



机器学习助力精准营销广告

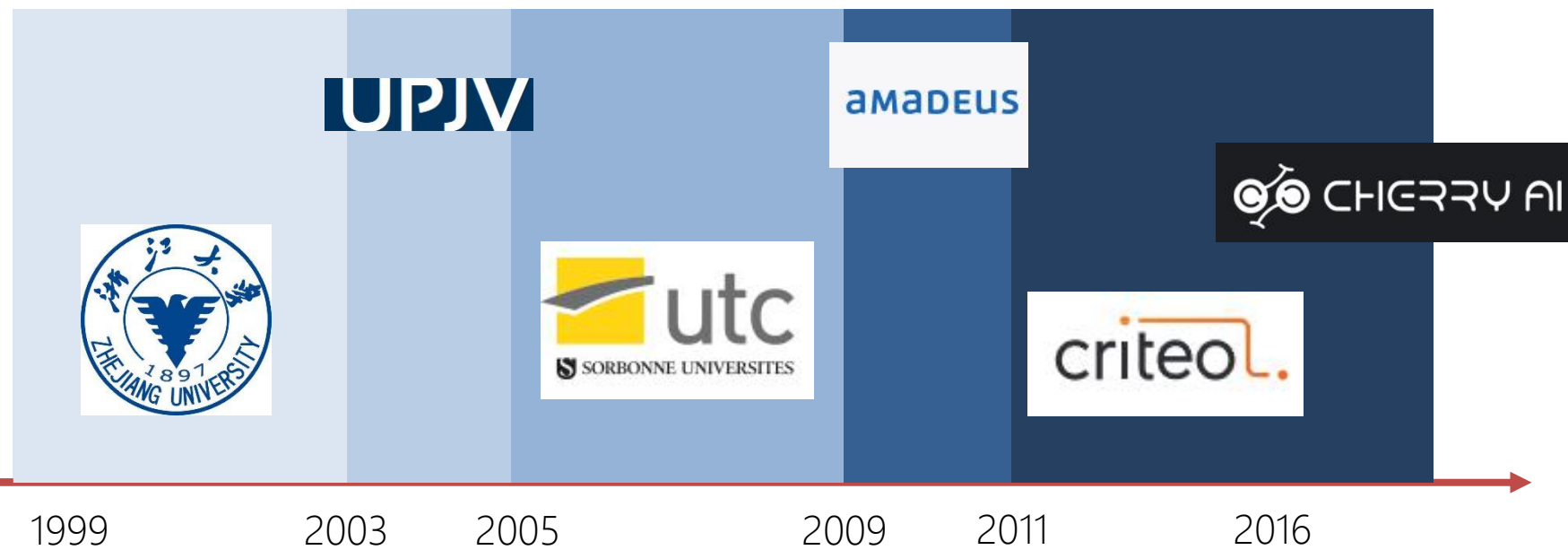
CherryAI

徐煌

About Me



▶ 徐煌



目录

CONTENTS

1. The Ads Bussiness
2. Machine Learning is everywhere
3. ML Life Cycle
4. Large Scale Machine Learning

The Bussiness

广告主

广告代理商

广告展示平台

最终用户



设置广告受众 复制已有受众人群

自定义人群包	不限	定向人群包	排除人群包
地域	不限	省市	区县
性别	不限	男	女
年龄	不限	指定年龄段	
兴趣分类	不限	添加兴趣	

兴趣分类

- 全选
- 游戏
 - 休闲时间
 - 跑酷竞速
 - 宝石消除
 - 网络游戏
 - 动作射击
 - 扑克棋牌

兴趣关键词	不限	自选关键词
app行为定向	不限	按分类 按app
用户首次激活时间	不限	指定时间段



Only Big Whales



Tencent 腾讯

Buy options

今日头条 投放管理平台

设置广告受众 复制已有受众人群

自定义人群包

不限	定向人群包	排除人群包
----	-------	-------

地域

不限	省市	区县	商
----	----	----	---

性别

不限	男	女
----	---	---

年龄

不限	指定年龄段
----	-------

兴趣分类

不限	添加兴趣
----	------

兴趣分类

- 全选
- 游戏
 - 休闲时间
 - 跑酷竞速
 - 宝石消除
 - 网络游戏
 - 动作射击
 - 扑克棋牌

兴趣关键词

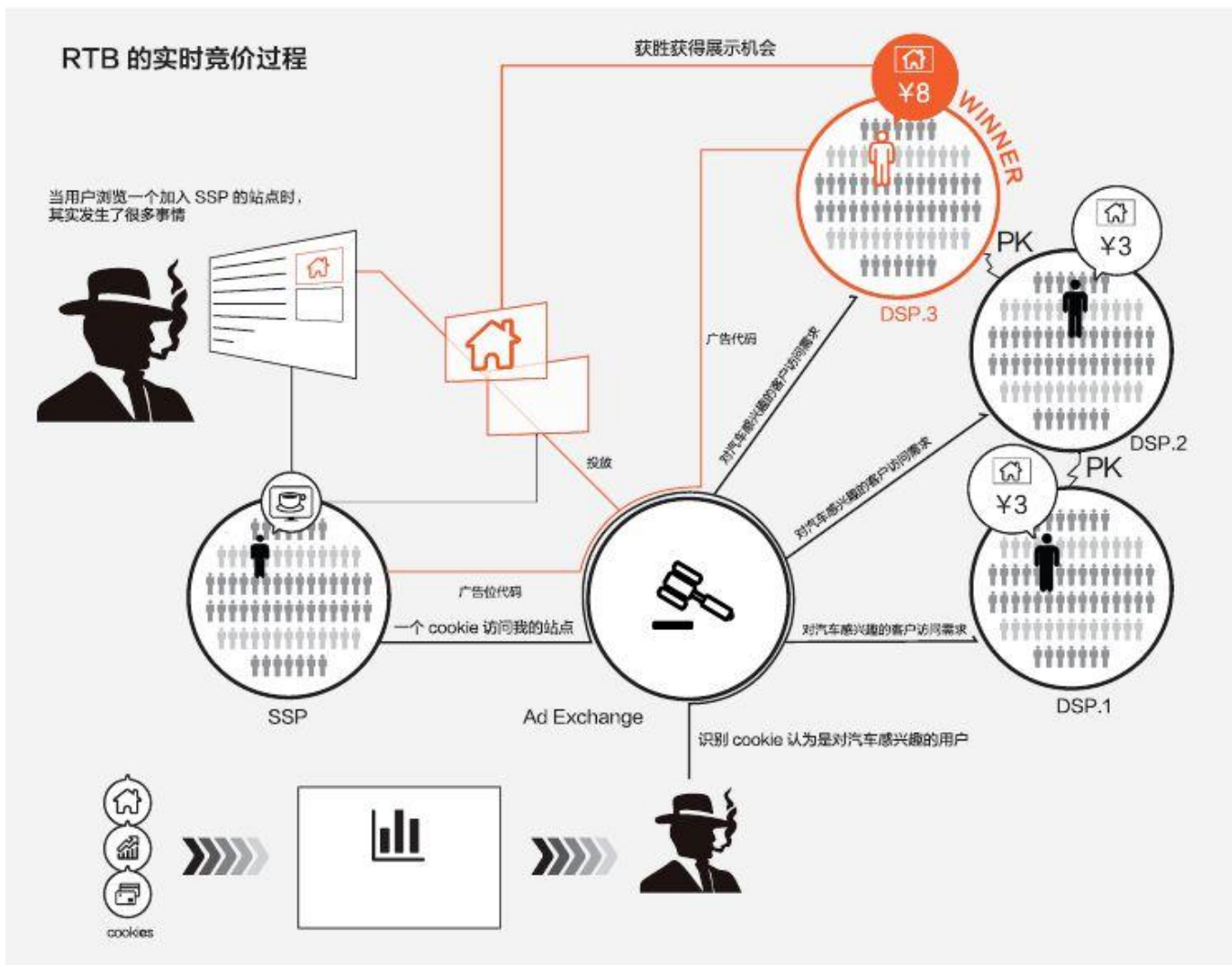
不限	自选关键词
----	-------

app行为定向

不限	按分类	按app
----	-----	------

用户首次激活时间

不限	指定时间段
----	-------



Problems for Advertisers

- Audience?
- Price to buy Ads?
- Content that shows to the audience?
- When?
- ...

Performance Advertising

« The right ad at the right time to the right user »



How we earn money?

- Clients pay us per click, sale etc.
- We buy advertisements from publishers (google, facebook, etc.) in cost of displays.
- We earn the difference: $\text{Click} * \text{Cost per click} - \text{Cost of displays}$



目录

CONTENTS

1. The Ads Bussiness
2. Machine Learning is everywhere
3. ML Life Cycle
4. Large Scale Machine Learning

ML is everywhere

- We use ML for:
 - Bidding
 - Campaign selection
 - Look&Feel optimization
 - Product recommendation

Campaign selection

- Choose the best client for current user



...

- Choose max estimated reward

Bidding

- Estimate the real value of a display
 - The estimated value (estimated cost per display) could be varied for different bussiness model (CPC/CRO/COS/Target COS)



© Can Stock Photo - csp9048331

Example: prediction CPM based on CPC in bidding



Buy ? $\mathbb{E}[CPM] > CPM$

$$\mathbb{E}[CPM] = \mathbb{E}[NbClicks] * CPC$$

Recommendation

- Choose the best (Click Rate, Conversion Rate, Estimated Sales Amount) products to show in the banner



CarGurus

2011 BMW 3 Series 328i SULEV
Save \$2,304. Space gray metallic, 328i...
Great Deal
\$23,900 [VIEW](#)

2011 BMW 3 Series 328i xDrive
Save \$4,181. Black sapphire metallic, ...
Great Deal
\$15,990 [VIEW](#)



Jules

Blouson mix m...
79,99 €
[Acheter](#)

Doudoune cot...
89,99 €
[Acheter](#)



intuit **QuickBooks** In partnership with **STAPLES**

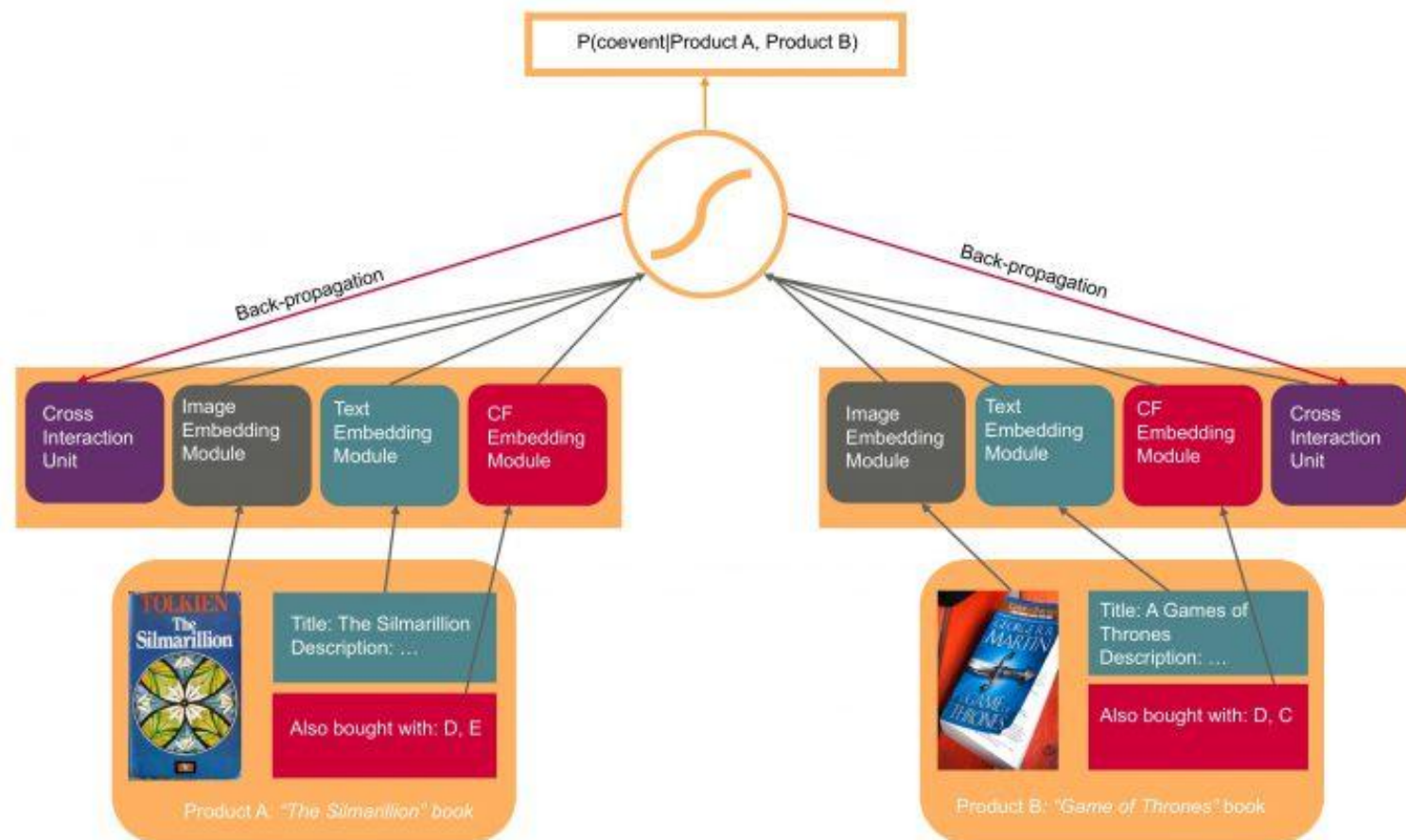
Asus, X551mav- Rcln06, 15.6" ... ~~\$299.99~~ [Shop](#)

Hp 950xl/951 High Yield... ~~\$349.99~~ [Shop](#)

Google Nexus 7 Lte 7", 32gb ... [Shop](#)

Content2vec: Unified product representation for recommender systems

- Collaborative Filtering and others informations
- Reach a unified product representation that gathers all information available on the products to enable us to do better recommendations.



Contextual Recurrent Neural Networks for Recommendation

- RNN with context information



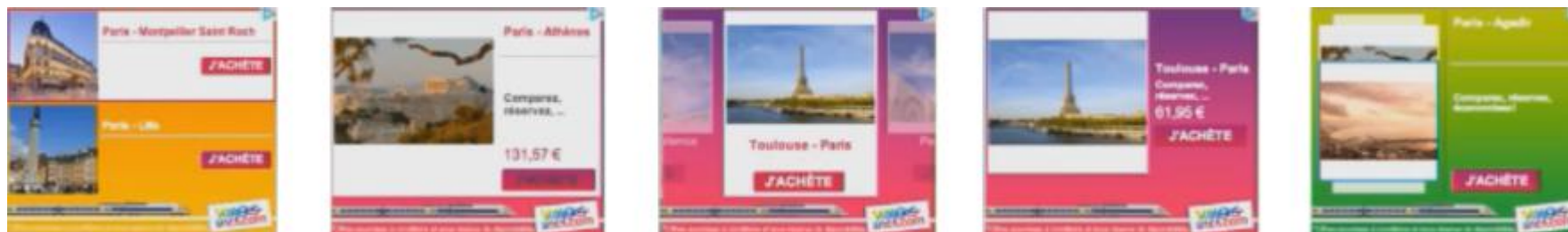
Dynamic rendering optimization

- Choose the look&feel of the banner

Layout:



ColorSet:

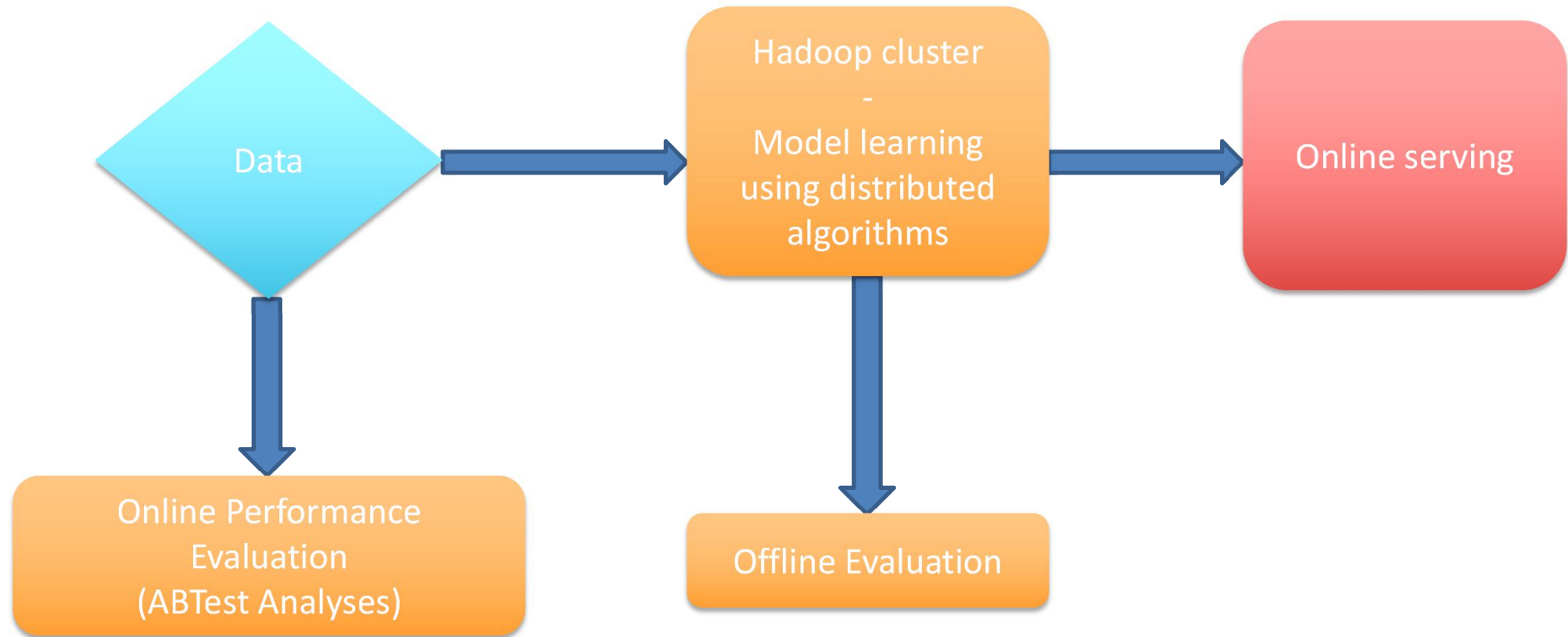


目录

CONTENTS

1. The Ads Bussiness
2. Machine Learning is everywhere
3. ML Life Cycle
4. Large Scale Machine Learning

Life Cycle



Data

- Stored on Hadoop Distributed File System
- Raw data:
 - Compressed json
 - Different data sources: displays, clicks, sales, etc..
- Refined data:
 - Produced by Hadoop jobs (cascading, scalding)
 - Combine different data types
 - Exported in Parquet (column based) to accelerate reading

Offline Evaluation

- An internal tool that replays Prod traffic from logs with different prediction models
- Target:
 - Evaluation of new models offline before going to AbTest
 - Advance investigation of production models

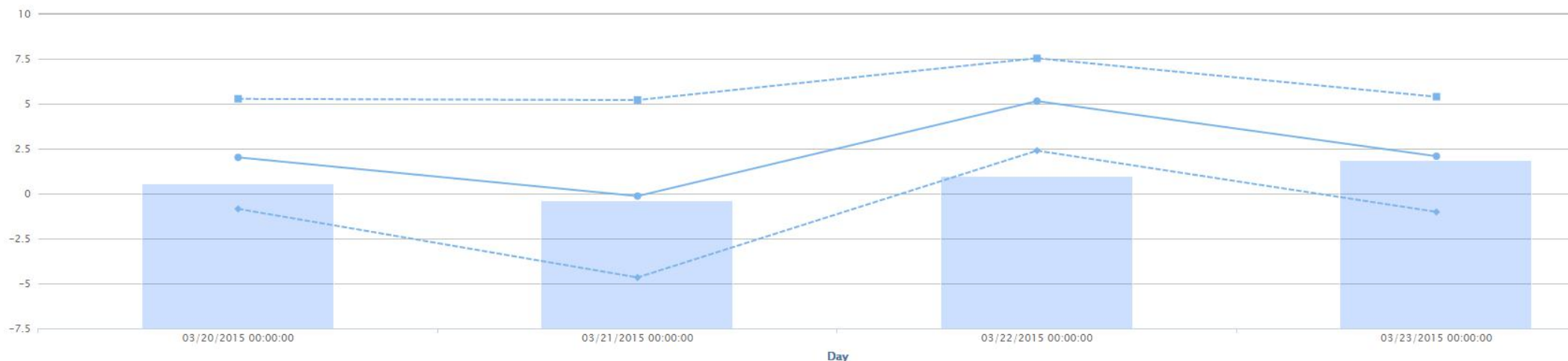


Offline Evaluation: In metrics we trust

My ABTest: +10% RevExTac on the first 2 hours
Could we trust this improvement or is it just noise?



➤ Computing confidence interval:



Offline Evaluation: Simulation

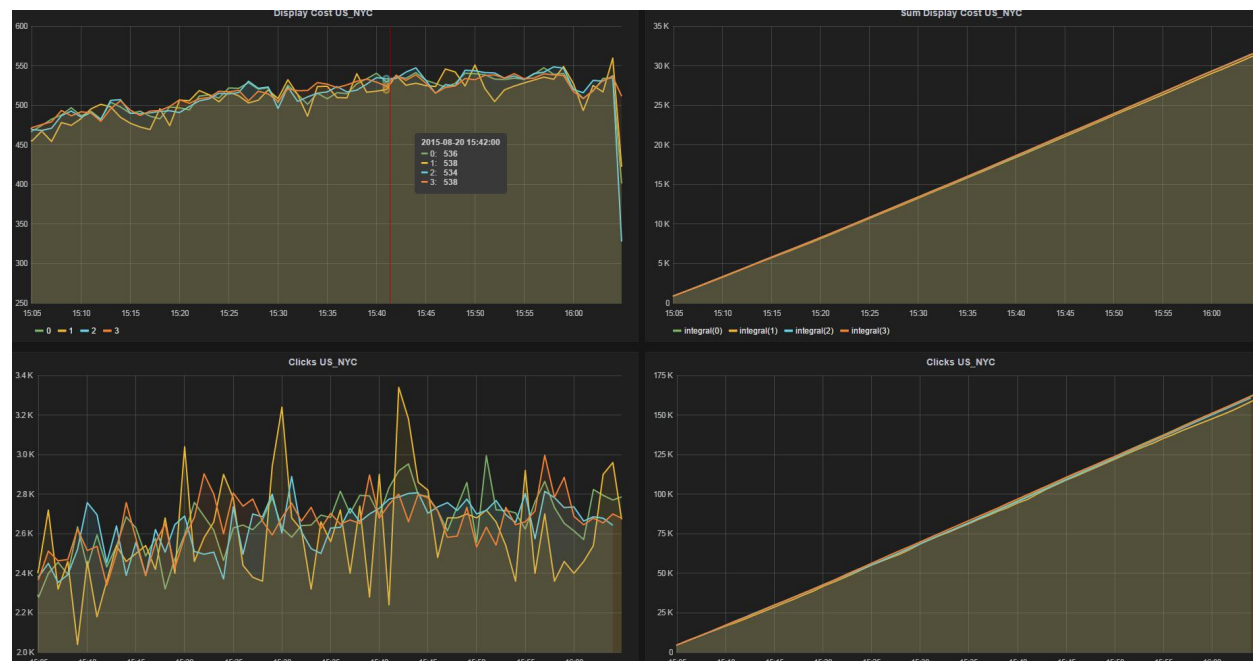
- As prediction will impact bidding, so the data we receive, simulation of acquired data should be done.
- Using counterfactual estimators
- Bid with a gaussian bias

Online Performance Evaluation (ABTest analyses)

- All Product changes are validated by an AbTest (performance is everything)
 - Realtime monitoring: to secure AbTest in realtime
 - AbTest analysis Framework: to validate AbTest with deep insight and confidential interval



AbTest analysis Framework



Realtime monitoring

Lesson learned

- Quality data is important
- Use offline tests to tune models
- Use abtest to secure and evaluate changes

目录

CONTENTS

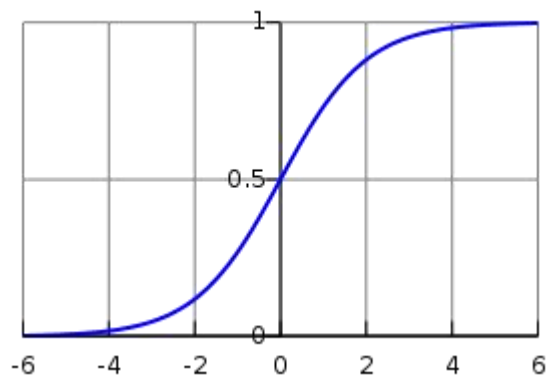
1. The Ads Bussiness
2. Machine Learning is everywhere
3. ML Life Cycle
4. Large Scale Machine Learning

Learning: Click prediction modelling

$$\mathbb{E}[NbClicks] = \mathbb{P}\{Click\} \mathbb{E}[NbClicks|Click]$$

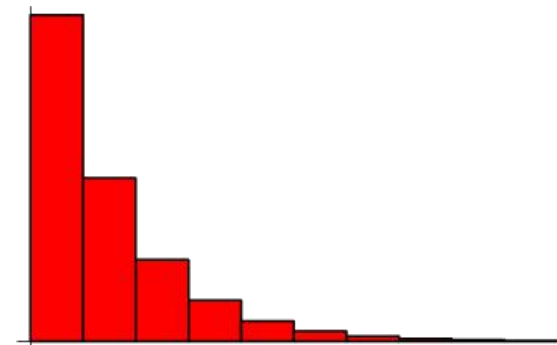
$$= \frac{1}{1 + e^{-\langle w.x \rangle}}$$

(logistic)



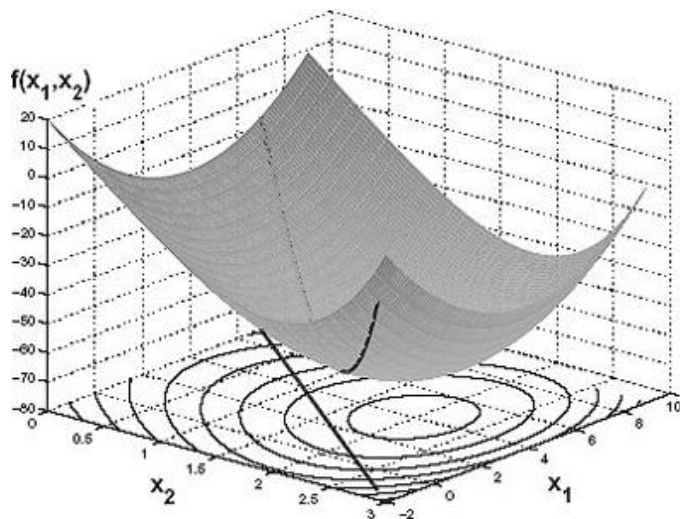
$$= 1 + e^{-\langle w.x \rangle}$$

(geometric)

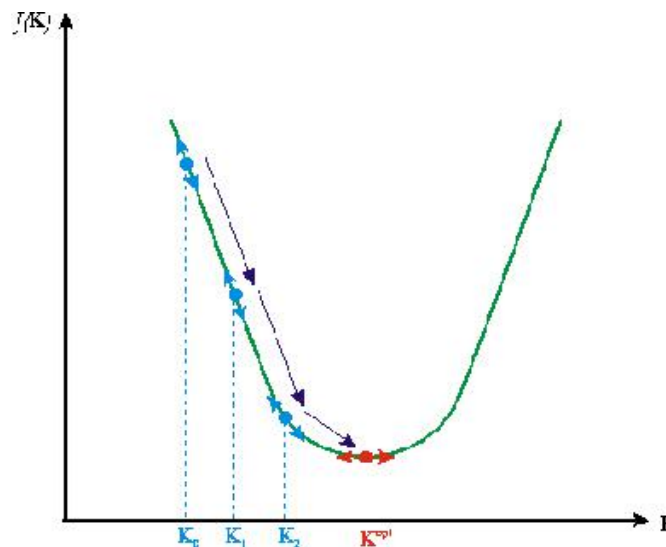


Learning: Logistic Regression

- The output model is a vector of weights : `double[]`



Convex Optimisation



Solvable with iterative Gradient Descent Algorithms (L-BFGS)



Fast Prediction at runtime

Learning: Need for scale

Learning a model (click prediction):

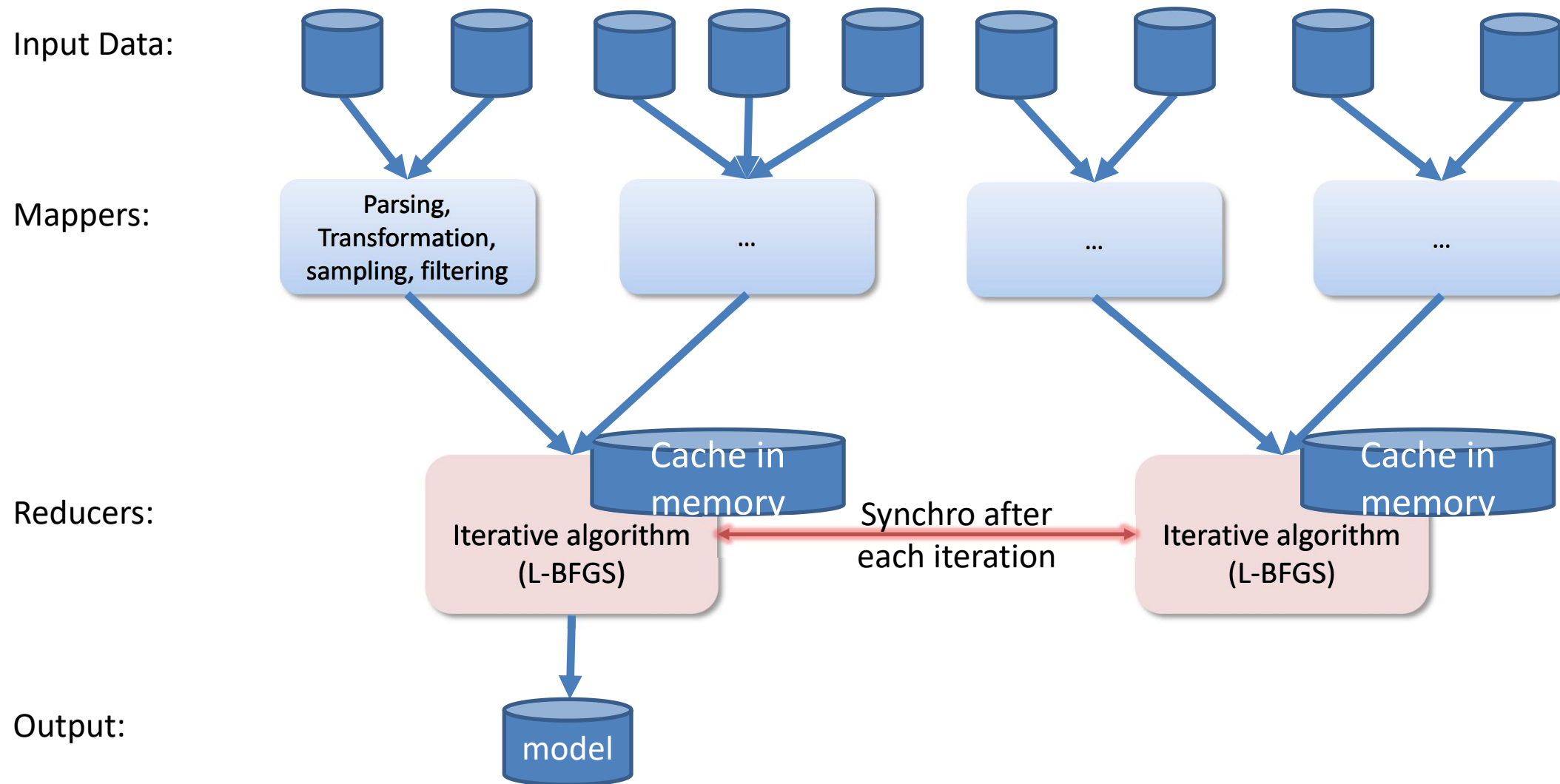
- Several days of data
- Billions of samples (after sampling)
- Millions of features

... and we have ~200 of them
(click/sales/... x DCs x ABTests)...

.. and we want to refresh as much as possible..



Learning with hadoop



Learning: Some numbers

- ~ 1300 models/day
- Ingesting 596 TB/day
- Consuming 6310 CPU day/day
- Learning time: [10min; 3h]
- Refresh rate: [3h; 6h]

From: 31-3-2015

To: 02-4-2015

DataId	Platform	Timeline																	
		31-3-2015					01-4-2015					02-4-2015							
793	EU	200478	200728	200980	300270	300404	300707	300920	301158	301377	301572	301702	301903	302138	302344	302541			
934	EU	200920	200883	300134	300400	300910	300840	301031	301270	301400	301581	301852	302040	302245	302448	302640			
935	EU	200007	200851	300107	300305	300555	300735	300927	301130	301330	301524	301702	301887	302051	302240	302420	302505		
943	EU	200572	200823	300100	300343	300525	300721	300943	301343	301558	301742	301900	302115	302207	302477	302552			
944	EU	200417	200451	200885	300131	300344	300530	300722	300913	301101	301204	301480	301655	301810	301902	302104	302334	302505	302670
999	EU	200518	200751	200990	300217	300401	300588	300700	300940	301138	301314	301500	301600	301835	302004	302165	302335	302500	302653

Learning: Lesson learned

- Balance your data
- Hash
- Tradeoff: reactive vs stable

Thanks!

Q & A