amadeus

SentimentTM Project





amadeus

The Team



Minh Phan Intern / DTI



Alejandro Mottini Data Scientist / R&D



Eoin ThomasData Scientist / DTI

2016 Amadeus S.A.S

Agenda

_ Project Introduction	-
Twitter API, Tweets Data	
_ Project Pipeline	
_ Classification Models	4
Product Prototype	!
Wrap-up	(



2016 Amadeus S.A.

Project Introduction

Travelers' sentiment?

- My direct to London Bridge is cancelled... going via Blackfriars ... if you don't see me by 2:40 I got lost :/
- Think the smell of microwaved fish is awful? The bus I'm on smells like that, first thing in the morning... #puke #tfl #London
- London st pancras train station is the most confusing place ever. Got lost about 4 times trying to find the way out



- Thing im most excited for in London? To walk across Abbey road
- Hello London !! It's been a while. Looking forward to playing at the Acklam Village Market at 5pm :) @Hot_Vox
- _ @LondonAsterino You are a beautiful Queen,London You sing like an angel You are awesome I will always believe in you I love you

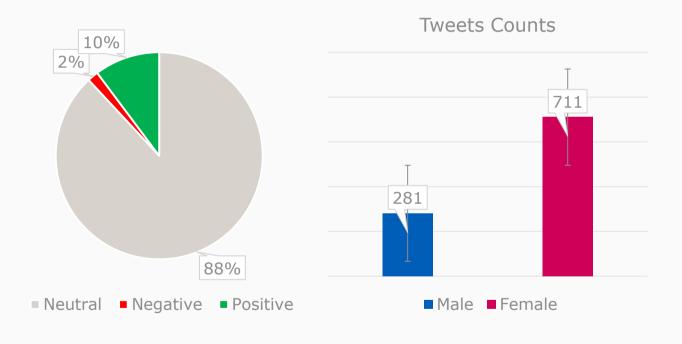
○ 2016 Amadeus S.A

Project Introduction

Tourism locations analysis?



Big Ben, London



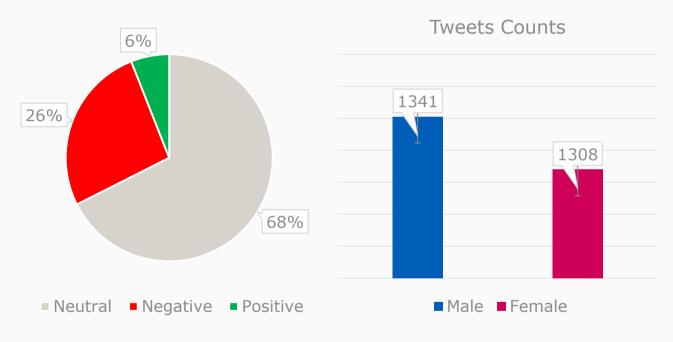
2016 Amadaile S.A

Project Introduction

Tourism locations analysis?



Royal Borough of Greenwich, London



© 2016 Amadeus S.A.

Project Introduction

Targets of the project

- Explore Twitter data (API, topics, reliability...)
- Predict and classify Age, Gender and Sentiment (feeling) of travelers specific to a destination based on the information extracted from Twitter
- Build a Business Intelligence prototype (dashboard) to monitor Twitter activity related to specific destination.





© 2016 Amadeus S.A

Twitter API, Tweets Data

Twitter API and Python package tweepy

_ Twitter API

- Free developer API vs. paid enterprise API (GNIP, https://gnip.com/)
- REST API vs. Streaming API
- REST API rate limit (Ref: https://dev.twitter.com/rest/public/rate-limits)
 - GET search/tweets: 450 requests / 15-min
 - GET users/lookup: 300 requests / 15-min
 - GET statuses/user_timeline: 1500 request / 15-min (max 3,200 tweets)
 - 200k tweets / 8 hours
 - 7-10 days in the past (not officially confirmed by Twitter)

Python package: tweepy v3.5.0

- Twitter account → Twitter app + key → Authorization → Query
- Automatically pause and wait when reaching request limit
- Support REST and Streaming APIs, return JSON format





Twitter API, Tweets Data

Twitter data

Data including

- Tweet (text, 140 characters)
- Profile picture (image)
- Some user information (text, location, selfdescription, etc.)
- Account information (number of followers, following, etc.)

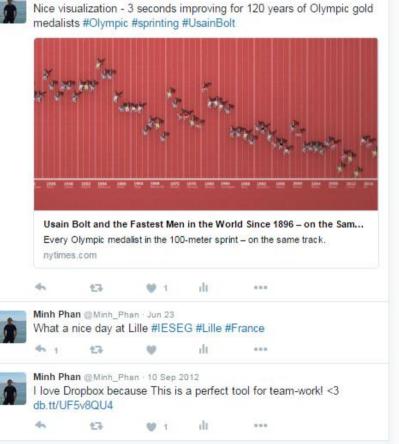
Pros

- Tons of free data
- Can access other accounts, no need to be their friends

_ Cons

- No gender, no birthday, no country
- Location is free text, some are unreal
- Spamming accounts







© 2016 Amadeus S.A.

Twitter API, Tweets Data

Twitter data, some statistics

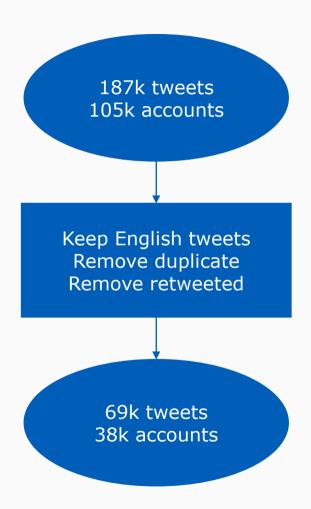
Tweets related to London (187k)

- Tweet language: en (94%), es (0.4%), de (0.4%), fr (0.3%)
- Retweeted: 53%
- Geo-location information: 3.1%
- Useful tweets: 37% (69k)

_ Accounts (105k)

- Account language: en (84%), es (4%), de (1%), fr (1.6%)
- Having user description: 12%
- User profile image changed: 98%
- Useful accounts: 36% (38k)

→ Losing 63% of tweets





Twitter API, Tweets Data

Twitter data, topics clustering

Vectorize methods

- TF-IDF
- Word2Vec (Google word2vec, Stanford GloVe)

Clustering methods

- K-means
- IDA

_ Evaluating methods

- Inertia (sum of distances)
- Silhouette score
- Gap statistic

→ Clusters are not stable

```
first time near work wharf lane centre palace people wall river hyde going home view street so ho south thames "so ho beautifulplease love east_end city #a good hammers mith covent tower_bridge england back eye west_end city #a greater greenwich place would station university on green lobkensing to not anary east angel take old man warfeld and market united day road jit's way see get work was people wall river hyde east_end city #a greater greenwich place eve west_end city #a greater greenwich place would station university on today hotel endthanks new #jobs#u day road jit's way around best mayfair cross north night
```



house still please it's eyelive one love wouldbest like exhibition irenanextexceltimeocttoday national lookingnov brixton show see take amazing november stadiumhyde_parkivecomehistory trafalgar_square listory_museumtrafalgar museumnew bridgeget theatre know natural_history wer chelsea square big weekend tower last first goodfestivalyear game academy videofashion look

guy security debate guy_filmed news_multiple

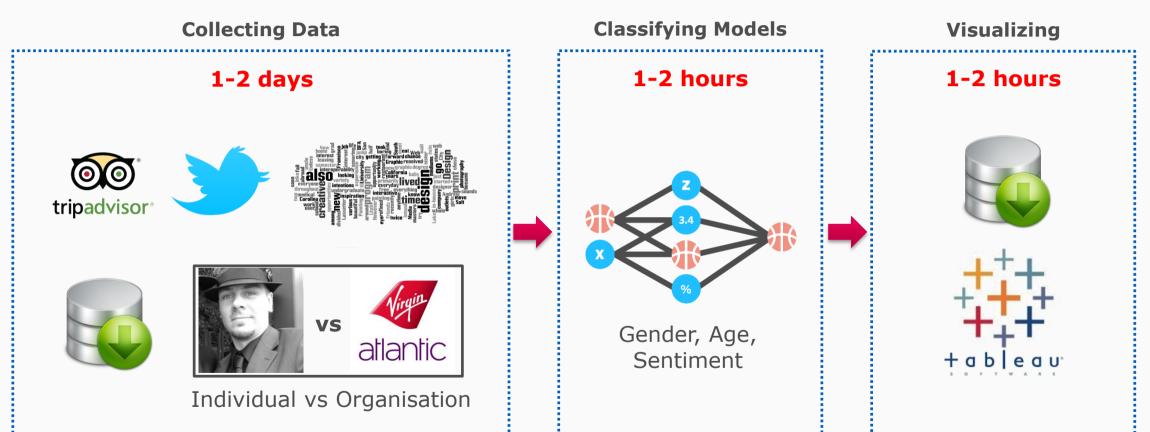
detectives believe_lone modern sure_couldamerican_politics
five_people_taken_hospitalgeek's_paradise surpsons-themed five_people_taken_hospitalgeek's_paradises_people_taken_topitalgeek's_paradises_people_taken_topitalgeek's_paradises_people_taken_topitalgeek's_paradises_people_taken_topitalgeek's_paradises_people_taken_topitalgeek's_paradises_people_taken_topitalgeek's_paradises_topitalgee



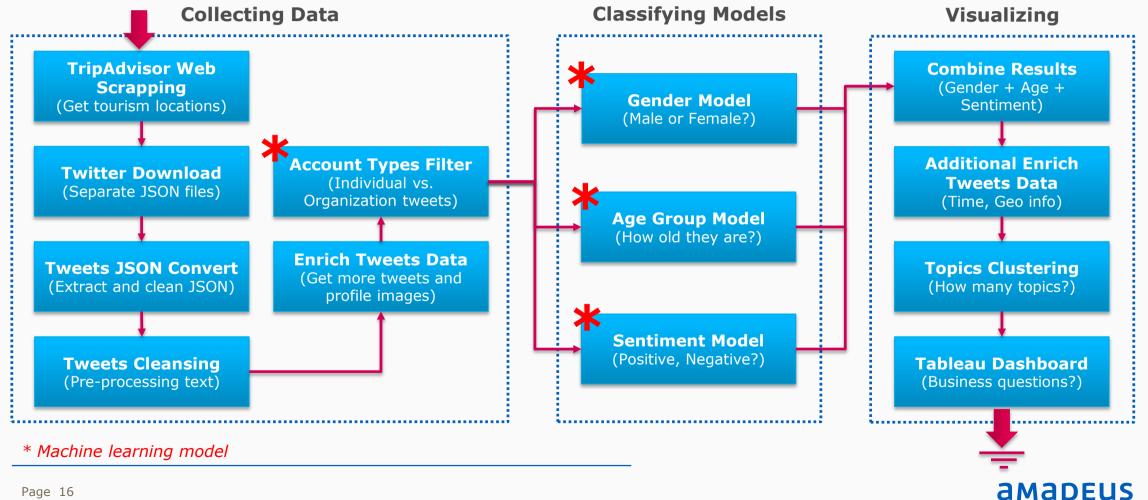
3_____ Project Pipeline

© 2016 Amadeus S A

Project Pipeline



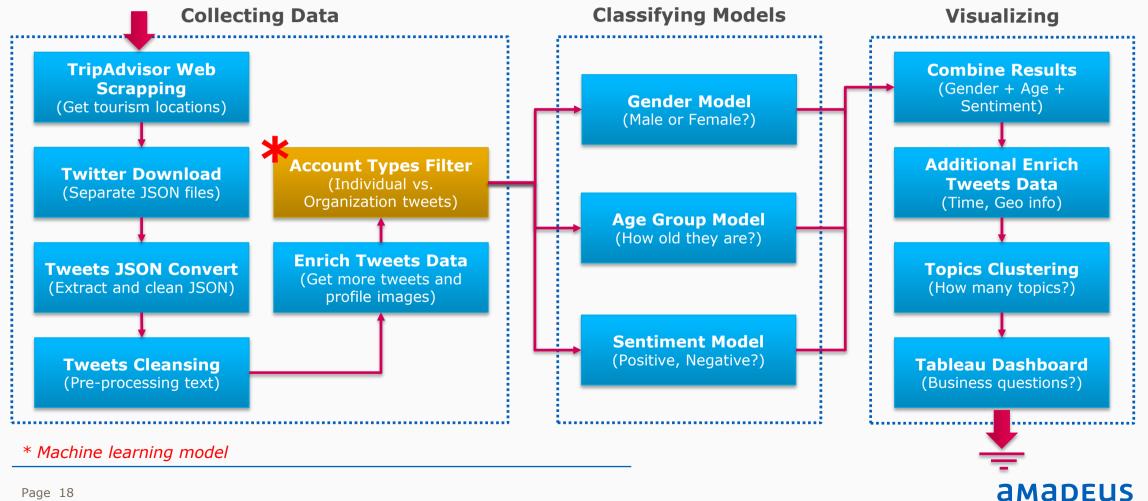
Project Pipeline







Project Pipeline



2016 Amadeus S.A.

Classification Models

Model 1: Account types classification model

Training data

 Uniform sampling Twitter accounts, hand labelling 2,783 accounts (i.e. individual, organization, unknown)

Features

- Image features: (43 features)
 - Texture
 - Luminance (contrast)
 - HSV (hue, saturation, brightness, pleasure, arousal, dominance
 - Itten's features
- Text features: (vector 300 dimension)
 - Google Word2Vec



Name: VirginAtlantic / Virgin Atlantic / 20626359

Location: Worldwide

Description: Hello Gorgeous! Follow us for news, banter & assistance 24/7. Visit our blog at https://t.co/ziOX3LRnnS or for official concerns visit https://t.co/ZQImKUNzyR

Tweet: A3: Secret Food Tours are a fun way to see London, just ask @cyneats @CerealKillerUK #ExpediaChat https://t.co/1Pnzy2L0tH



Name: JimWrongUn / Jim / 769480166338854912

Location:

Description: Scientist incognito. The left think I'm right; the right think I'm left. I'm stuck in the middle with you lot...

Tweet: @TappnRay @LBC Oxford street is full of hare khrishnas #NottingHillCarnival #london



1.00

1.00

© 2016 Amadeus S.A.9

Classification Models

Model 1: Account types classification model

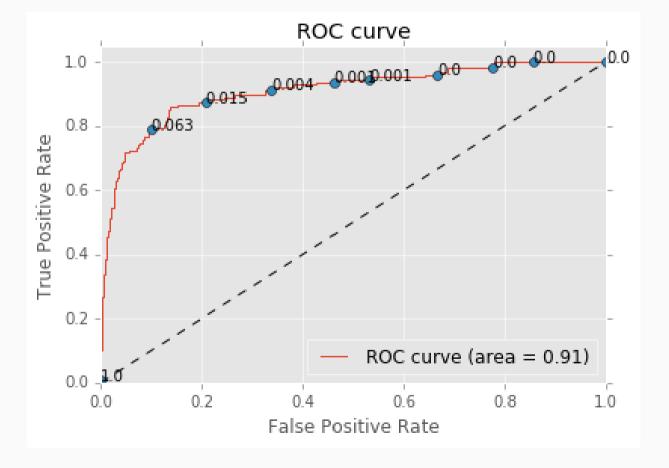
ML model

XGBoost + GridSearchCV

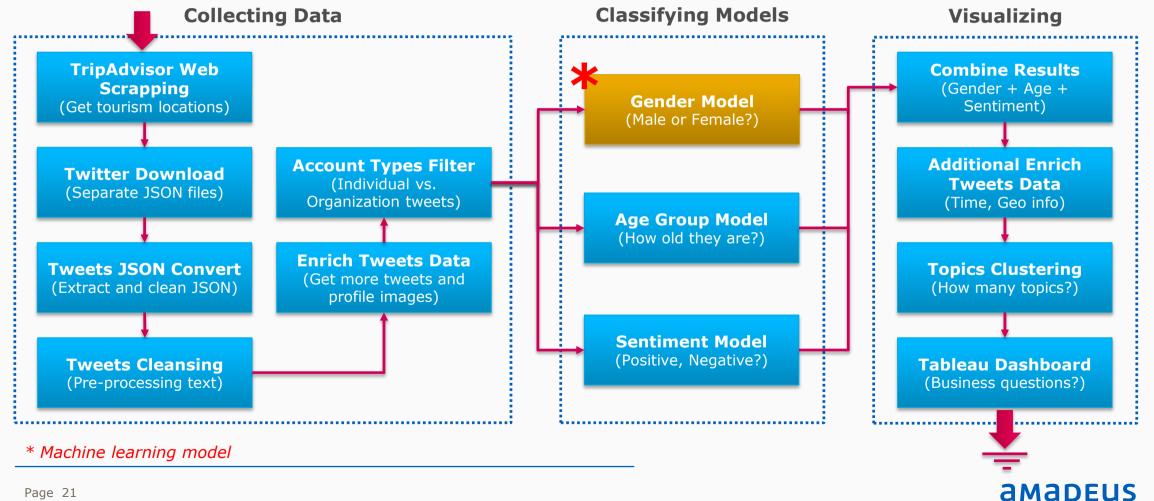
Evaluation

• AUC: 0.91

F1-score: 0.89Accuracy: 90%



Project Pipeline



© 2016 Amadeus S.A.S.

Classification Models

Model 2: Gender classification model

_ Training data

- List of Male and Female names
- Look up on Twitter: 1,901 accounts

Features

- Gender lexicon, Penn State University
- Text features:
 - Is retweeted
 - Word length
 - · Count users replies, links, hashtags, emoticons
 - Count all-capital-words, e.g. AMADEUS
 - Count long-end words: e.g. greatttttt
 - Count number of Pronouns, e.g. I, he, she...

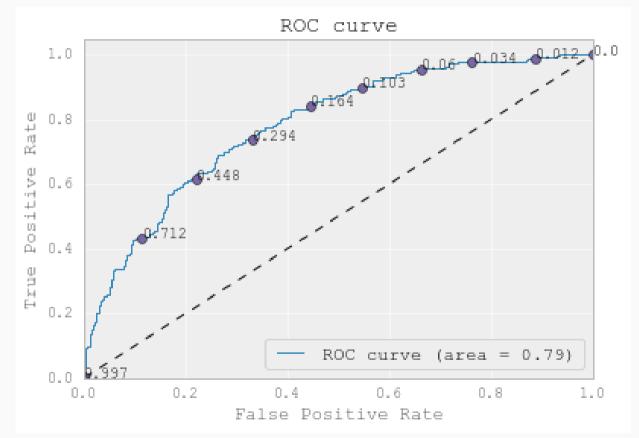
ML model

XGBoost + GridSearchCV

Evaluation

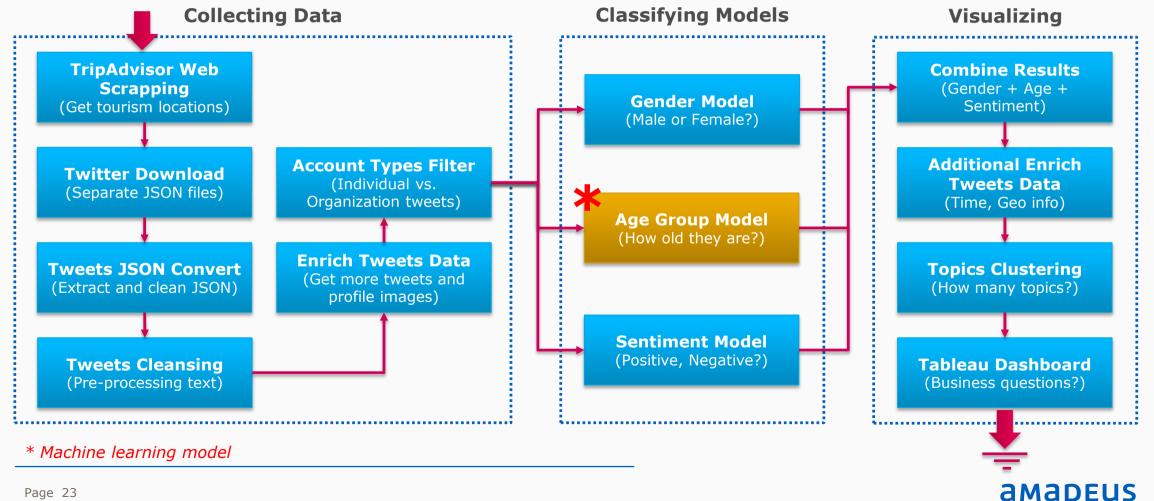
• AUC: 0.79

F1-score: 0.71Accuracy: 72%





Project Pipeline



2016 Amadeus S.A.S

Classification Models

Model 3: Age group classification model

_ Training data

- Look up the term "Happy birthday [age] to me"
- 1,477 Twitter accounts

Features

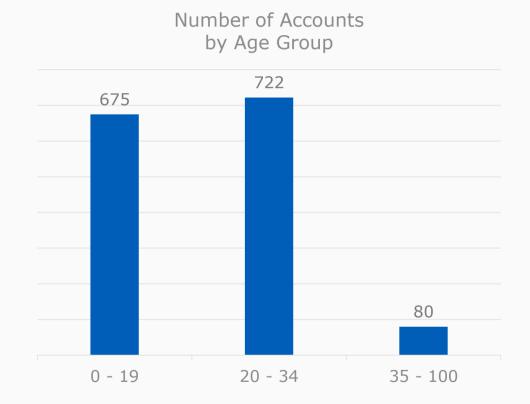
- · Age lexicon, Penn University
- Text features:
 - Is retweeted
 - Word length
 - · Count users replies, links, hashtags, emoticons
 - Count all-capital-words, e.g. AMADEUS
 - Count long-end words: e.g. greatttttt
 - Count number of Pronouns, e.g. I, he, she...

ML model

XGBoost + GridSearchCV

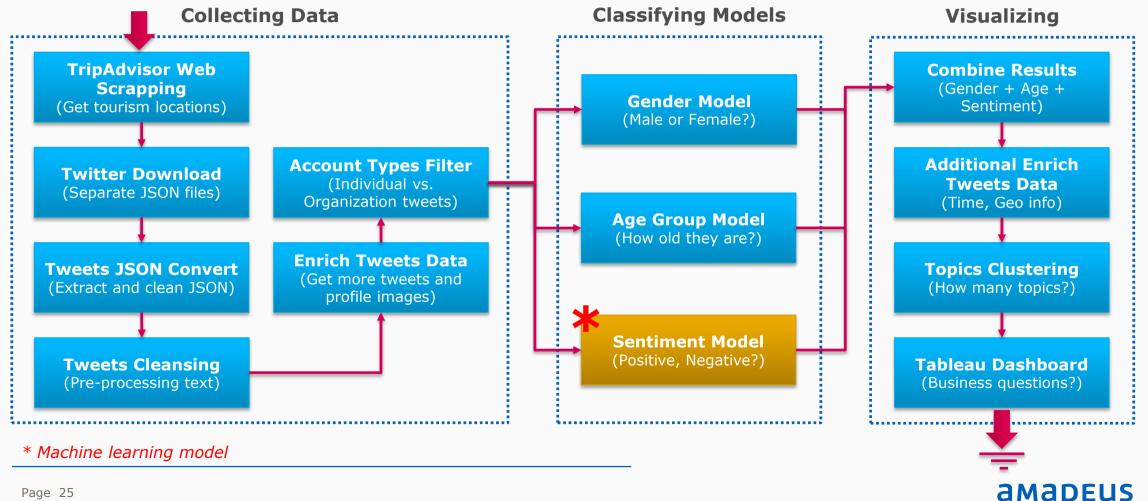
Evaluation

F1-score: 0.63Accuracy: 62%





Project Pipeline



2016 Amadeus S.A.

Classification Models

Model 4: Sentiment classification model

Training data

- SemEval Workshop 2013 (5,670 train, 851 test)
- SESAMm sentiment dataset (6,588 train, 1,044 test)

Features

- Lexicons:
 - NRC Hashtag Sentiment Lexicon
 - Sentiment140 Context Lexicon
 - MPQA Lexicon
 - NRC Word-Emotion Association Lexicon
 - Opinion Lexicon English
- Counts:
 - Long-end words, e.g. soooooo
 - Continuous exclamation marks, e.g. !!!???
 - Number of hashtags
 - Upper-case words, e.g. AMADEUS
- Word2Vec:
 - Twitter Word Clusters Lexicon (CMU)
 - Google Word2Vec model



Ironic or Neutral

"The view is bad but the food is great."

"Excellent! This day couldn't start off any better!"

"See London's reflection shimmers on its dark waters! Paddle down river to see the Houses of Parliament in all..."



© 2016 Amadeus S.A

Classification Models

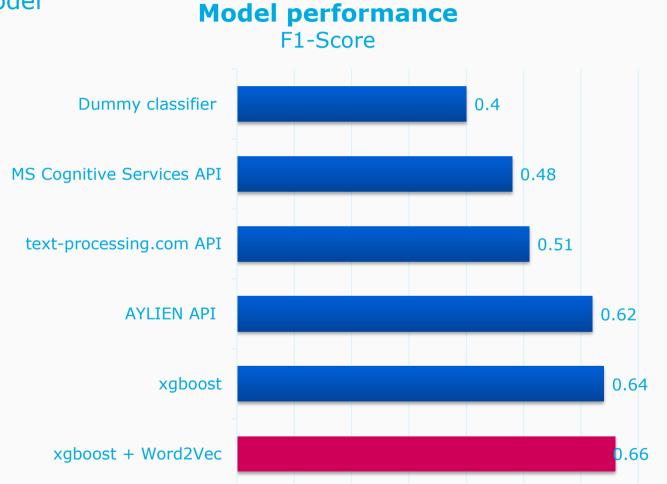
Model 4: Sentiment classification model

ML model

XGBoost + GridSearchCV

Evaluation

- F1-score 0.66
- Accuracy: 67%

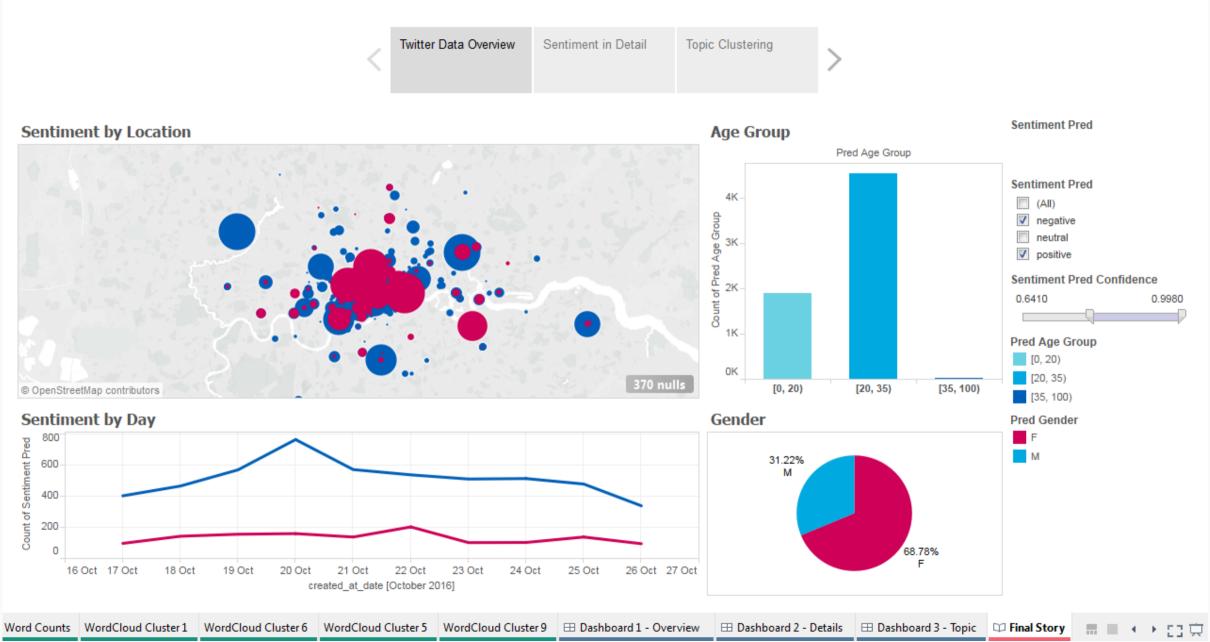




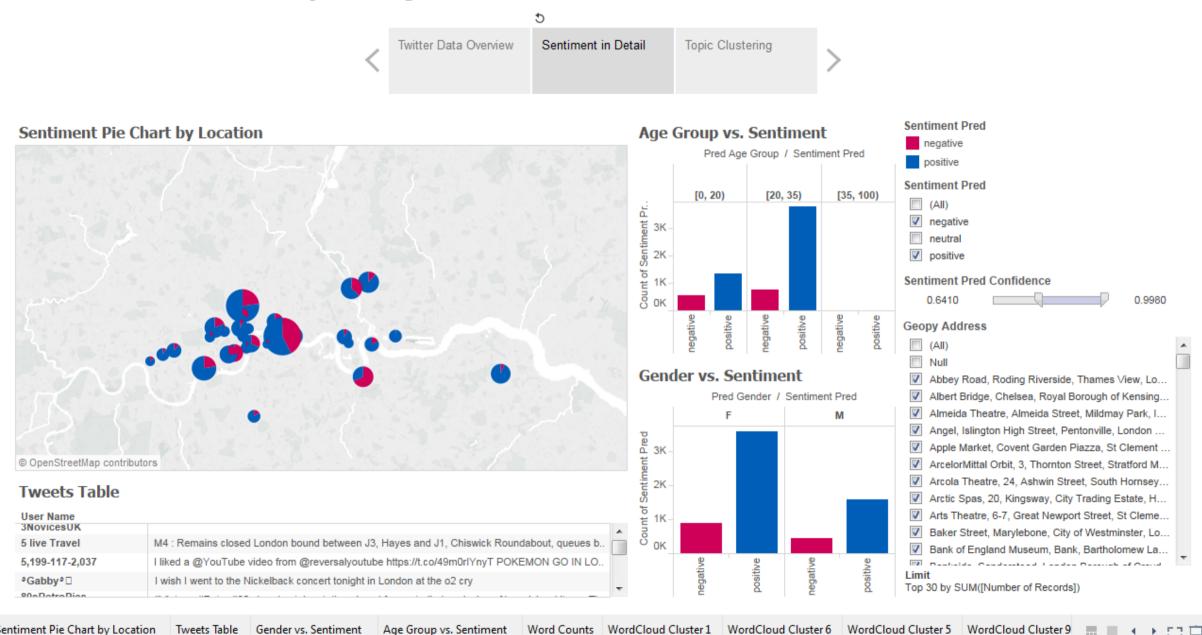
5

Product Prototype

Twitter Sentiment Analysis Project



Twitter Sentiment Analysis Project

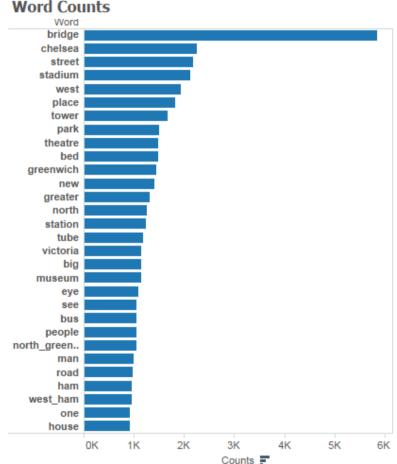


Twitter Sentiment Analysis Project

Twitter Data Overview

Sentiment in Detail

Topic Clustering



WordCloud Cluster 1

take roadriver centreend like back cross view #chelsea first house east_end theatre eye beautiful see station soho one city please wharftower_bridge garden greenwichbigtoday greater near south great west_endpark place eaststadiumtower university #jobs chelsea westminstergood ill school would hammersmith victoria morning people West wallcovent_garden_time .canary_wharflooking workaround palace manager #soho hotel man day central live covent market loveway bank best green

WordCloud Cluster 6

would house onegallery palladium lastgreat back love sunday tower new show tomorrowgoing eye tonight arena exhibition IVE britishchelseawntickets @youtube trafalgar_{year}time hyde parkstadium place novemberplease museum come bridge see theatre victoria night trafalgar_square history visitsquare natural_history get national today natural lookingcan't weekstill october fashion first day next park brixton home big view royal falling festival history visitsquare natural_history get national lookingcan't weekstill october amazing weekend olympia olympia olympia want event good people want

WordCloud Cluster 5

close avenue bed_terraced bath_rec canary_wharf se10 lohn's fees_applyhampstead se1chlswick recwharf_e14 class_escorts west_hampstead camberwell terraced bed_bath westminsterhigh_street nw6 rec wharf_e14 class_escorts terraced_house_kensington_high_canary_road_bethnal oxford_street studio wimbledononthemarket road #hammersmith wharf zoopta_camberweii #hammersmith_bed court house_rightmove bed_apartmentescorts apartment_rightmove flat_onthemarket bed flat_zoopla rightmove bedroom_flat gardens bed_flat bed_studio Zoopla apartment kensington bedroom lane flat_rightmovebridge high_classapply park_lane oxford south hill

WordCloud Cluster 9

could_fivetweet_american_cat_ears security_newsinjured_enhanced_railway
tate_could_stick_guaranteed_totage_remember_costume_geek's_footage_paradise_sure
tate_five_people_roof_tate_simpsons-themed_tweet_politics_injuring_crashes_geek's_paradise_ hospitaltottenhambus_hits people nail-biting heading_westminster_sure stick gopro_footageguy double-decker left multiple_injuries Inside_welrd_nail-biting_gopro modern_nail-bitingsexy_cattottenham_injuring hits_bridge_bridge_tottenham gopro injuring_people bridgebus man_climbs double-decker_bussexy taken injuries_double-decker enhanced_security manclimbs_roof injuries cheer_debate guaranteed_cheer tate_modern_debate_americancostume.simpsons-themed people_five paradise hospital_injured _inside_wonderful_costumecrashes_rallway_rallway_bridge _remember_guy_debate sure_could modernhits tweet_guaranteedremember_taken_hospital_bus_crashes_climbs wonderful people_taken roof multiple might believe_lone crude welrd oeek's filmed



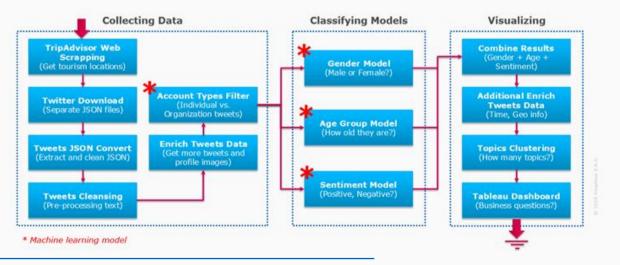


© 2016 Amadeus S.A

Wrap-up

Things we did

- __ Built a complete pipeline for the project, easy to add or improve modules in the future
- Built an Account Types classification model, AUC=0.91, Accuracy=90%
- Built a Gender classification model, AUC=0.79, Accuracy=72%
- Built a Age Group classification model, F1-Score=0.63, Accuracy=62% (21% better than a dummy model)
- Built a Sentiment Classification model, F1-Score=0.66, Accuracy=67%, better than commercial APIs (e.g. Microsoft Cognitive API, AYLIEN API), free API (e.g. text-processing.com API)





Wrap-up

Things to improve

- More data to improve the 4 models
 - More ways to collect data, e.g. "London":)
 - Buy data from Twitter enterprise API (GNIP)
 - Other training data sources, e.g. Sentiment140 http://help.sentiment140.com/for-students/
- Upgrade Machine Learning model
 - Deep learning model to improve Age and Gender classification
 - Ensemble 16 models can increase Sentiment classification model's F1-Score up to 0.70
- Add company/user history database to save time and avoid duplicated work
- _ Tweak decision threshold to fit with the local data distribution, e.g. percentage of organization vs. individual accounts is different by city



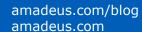


#ThankYou! and #MerryChristmas!



You can follow us on: AmadeusITGroup





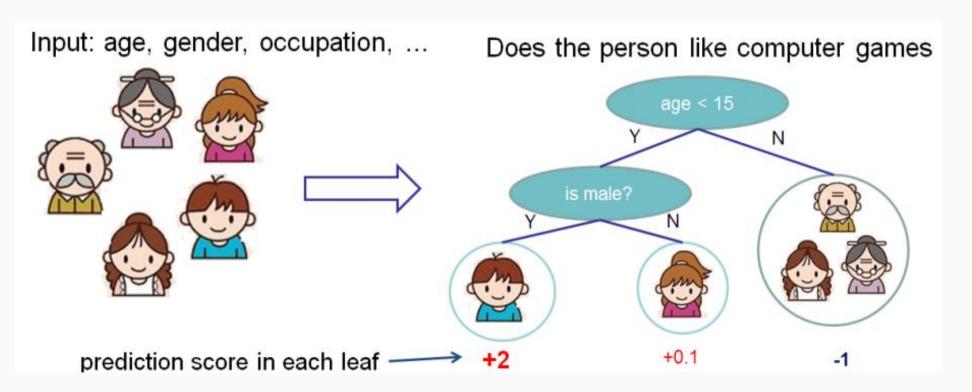


2016 Amadeus S./

Appendix 1

XGBoost - Extreme Gradient Boosting

Tree Ensemble



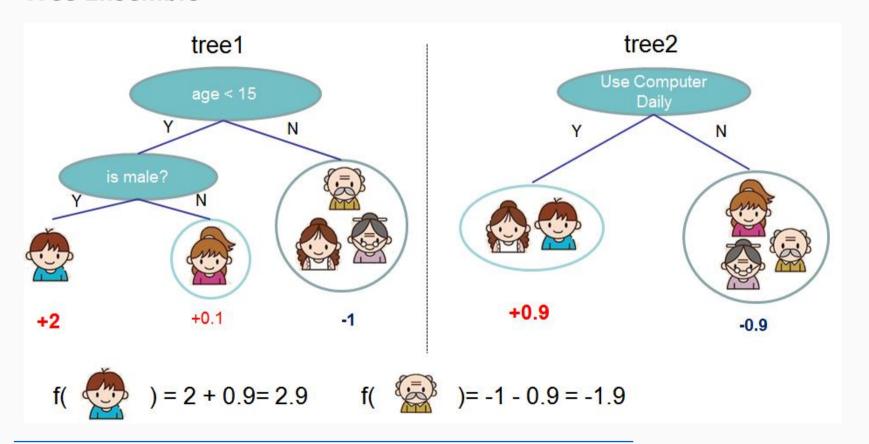


© 2016 Amadeus S

Appendix 1

XGBoost - Extreme Gradient Boosting

Tree Ensemble





2016 Amadeus S.,

Appendix 1

XGBoost - Extreme Gradient Boosting

Gradient Boosting



1



g1, h1

2



g2, h2

3



g3, h3

4

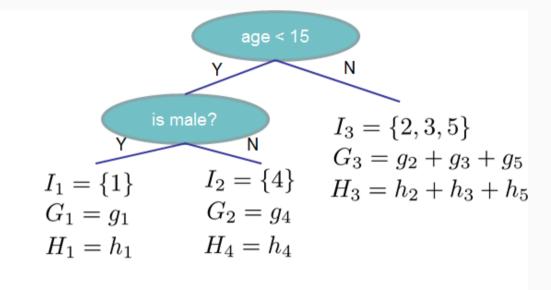


g4, h4

Ę



g5, h5



The smaller the score is, the better the structure is

 $Obj = -\sum_{j} \frac{G_j^2}{H_i + \lambda} + 3\gamma$



2016 Amadeus S A S

Appendix 2

Word2Vec

- Represent a word by a vector (vectorize)
- Reduce the data dimension compare with other word frequency method (TF-IDF)
- Better results most of the time

