# Quantmetry

## Data Science Consulting

## Project use case

Improve your marketing reach using large scale machine learning on Spark

Matthieu Vautrot
mvautrot@quantmetry.com
Nina Bertrand
nbertrand@quantmetry.com

June, 6th

BERLIN BUZZWORDS 2016 JUNE 05-07

**Quantmetry**
Data Science Consulting

- Founded in 2011
- 25 consultants - Data Scientists & engineers

- @bertrand_nina
- @matthieuvautrot

**Share of experience of an on-going project**

- Our client, a insurance and bank …

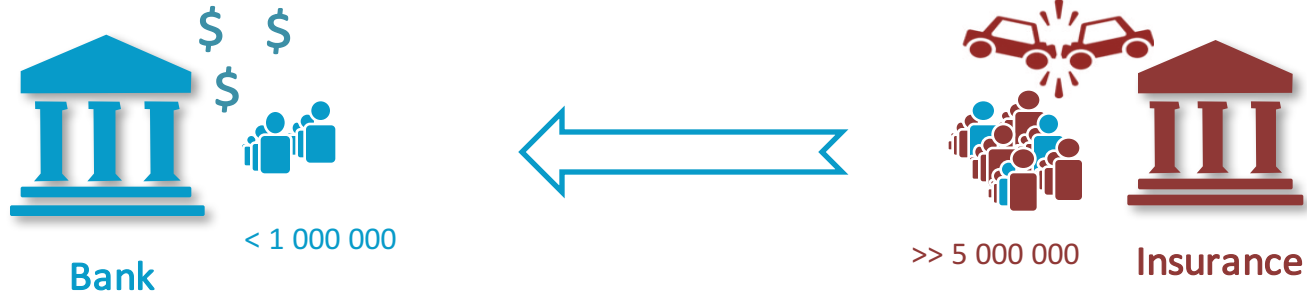- … frequently build and run Marketing Campaigns to sell their products

0,15%   0,76%

0,08%

0,83%   0,23%

- Can you do better than just "scoring" approaches with "Big Data" ?

**Bank**  
< 1 000 000

>> 5 000 000  **Insurance**

- Goal of the campaign : "Bancarize" as much Insurance clients as possible

- Can use multiple canals for it :

  - Display (web adds through DMP and / or client web site tagging)
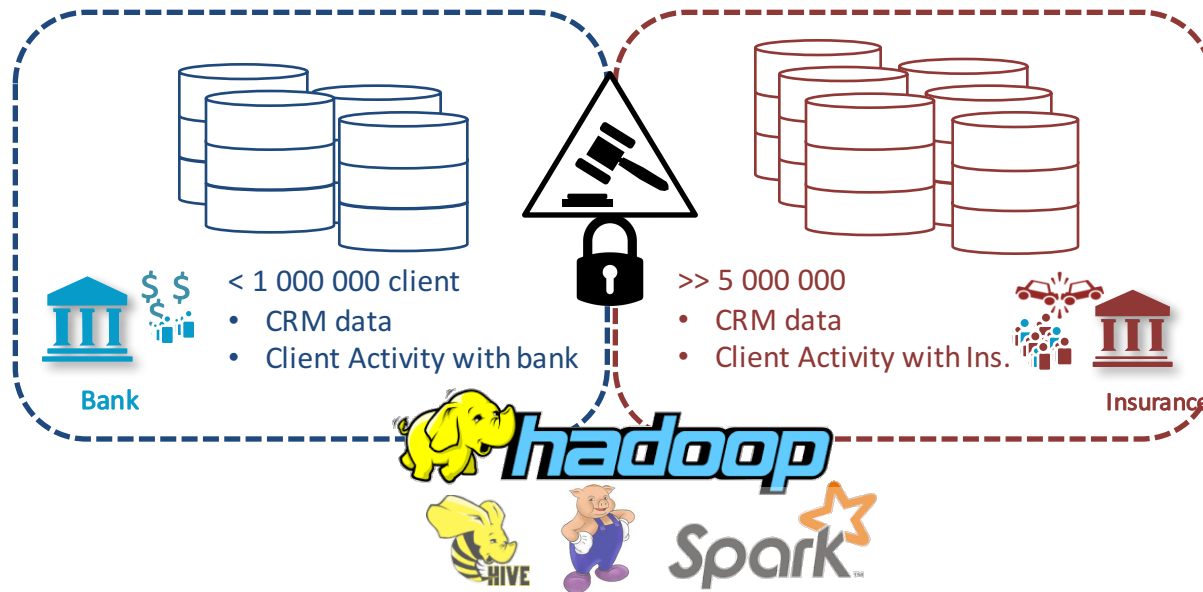
  - Email

  - Direct call

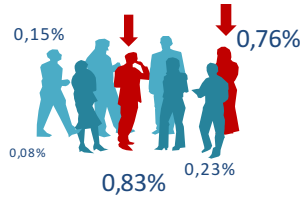- Can you do better than just "scoring" approaches with "Big Data" ?

**Bank**

< 1 000 000 client
- CRM data
- Client Activity with bank

**Insurance**

>> 5 000 000
- CRM data
- Client Activity with Ins.

- A same Hadoop Cluster, 2 distinct tenants (usual legal stuff)

0,15%

0,76%

0,08%

0,83%    0,23%

- Can you do better than just "scoring" approaches with "Big Data" ?

Sure we can !!

0,15%

0,02%

**0,76%**

0,08%

0,23%

**0,83%**

**A typical scoring campaign is generally built like that :**

- Train predictive **model** on the client base : observed buyers of the product (or similar)

- Apply the **model** and score the whole eligible client data base

- Send to the **N top score** a marketing message

0,15%

0,02%

**0,76%**
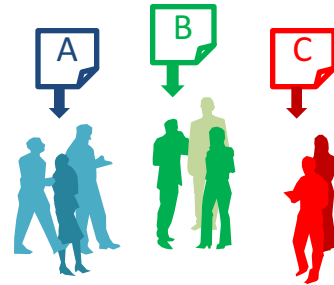
0,08%

0,23%

**0,83%**

**Two major limitations to this approach :**

- **Lack of personalisation** : Same message is sent to the top scored group

- **Scoop natural noise :** Target who would buy the product anyway

0,15%  0,02% **0,76%**

0,08%  **0,83%**  0,23%

**Simple scoring campaign**

A
B
C

**Message personalisation campaign**

**Come because you can :**

A : get money

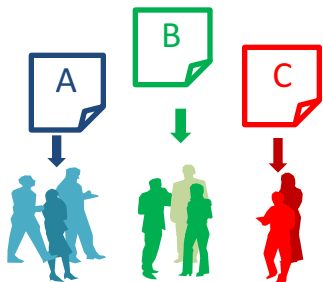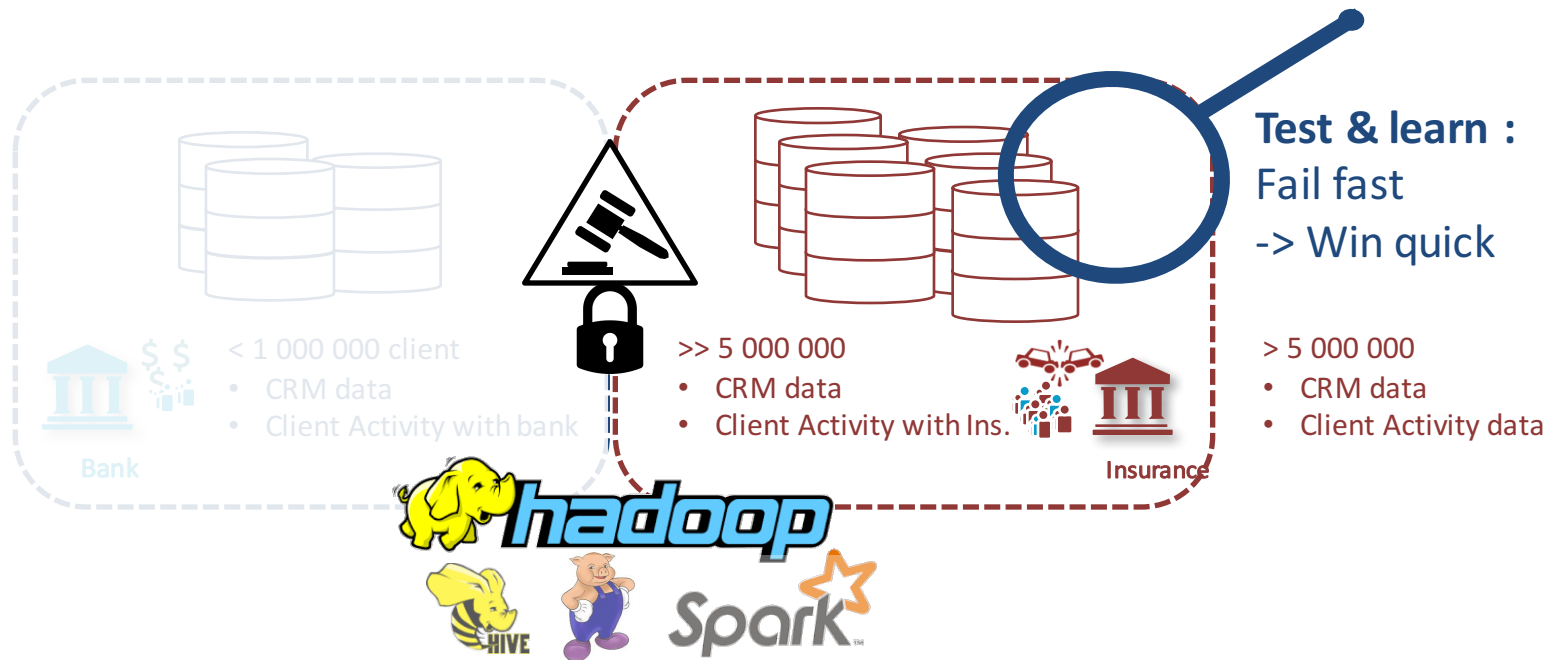B : we are better than the other Banks

C : win an ipad !!

**Two major limitations to this approach :**

- **Lack of personalisation** : Same message is sent to the top scored group

  **-> Use multiple messages and find out who likes which one with ML**

- **Scoop natural noise :** Target who would buy the product anyway

-> **Use multiple messages and find out who likes which one with ML**

**Test & learn :**
Fail fast
-> Win quick

< 1 000 000 client
• CRM data
• Client Activity with bank

**Bank**

>> 5 000 000
• CRM data
• Client Activity with Ins.

**Insurance**

> 5 000 000
• CRM data
• Client Activity data

Test & learn the different messages directly from the insurance client data base

>> 5 000 000

Insurance

**1**

**Every day :**
- List of cookies to target (DMP)

**Every couple weeks**
- List of people to contact (email / tel)

>> 5 000 000

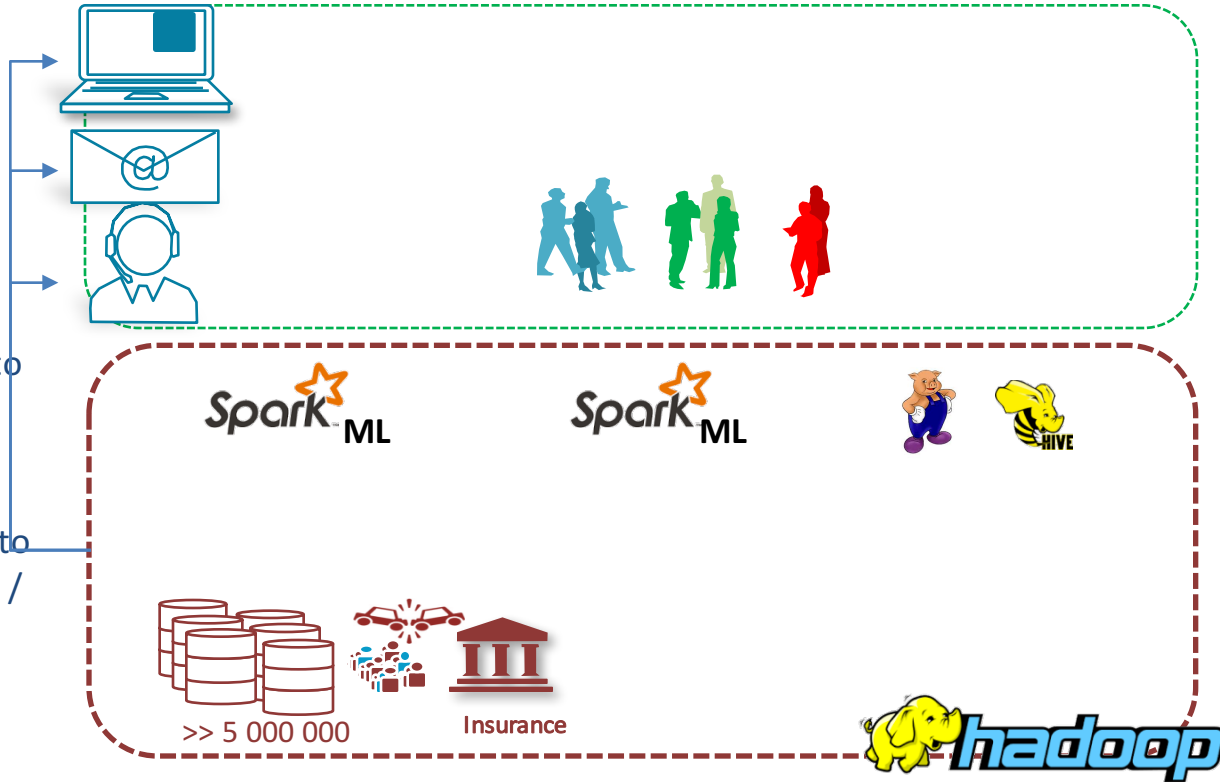Insurance

- **At day 0** : random client sample / random message

**1**

**Every day :**
- List of cookies to target (DMP)

**Every couple weeks**
- List of people to contact (email / tel)

**2**

A

B

C

Message Display to the target pop.

Spark ML

Spark ML

>> 5 000 000

Insurance

hadoop

- **At day 0** : random client sample / random message

1

**Every day :**
- List of cookies to target (DMP)

**Every couple weeks**
- List of people to contact (email / tel)

2

A

B

C

**Message Display to the target pop.**

3

**Every day :**
Who responded (or not) to which message

*Spark* ML

*Spark* ML

HIVE

>> 5 000 000

Insurance

hadoop

- **At day 0** : random client sample / random message

**2** Message Display to the target pop.

**1** **Every day :**
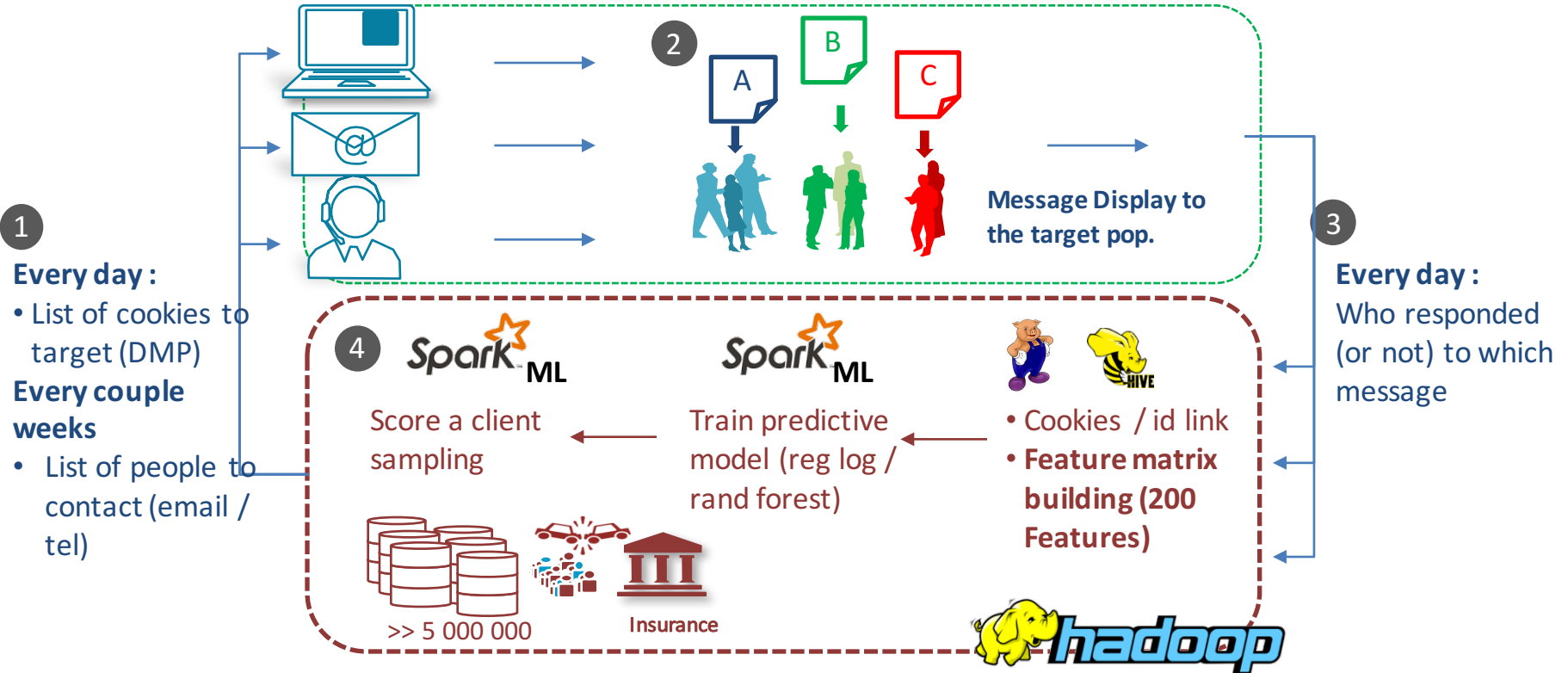- List of cookies to target (DMP)

**Every couple weeks**
- List of people to contact (email / tel)

**3** **Every day :** Who responded (or not) to which message

**4** Spark ML — Score a client sampling

Spark ML — Train predictive model (reg log / rand forest)

- Cookies / id link
- **Feature matrix building (200 Features)**

>> 5 000 000

Insurance

hadoop

- **At day 0** : random client sample / random message

# Test & learn multiple messages



**2** Message Display to the target pop.

**1**

**Every day :**
- List of cookies to target (DMP)

**Every couple weeks**
- List of people to contact (email / tel)

**3**

**Every day :**
Who responded (or not) to which message

**4** Spark ML

Score a client sampling

Spark ML

Train predictive model (reg log / rand forest)

- Cookies / id link
- **Feature matrix building (200 Features)**

>> 5 000 000

Insurance

hadoop

- **At day 0** : random client sample / random message
- **At day N+1** : use **N** days results from learning
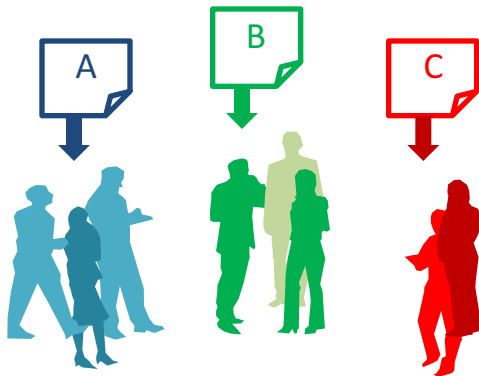
16

**Two major limitations to this approach :**

- **Lack of personalisation** : Same message is sent to the top scored group

    -> Use Multiple messages and find out who likes which one with ML

- **Scoop natural noise :** Target who would buy the product anyway

**Two major limitations to this approach :**

- **Lack of personalisation** : Same message is sent to the top scored group

  -> Use multiple messages and find out who likes which one with ML

- **Scoop natural noise :** Target who would buy the product anyway

  -> Use ML models that get rid off natural noise : uplift models

## Idea

✓ Describe **message effect** on target

## Motivation

- Do not call self-converted-people

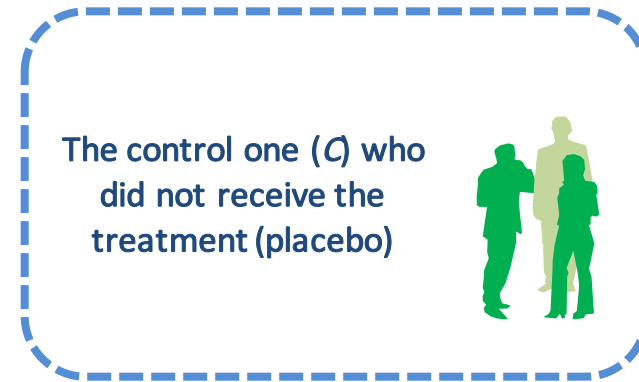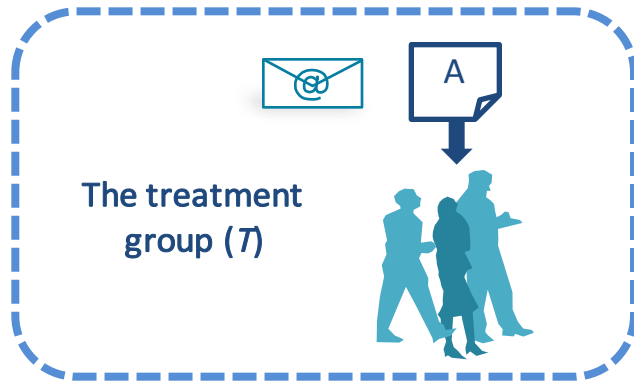- Some customers are liable to buy but marketing phone call have a negative influence on them

**Methodology :**

- Two samples should be regarded :



The treatment group ($T$)

The control one ($C$) who did not receive the treatment (placebo)

**Different possible implementations :**

- Independent models

- Regression with tuning parameters

- Sequential models

**Methodology :**

- Two samples should be regarded :



The treatment group ($T$)

The control one ($C$) who did not receive the treatment (placebo)

**Different possible implementations :**

- **Independent models**

- Regression with tuning parameters

- Sequential models

**Spark**

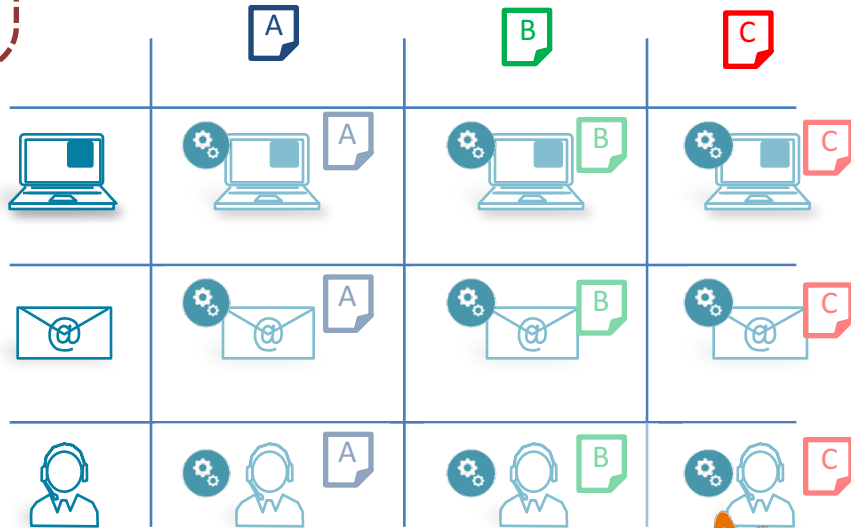Train predictive model (reg log / rand forest)

- Requirement of Independent models



Random Forest / Log reg.
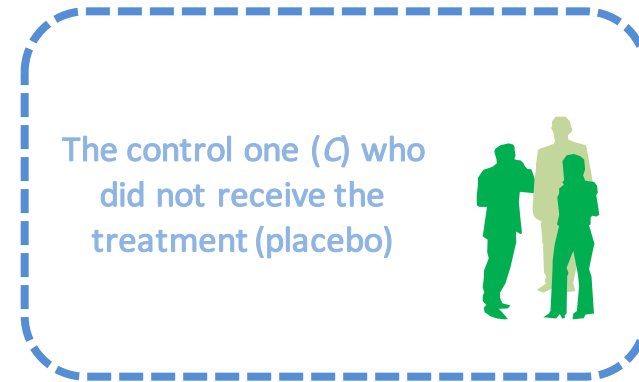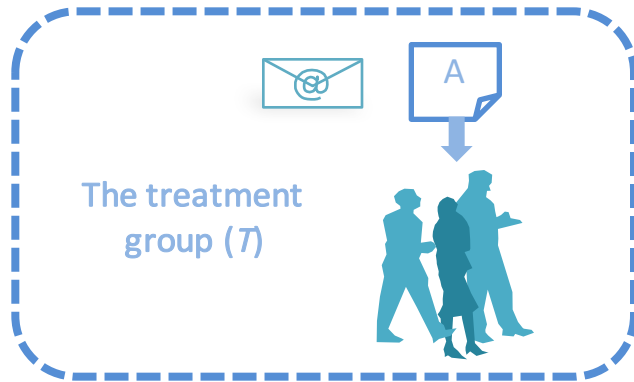Controlled by Uplift

**Spark** ML

**Every day :**

- One predictive model is calculated for every Message X Canal

- Models as usual : random forest or logistic regression

**Methodology :**

- Two samples should be regarded :

The treatment group (*T*)

The control one (*C*) who did not receive the treatment (placebo)

**Different possible implementations :**

- **Independent models:**

- Regression with tuning parameters

- Sequential models

$$Uplift(x) = P(\ S\ |\ x,\ T=1) - P(\ S\ |\ x,\ T=0)$$

*S the subscription event*

$$\text{Uplift}(x) = P( S \mid x, T=1) - P( S \mid x, T=0)$$

## Difficulty :

- There is a predicted uplift by customer but no individual real uplift → no individual target..
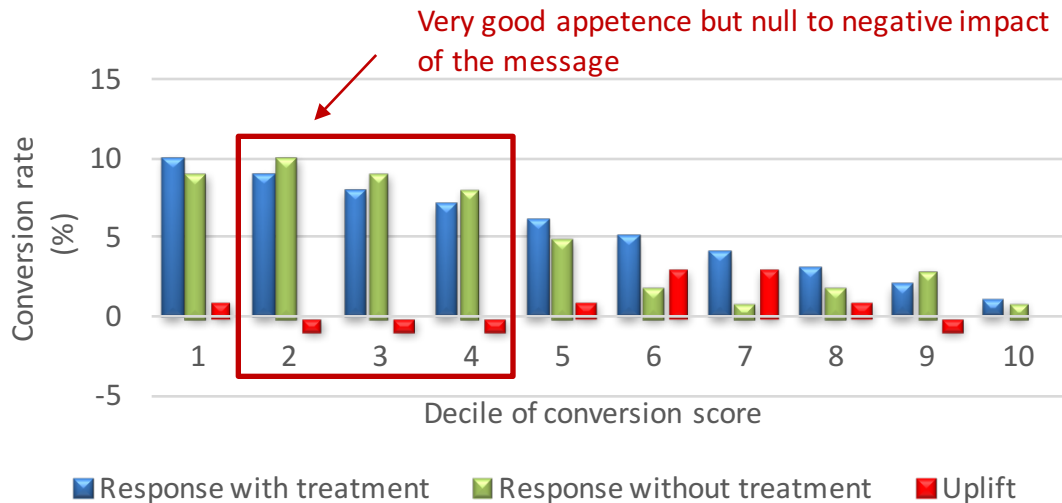
## Solution :

- Sort customers by their uplift score in decreasing order
- Focus on quantile of customers
- Calculate difference between conversion rate of treated group and natural conversion rate

## **Appetence** sorted by conversion probability

Very good appetence but null to negative impact of the message



Conversion rate (%)

Decile of conversion score

■ Response with treatment   ■ Response without treatment   ■ Uplift
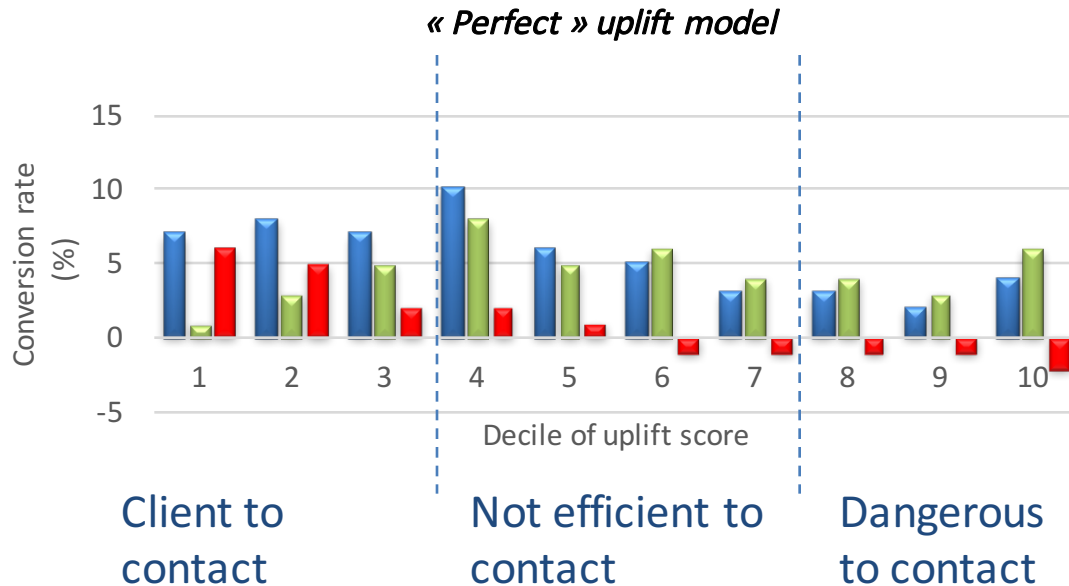
- Groups with highest conversion score has not necessarily been scored with the highest uplift.

- This people may have converted without any treatment.

## **Uplift model** sorted by predicted uplift

*« Perfect » uplift model*



Client to contact

Not efficient to contact

Dangerous to contact

- … What about the real uplift ?

- How do you assess the performance ?

Train predictive model (reg log / rand forest)

A

Come you get money $$$

B

Come we simple your life !!
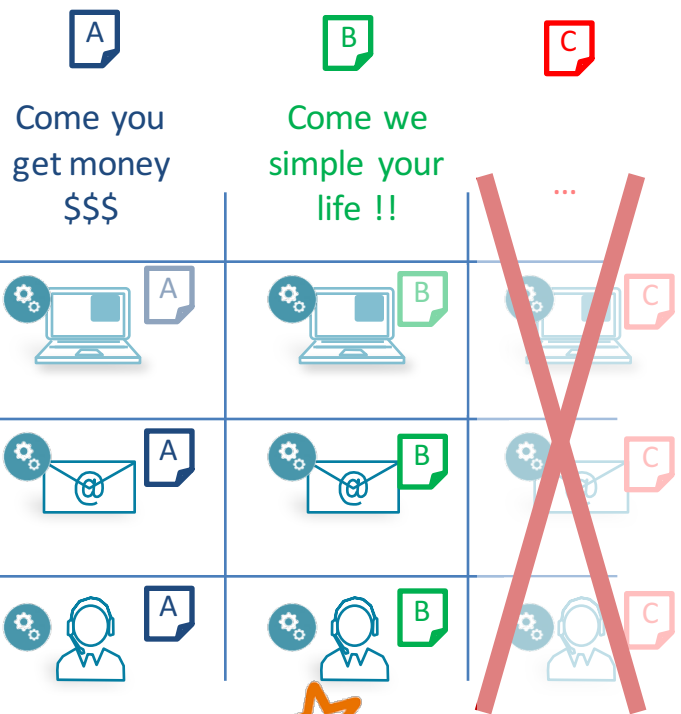
C

...

Contact canal

Spark ML

**Need POC**

- Quick agile POC iterations

- Limited to 2 messages to push

**For all 3 canals**

- Data preparation (Pig Hive)  done

- Predictive Algorithms : done

Train predictive model (reg log / rand forest)

Come you get money $$$

Come we simple your life !!

...

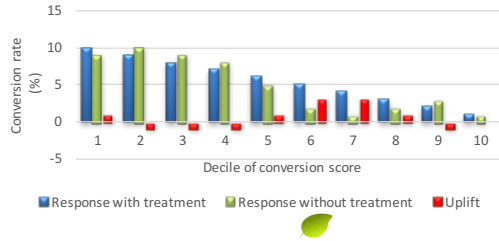Contact canal

Spark ML

**Need POC**

- Quick agile POC iterations

- Limited to 2 messages to push

**For all 3 canals**

- Data preparation (Pig Hive) done

- Predictive Algorithms : done

- 2 waves already achieved in mail and tel
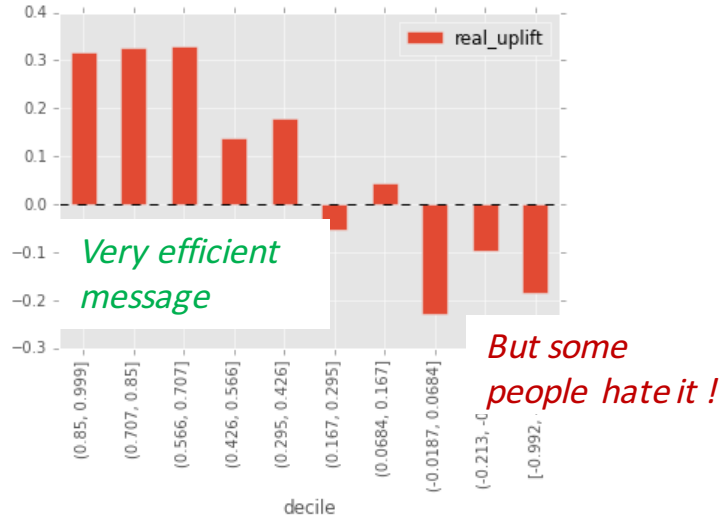
- DMP results analysis is on going

# Uplift model : improve an effect treatment
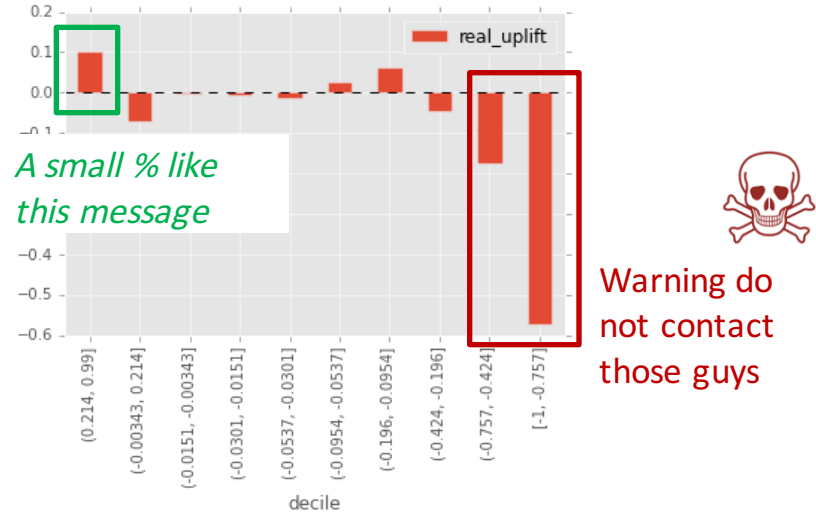## Use case observed uplift and marketing insights

**Observed uplift** : for mail canal after 1rst wave

**1st message : « Come you get money $$$ »**

*Very efficient message*

*But some people hate it !*

**2nd message : « Come we simple your life »**

*A small % like this message*

Warning do not contact those guys

**We just have to take best score between the 2 models**
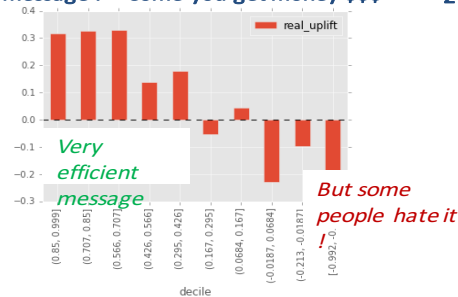
## Data engineering the Marketing campaign

- Easy on paper but watch out to business and IT organization constraints (eg : DMP and Hadoop Cluster not easly linkable)

- Spark is good but sometimes Scikit learn can do the trick for first quicker ML iteration
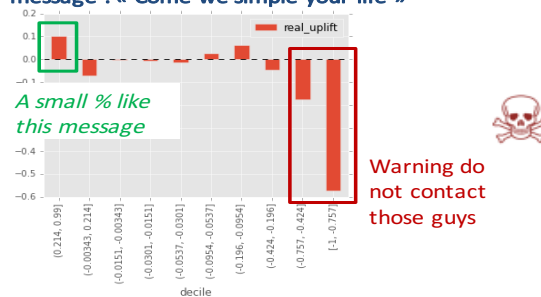
## Uplift modeling



POC → Spark Indus. ML

- Very efficient for marketing insight already on first waves -> Promising for the following up of the project !



1st message : « Come you get money $$$ »

Very efficient message

But some people hate it !

2nd message : « Come we simple your life »

A small % like this message

Warning do not contact those guys

# Quantmetry
## Data Science Consulting

BERLIN
BUZZWORDS
2016 JUNE 05-07

## Q & A ?

**Thank you**

**Nina Bertrand :** nbertrand@quantmetry.com

**Matthieu Vautrot :** mvautrot@quantmetry.com