

# Advanced Analytics

## Predictive Analytics and AI: the future of actuaries?

Antoine LY  
antoine.ly@milliman.com

December 2018 – Tel Aviv

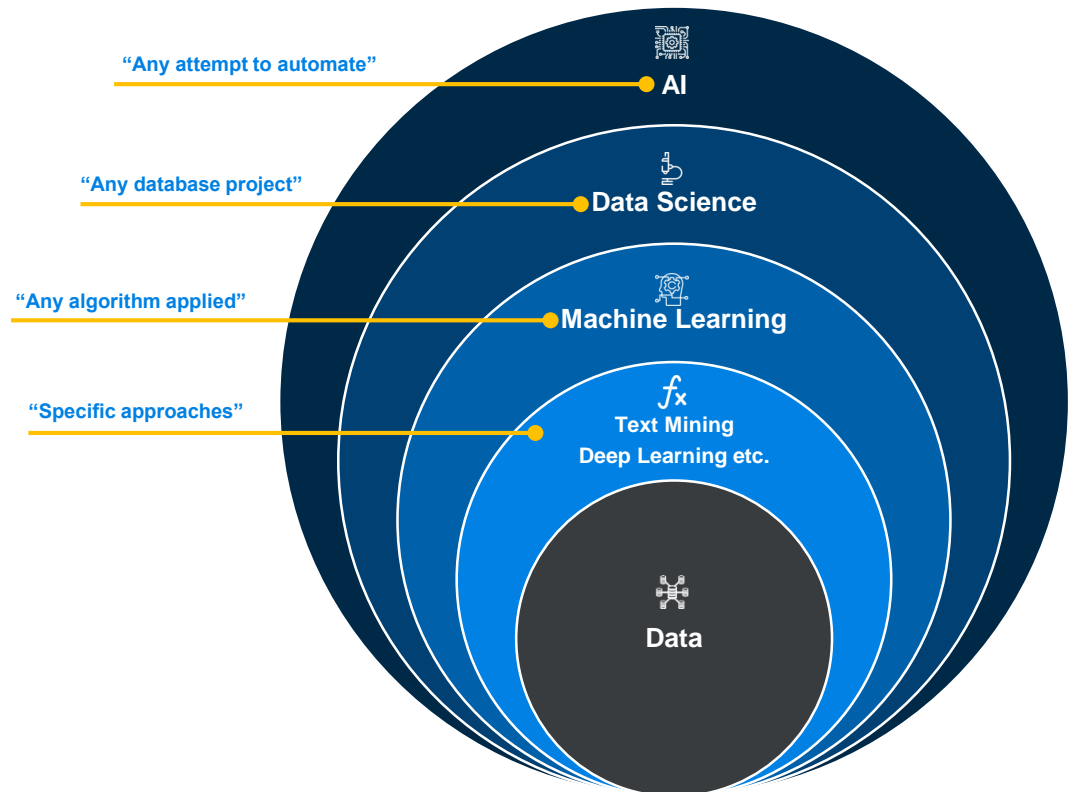


Copyright © Milliman 2018. All rights reserved. The information contained in this document ("Information") is for informational purposes only and Milliman assumes no liability or advisory duty in relation to said Information. The Information must not be modified, reproduced or distributed without the express consent of Milliman.

# Advanced Analytics

## Artificial Intelligence?

- Origin: Alan Turing (1950).
- Today, **many definitions**:
  - Big Data + Machine Learning = AI,
  - AI = Handle data with intelligence,
  - Etc.
- This covers **algorithms, automation, data visualization**,...
- Recent **buzz words**: deep learning, GDPR, discriminative algorithm, ...



# Advanced Analytics

## General presentation

- Advanced analytics combines **multiple fields of expertise** (Predictive modelling, IA, IT) to tackle complex issues which involve diverse sources of data and strong structural constraints
- The diversity of profiles and experience of the Milliman Analytics team allows us to work on **different topics** and ensures the success of projects.



### Predictive modelling

Projects with **machine learning** and **actuarial** modelling



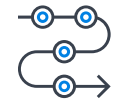
### Process automation

Design of tools and dashboards in R (**Shiny**), Python (**Dash**, **Flask**)



### Data management

Management of different data sources, **big data** processing, use of structured and unstructured data



### Software development

Development of **SaaS softwares**, use of best practices (CI, unit tests, git), local or cloud deployment

# Advanced Analytics

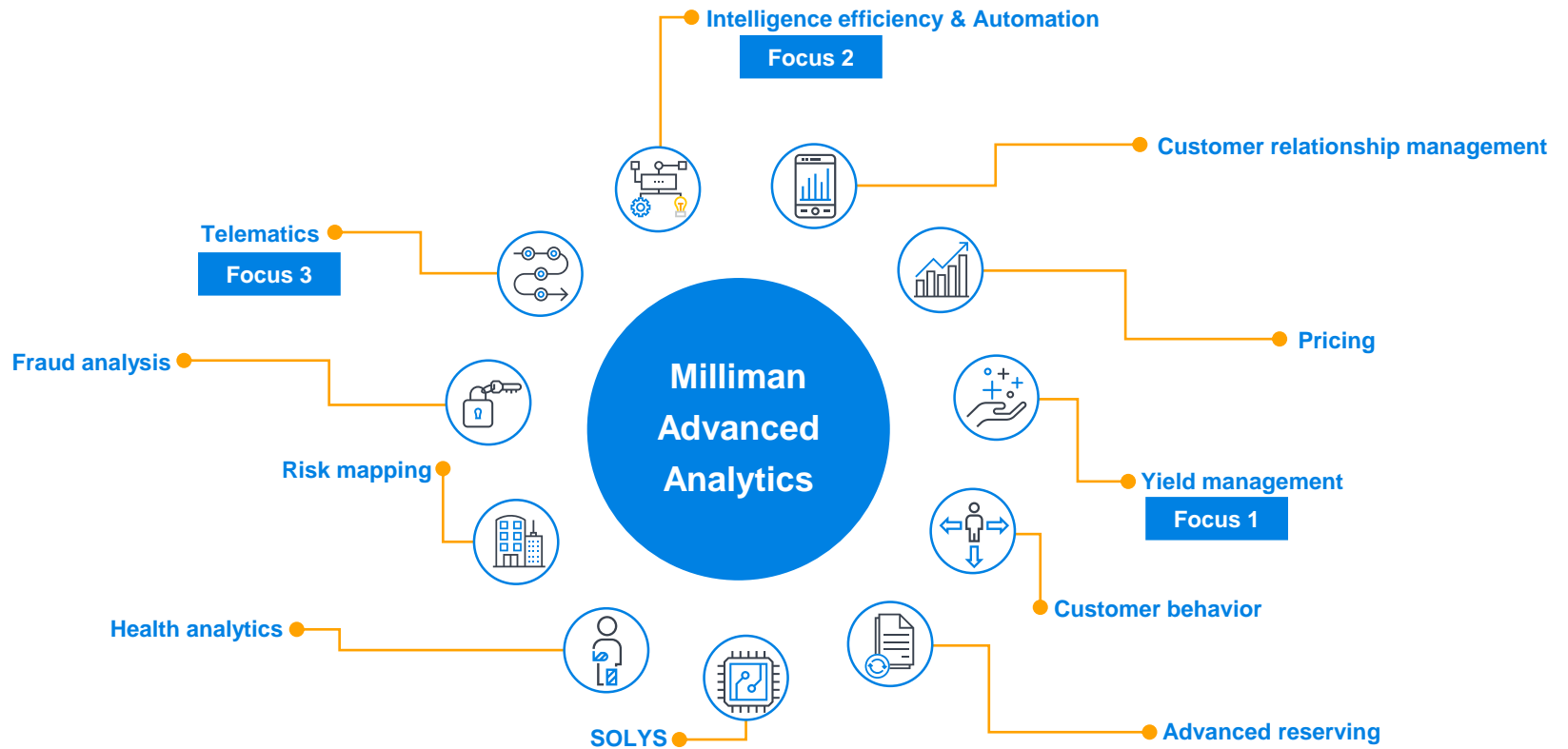
Our vision: Key points

- Milliman has been working on different analytics projects for 5 years. From our experience ranging from studies to tool developments, we believe that successful analytics in a large company depends on these **key points**:



# Advanced Analytics

Different areas



The background features a complex network of thin, light grey lines crisscrossing across the frame. Overlaid on this are several prominent, thicker lines in shades of blue and orange. Scattered throughout are circular nodes, some solid and some with a double-ring effect, in matching blue and orange colors. A large, solid blue shape with a diagonal cutout is positioned on the left side, serving as a backdrop for the text.

# Focus 1 Yield Management

# Artificial efficiency applied to yield management

Traffic and Budget forecasting for non-insurance company (1/2)

## Description:

A European transportation company wants to deploy an open source solution to lead their budget and predict traffic .

## Constraints:

- Use only open source technologies;
- Testing and integrating a machine learning model;
- Developing a Graphical User Interface;
- Respecting current IT infra.



# Artificial efficiency applied to yield management

Traffic and Budget forecasting for non-insurance company (2/2)

## Highlights

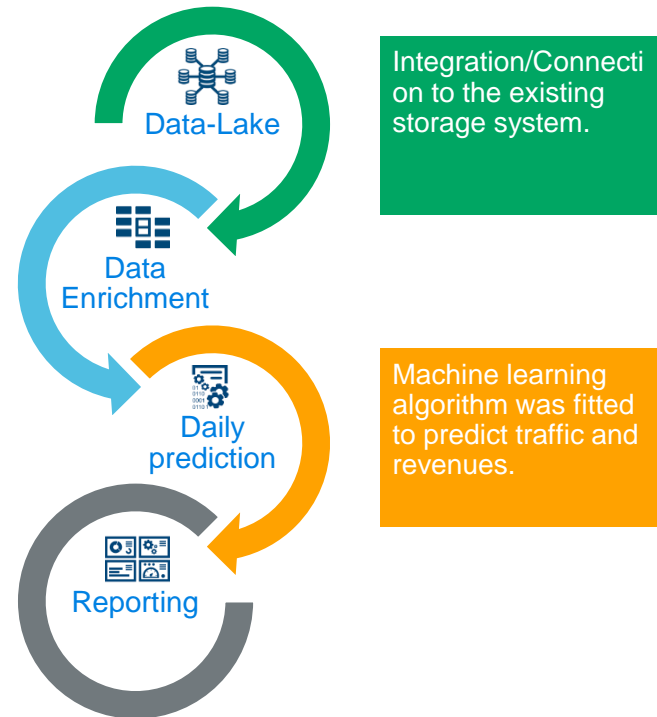
We demonstrated that Milliman could successfully achieve a data science project. All the expected skills of a data scientist were gathered into one project: data management, machine learning modelling, web reporting, IT knowledge.

## Duration of the project

2month + 6 month  
(PoC + full project)

Extraction of relevant data and joining with external data (calendar, targets etc.)

Indicators are generated after the predictions and reported through an interactive dashboard.





# Focus 2 Text Mining



# Comment analysis

# Comment analysis on insurance companies

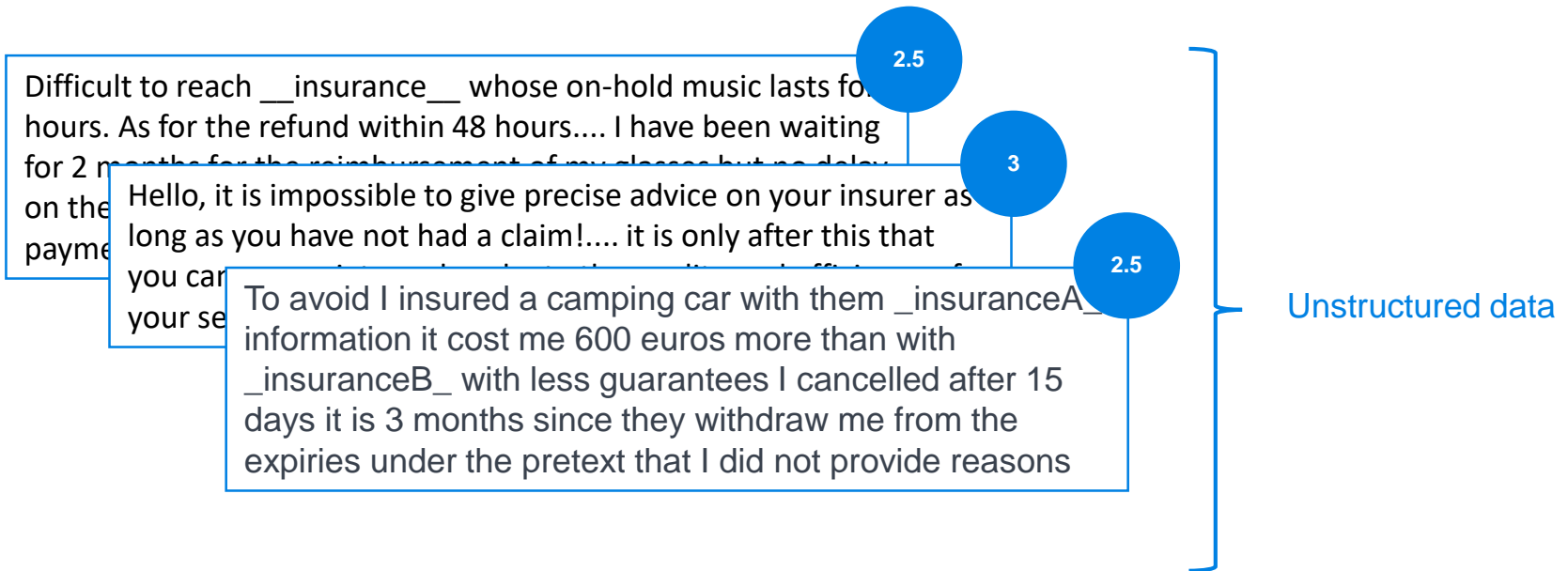
## Context (1/3)

- Why analyzing comments could be useful?
  - **Clients share experience:** comments and ratings influence underwriting because people like getting feedback.
  - **Understanding user experience:** unsatisfied clients make more noise... understanding main reasons why someone can move to another insurance company is valuable.
  - **Competitive market:** new digital actors with digital services.
- How do we proceed?
  - **Lots of information:** thousands of comments are published, reading them one by one could be time consuming and not efficient to extract main topics.
  - **Automating analysis and reporting:** the programming languages used in machine learning provide a large amount of tools to ease exploration.

# Comment analysis on insurance companies

## Context (2/3)

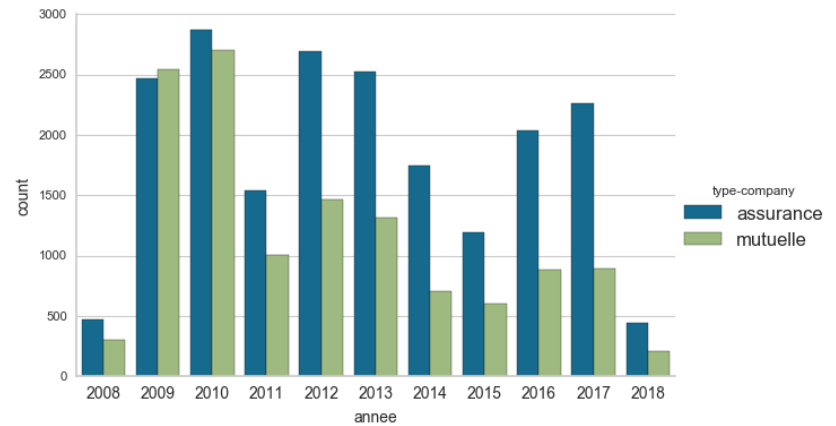
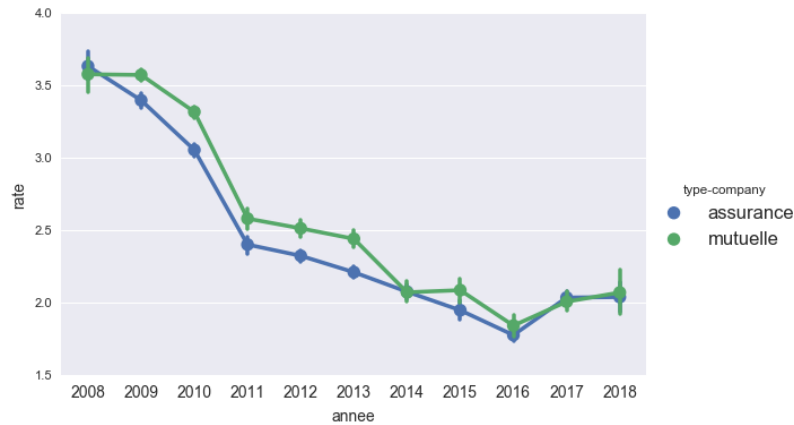
- A study lead on almost 40k comments about a mix of 180+ French insurance companies
- Each comment is rated with a mark of five points



# Comment analysis on insurance companies

## Context (3/3)

- General trend: rate decreases 1.5 points lost in average in 10 years
- No specific distinction between mutual or insurance companies



- Can we learn more?

# Comment analysis on insurance companies

Data preprocessing (1/2)

## How to deal with text information?

no refunds from September to March despite multiple registered letters left unanswered. the people with whom we are connected by phone are never able to answer requests due to computer problems, holidays, illness... NO customer services! more than dubious practice! do not subscribe!!!!!!



[ no , refund , september , multiple , letter , left , despite , unanswer , people , whom , we , be , connect , phone , never , able , request , due , computer , problem , holidays , demand , illness , service , more , dubious , practice , not , subscribe ]

- Convert string sequence into list of words
- Stemming
- Remove “stopwords”
- Remove special characters
- Apply custom scrubbing

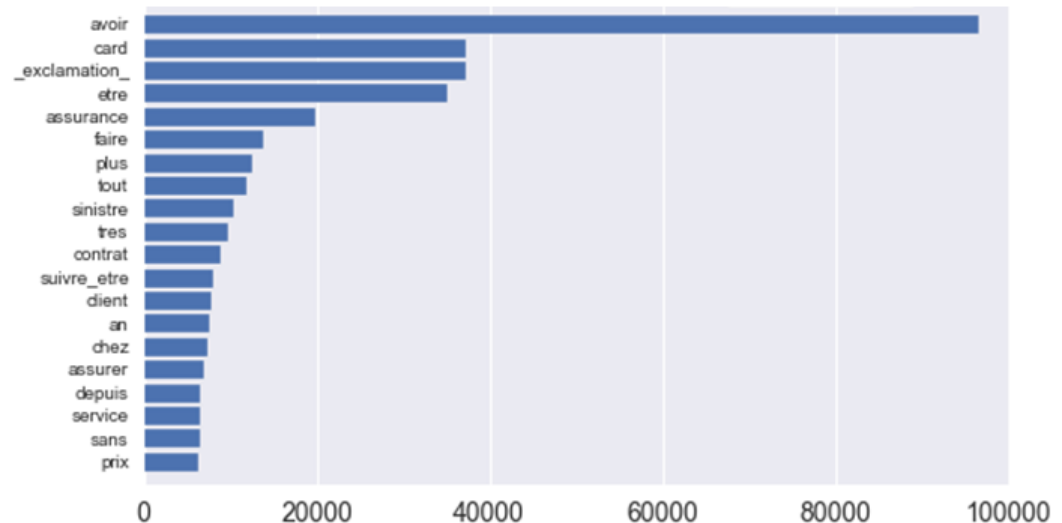
# Comment analysis on insurance companies

Data preprocessing (2/2)

## First analysis

- People who comment are **very expressive**: 1 exclamation point in average by comment
- People are **very descriptive**, large a amount of qualifying verbs and numbers to describe prices and times
- Different perspectives for the analysis:
  - Focus on **figures** (time or prices)
  - Define a **normed rating** (retreated from exposure, etc.)
  - Extract **topics** from text
  - etc.

Frequency Distribution of Top 20 tokens



# Comment analysis on insurance companies

Text representation

## Methodologies overview

[ aucun , remboursement , septembre , a , mars , malgres , multiple , courrier , recommande , laisser , sans , reponse , personne , sommer\_etre , mettre , relation , telephone , etre , jamais , capable , repondre , demande , suite , avoir , problem , informatique , conge , maladie , aucun , service , client , pratiquer , plus , douteur\_douteux , avoir , surtout , souscrire ]

Frequency

### Bag of Words

dictionary aucun, avoir, aller, ..., jamais, ..., problem, ...

Sentence

1	2	0	...	1	...	1	...
---	---	---	-----	---	-----	---	-----

Each sentence is represented by a vector of length 37896 counting the occurrence of each word



### Vector representation: « Word2Vec »

Context

dictionary aucun, avoir, aller, ..., jamais, ..., problem, ...

word

0	0.01	$2^{e-4}$	...	$w_k$	...	...	...
---	------	-----------	-----	-------	-----	-----	-----

Each word is represented by a vector but there, each value indicates a score (probability like) that a word of the dictionary co-appears with the word of interest



# Comment analysis on insurance companies

Focus on topic detection (1/2)

- Unsupervised methodology to detect topic

- Ratings are biased and subjected to personal appreciation
- Extracting topics helps in understanding reason of (dis)satisfaction
- It provides more insights to enhance client experience and improve process

→ Example with **Latent Dirichlet Allocation**: Bayesian model to cluster words/sentences

Topic 0

tres bon  
rapport a  
bon assurance  
service client  
card card  
tres cher  
bon assureur  
moins cher  
etre tres  
concurrence  
aucun suivi  
tres mauvais  
a deconseiller  
suivi dossier  
card juillet

Topic 3

card an  
avoir avoir  
depuis card  
card euro  
avoir card  
non responsable  
avoir etre  
avoir sinistre  
jamais avoir  
tout risque  
avoir jamais  
card avoir  
an avoir  
avoir accident  
moins cher

Topic 1

avoir ete  
avoir avoir  
card moi\_mois  
suite a  
assurance a  
card jour  
service client  
card card  
a eviter  
apres avoir  
avoir dire  
a recevoir  
avoir sinistre  
avoir envoye  
a jamais

Topic 4

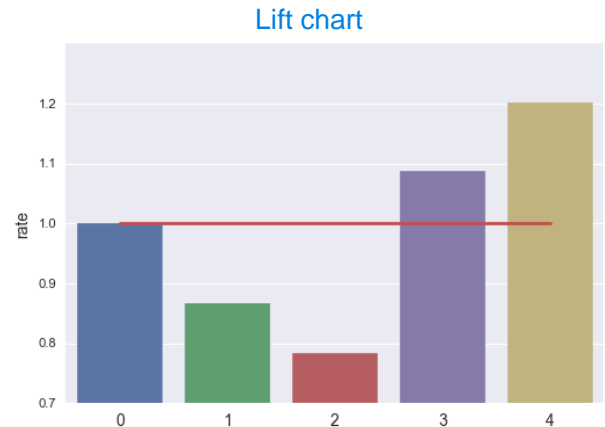
cas sinistre  
service client  
avoir ecoute  
etre tres  
prise charge  
conseiller etre  
qualite service  
prix etre  
tres rapide  
intervention cas  
assurance tres  
conseiller clientele  
disponibilite conseiller  
qualite intervention

Topic 2

\_exclamation\_  
qualite prix  
a fuir  
rapport qualite  
fuir \_exclamation\_  
tout aller  
\_exclamation\_ etre  
aller bien  
card moi\_mois  
a eviter  
rien \_exclamation\_  
client \_exclamation\_  
tres mauvais  
\_exclamation\_ plus  
\_exclamation\_ alors

# Comment analysis on insurance companies

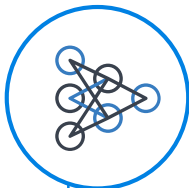
Focus on topic detection (2/2)



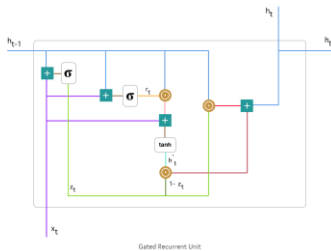
# Use case: Net promoter score automation

An example to apply text-mining

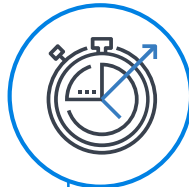
## Understand and design models



- Understanding models
- Programming skills
- Business knowledge



## API



- Capitalizing on tool development experience
- Risk modeling
- Optimizing production process with **open source**

## Dashboarding Rshiny / Dash



- Monitoring results
- Optimize long term ROI



# Optical Character Recognition (OCR)

### Pdf / word document

**Rechercher**

Sequential

Précédent      Suivant

GMAT / Princeton Review  
Vol. 6, pp. 748-775  
Sequential Design for C  
Robert B. Gramacy<sup>1</sup> and Michael Ludkovski<sup>2</sup>

**Abstract.** We propose a new approach to solving optimal stopping problems via simulation. Working within the backward dynamic programming/Shell envelope framework, we augment the methodology of Longstaff and Schwartz that focuses on approximating the stopping strategy. Namely, we introduce adaptive generation of the stochastic grids anchoring the simulated sample paths of the underlying state process. This allows for active learning of the classifier partitioning the state space into the continuation and stopping regions. To this end, we examine exponential domain schemes that adaptively place new domain points close to the stopping boundaries. We then discuss dynamic regression algorithms that can implement such recursive estimation and local refinement of the classifiers. The new algorithm is illustrated with a variety of numerical experiments, showing that an order of magnitude savings in terms of design size can be achieved. We also compare with existing benchmarks in the context of pricing multidimensional Bermudan options.

**Key words.** optimal stopping, regression Monte Carlo, dynamic time, active learning, expected improvement

**AMS subject classifications.** 91G60, 62L65, 65G40

**DOI.** 10.1137/10980898

1. Introduction. Numerical solution of optimal stopping problems remains a fertile area of research with applications in derivatives pricing, optimization of trading strategies, and options, and algorithmic trading. As the underlying models continue to get more and more complex, the computational holy grail of robust, fast, and accurate solvers remains elusive. Essentially all analytic methods deteriorate in high-dimensional problems where geometric intuition vanishes. Thus, recent attention has turned to probabilistic approaches, based on the Shell envelope representation. These methods reduce to recursive estimation of conditional expectations,

$$(1.1) \quad f_t(x) := \mathbb{E}[Y_t | X_t = x], \quad t = T-1, T-2, \dots, 0,$$

where the response is real-valued,  $Y_t \in \mathbb{R}$ , and the state process  $X_t \in \mathbb{R}^d$  is a multidimensional Markov process, typically with moderate dimension  $d \in [1, 10]$ .

The estimation problem in (1.1) comes from a dynamic programming (DP) argument. Namely, the process  $(Y_t)$  is the Shell envelope of the payoff process  $(h_t(X_t))$ , and its expected value  $f_t(x)$  is interpreted as the price of the corresponding claim, given initial condition  $X_t = x$ . The envelope  $Y_t$  is defined recursively via  $Y_t = h_t(X_t)$ , where  $\tau_{t+1} := \inf\{s \geq t+1 : h_s(X_s) \geq f_t(X_s)\}$  depends on the future expectations  $\{f_s(x) : s \geq t, x \in \mathbb{R}^d\}$ .

Received by the editors July 30, 2014; accepted for publication (in revised form) June 19, 2015; published electronically August 18, 2015.  
<http://www.siam.org/journals/afm/9/10980898.html>  
 Booth School of Business, The University of Chicago, Chicago, IL 60637 (rgramacy@chicagobooth.edu)  
 Department of Statistics and Applied Probability, University of California, Santa Barbara, CA 93106-3110 (ludkovski@pstat.ucsb.edu). This author's research was partially supported by NSF grant ATD-1222282.

748

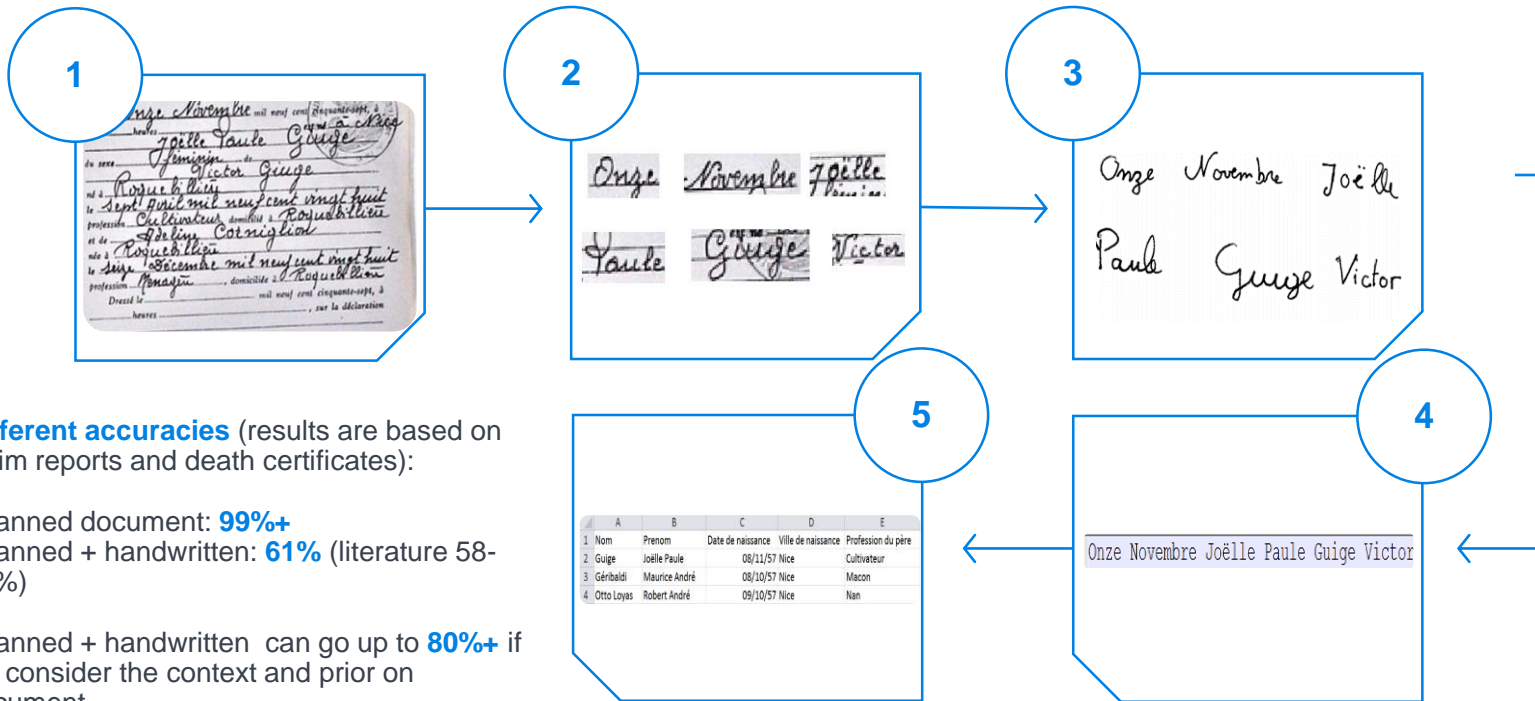
Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.



### Scanned + Handwritings

# OCR algorithms with handwritten writing style

Text-Mining and document analysis



**Different accuracies** (results are based on claim reports and death certificates):

Scanned document: **99%+**

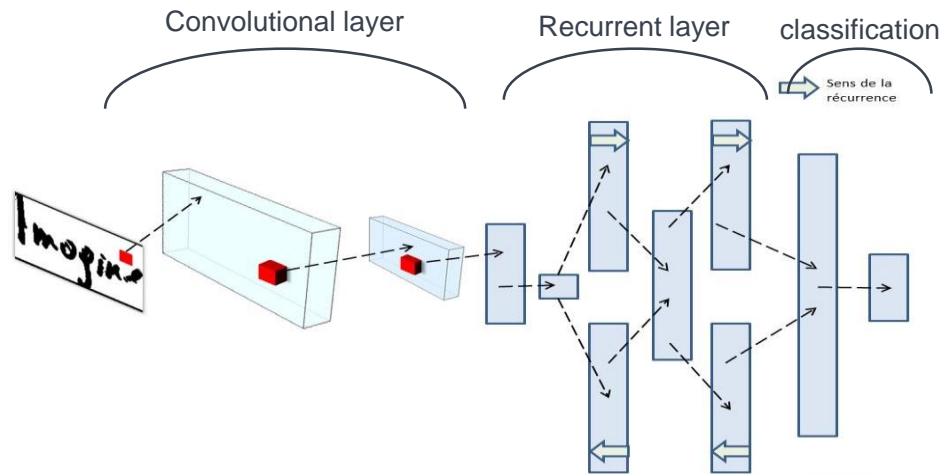
Scanned + handwritten: **61%** (literature 58-60%)

Scanned + handwritten can go up to **80%+** if we consider the context and prior on document.

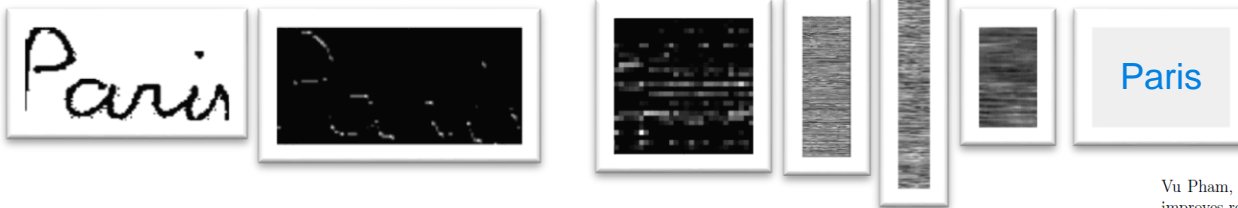
	A	B	C	D	E
1	Nom	Prenom	Date de naissance	Ville de naissance	Profession du père
2	Guige	Joëlle Paule	08/11/57	Nice	Cultivateur
3	Géribaldi	Maurice André	08/10/57	Nice	Macon
4	Otto Loyas	Robert André	09/10/57	Nice	Nan

# OCR algorithms with handwritten writing style


## Training and results



- Data base (Rimes)
  - Training on **70 000 pictures**
  - Testing 7000
- Use of AWS GPU
- Error metric WER: % of words badly detected (i.e. at least one character is wrong)
- Literature best score (LSTM 100): **44.37 %**
- Our OCR model: **46.5 %**



Vu Pham, Théodore Bluche, Christopher Kermorant, and Jérôme Louradour. Dropout improves recurrent neural networks for handwriting recognition. In *Frontiers in Handwriting Recognition (ICFHR), 2014 14th International Conference on*, pages 285-290. IEEE, 2014.



# Focus 3 Telematics



# General overview

# Telematics

At a glance (1/3)

- **Key words** within the insurance business:

- Black Box Insurance
- Usage Based Insurance
- Pay As You Drive
- Pay How You Drive
- Driving pattern
- Driving styles
- Pattern recognition
- Data provider



- The information is called driving behavior data (DBD). It covers the typical questions: when, where, how.

- **Disruption**

- **Customers:** demand personalized relationships, reactivity, transparency, valuable services
- **Insurers:** looking to boost profitability, anticipate and understand clients needs
- Other **players**

# Telematics

At a glance (2/3)

- Applications:
  - Pricing,
  - Scoring,
  - Claims (first notice of loss),
  - Fraud,
  - etc.
  - **Underwriting & Services.**
- Involve potential causative variables in the equation that could replace **traditional proxies**:
  - **Accuracy**
  - **Incentives**
- **Telematics is no longer in early stages:**
  - Italy, UK, US
  - French market

# Telematics

At a glance (3/3)



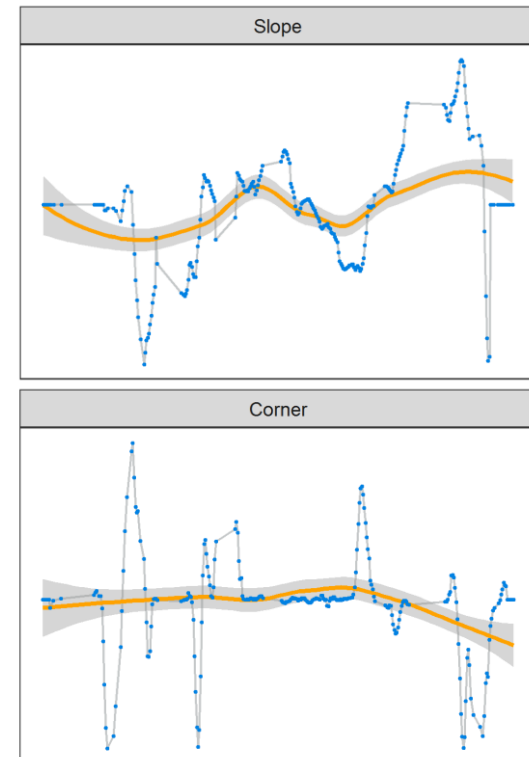
**In practice**

# Telematics – in practice

## Data preparation

- Provided by (at least) **GPS** latitude, longitude and **timestamp**.
- **Data management**: filtering techniques (causal moving mean or median, Kalman, etc.)
- Compute **indicators**: speed, acceleration, etc.

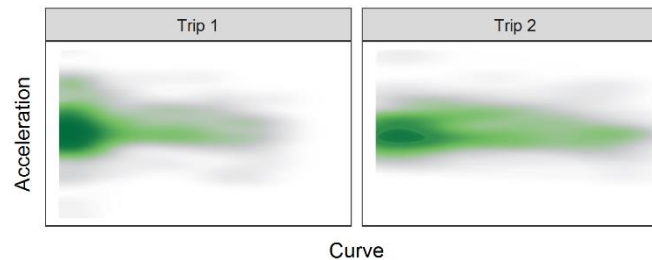
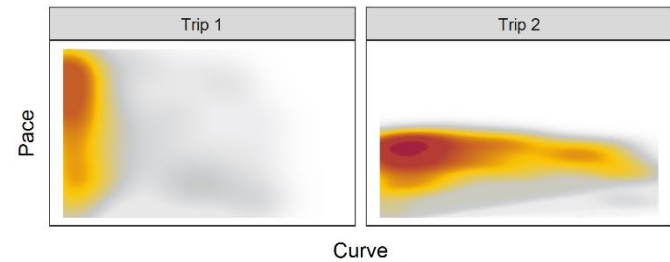
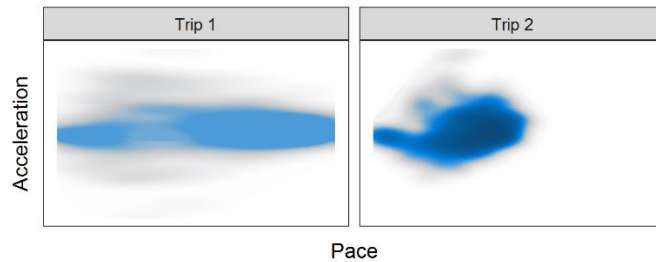
Latitude	Longitude	Elevation	Time
degree	degree	meters	dd/mm/yyyy hh:mm:ss
48.8755	2.2848	50	21/05/2018 20:30:00
48.8760	2.2868	55	21/05/2018 20:30:05
48.8757	2.2884	60	21/05/2018 20:30:10
48.8755	2.2905	65	21/05/2018 20:30:15
48.8753	2.2919	75	21/05/2018 20:30:20



# Telematics – in practice

## Data analysis – overview

- Build your scores to answer questions:
  - How **far** do you drive?
  - How **aggressive** are you?
  - How **fast** do you drive?
  - **When** do you drive?
  - ...



# Telematics – Case study

Get insights from the data

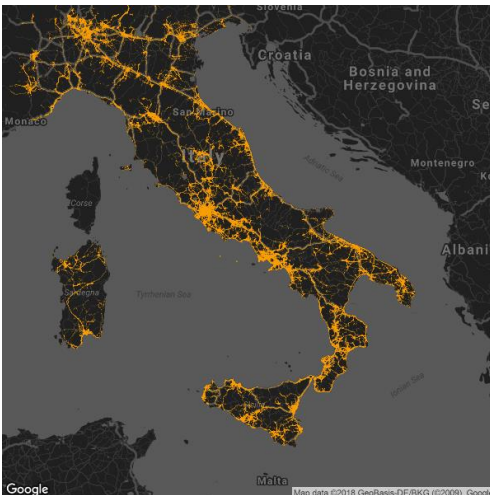
- From very raw data, we can enrich/challenge the KPIs reported by the data provider. Information is being controlled with a lot of applications (**pricing** and **services**).
  - More: White paper 2018 <http://www.milliman.com/uploadedFiles/insight/2018/raw-telematics-data-driving-behaviour.pdf>



## Data assessment

Take control

- Explore** with simplified analysis
- Scrub** the data
- According to the frequency we can **compute information**



## Visualization

View the data

- Use custom graphics** to take decisions for the modelling part



## Use information

Build profiles, KPIs, services

Use the information for:

- Trips and drivers analysis:** unsupervised techniques to build objective clusters
- Crash analysis:** get a better understanding of how an accident happened
- Build a driving risk score** the driver can use to optimize the insurance premium
- Fraud detection:** check if claims report matches with the telematics data
- Predictive maintenance:** considering the data available, predict and anticipate the maintenance of the vehicle
- Etc.:** other roads conditions, etc.

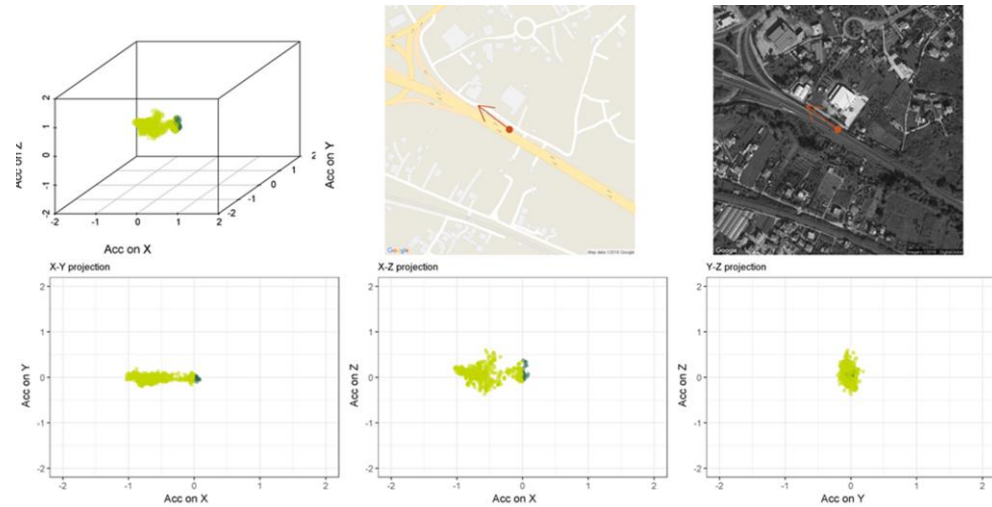


# Telematics – Case study

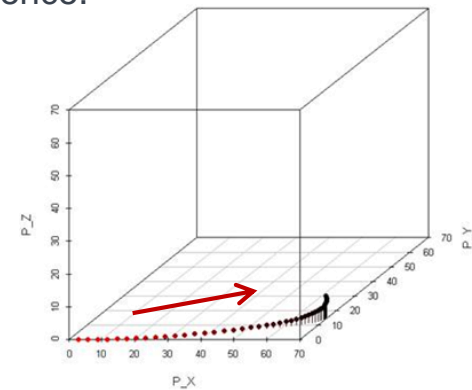
Focus: How to analyze crash/event?



- Analysis of detected events.



- Trajectory inference.



The background of the slide features a complex network diagram with numerous nodes and connecting lines, rendered in a light blue and grey color scheme. A large, semi-transparent blue rectangle is overlaid on the left side of the image, containing the word "Conclusion" in white text.

# Conclusion

Thank you. Questions?

