

# Projet de Stage : Détection d'Anomalies

Wensi DING

August 25, 2016

# Outline

- 1 Introduction
- 2 Modélisation
- 3 Test et Conclusion
- 4 Référence

## Données :

- données brutes : imsi, event, cgi, timestamp
- données stop-segment : imsi, segmentType, tsStart, tsEnd, cuidStart, cuidEnd, speed
- données des cellules : cgi, cuid, longitude, latitude, rayon

## Objective : identifier les **anomalies**

- comportements inhumains?
- comportements particuliers?
- imprécision ou erreurs dans les données ?
- etc...

# Présentation des Modèles

## Modèle Gaussien :

- chaque variable suit la loi gaussienne
- toutes les variables sont indépendantes
- anomalies sont les points avec les plus petites probabilités

## Modèle LOF :

- lof : rapport entre la moyenne des densités des points voisins et la densité locale
- anomalies sont des points qui ont une densité relativement petite par rapport leurs voisins

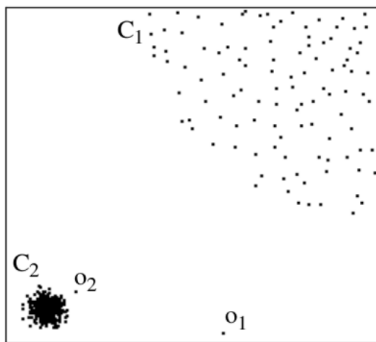
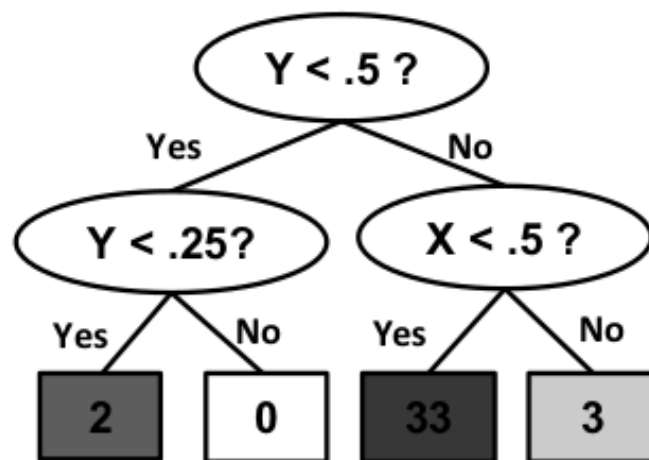


Figure 1: 2-d dataset DS1

# Présentation des Modèles

Modèle Half Space Tree :

- Arbre de décision
- Chaque nœud représente les intervalles des toutes les dimensions
- anomalies sont des points qui entrent dans les feuilles avec les plus petits nombres de points



Avec les **données brutes**:

- **v\_max** : le maximum des vitesses entre deux événements consécutifs
- **v\_mean** : la moyenne sur l'ensemble des valeurs moyennes des vitesses entre deux événements consécutifs chaque jour
- **n\_events\_mean** : la moyenne du nombre d'événements chaque jour
- **n\_events\_sd** : la variance du nombre d'événements chaque jour
- **duree\_stop\_mean** : la moyenne des sommes de durée d'être immobile chaque jour
- **duree\_stop\_sd** : la variance des sommes de durée d'être immobile chaque jour

Avec les **données stop-segment**:

- **v\_max\_ss, v\_mean\_ss, duree\_stop\_mean\_ss, duree\_stop\_sd\_ss**
- **n\_stop\_mean** : la moyenne du nombre d'événements de type "STOP" chaque jour
- **n\_stop\_sd** : la variance du nombre d'événements de type "STOP" chaque jour
- **dis\_ss** : la médiane des distances entre les endroits où la personne a fait le dernier stop chaque jour

sur les variables :

- $v_{max}$  ✓ : capable de sortir deux types d'anomalies: cellules mal positionnées, oscillations entre deux cellules avec une distance grande
- $v_{max\_ss}$  : capable d'enlever les événements anormaux à cause de cellules très mal placées; pour les cellules pas trop mal positionnés ou les cellules avec rayons grands, le résultat peut varier
- $v_{mean}$  : écrasé par les  $v_{max}$  très grands, pas très performant
- $v_{mean\_ss}$  ✓ : capable de représenter une sorte de mobilité d'une personne en moyenne



sur les variables :

- `n_events_mean` ✓: capable de sortir des imsi avec le nombre d'oscillations entre plusieurs cellules très grand
- `n_events_sd` : écrasé par les `n_events_mean` très grands, pas très performant
- `n_stop_mean` ✓: capable de représenter une sorte de mobilité d'une personne, exemple: un voyageur peut faire plein de stop pendant une journée
- `n_stop_sd` ✓: capable de représenter une sorte de mobilité d'une personne, exemple: un voyageur qui visite plusieurs villes peut faire plein de stop dans une journée de visite, très peu de stop dans une journée de déplacement

sur les variables :

- `duree_stop_mean` : pas très précis pour représenter la vraie durée de stop
- `duree_stop_sd` : pas très précis à cause d'imprécision de `duree_stop_mean`
- `duree_stop_mean_ss` : capable de représenter une sorte de mobilité d'une personne, mais cette variable a une relation linéaire avec `n_stop_mean`
- `duree_stop_sd_ss` : pas de typique comportement sorti
- `dis_ss` : capable de représenter une sorte de mobilité dans le sens global, exemple: les personnes qui bougent tout le temps avec une grande distance chaque jour

variables choisies :

- E1 : v\_max, v\_mean, n\_events\_mean, n\_events\_sd, duree\_stop\_mean, duree\_stop\_sd
- E2 : v\_max, v\_mean\_ss, n\_events\_mean, n\_stop\_mean, n\_stop\_sd

sur les modèles :

1. Courbe ROC et AUC Score(référence : modèle gaussien)
  - LOF : 0.78 (E1) ; 0.82 (E2)
  - HS\_Tree : 0.94 (E1) ; 0.93 (E2)

# Test et Conclusion

2. Matrice de Confusion : 679 anomalies sur 157,107 imsi pour tous les trois modèles

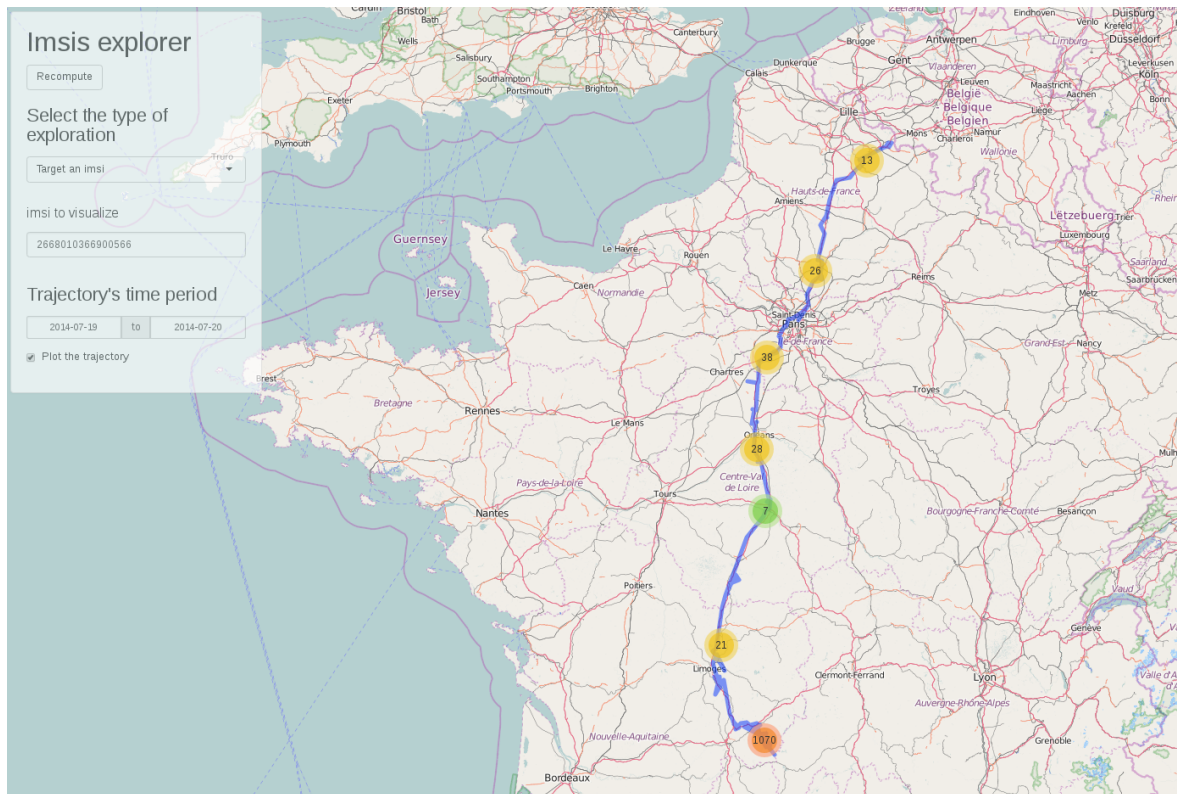
<b>Modèle Réfèrent :</b>	Gaussien	Gaussien	LOF
<b>Modèle Testé :</b>	LOF	HS_Tree	HS_Tree
<b>Accuracy Score :</b>	0.99	0.99	0.99
<b>F1 Score :</b>	0.23	0.22	0.11
<b>Matrice de Confusion :</b>	TN : 155906; FP, FN : 522; TP : 157	TN : 155897; FP, FN : 531; TP : 148	TN : 155825; FP, FN : 603; TP : 76

Table 1: Résultat matrice de confusion

- Pour comparer ou évaluer les résultats des trois modèles, il vaut mieux utiliser F1 Score ou Similarité de Jaccard.
- Pour expliquer les grandes différences entre les trois modèles, une hypothèse possible : la transformation à la distribution gaussienne sur les variables dans le modèle gaussien a un effet important.

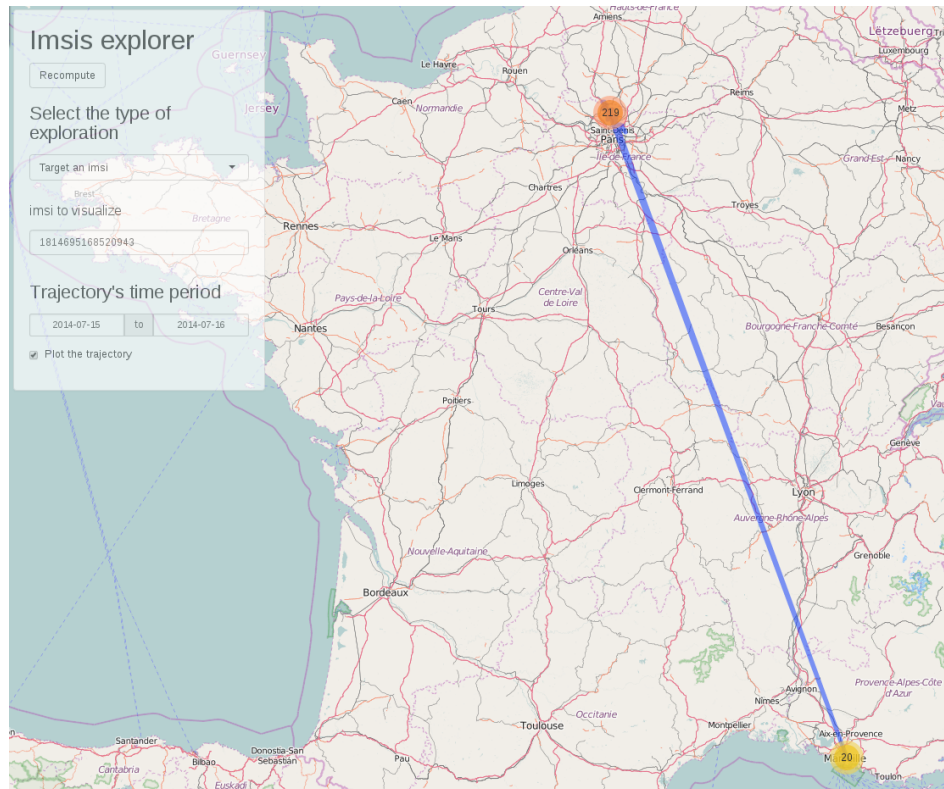
# Quelques types d'anomalies

## 1. profile spécial



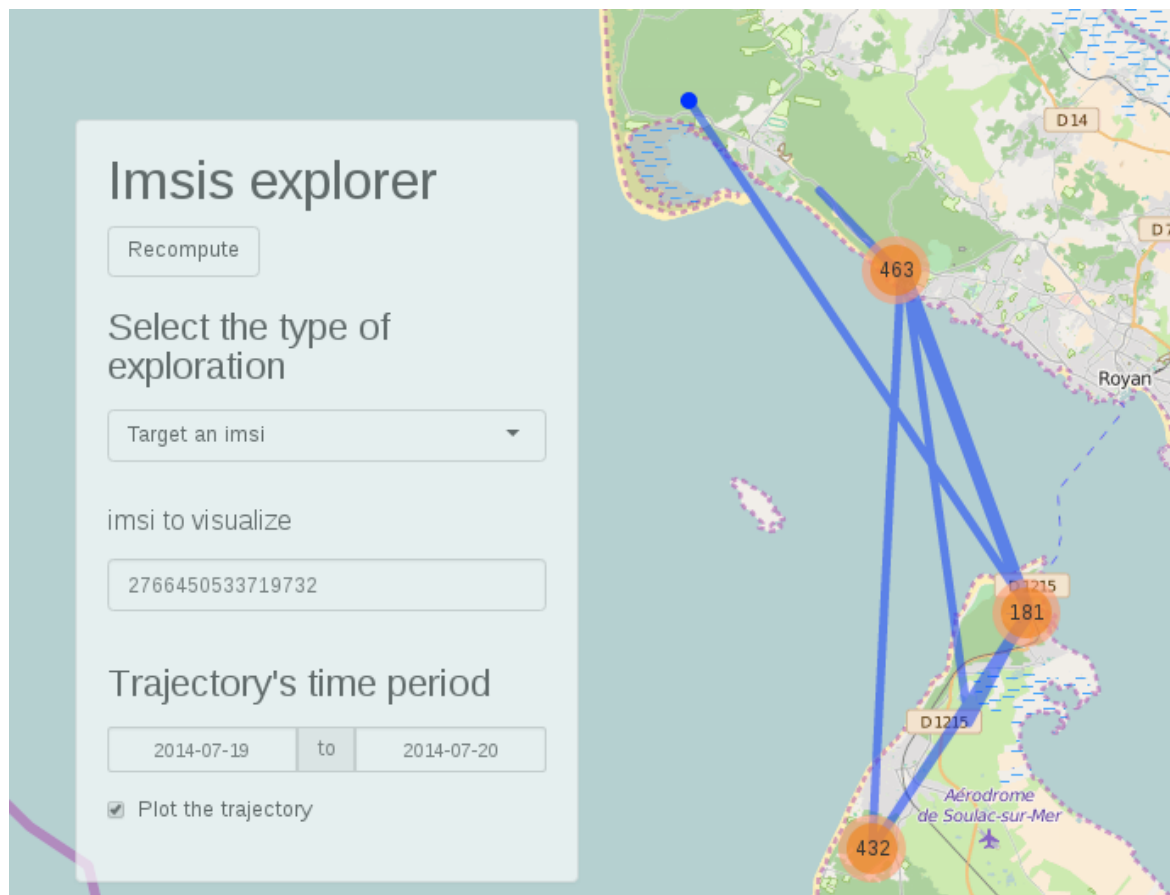
# Quelques types d'anomalies

## 2. cellule très mal placée



# Quelques types d'anomalies

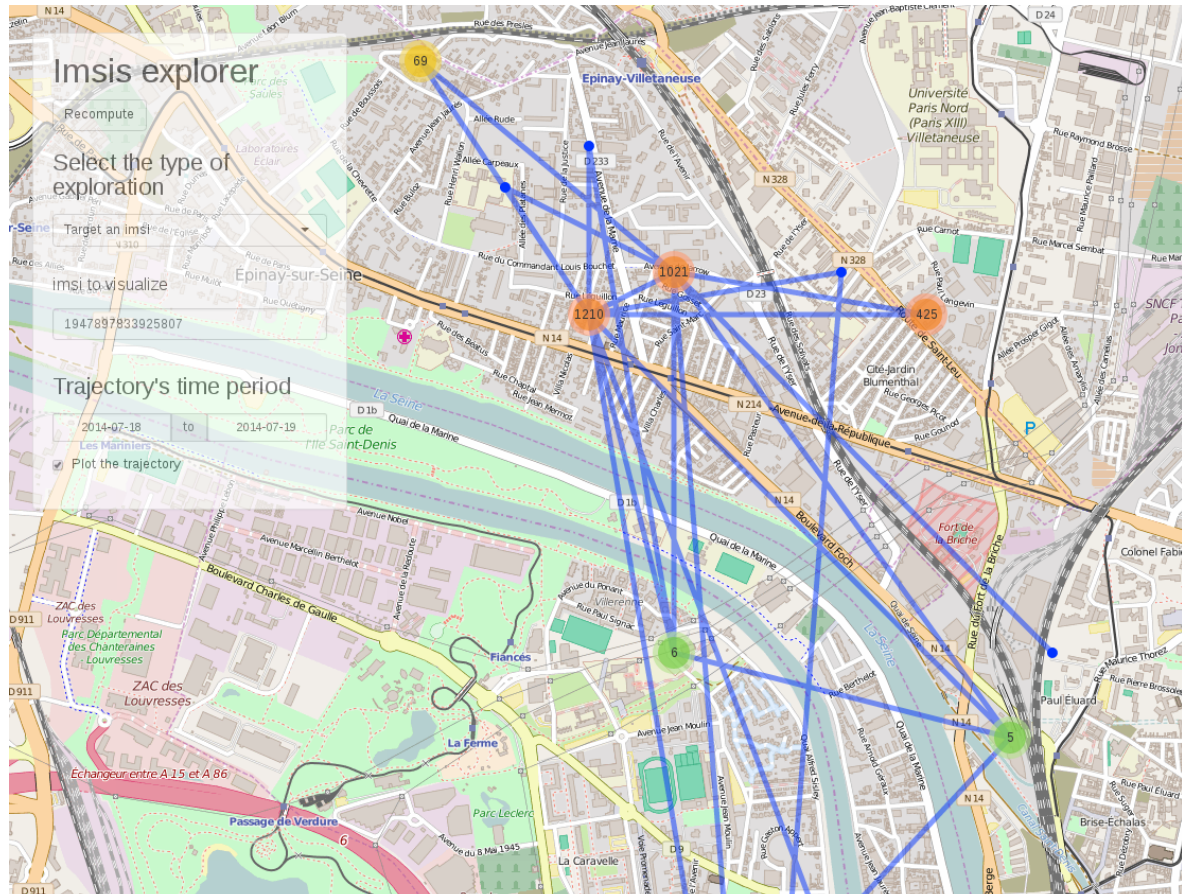
## 3. nombre d'événement expose





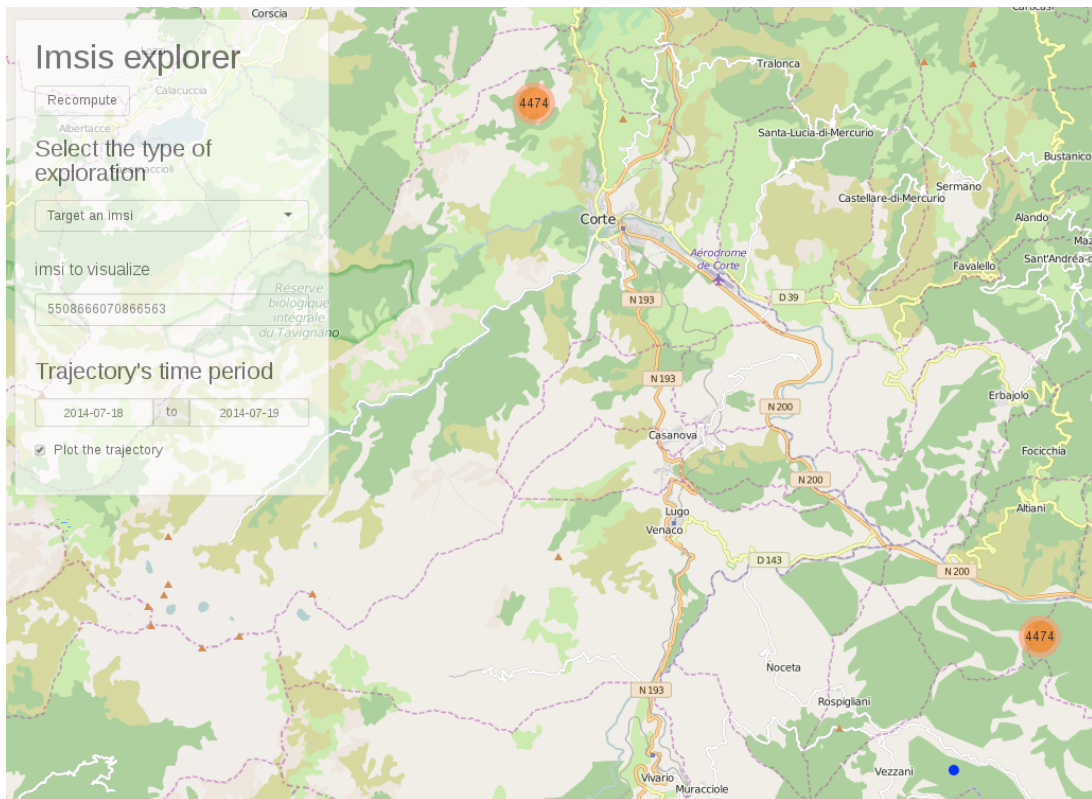
# Quelques types d'anomalies

## 3. nombre d'événement explose



# Quelques types d'anomalies

## 3. nombre d'événement explose



# Proposition

- modèle : comprendre les différences des modèles
- variables : créer et examiner des nouvelles variables (par exemple: différencier les jours dans la semaine et les weekends; utiliser les données avec d'autre type d'événement; etc...)

- Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, Jörg Sander (2000). LOF: Identifying Density-Based Local Outliers. ACM, 1-58113-218-2/00/05
- Swee Chuan Tan, Kai Ming Ting, Tony Fei Liu. Fast Anomaly Detection for Streaming Data. Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence