



## Conception de logiciels Adaptés Contexte & téléphone intelligent

UQÀM | Département d'informatique

Sébastien Mosser  
INF600G - E20 - Séquence 3- Partie 3



## Remerciements



Pr. Arthur Charpentier  
Département de Maths



Pr. Marie-Jean Meurs  
Département Informatique



Dr. Yonatan Vaizman  
UCSD

Ceci n'est TOUJOURS PAS un cours d'IA

## Bibliographie



### Machine Learning with R

Learn how to use R to apply powerful machine learning methods and gain an insight into real-world applications

Brett Lantz

[PACKT] open source

[Lantz, 2013]

### Recognizing Detailed Human Context In-the-Wild from Smartphones and Smartwatches

Yonatan Vaizman<sup>1</sup>, Katherine Ellis<sup>2</sup> and Gert Lanckriet, senior member, IEEE  
This version was accepted to IEEE Pervasive Computing magazine.  
Please cite the official (edited and stylized) paper, which appears in IEEE Pervasive Computing, vol. 16, no. 4, October–December 2017, pp. 62–74.

**Abstract**—The ability to automatically recognize a person's behavior context can contribute to health monitoring, aging care and many other applications. Most existing work focuses on laboratory settings, which are not representative of real-world usage. We collected over 300k minutes of sensor data with context labels from 62 subjects. Unlike previous studies, our subjects used their phones in their natural environments, including both constrained and unconstrained scenarios. We show that unconstrained phone usage resulted in situations that are harder to recognize. We demonstrate how fusion of multi-modal sensors can help to improve context recognition in the wild.

**Index Terms**—Contact awareness, Mobile sensors, Artificial intelligence, Machine learning, Human activity recognition.

#### 1 INTRODUCTION

The biomedical research community acknowledges the value of mobile sensors for health monitoring [1]. These sensors can be used to support aging at home [2], and to enable for such applications to succeed on a large scale, the sensors and the algorithms have to be accurate and to work smoothly, without requiring the person to change his/her behavior. This is challenging in the real-world settings, where such applications will eventually run. In this paper, we focus on the problem of recognizing in-the-wild, meaning capturing people's authentic behavior in their natural environments, with the use of everyday mobile phones. This is challenging because it increases the difficulty that in-the-wild conditions add, and show how multi-modal sensors can help.

#### Related work

It is common for people to have their phone close to them all day long. This is due to the fact that the variety of built-in sensors, make phones popular agents for recognizing people's behavior, especially in the home setting. While capturing informative signals about hand gestures, it is also important to keep the sensor very small so as not to add any burden to the user.

Pervasive works have shown the advantage of fusing multiple sensors from phones, tablets and smartwatches, to improve recognition of basic movement patterns, such as walking, running, sitting, standing or drinking or drinking coffee [8]. However, most past works have collected data under heavily controlled conditions, with subjects performing specific tasks, such as walking or running. Fitting models to recognize described activities may result in poor performance in real-life scenarios [9]. To generate real-life working applications, we argue that research has to be done in natural and realistic settings, satisfying four

[Vaizman et al, 2017]

## Rétrospective L2

## Détection de contexte

## Application au téléphone intelligent

## Travail à faire pour L3

1 2

3 4

# Jeux de données : ExtraSensory

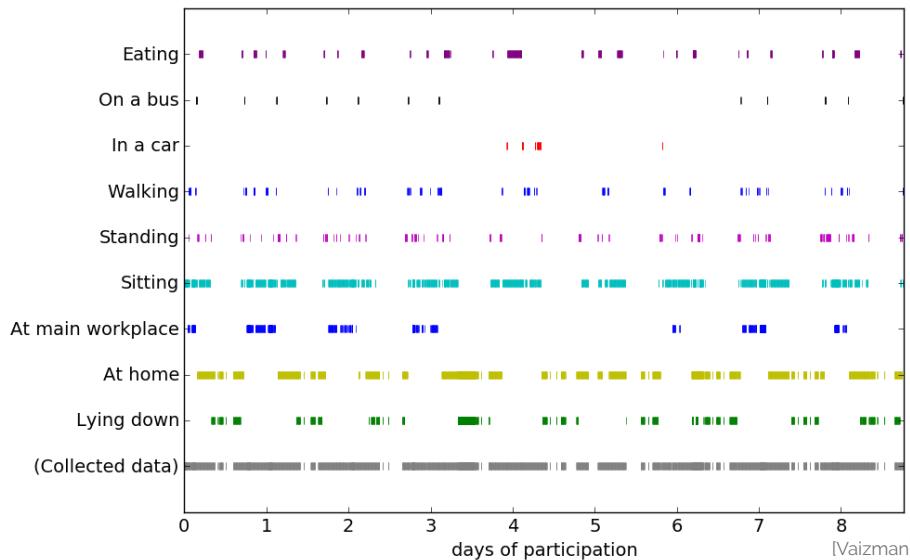


Name	Date Modified	Size	Kind
# 0A986513-7828-4D53-AA1F-E02D6DF9561B.features_labels.csv	May 15, 2017 at 6:17 PM	8.6 MB	CSV Document
# 0BFC35E2-4817-4B65-BF47-76472302AD.features_labels.csv	May 15, 2017 at 6:15 PM	6.7 MB	CSV Document
# 0E6184E1-90C0-45B4-B25A-F1ECB79714E.features_labels.csv	May 15, 2017 at 6:16 PM	15.6 MB	CSV Document
# 00EABED2-271D-49D8-B599-1DA09240601.features_labels.csv	May 15, 2017 at 6:14 PM	4.5 MB	CSV Document
# 1DBB0F6F-1F81-A5B4-0D94-C062ACFA4842.features_labels.csv	May 15, 2017 at 6:14 PM	15.9 MB	CSV Document
# 2C32C23E-E300-498A-8D04-C069150A02E.features_labels.csv	May 15, 2017 at 6:13 PM	18.5 MB	CSV Document
# 4E99F91F-4654-42E9-9900-433B94A27E7.features_labels.csv	May 15, 2017 at 6:15 PM	6.7 MB	CSV Document
# 4FC32141-E8B8-4BF-8804-1259A4910C8.features_labels.csv	May 15, 2017 at 6:13 PM	10.4 MB	CSV Document
# 5EF4122-8513-46A4-BCF1-E62AAC28502C.features_labels.csv	May 15, 2017 at 6:14 PM	8.1 MB	CSV Document
# 7CE37510-56D0-4120-AC1F-023351428D2Z.features_labels.csv	May 15, 2017 at 6:15 PM	20.6 MB	CSV Document
# 7D9B8B102-A612-4E2A-BE22-3159752F5D58.features_labels.csv	May 15, 2017 at 6:16 PM	3.1 MB	CSV Document
# 90C3BD04-E8B2-4F29-AB52-B476535226F2.features_labels.csv	May 15, 2017 at 6:16 PM	20.4 MB	CSV Document
# 11B5EC4D-A133-42B9-B475-4E737182A406.features_labels.csv	May 15, 2017 at 6:13 PM	19.2 MB	CSV Document
# 24E40C4C-A349-4F09-03AB-01D00FB994AF.features_labels.csv	May 15, 2017 at 6:15 PM	10.2 MB	CSV Document
# 27E04243-B138-4F40-A164-F40860165CF3.features_labels.csv	May 15, 2017 at 6:16 PM	10.2 MB	CSV Document
# 33A85C34-CFEE-4732-9E73-0A7ACB61B27A.features_labels.csv	May 15, 2017 at 6:14 PM	13.3 MB	CSV Document
# 40E170A7-607B-4578-A04-F021C3B0384A.features_labels.csv	May 15, 2017 at 6:17 PM	14.9 MB	CSV Document
# 59EEFAE0-DEB0-4FF-9250-54D2A3D0CF2.features_labels.csv	May 15, 2017 at 6:17 PM	15.9 MB	CSV Document
# 74BB6607-5D4B-43CF-82C2-341B76BEA0F4.features_labels.csv	May 15, 2017 at 6:14 PM	15.1 MB	CSV Document
# 78A91A6E-A451-40E5-8D47-94755F0B83B.features_labels.csv	May 15, 2017 at 6:16 PM	25.9 MB	CSV Document
# 83CF687B-7CEC-434B-9FE8-00CD3D5799BE6.features_labels.csv	May 15, 2017 at 6:16 PM	20.5 MB	CSV Document
# 86A4F379-B305-473D-9D83-FCTD00180EF.features_labels.csv	May 15, 2017 at 6:13 PM	20.4 MB	CSV Document
# 96A35BAA-FF2-4223-B93-C7425B901B47.features_labels.csv	May 15, 2017 at 6:16 PM	12 MB	CSV Document
# 098A72A5-E3E5-4F54-1B52-BBDAA0D7B694.features_labels.csv	May 15, 2017 at 6:13 PM	13.7 MB	CSV Document
# 99B204C0-DD5C-4B5-4B87-83E8-A372B1B0D769.features_labels.csv	May 15, 2017 at 6:15 PM	12.1 MB	CSV Document
# 48F14D02-7689-43B9-A2AA-C87722271628.features_labels.csv	May 15, 2017 at 6:17 PM	13.8 MB	CSV Document
# 797D145F-3858-4A7F-47C2-4AEB721E133C.features_labels.csv	May 15, 2017 at 6:17 PM	7.5 MB	CSV Document
# 1155FF54-63D3-4A82-9863-8385D0B00A13.features_labels.csv	May 15, 2017 at 6:17 PM	5.6 MB	CSV Document
# 1538C99F-BA1E-4EFA-9A94-6C7C47701820.features_labels.csv	May 15, 2017 at 6:15 PM	12.7 MB	CSV Document

Macintosh HD ▾ Users ▾ mosser ▾ work ▾ teaching ▾ INF600G ▾ behavioral-context-recognition ▾ datasets ▾ ExtraSensory

<http://extrasensory.ucsd.edu/>

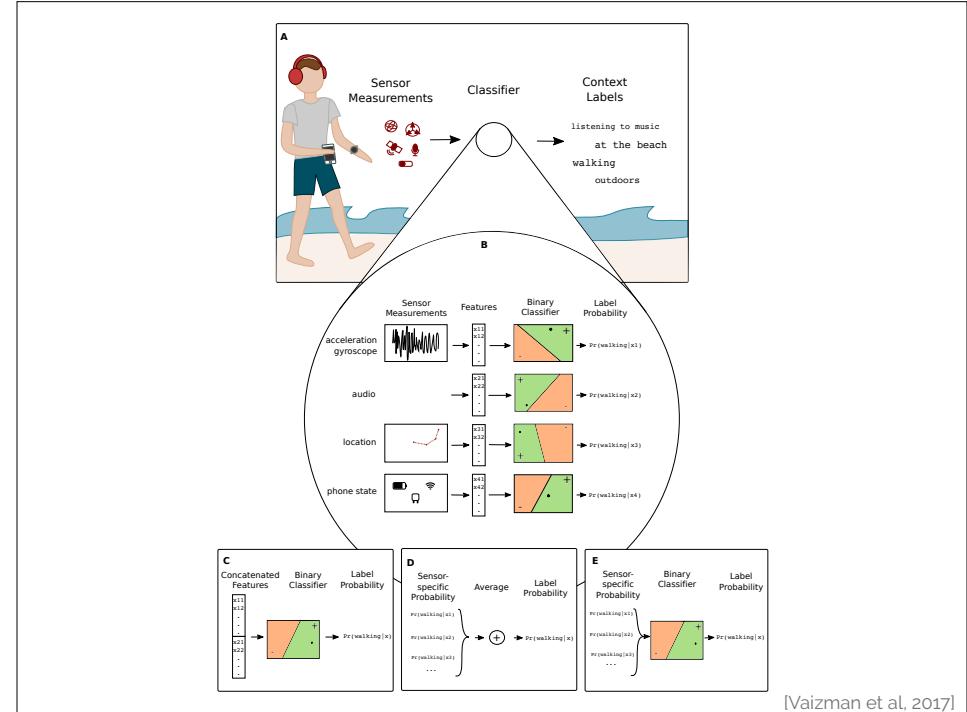
## Durée de la collecte



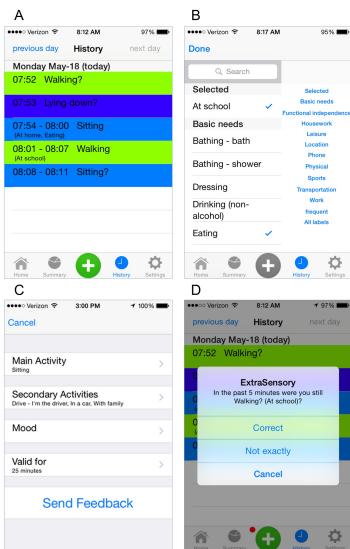
## Exemple de données collectées

0A986513-7828-4D53-AA1F-E02D6DF9561B.features_labels						
	A	B	C	D	E	
1	timestamp	raw_acc:magnitude_stats:mean	raw_acc:magnitude_stats:std	raw_acc:magnitude_stats:moment3	raw_acc:magnitude_stats:moment4	raw
2	1449601597	1.000371	0.007671	-0.016173	0.027860	
3	1449601657	1.000243	0.003782	-0.002713	0.007046	
4	1449601717	1.000811	0.002082	-0.001922	0.003575	
5	1449601777	1.001245	0.004715	-0.002895	0.008881	
6	1449601855	1.001354	0.065186	-0.096520	0.165298	
7	1449601918	1.002472	0.067833	0.164222	0.289580	
8	1449601975	1.001820	0.002009	-0.001022	0.003073	
9	1449602035	1.001889	0.001781	0.000630	0.003392	
10	1449602095	1.004284	0.105347	0.163315	0.295795	
11	1449602155	1.001456	0.003742	-0.004386	0.013095	
12	1449602215	1.001763	0.004617	0.011668	0.020438	
13	1449602275	1.001567	0.000635	0.000635	0.012204	

60 participants



# Application mobile de collecte



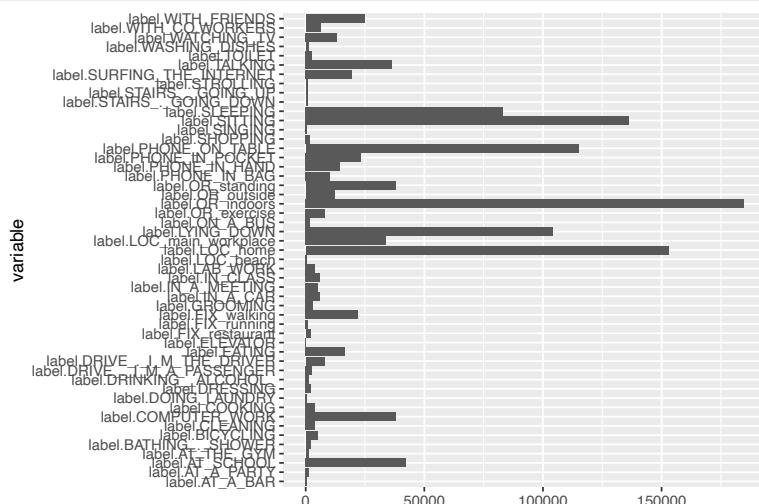
Surveillance du comportement

Annotation par l'utilisateur

Validation des prédictions

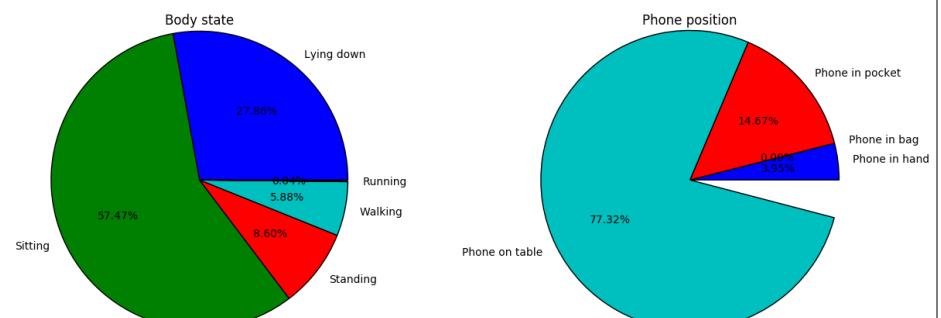
[Vaizman et al. 2017]

# Analyse quantitative du jeux de données

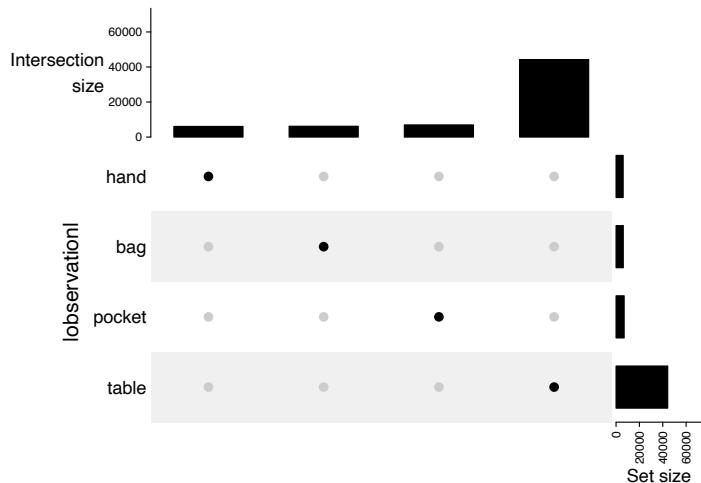


377,346 observations. Espace des données: 278 dimensions

# Représentation des observations



## Exclusivité des observations



## Fabrication d'un jeu de données initial

```
phone_observations <- dataset %>% select(starts_with("raw_acc"),
                                         starts_with("proc_gyro"),
                                         starts_with("raw_magnet"),
                                         starts_with("label:PHONE_"))
phone_observations <- phone_observations[complete.cases(phone_observations),]
```

- Sur les 377.348 observations et 278 dimensions, on garde :
  - Les données de l'accéléromètre, du gyroscope et du magnétomètre
  - Les étiquettes de reconnaissance du téléphone
  - Uniquement les observations complètes pour ces variables là
- Cela nous donne 59,982 observations

## Normalisation des données [1/2]

```
## raw_acc:3d:mean_x raw_acc:3d:mean_y raw_acc:3d:mean_z
## Min. :-1.134850 Min. :-1.022699 Min. :-1.08164
## 1st Qu.:-0.020165 1st Qu.:-0.028909 1st Qu.:-0.99016
## Median :0.012589 Median :0.004741 Median :-0.11415
## Mean :0.005651 Mean :0.028688 Mean :-0.05196
## 3rd Qu.:0.064781 3rd Qu.:0.082794 3rd Qu.:0.96098
## Max. :1.289129 Max. :1.290554 Max. :1.05319
```

- Les données ne sont pas normalisées
  - et donc difficilement comparables
- On va ramener les données dans un intervalle [0,1]

```
norm_variable <- function(v) { return ((v - min(v)) / (max(v) - min(v))) }
```

## Normalisation des données [2/2]

```
## raw_acc:3d:mean_x raw_acc:3d:mean_y raw_acc:3d:mean_z
## Min. :-1.134850 Min. :-1.022699 Min. :-1.08164
## 1st Qu.:-0.020165 1st Qu.:-0.028909 1st Qu.:-0.99016
## Median :0.012589 Median :0.004741 Median :-0.11415
## Mean :0.005651 Mean :0.028688 Mean :-0.05196
## 3rd Qu.:0.064781 3rd Qu.:0.082794 3rd Qu.:0.96098
## Max. :1.289129 Max. :1.290554 Max. :1.05319
```

```
normalize <- function(dataset) {
  ds_n <- as.data.frame(lapply(dataset[1:ncol(dataset)-1],
                                norm_variable))
  ds_n$status <- dataset$STATUS
  return(ds_n)
}
```

```
## raw_acc:3d:mean_x raw_acc:3d:mean_y raw_acc:3d:mean_z
## Min. :0.0000 Min. :0.0000 Min. :0.0000
## 1st Qu.:0.4599 1st Qu.:0.4296 1st Qu.:0.04285
## Median :0.4734 Median :0.4442 Median :0.45319
## Mean :0.4705 Mean :0.4545 Mean :0.48233
## 3rd Qu.:0.4949 3rd Qu.:0.4779 3rd Qu.:0.95680
## Max. :1.0000 Max. :1.0000 Max. :1.0000
```

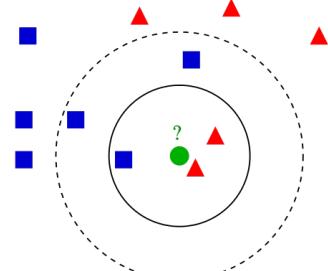
# Entrainement & Test

```
prepare_datasets <- function(complete, random.seed=42) {  
  set.seed(random.seed) # Pour rendre les expériences reproductibles  
  training <- complete %>% sample_n(size = 0.8*nrow(complete))  
  test <- setdiff(complete, training)  
  list(train = training, test = test)  
}
```

- On va garder 20% des données pour mesurer les performances
- Et entraîner les modèles avec 80%
- *On utilise une "graine" pour rendre les expériences reproductibles*

# Un premier prédicteur : kNN

- On applique la méthode des *k-nearest neighbours*:
- Pour une nouvelle observation :
  - On regarde les classes de ses *k*-plus proches voisins
  - On vote pour la classe majoritaire
- Simple à mettre en place :
  - Distance euclidienne

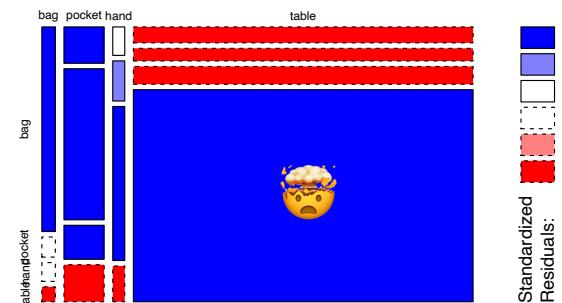


[https://en.wikipedia.org/wiki/K-nearest\\_neighbors\\_algorithm](https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)

# Une version de kNN en R

```
run_knn <- function(training, test, k.value=42) {  
  # Preparing the datasets for training  
  training_data <- training %>% select(-status)  
  training_classes <- training %>% select(status)  
  # Same, for testing  
  test_data <- test %>% select(-status)  
  test_classes <- test %>% select(status)  
  # Running the classifier  
  predictions <- knn(train = training_data,  
                     test = test_data,  
                     cl = training_classes[,1],  
                     k = k.value)  
  # Building the confusion matrix  
  cmat <- confusionMatrix(reference = test_classes[,1],  
                           data = predictions)  
  return(cmat)  
}
```

# Visualisation de la matrice de confusion



Standardized Residuals:  
Legend:  
-4 < -2.4 < 0 < 2.4 < 4

	Balanced Accuracy	Precision	Recall
Class: bag	0.63	0.79	0.28
Class: pocket	0.75	0.58	0.55
Class: hand	0.59	0.60	0.19
Class: table	0.73	0.82	0.97

## Un prédicteur encore plus simple : Constant

```
constant_table <- function(training, test) {  
  test_classes <- test %>% select(status)  
  # On répond toujours table ... facile !  
  predictions <- as.factor(rep('table', nrow(test)))  
  cmat <- confusionMatrix(reference = test_classes[,1],  
                           data = predictions)  
  return(cmat)  
}  
cmat_constant <- with(ds <- prepare_datasets(normalize(phone_observations)),  
                      constant_table(ds$train, ds$test))
```

Ce prédicteur à une accuracy de ... 70.80% !!!!

## On doit nettoyer les données !!

```
summary(phone_observations$STATUS)  
  
##      bag   pocket    hand   table  
##      5668     6533     5465   42216
```

- Les observations "téléphone sur la table" sont déséquilibrées
- On va ré-équilibrer le jeu de données :
  - On garde 'bag', 'pocket' et 'hand', environ 6,000 obs/classes
  - On échantillonne (*random*) 6,000 observations pour 'table'

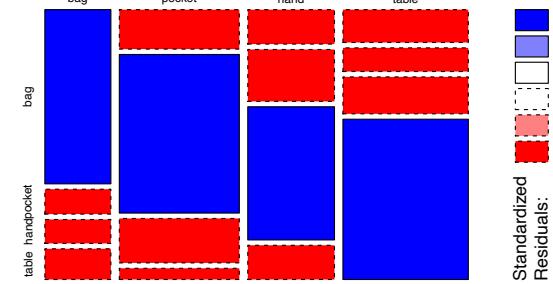
## Équilibrage, et ré-apprentissage

```
balance <- function(dataset, random.seed=42) {  
  set.seed(random.seed)  
  res <- dataset[dataset$STATUS == 'bag',]  
  res <- rbind(res, dataset[dataset$STATUS == 'pocket',])  
  res <- rbind(res, dataset[dataset$STATUS == 'hand',])  
  res <- rbind(res, sample_n(dataset[dataset$STATUS == 'table'], 6000))  
  return(res)  
}
```

On parle de "pipeline" de traitement de données

```
cmat_balanced <-  
  with(ds <- prepare_datasets(normalize(balance(phone_observations))),  
       run_knn(ds$train, ds$test))
```

## Résultats sur le jeu équilibré



ACC = 61.75%

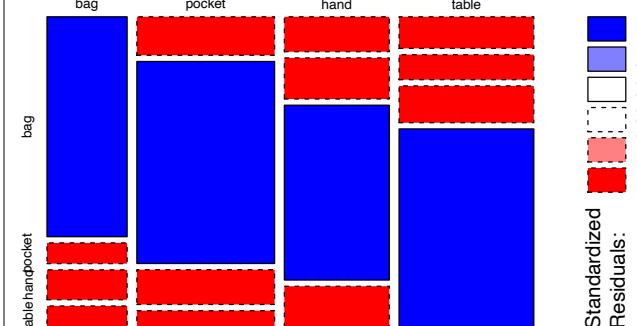
	Balanced.Accuracy	Precision	Recall	F1
Class: bag	0.71	0.69	0.49	0.57
Class: pocket	0.76	0.63	0.68	0.65
Class: hand	0.68	0.52	0.50	0.51
Class: table	0.80	0.63	0.76	0.69

# Problème de normalisation

- Notre normalisation "écrase" les valeurs extrêmes
- Ce n'est pas très adapté à des valeurs de capteurs
  - Changement brusque de cap par exemple
- On va utiliser une autre méthode de normalisation : le Z-score
  - Inconvénient : pas dans [0,1] (pas de min ni de max prédéfini)

```
z_normalize <- function(dataset) {  
  ds_n <- as.data.frame(scale(dataset[1:ncol(dataset)-1]))  
  ds_n$status <- dataset$STATUS  
  return(ds_n)  
}
```

# Matrice de confusion avec Z-score

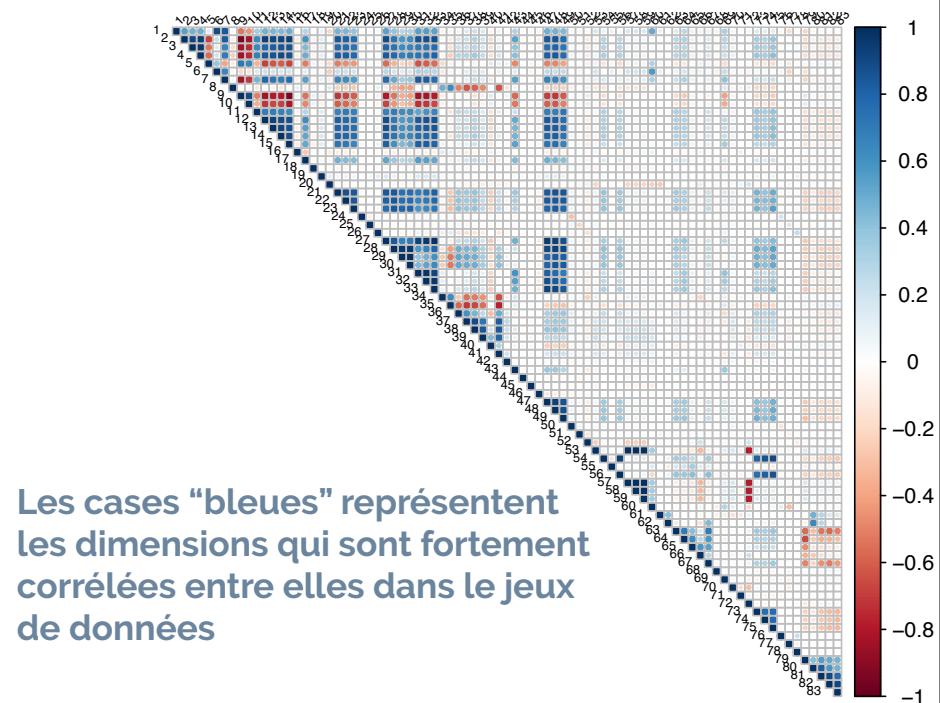


	Balanced Accuracy	Precision	Recall	F1
Class: bag	0.76	0.75	0.57	0.65
Class: pocket	0.81	0.69	0.75	0.72
Class: hand	0.74	0.60	0.60	0.60
Class: table	0.82	0.68	0.76	0.72

# Et si on regardait les dimensions ?

- On travaille avec 83 dimensions sur les observations
- L'algorithme kNN repose sur une distance euclidienne
  - Il est sensible à la présence de "trop" de dimensions
- Peut-on se débarrasser de dimensions corrélées ?

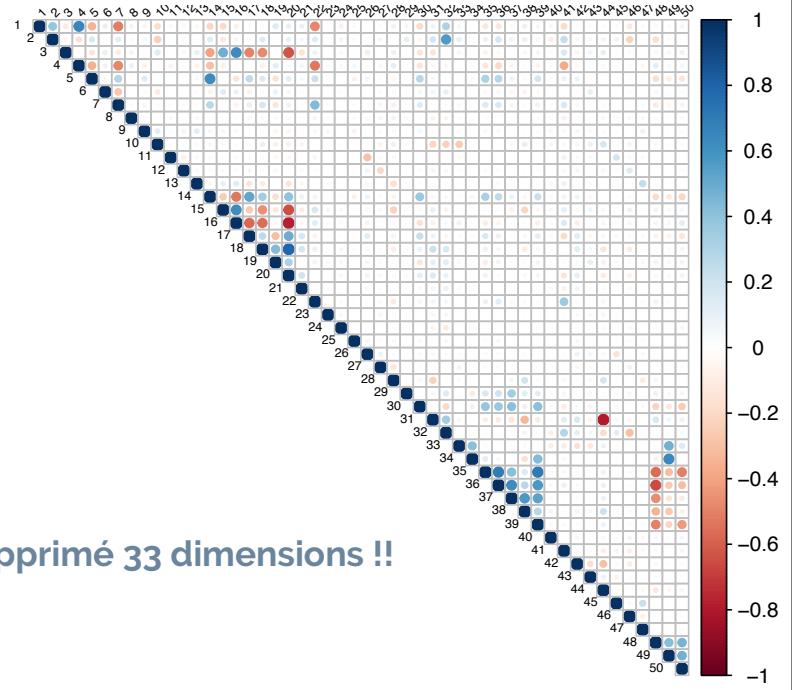
```
print_corplot <- function(dataset) {  
  corr_obs <- dataset %>% select(-STATUS)  
  colnames(corr_obs) <- 1:ncol(corr_obs)  
  corr_matrix <- cor(corr_obs)  
  corrpplot(corr_matrix, type = "upper",  
            tl.col = "black", tl.srt = 45, tl.cex=0.5)  
}
```



# Épuration du jeu de données

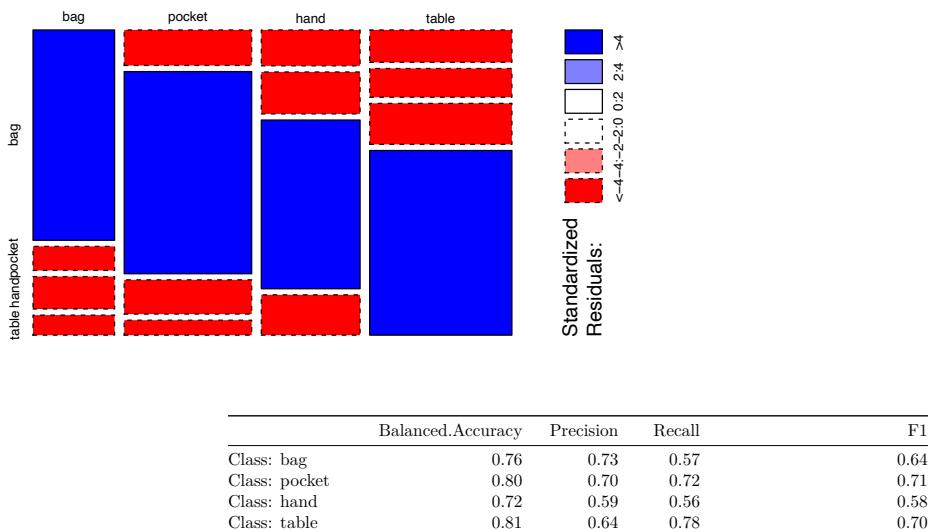
- On va garder uniquement les dimensions "pas trop" corrélées
  - On fixe un seuil arbitraire à 80%
- Cela va grandement accélérer le temps d'entraînement
  - On part d'un espace initial à 83 dimensions ...

```
slice_relevant <- function(dataset, threshold = 0.8) {  
  corr_obs <- dataset %>% select(-STATUS)  
  corr_matrix <- cor(corr_obs)  
  high <- findCorrelation(corr_matrix, cutoff = threshold)  
  res <- dataset[, -c(high)]  
  res$STATUS <- dataset$STATUS  
  return(res)  
}  
print_corplot(slice_relevant(phone_observations))
```



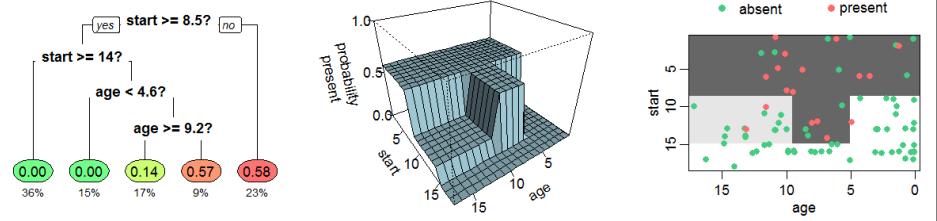
On a supprimé 33 dimensions !!

## Résultat sur le jeu épuré



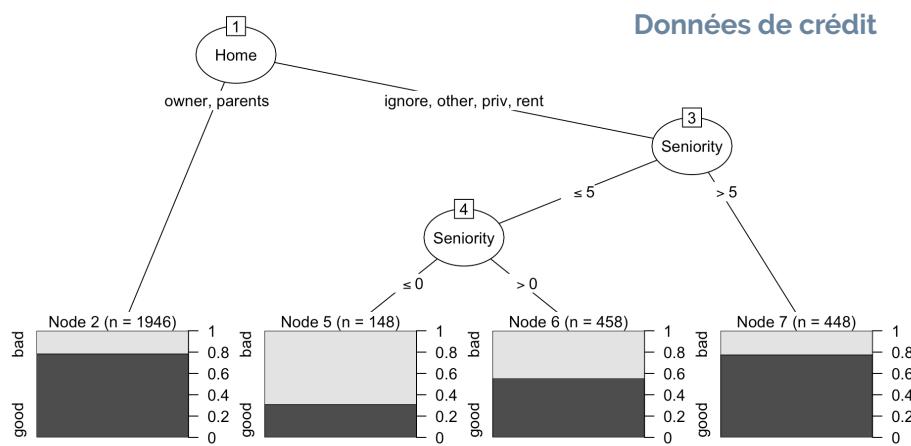
## Changement de paradigme : Decision Trees

- Plutôt qu'un seul choix, on va passer par une succession de décisions qui vont nous amener vers la classe de l'élément



[https://en.wikipedia.org/wiki/Decision\\_tree\\_learning](https://en.wikipedia.org/wiki/Decision_tree_learning)

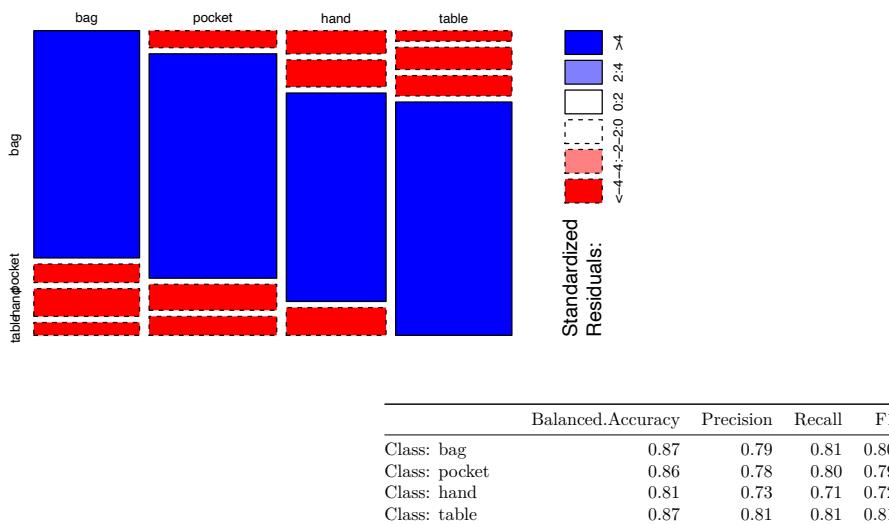
## Exemple : l'algorithme C5.0



## Entraînement d'un modèle C5.0

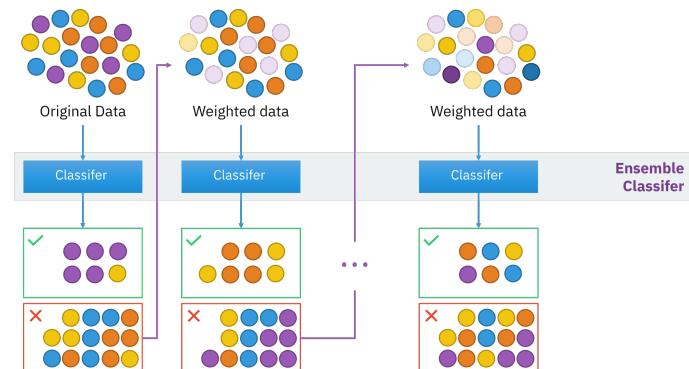
```
run_dt <- function(training, test, nb.trials=1) {  
  training_data <- training %>% select(-status)  
  training_classes <- training %>% select(status)  
  test_data <- test %>% select(-status)  
  test_classes <- test %>% select(status)  
  model <- C5.0(training_data,  
                 training_classes[,1],  
                 trials = nb.trials)  
  predictions <- predict(model, test_data, type = "class")  
  cmat <- confusionMatrix(reference = test_classes[,1],  
                           data = predictions)  
  return(list(cmat=cmat, tree=model))}
```

## Résultat avec un arbre de décision

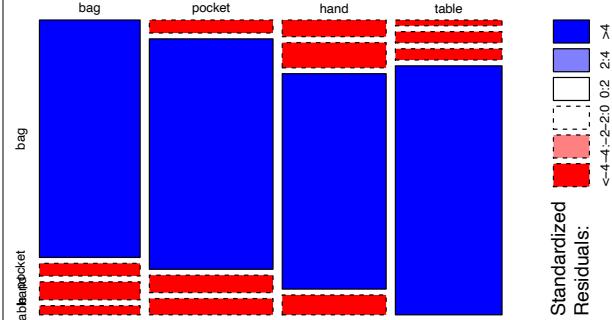


## Méta-algorithme: Adaptive Boosting

- Plutôt qu'un seul arbre, on va en fabriquer plusieurs
- Les arbres vont se "spécialiser" pour trouver les cas "spéciaux"



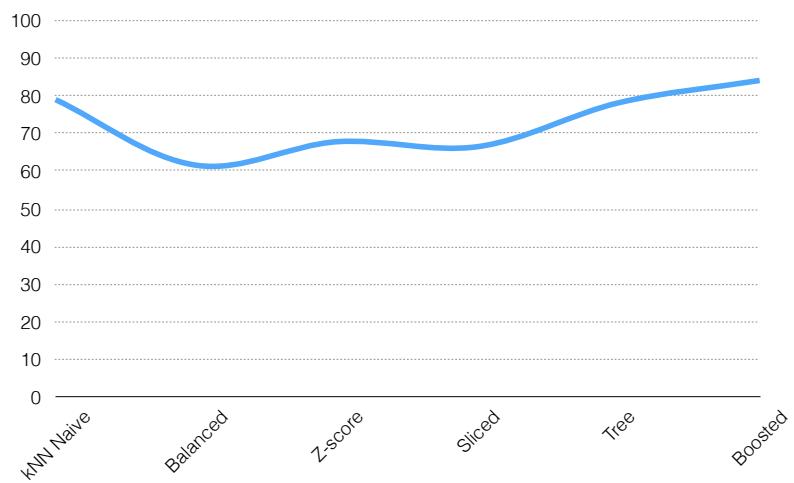
## Résultat avec une accélération de 10



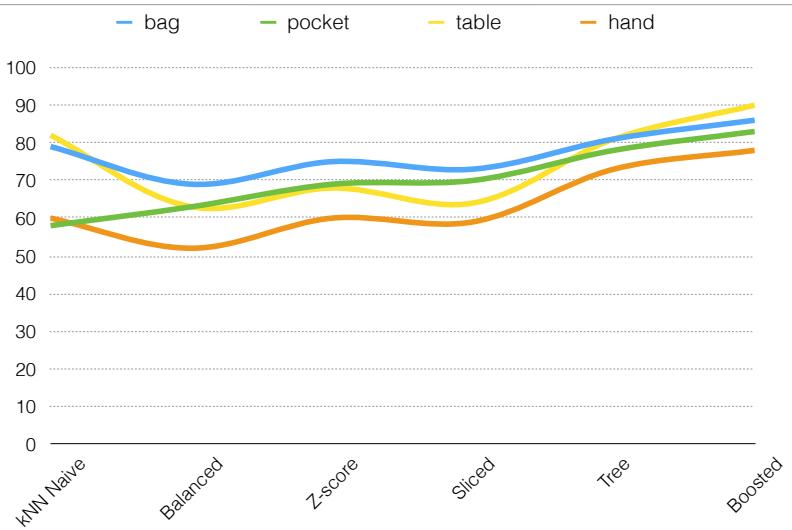
	Balanced.Accuracy	Precision	Recall	F1
Class: bag	0.91	0.86	0.86	0.86
Class: pocket	0.89	0.83	0.85	0.84
Class: hand	0.87	0.78	0.81	0.79
Class: table	0.90	0.90	0.84	0.87



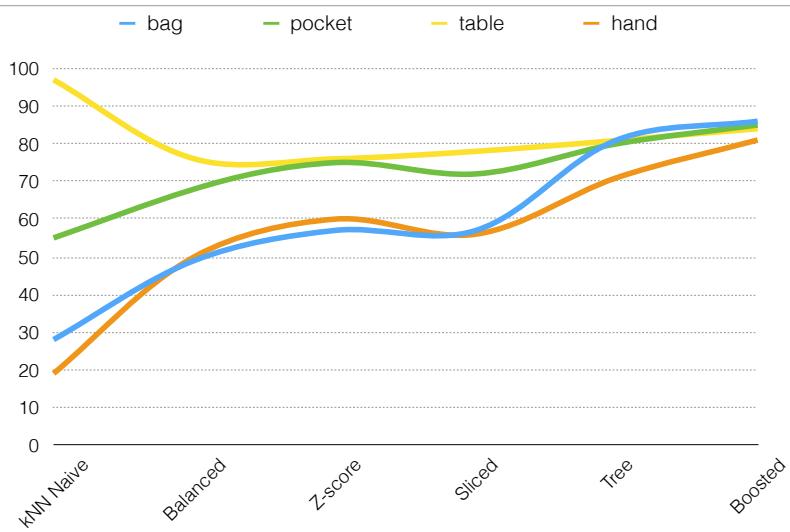
## Évolution de l'Accuracy générale



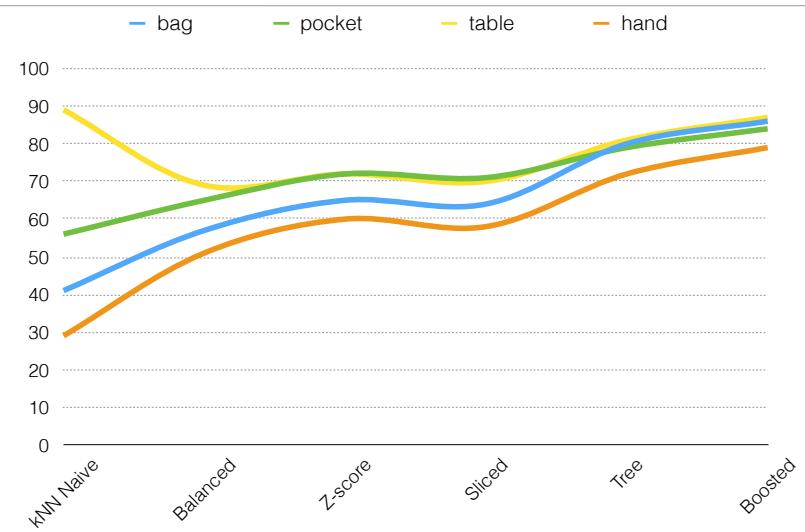
## Évolution de la Précision



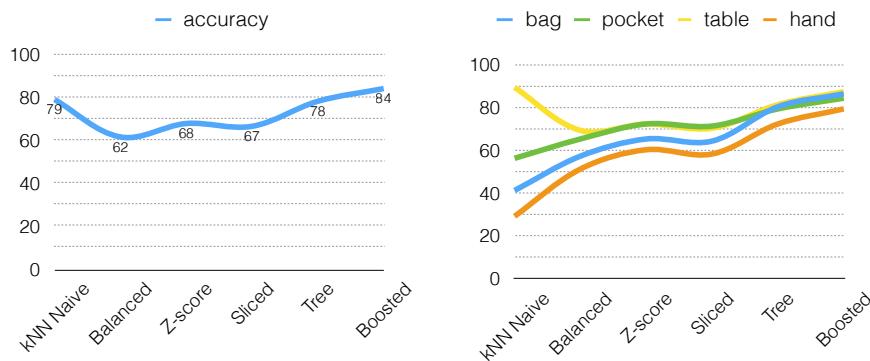
## Évolution du Rappel



## Évolution de la F-mesure



## Accuracy versus F-mesure



Faire un prédicteur, c'est simple.  
Bien le faire, c'est compliqué ...

## Publicité (e.g., INF6200)

- En tant que développeur, je fais un commit
- Des modifications concurrentes existent



Quel algorithme de fusion utiliser pour avoir le moins de chances possible de devoir résoudre un “merge conflict” ?

