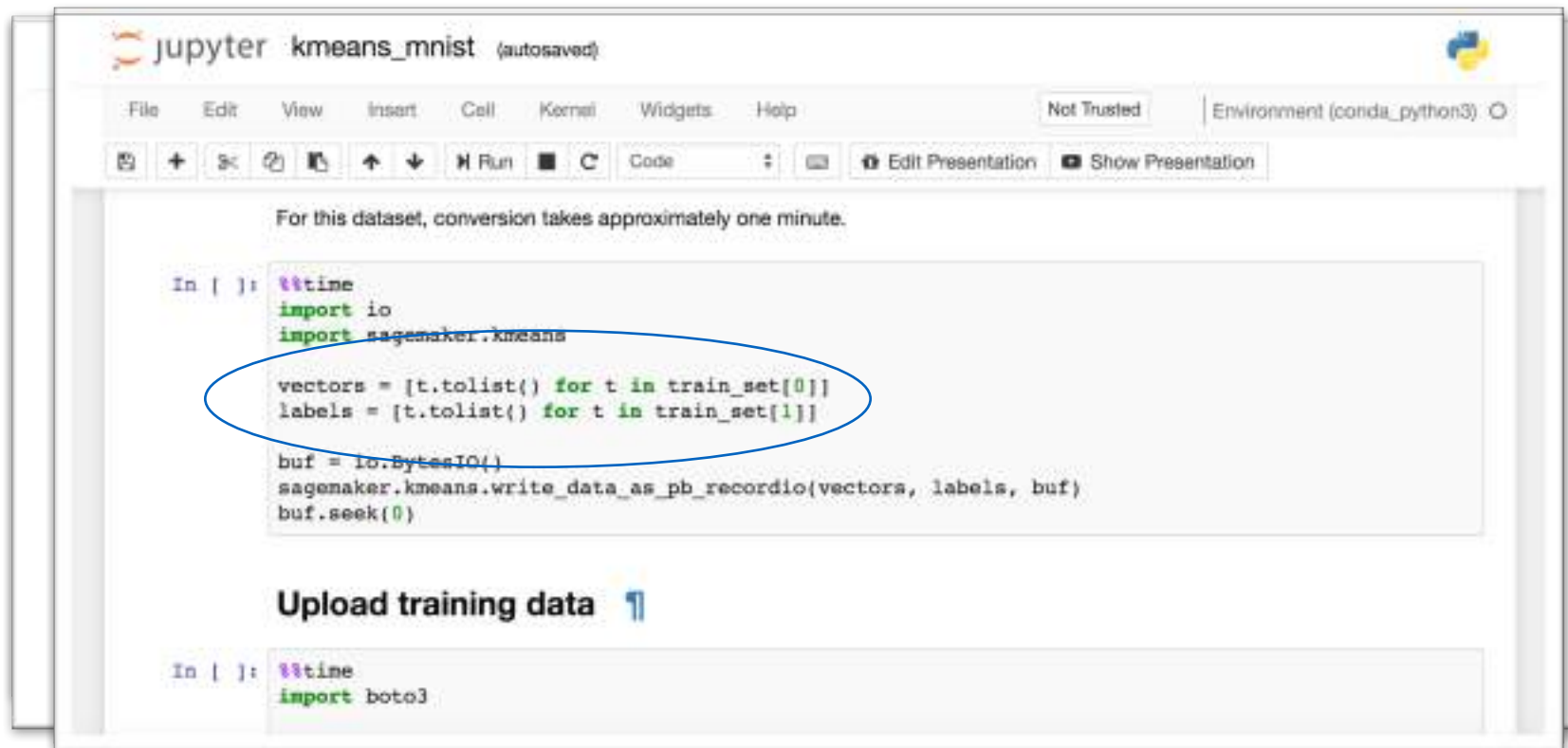
Amazon SageMaker

如何使用Amazon SageMaker

从Notebook样例开始



配置训练数据集



The screenshot shows a Jupyter Notebook titled "kmeans_mnist (autosaved)". The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), a "Not Trusted" status indicator, and an environment selector "Environment (conda_python3)". Below the menu is a toolbar with icons for file operations and execution. The main content area contains a text prompt: "For this dataset, conversion takes approximately one minute." Below this is a code cell with the following Python code:

```
In [ ]: %%time
import io
import sagemaker.kmeans

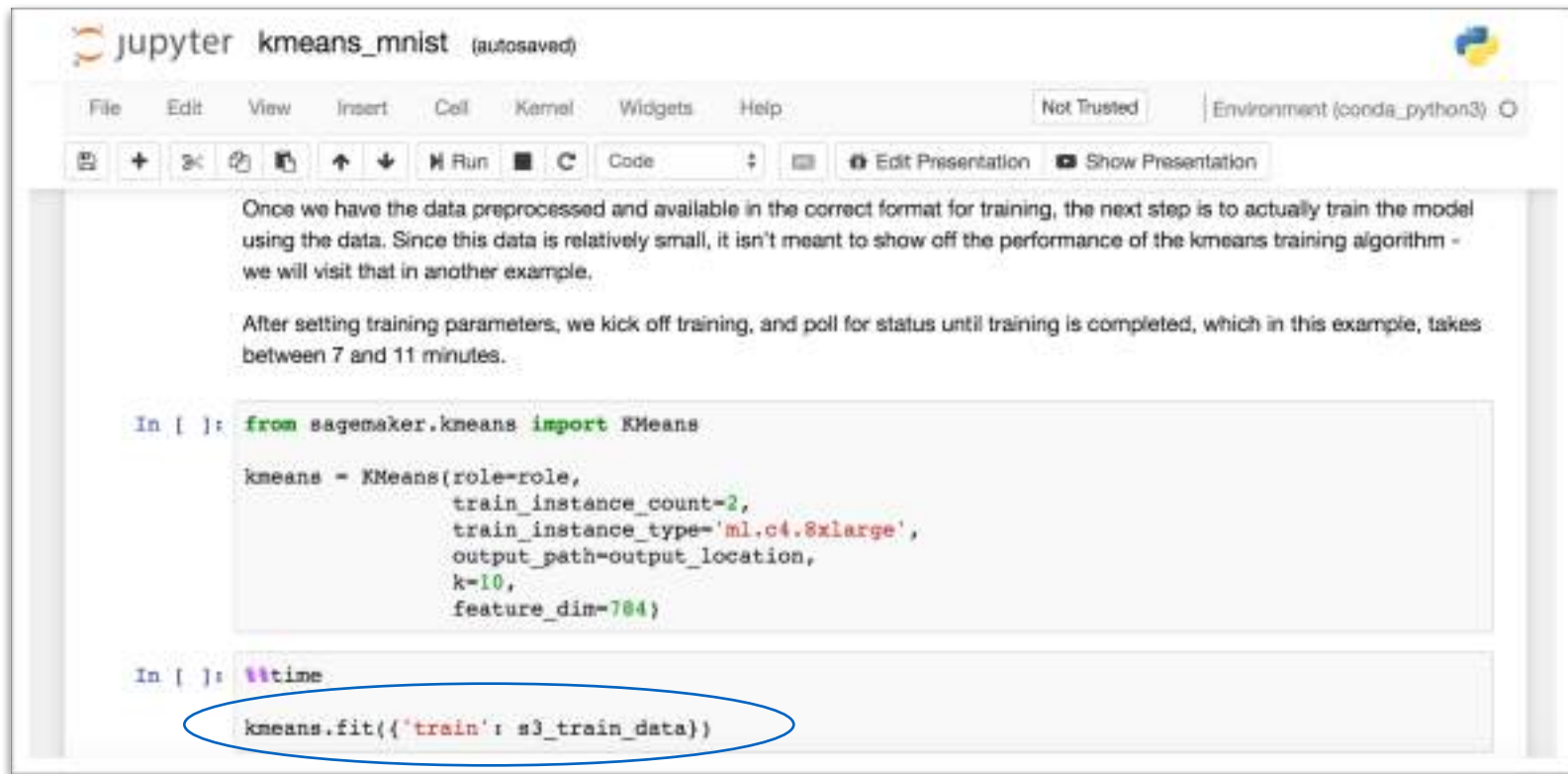
vectors = [t.tolist() for t in train_set[0]]
labels = [t.tolist() for t in train_set[1]]

buf = io.BytesIO()
sagemaker.kmeans.write_data_as_pb_recordio(vectors, labels, buf)
buf.seek(0)
```

The code cell is followed by a section titled "Upload training data" with a blue arrow icon. Below this is another code cell with the following Python code:

```
In [ ]: %%time
import boto3
```

训练你的模型



The screenshot shows a Jupyter Notebook titled "kmeans_mnist (autosaved)". The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), a status bar (Not Trusted, Environment (conda_python3)), and a toolbar with icons for file operations, running, and presentation. The notebook content consists of two text blocks and two code cells. The first text block explains that the data is preprocessed and that training the model is the next step. The second text block notes that training takes between 7 and 11 minutes. The first code cell imports KMeans from sagemaker.kmeans and defines a KMeans object with specific parameters. The second code cell starts with a time measurement and then calls the fit method on the KMeans object, which is circled in blue.

Once we have the data preprocessed and available in the correct format for training, the next step is to actually train the model using the data. Since this data is relatively small, it isn't meant to show off the performance of the kmeans training algorithm - we will visit that in another example.

After setting training parameters, we kick off training, and poll for status until training is completed, which in this example, takes between 7 and 11 minutes.

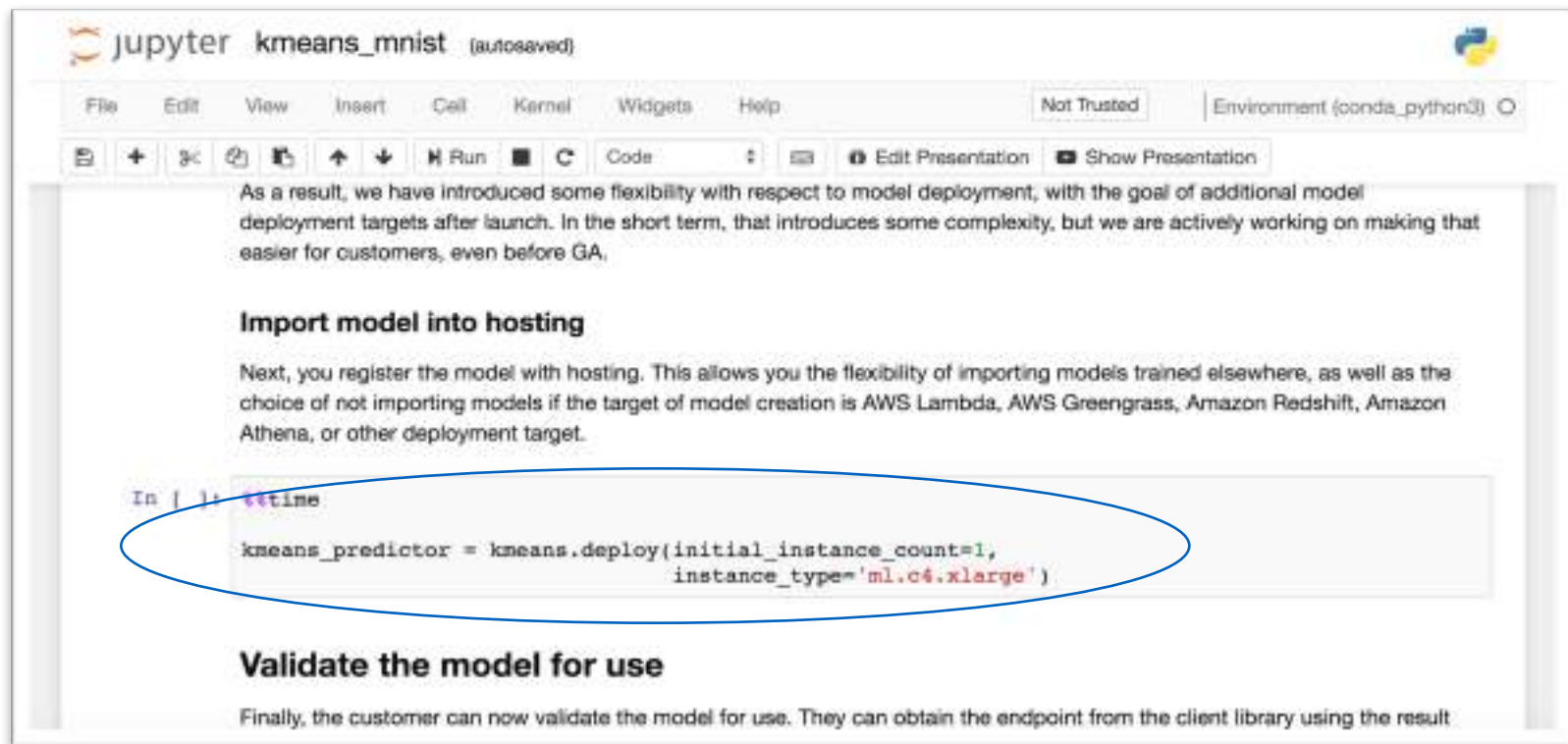
```
In [ ]: from sagemaker.kmeans import KMeans

kmeans = KMeans(role=role,
                 train_instance_count=2,
                 train_instance_type='ml.c4.8xlarge',
                 output_path=output_location,
                 k=10,
                 feature_dim=784)
```

```
In [ ]: %time

kmeans.fit({'train': s3_train_data})
```

部署你的模型



The screenshot shows a Jupyter Notebook window titled "jupyter kmeans_mnist (autosaved)". The interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), a "Not Trusted" warning, and an environment selector "Environment (conda_python3)". Below the menu is a toolbar with icons for file operations, running, and presentation. The main content area contains text explaining model deployment flexibility, followed by a section titled "Import model into hosting" which describes registering a model for various AWS targets. A code cell is highlighted with a blue oval, containing the following Python code:

```
In [ ]: %%time
kmeans_predictor = kmeans.deploy(initial_instance_count=1,
                                 instance_type='ml.c4.xlarge')
```

Below the code cell is a section titled "Validate the model for use" with introductory text.

检验你的模型

jupyter kmeans_mnist (autosaved)

File Edit View Insert Cell Kernel Widgets Help Not Trusted Environment (conda_python3)

instance_type='ml.c4.xlarge')

Validate the model for use

Finally, the customer can now validate the model for use. They can obtain the endpoint from the client library using the result from previous operations, and generate classifications from the trained model using that endpoint.

```
In [ ]: result = kmeans_predictor.predict(train_set[0][30:31])
        print(result)
```

OK, a single prediction works.

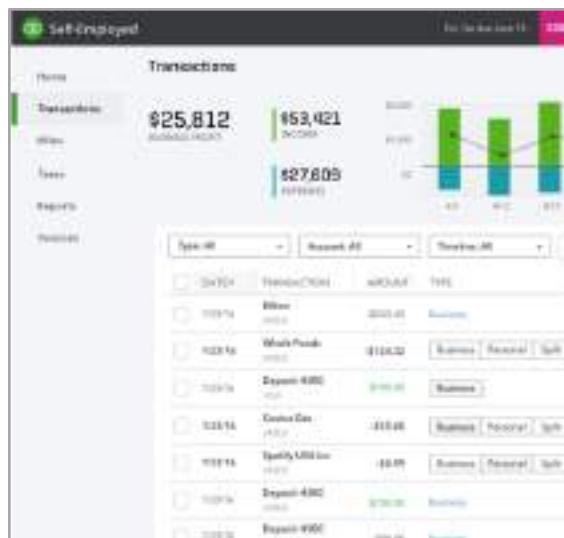
Let's do a whole batch and see how well the clustering works.

```
In [ ]: %time
        result = kmeans_predictor.predict(valid_set[0][0:100])
        clusters = result['labels']
```

Amazon SageMaker案例分享

AI and ML at Intuit

产品推荐



欺诈检测与防止



用户关怀与建议



intuit

Intuit使用SageMaker获得的好处

From

To

需要临时设置
和管理notebook环境



使用SageMaker notebook轻松完
成数据探索工作

有限的模型部署选择



通过虚拟化手段
达成极强的灵活性

团队之间需要争抢计算资源

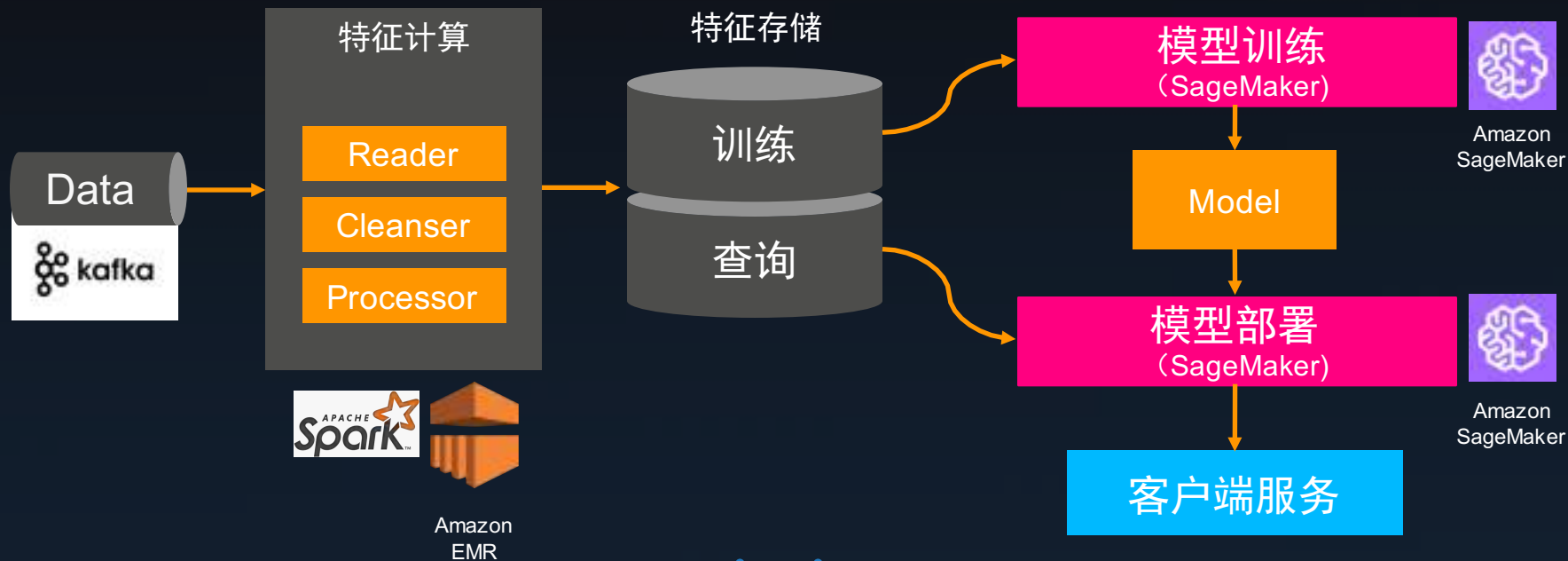


自动扩展的模型部署环境

intuit.



使用 SageMaker 在 AWS 上构建近实时的欺诈检测



intuit.

Amazon SageMaker – 下一步动作

- Amazon SageMaker入门文档: <https://aws.amazon.com/sagemaker/>
- 使用Amazon SageMaker SDK:
 - Python: <https://github.com/aws/sagemaker-python-sdk>
 - Spark: <https://github.com/aws/sagemaker-spark>
- SageMaker 样例: <https://github.com/aws-labs/amazon-sagemaker-examples>
- 把你构建的应用告诉我们!



Amazon SageMaker

端到端的托管机器学习平台

!GO BUILD!

GMTC 2018

全球大前端技术大会

—— 大前端的下一站 ——



<<扫码了解更多详情>>

关注 ArchSummit 公众号
获取国内外一线架构设计
了解上千名知名架构师的实践动向



Apple • Google • Microsoft • Facebook • Amazon 腾讯 • 阿里 • 百度 • 京东 • 小米 • 网易 • 微博

深圳站：2018年7月6-9日

北京站：2018年12月7-10日

QCon

全球软件开发大会【2018】

上海站

2018年10月18-20日

7折

预售中, 现在报名立减2040元

团购享更多优惠, 截至2018年7月1日



扫码关注
获取更多培训信息





关注QCon微信公众号，
获得更多干货！

Thanks!



Geekbang > InfoQ
极客公园