

Optimizing Sports Performance:

Activity Monitoring focused on vacation period



Data Science Fullstack (ds-ft-od-03)



Sommaire:

- → Problématique & Objectifs
- → Contexte & Dataset
- → Collecte des données
- → Exploratory Data Analysis (EDA)
- → Modèles d'IA
- → Résultats et challenges





Problématique & Objectifs





Problématique

- Maintenir un niveau optimal de performance et condition physique pendant les périodes de repos
- Maximisation du temps de récupération.
- Les sportifs oublient ou évitent de partager toutes les informations.

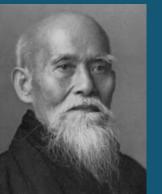


Objectif

Créer un modèle pour pouvoir identifier **passivement** le activités réalisées via le senseur (accéléromètre) intégré du smartphone 🜟



"Huit forces soutiennent la Création : Le mouvement et l'immobilité, La solidification et la fluidité , L'extension et la contraction, L'unification et la division"



Morihei Ueshiba

Morihei Ueshiba est le fondateur de l'aïkido.

14 décembre 1883 - 26 avril 1969



Contexte



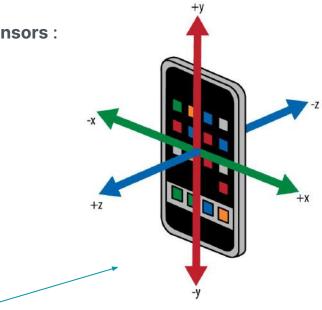


Sans y prêter attention, nos smartphones regorgent de capteurs / sensors :

- Capteurs GPS
- Capteurs de vision (cameras)
- Capteurs audio (microphones)
- Capteurs de lumière
- Capteurs de température
- Capteurs de direction (compas magnétiques)
- Capteurs d'accélération (accéléromètres)









- Et si on utilisait ce capteur (accéléromètres téléphoniques) afin d'exploration ses données ?
- Dans quel but?
 - Reconnaissance d'activité V
 - Identifier l'activité physique pratiquée par un utilisateur
 - Marche
 - Jogging
 - Descente escaliers
 - Monter escaliers
 - Assis
 - Debout







- Time series data : série chronologique agrégé
 - résument l'activité de l'utilisateur sur 10-secondes intervalles.



- <u>Et si ces données résultantes nous permettait d'induire un modèle prédictif de reconnaissance d'activité</u>?
 - o L'utilisateur *n'aurait rien à faire* 😎
 - Seulement porter son smartphone dans sa poche
 - Récolte passive de données
 - Reconnaissance d'activité physique 🔽



- Dans un but d'optimisation des performances sportives :
 - Suivi d'activité du sportif centré sur les périodes off (vacances, week-end)
- Permettrait de :
 - Générer un profil d'activité quotidien/hebdomadaire **
 - Déterminer si le sportif effectue une quantité d'activité saine pré compétition
 - Envoyer notification téléphonique →
 - Si trop d'activité (non respect temps de repos) X
 - Si pas assez d'activité (avant compétition importante par exemple)





Dataset





Laboratoire WISDM

(Wireless Sensor Data Mining)

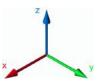


Les variables

- User : 36 utilisateurs différents
- Activity: 6 activités différentes
- Timestamp : nanosecondes
- X-axis
- Y-axis

Z-axis

3D variables (accéléromètre)





Collecte de la donnée





Contrôlée par une application créée par le Laboratoire WISDM

Application & Smartphone basé sur Android (car gratuit)

L'application :

Grâce à une simple interface utilisateur graphique, a permis d'enregistrer

- Nom de l'utilisateur,
- Démarrer/Arrêter la collecte de données
- Etiqueter l'activité en cours.

Les sujets:

Port téléphone Android : poche avant du pantalon

Doivent pratiquer les 6 activités (périodes de temps spécifiques)

Données collecté de l'accéléromètre

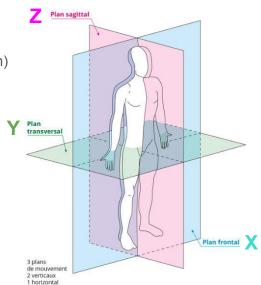
Toutes les 50 ms = **20 échantillons par seconde.**





Comment l'accéléromètre téléphonique analyse les mouvements?

- Analyse linéaire tridimensionnelle (3D)
 - Détecte et analyse les changement de vitesse
 - Mesure la force résultante exercée (composante sensible à l'accélération)
- X axis
 - Représente accélération dans le plan horizontal
 - De droite à gauche et vice versa
- Y axis
 - Représente accélération dans le plan vertical
 - De haut en bas or vice versa
- Z axis
 - Représente accélération dans le plan axial
 - **Devant / derrière** par rapport au centre de l'appareil



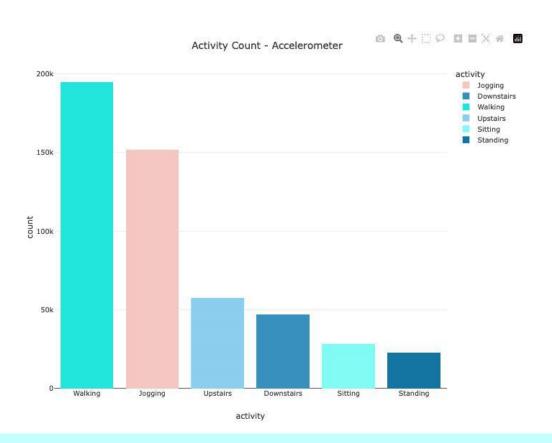


Analysons nos données







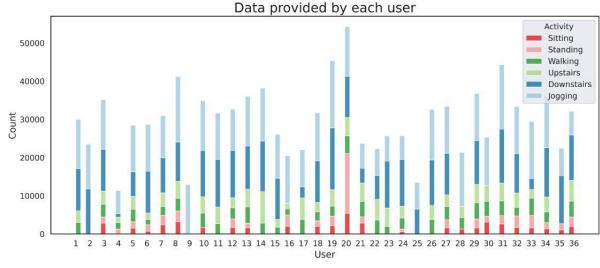


```
# Graph of acitivity frequency with matplotlib
import plotly.express as px # importations of plotly librairies
import plotly.graph_objects as go
import plotly.io as pio
pio.templates["doriane"] = go.layout.Template( # building a color template)
    layout_colorway=[
        "#f3cec9".
        "#4B9AC7",
        "#4BE8E0",
        "#9DD4F3",
        "#97FBF6",
        "#2A7FAF",
        "#23B1AB",
        "#0E3449",
        "#015955".1)
pio.templates.default = "doriane"
# Graph with diffrent colors for each activities
fig = px.bar(df sample, x='activity', color='activity')
# Setting title and figsize
fig.update layout(title="Activity Count - Accelerometer",
                  autosize=False, width=900, height=700)
fig.update_xaxes(categoryorder='array',
                 categoryarray= ['Walking', 'Jogging', 'Upstairs',
                                 'Downstairs', 'Sitting', 'Standing'])
fig.show('colab') # after ordering values, show the graph
```

What are the most frequent activities?

- Top 3 most frequent activities:
 - Walking
 - Largely above all other subclasses
 - Jogging
 - In second position
 - **Upstairs**
 - In third position



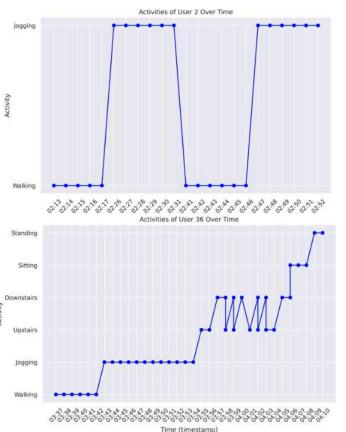


Are activities frequency different between users?

- Almost all users performed 6 different activities
- Let's take a closer look at the behavior over a timeline of some user

```
# Counting frequencies of each activity for each user
activity_counts = df_acc.groupby(['user', 'activity']).size().unstack(fill_value=0)
activity counts = activity counts[activity order]
plt.figure(figsize=(15, 6))
plt.title('Data provided by each user', fontsize=20)
sns.set_style("white")
sns.set_theme(rc={'figure.dpi': 600})
sns.histplot(data=df_acc, x='user', hue='activity', multiple='stack', palette='Paired'
# Manually create a legend based on actity order
legend labels = activity order
plt.legend(legend_labels, title='Activity', fontsize=12)
# Setting titles and showing graph
plt.xlabel('User', fontsize=14)
plt.ylabel('Count', fontsize=14)
plt.xticks(ticks=df acc['user'].unique(), fontsize=12)
plt.yticks(fontsize=12)
```





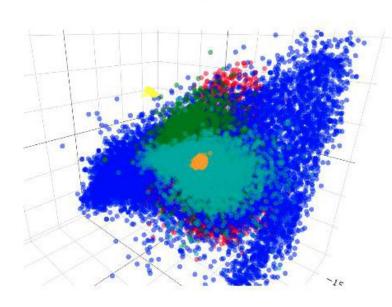
```
_ _ _ _ _ _ _ _ _ _ _ _ _ _ _
# Timestamps per user: taking the example of user number 2 and user number 36
df_acc['timestamp_str'] = df_acc['timestamp'].apply(lambda x: timedelta(microseconds=round(x, -3) // 1000))
# Filter data for user 2 and user 36
user 2 data = df acc[df acc['user'] == 2]
user 36 data = df acc[df acc['user'] == 36]
#user 2 data['timestamp hour'] = user 2 data['timestamp str'].apply(lambda x: str(x)[6:12])
#user 36 data['timestamp hour'] = user 36 data['timestamp str'].apply(lambda x: str(x)[6:12])
user_2_data.loc[:, 'timestamp_hour'] = user_2_data['timestamp_str'].apply(lambda x: str(x)[6:12])
user 36 data.loc(: 'timestamp hour' = user 36 data['timestamp str'].apply(lambda x: str(x)[6:12])
user 2 data grouped = user 2 data.groupby(['timestamp hour', 'activity']).size().reset index(name='count')
user 36 data grouped = user 36 data.groupby(['timestamp hour', 'activity']).size().reset index(name='count')
# Set up the figure with two subplots
fig, axs = plt.subplots(1, 2, figsize=(18, 6))
# Plot the activities of user 2
axs[0].plot(user 2 data grouped['timestamp hour'], user 2 data grouped['activity'], marker='o', color='blue')
axs[0].set xlabel('Time (timestamp)')
axs[0].set_ylabel('Activity')
axs[0].set_title('Activities of User 2 Over Time')
axs[0].tick params(axis='x', labelrotation=45)
axs[0].grid(True)
# Plot the activities of user 36
axs[1].plot(user 36 data grouped['timestamp hour'], user 36 data grouped['activity'], marker='o', color='blue')
axs[1].set xlabel('Time (timestamp)')
axs[1].set_ylabel('Activity')
axs[1].set_title('Activities of User 36 Over Time')
axs[1].tick params(axis='x', labelrotation=45)
axs[1].grid(True)
# Adjust layout
plt.tight layout()
# Display the plot
```

Let's have a look of activities evolution over time (2 users example)

- **b** User 2: we can see this user alternatively jogging and walking
- b User 36: this user is going through all types of the 6 activities



Activities in 3D Space - User 36



activity

- Walking
- Jogging
- Upstairs
- Downstairs Sitting
- Standing

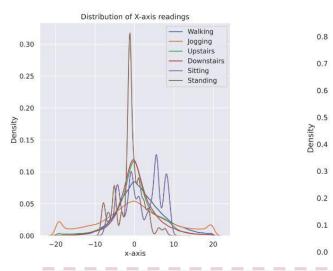
```
_____
# Let's represent in 3D all axis in relation with activities (example of user number 36)
df user36 = df acc[df acc['user'] == 36] # we select on user going through all activities
activity_colors = {'Walking': 'red', # Defining a unique color for each activities
                  'Jogging': 'blue',
                  'Upstairs': 'green'.
                  'Standing': 'orange',
                  'Upstair': 'purple',
                  'Sitting': 'yellow'}
fig = px.scatter_3d(df_user36, x='x-axis', y='y-axis', z='z-axis', # 3D graph with plotly
                   color='activity', color discrete map=activity colors,
                   title='Activities in 3D Space - User 36')
fig.update traces(marker=dict(size=4)) # Reducing size of points for better visualization
fig.update_traces(opacity=0.5) # Ajust opacity of points
fig.update_layout(legend=dict(font=dict(size=16)))
fig.show() # Show the interactive graph
```

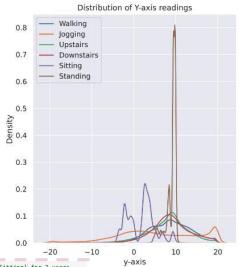
Dispersion of activities in a tri-axial plan (3D)

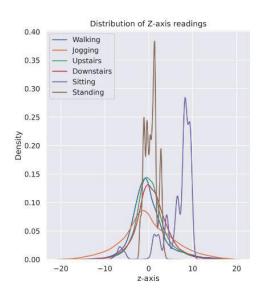
- While Jogging is the most spread variable (Walking also but less)
- Upstairs and Downstairs are spread in 3 plan but more concentrated











Let's represent in 3D all axis in relation with 2 activities (Jogging/Sitting) for 2 users x acc cols = ['x-axis'] y_acc_cols = ['y-axis'] # List of columns representing X, Y and Z axis z_acc_cols = ['z-axis'] axis_labels = ['X', 'Y', 'Z'] # Naming the axis fig, axes = plt.subplots(1, 3, figsize=(16, 6)) # Creation of subgraphs axes = axes.flatten() # For each axis we generate a plot in relation with activities (for loop) for label, cols, ax in zip(axis_labels, [x_acc_cols, y_acc_cols, z_acc_cols], axes): for activity in df acc['activity'].unique(): # Kernel density for each activities (per axis) sns.kdeplot(data=df acc[df acc['activity'] == activity][cols[0]], ax=ax, label=activity) ax.set_title(f'Distribution of {label}-axis readings') # Setting the title ax.legend() # Aggregating legend plt.tight_layout() plt.show() # Display of graphs

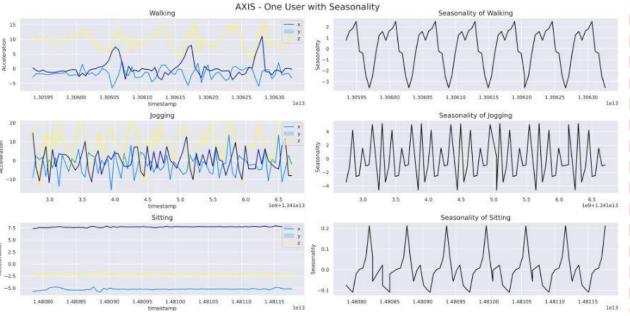
Representation of activities for each axis?

- Standing has the highest kernel density in all 3 plans
- Activities for Y axis (transversal plan) are NOT centered :
 - pant : creating a deviation V





Exemple 3 activités





- Accélérations plus élevés dans les valeurs Y dans la plupart des activités
- la marche et la course peuvent être décrits en terme du temps entre les pics
- S'asseoir présent des différences dans les valeurs des axes en raison de l'orientation du dispositif par rapport à la Terre



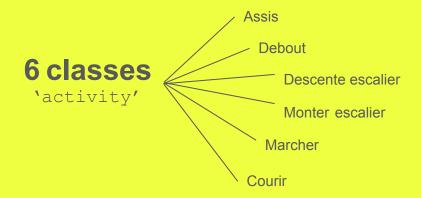
Les modèles Analyse prédictive



Type de problématique

(intelligence artificielle)

Problème de classification multiclasse



Variable cible (target)

'activity'



Recurrent Neural Network

Modèle de Deep Learning

LSTM Network





Pourquoi choisir un modèle LSTM?

- Capable de capturer et d'apprendre dans les données
- Capables de se souvenir des informations sur de longue durée
- Il n'est pas nécessaire de réaliser une extraction manuelle des caractéristiques des données des capteurs

<u>Recurrent Neural Network</u> 🧠 **LSTM Network**

Extraction des segments de données et leurs étiquettes



- **Segments**: Liste vide pour stocker les **segments** de données.
- <u>labels</u>: Liste vide pour stocker les étiquettes des activités.
- 🔔 ძ 🔭 acc: DataFrame contenant des **données d'accélération**.
- N TIME STEPS: Taille de chaque segment de données.
- **step**: Pas entre les segments Avec un chevauchement de 80 % dans ce cas).
- <u>xs, ys, zs: Valeurs d'accélération dans les segments.</u> label. L'étiquette la plus commune dans chaque segment.



Data preparation

Méthode Windowing:

- Garantie que les informations pertinentes ne soient pas perdues lors de la segmentation
- Utile quand les segment adjacents sont importantes pour l'analyse et prédiction

Recurrent Neural Network LSTM Network

🔹 The model 🚀

```
# THE MODEL (Long Short-Term Memory)

# epochs = 20
batch_size = 1024

model = Sequential()|

# RNN layer
model.add(LSTM(units = 128, input_shape= (X_train_m.shape[1], X_train_m.shape[2]))) # input_shape = 50, 3

# Dropout layer
model.add(Dropout(0.4))

# Dense layer with ReLu
model.add(Dense(units = 64, activation='relu'))

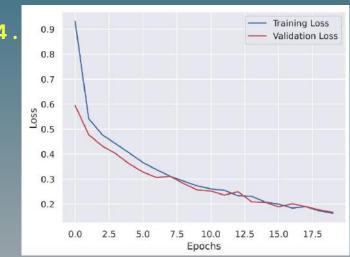
# Softmax layer
model.add(Dense(y_train_m.shape[1], activation = 'softmax'))

# Compile model
model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
```

Recurrent Neural Network 🧠

LSTM Network

- Performances
 - → Training and Validation loss over epochs



Visualize the loss during training
#
plt.plot(history.history["loss"], color="b", label="Training Loss")
plt.plot(history.history["val_loss"], color="r", label="Validation Loss")
plt.ylabel("Loss")
plt.xlabel("Epochs")
plt.tegend()
plt.show()

→ Training and Validation accuracy over epochs



```
# Use Seaborn to visualize traing and validation curves

#
plt.figure(figsize=(10, 6))
sns.lineplot(x='Epochs', y='value', hue='variable', data=pd.melt(accuracy_data, ['Epochs']))
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.title('Training and Validation Accuracy over Epochs')
plt.show()
```



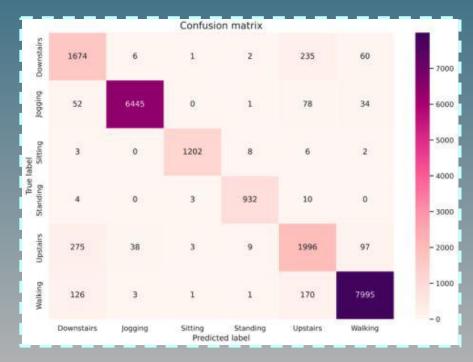
Performances

→ Weighted F1 Score obtained

F1-score for Downstairs: 0.8042 F1-score for Jogging: 0.9826 F1-score for Sitting: 0.9905 F1-score for Standing: 0.9825 F1-score for Upstairs: 0.8092 F1-score for Walking: 0.9674

Weighted F1-score: 0.9412

→ Confusion Matrix of the model





Logistic Regression

Modèle de Machine Learning

Supervisé





Train and test sets

- → We do not divide with Train Test Split
- → We do it manually



X_train & X_test preprocessing Q

Statistical characteristics in the time domain for x, y and z

Averages (classic & absolute)

Standard deviation

Minimum

Maximum

```
# Function to detect peaks in a time series
def detect peaks(data):
    peaks_indices = find_peaks(data)[0]
   if len(peaks indices) < 2:
       return np.nan
   time between peaks = np.diff(peaks indices) * 10 #Convert from indices to milliseconds (each index represents 10 ms)
   return np.mean(time between peaks) # Returns average time between peaks
# Lists to store features and labels
x list = []
y list = []
z list = []
train labels = []
time between peaks x = []
time between peaks v = 11
time between peaks z = []
window_size = 100
step size = 25
# Creation of window-size (200)
for i in range(0, df train.shape[0] - window size, step size):
   xs = df_train['x-axis'].values[i: i + window_size]
   vs = df train['v-axis'].values[i: i + window size]
   zs = df train['z-axis'].values[i: i + window size]
   label = df train['activity'][i: i + window size].value counts().idxmax() # Obtener la actividad más común
   # Calculate time between peaks for each axis
   time between peaks x.append(detect peaks(xs))
   time_between_peaks_y.append(detect_peaks(ys))
   time between peaks z.append(detect peaks(zs))
   x list.append(xs)
   y list.append(ys)
   z list.append(zs)
   train labels.append(label)
# Statistical characteristics in the time domain for x, y and z
X_train = pd.DataFrame()
# Average calculation for each axis
X train['x mean'] = pd.Series(x list).apply(lambda x: x.mean())
X_train['y_mean'] = pd.Series(y_list).apply(lambda x: x.mean())
X_train['z_mean'] = pd.Series(z_list).apply(lambda x: x.mean())
```

Logistic Regression ()

Train & test sets répartition

```
[56] # Y train and test definition
# ------
y_train = np.array(train_labels)
y_test = np.array(test_labels)
```

→ Drop of most common activity (created col) to not bias data





Preprocessing values

- → No missing values so no need to do imputations
- → We just instantiate OneHotEncoder() & StandardScaler()





- The model 🚀
 - → Classical model of Logistic Regression

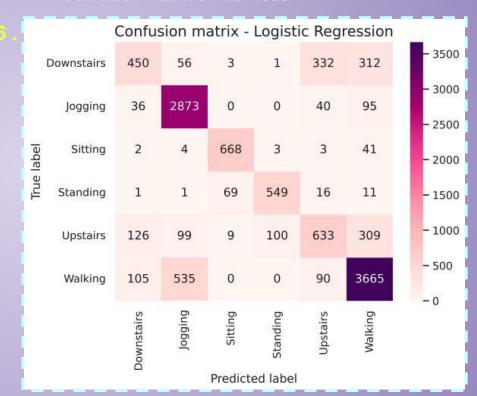
```
# Logistic regression with scores and classification reports
from sklearn.metrics import accuracy_score, f1_score, confusion_matrix
from sklearn.model_selection import GridSearchCV
lr = LogisticRegression()
lr = LogisticRegression(random_state = 21,max_iter=2000)
lr.fit(X_train_data_lr, y_train)
y_pred = lr.predict(X_test_data_lr)
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\n -----\n")
print(classification_report(y_test, y_pred))
```





Supervised model

- Performances Q
- Accuracy: 0.7865088546765151 ----- Classification Report precision recall f1-score support Downstairs 0.62 0.39 0.48 1154 0.81 0.94 0.87 3044 Jogging 0.89 0.91 Sitting 0.93 721 647 Standing 0.84 0.85 0.84 Upstairs 0.57 0.50 0.53 1276 0.83 Walking 0.83 0.83 4395 0.79 11237 accuracy macro avo 0.76 0.74 0.74 11237 weighted avg 0.78 0.79 0.78 11237





Adaboost Classifier

Combined with Logistic Regression

Ensemble





Adaboost classifier w/ Logistic Regression 🥸

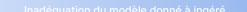


Ensemble technic

- The model 🚀 & Performances 🔍
 - Code for launching the algorithm

```
from sklearn.metrics import accuracy score, f1 score, confusion matrix
    from sklearn.ensemble import AdaBoostClassifier
    from sklearn.model selection import GridSearchCV
    print("Grid search...")
    print()
    logistic_regression = LogisticRegression(max_iter = 1000)
    adaboost logreg = AdaBoostClassifier(logistic regression)
    # Grid of values to be tested
    adaboost_logreg_params = {
        'base_estimator_C': [0.01, 0.05, 0.1, 0.5, 1.0, 5.0, 10.0, 50.0],
        'n_estimators': [5, 10, 20, 40, 60, 80, 100]
    gridsearch = GridSearchCV(adaboost_logreg, param_grid = adaboost_logreg_params, cv = 3,
    gridsearch.fit(X_train_data_lr, y_train)
    print("...Done.")
    print()
    print("Best hyperparameters : ", gridsearch.best params )
    print("Best validation accuracy : ", gridsearch.best_score_)
    print()
    GS Y test pred = gridsearch.predict(X test data lr)
    print('Predictions on the test set : ', GS Y test pred)
```

Raison décroissance des performances d'AdaBoost?



- Plutôt qu'une Logistic Regression : un RNN / CNN aurait été
- Classification report of the model

	precision	recall	f1-score	support
Downstairs	0.73	0.22	0.34	1154
Jogging	0.82	0.92	0.87	3044
Sitting	0.99	0.90	0.94	721
Standing	0.92	0.93	0.92	647
Upstairs	0.53	0.32	0.40	1276
Walking	0.71	0.87	0.78	4395
accuracy			0.76	11237
macro avg	0.78	0.69	0.71	11237
weighted avg	0.75	0.76	0.73	11237



Axes d'amélioration potentiels





Le positionnement de l'accéléromètre

Positionner l'accéléromètre sur la hanche (téléphone dans la poche) est ce suffisant ?

Mettre en relation une Montre connectée pourrait améliorer les prédictions

0

Montre connectée :

- Grâce à la mesure de la fréquence cardiaque
 - Mise en relation entre application téléphonique + la montre connectée
 - Communication entre les deux outils :
 - Ajouterait une variable
 - Augmenterait les performances algorithmique
 - Meilleure reconnaissance d'activité physique

0

Reconnaître plus d'activités que les 6 mentionnées :

On pourrait également apprendre à l'algorithme à :

- Détecter d'autres activités physique
 - Velo / Bicyclette
 - Jardinage
 - Rameur / Aviron





Obtenir un échantillon plus grand

- En statistique on sait bien que :
 - Plus l'échantillon est grand meilleur est la qualité des données exploitées
 - . Avoir plus que 36 'users' = meilleurs performances de l'algorithme

Tester different positions du téléphones lors de la récolte de données

- Avoir le téléphone dans la poche avant apporte un biais de mesure :
 - . L'utilisateur le met du côté droit ou gauche :
 - Si une asymétrie corporelle est présente alors les données varient
- Tester un positionnement au niveau du centre de gravité
 - . Pourrait potentiellement augmenter la qualité des mesures prises





- BEST MODEL → F1 score = 0.94
- From the confusion matrix :
 - Model hardly detect ascending or descending stairs
 - For the rest of activities : excellent performances

____ <u>Logistic Regression</u> 🔛

- Requires much more significant feature engineering
- Compared to the LSTM model: do not achieve as good results. F1 score 0.78

AdaBoost Classifier with LogReg (ensemble)

- Takes a lot of computation time
- Do not achieve significant results. F1 Score = 0.73
- Less interesting performances from all the above models.



Merci beaucoup!



Nous vous remercions pour votre attention.

