

HumanActivityRecognition

This project is to build a model that predicts the human activities such as Walking, Walking_Upstairs, Walking_Downstairs, Sitting, Standing or Laying.

This dataset is collected from 30 persons(referred as subjects in this dataset), performing different activities with a smartphone to their waists. The data is recorded with the help of sensors (accelerometer and Gyroscope) in that smartphone. This experiment was video recorded to label the data manually.

How data was recorded

By using the sensors(Gyroscope and accelerometer) in a smartphone, they have captured '3-axial linear acceleration'(tAcc-XYZ) from accelerometer and '3-axial angular velocity' (tGyro-XYZ) from Gyroscope with several variations.

prefix 't' in those metrics denotes time.

suffix 'XYZ' represents 3-axial signals in X , Y, and Z directions.

Feature names

1. These sensor signals are preprocessed by applying noise filters and then sampled in fixed-width windows(sliding windows) of 2.56 seconds each with 50% overlap. ie., each window has 128 readings.
2. From Each window, a feature vector was obtained by calculating variables from the time and frequency domain.

In our dataset, each datapoint represents a window with different readings

3. The accelertion signal was saperated into Body and Gravity acceleration signals(**tBodyAcc-XYZ** and **tGravityAcc-XYZ**) using some low pass filter with corner frequency of 0.3Hz.
4. After that, the body linear acceleration and angular velocity were derived in time to obtain *jerk signals* (**tBodyAccJerk-XYZ** and **tBodyGyroJerk-XYZ**).
5. The magnitude of these 3-dimensional signals were calculated using the Euclidian norm. This magnitudes are represented as features with names like *tBodyAccMag*, *tGravityAccMag*, *tBodyAccJerkMag*, *tBodyGyroMag* and *tBodyGyroJerkMag*.
6. Finally, We've got frequency domain signals from some of the available signals by applying a FFT (Fast Fourier Transform). These signals obtained were labeled with **prefix 'f'** just like original signals with **prefix 't'**. These signals are labeled as **fBodyAcc-XYZ**, **fBodyGyroMag** etc.,.
7. These are the signals that we got so far.

- tBodyAcc-XYZ
- tGravityAcc-XYZ
- tBodyAccJerk-XYZ
- tBodyGyro-XYZ
- tBodyGyroJerk-XYZ
- tBodyAccMag
- tGravityAccMag
- tBodyAccJerkMag
- tBodyGyroMag
- tBodyGyroJerkMag
- fBodyAcc-XYZ
- fBodyAccJerk-XYZ
- fBodyGyro-XYZ
- fBodyAccMag
- fBodyAccJerkMag
- fBodyGyroMag
- fBodyGyroJerkMag

8. We can esitmate some set of variables from the above signals. ie., We will estimate the following properties on each and every signal that we recoreded so far.

- **mean()**: Mean value
- **std()**: Standard deviation
- **mad()**: Median absolute deviation
- **max()**: Largest value in array

- **max()**: Largest value in array
- **min()**: Smallest value in array
- **sma()**: Signal magnitude area
- **energy()**: Energy measure. Sum of the squares divided by the number of values.
- **iqr()**: Interquartile range
- **entropy()**: Signal entropy
- **arCoeff()**: Autoregression coefficients with Burg order equal to 4
- **correlation()**: correlation coefficient between two signals
- **maxInds()**: index of the frequency component with largest magnitude
- **meanFreq()**: Weighted average of the frequency components to obtain a mean frequency
- **skewness()**: skewness of the frequency domain signal
- **kurtosis()**: kurtosis of the frequency domain signal
- **bandsEnergy()**: Energy of a frequency interval within the 64 bins of the FFT of each window.
- **angle()**: Angle between two vectors.

9. We can obtain some other vectors by taking the average of signals in a single window sample. These are used on the `angle()` variable`

- gravityMean
- tBodyAccMean
- tBodyAccJerkMean
- tBodyGyroMean
- tBodyGyroJerkMean

Y_Labels(Encoded)

- In the dataset, Y_labels are represented as numbers from 1 to 6 as their identifiers.
 - WALKING as 1
 - WALKING_UPSTAIRS as 2
 - WALKING_DOWNSTAIRS as 3
 - SITTING as 4
 - STANDING as 5
 - LAYING as 6

Train and test data were saperated

- The readings from **70%** of the volunteers were taken as **training data** and remaining **30%** subjects recordings were taken for **test data**

Data

- All the data is present in 'UCI_HAR_dataset/' folder in present working directory.
 - Feature names are present in 'UCI_HAR_dataset/features.txt'
 - **Train Data**
 - 'UCI_HAR_dataset/train/X_train.txt'
 - 'UCI_HAR_dataset/train/subject_train.txt'
 - 'UCI_HAR_dataset/train/y_train.txt'
 - **Test Data**
 - 'UCI_HAR_dataset/test/X_test.txt'
 - 'UCI_HAR_dataset/test/subject_test.txt'
 - 'UCI_HAR_dataset/test/y_test.txt'

Data Size :

27 MB

Quick overview of the dataset :

- Accelerometer and Gyroscope readings are taken from 30 volunteers(referred as subjects) while performing the following 6 Activities.

1. Walking
2. WalkingUpstairs
3. WalkingDownstairs
4. Standing
5. Sitting
6. Lying.

- Readings are divided into a window of 2.56 seconds with 50% overlapping.
- Accelerometer readings are divided into gravity acceleration and body acceleration readings, which has x,y and z components each.
- Gyroscope readings are the measure of angular velocities which has x,y and z components.
- Jerk signals are calculated for BodyAcceleration readings.
- Fourier Transforms are made on the above time readings to obtain frequency readings.
- Now, on all the base signal readings., mean, max, mad, sma, arcoefficient, engerybands,entropy etc., are calculated for each window.
- We get a feature vector of 561 features and these features are given in the dataset.
- Each window of readings is a datapoint of 561 features.

Problem Framework

- 30 subjects(volunteers) data is randomly split to 70%(21) test and 30%(7) train data.
- Each datapoint corresponds one of the 6 Activities.

Problem Statement

- Given a new datapoint we have to predict the Activity

In [1]:

```
import numpy as np
import pandas as pd

# get the features from the file features.txt
features = list()
with open('UCI_HAR_Dataset/features.txt') as f:
    features = [line.split()[1] for line in f.readlines()]
print('No of Features: {}'.format(len(features)))
```

No of Features: 561

Obtain the train data

In [2]:

```
# get the data from txt files to pandas dataffame
X_train = pd.read_csv('UCI_HAR_dataset/train/X_train.txt', delim_whitespace=True, header=None,
names=features)

# add subject column to the dataframe
X_train['subject'] = pd.read_csv('UCI_HAR_dataset/train/subject_train.txt', header=None,
squeeze=True)

y_train = pd.read_csv('UCI_HAR_dataset/train/y_train.txt', names=['Activity'], squeeze=True)
y_train_labels = y_train.map({1: 'WALKING', 2: 'WALKING_UPSTAIRS', 3: 'WALKING_DOWNSTAIRS', \
4: 'SITTING', 5: 'STANDING', 6: 'LAYING'})

# put all columns in a single dataframe
train = X_train
train['Activity'] = y_train
train['ActivityName'] = y_train_labels
train.sample()
```

C:\Users\krush\Anaconda3\lib\site-packages\pandas\io\parsers.py:678: UserWarning: Duplicate names specified. This will raise an error in the future.
return _read(filepath_or_buffer, kwds)

Out[2]:

	tBodyAcc-mean()-X	tBodyAcc-mean()-Y	tBodyAcc-mean()-Z	tBodyAcc-std()-X	tBodyAcc-std()-Y	tBodyAcc-std()-Z	tBodyAcc-mad()-X	tBodyAcc-mad()-Y	tBodyAcc-mad()-Z	tBodyAcc-max()-X	..
732	0.27407	0.001963	-0.119016	-0.975956	-0.907987	-0.93833	-0.980617	-0.91002	-0.941729	-0.909069	..

1 rows × 564 columns

◀		▶
---	--	---

In [3]:

```
train.shape
```

Out[3]:

(7352, 564)

Obtain the test data

In [4]:

```
# get the data from txt files to pandas dataffame
X_test = pd.read_csv('UCI_HAR_dataset/test/X_test.txt', delim_whitespace=True, header=None, names=features)

# add subject column to the dataframe
X_test['subject'] = pd.read_csv('UCI_HAR_dataset/test/subject_test.txt', header=None, squeeze=True)

# get y labels from the txt file
y_test = pd.read_csv('UCI_HAR_dataset/test/y_test.txt', names=['Activity'], squeeze=True)
y_test_labels = y_test.map({1: 'WALKING', 2: 'WALKING_UPSTAIRS', 3: 'WALKING_DOWNSTAIRS', \
                             4: 'SITTING', 5: 'STANDING', 6: 'LAYING'})

# put all columns in a single dataframe
test = X_test
test['Activity'] = y_test
test['ActivityName'] = y_test_labels
test.sample()
```

C:\Users\krush\Anaconda3\lib\site-packages\pandas\io\parsers.py:678: UserWarning: Duplicate names specified. This will raise an error in the future.
return _read(filepath_or_buffer, kwds)

Out[4]:

	tBodyAcc-mean()-X	tBodyAcc-mean()-Y	tBodyAcc-mean()-Z	tBodyAcc-std()-X	tBodyAcc-std()-Y	tBodyAcc-std()-Z	tBodyAcc-mad()-X	tBodyAcc-mad()-Y	tBodyAcc-mad()-Z	tBodyAcc-max()-X
2040	0.274196	-0.023026	-0.118437	-0.991049	-0.92557	-0.960282	-0.991526	-0.922057	-0.959172	-0.937224

1 rows × 564 columns

◀		▶
---	--	---

In [5]:

```
test.shape
```

Out[5]:

(2947, 564)

Data Cleaning

1. Check for Duplicates

In [6]:

```
print('No of duplicates in train: {}'.format(sum(train.duplicated())))
print('No of duplicates in test : {}'.format(sum(test.duplicated())))
```

No of duplicates in train: 0
No of duplicates in test : 0

2. Checking for NaN/null values

In [7]:

```
print('We have {} NaN/Null values in train'.format(train.isnull().values.sum()))
print('We have {} NaN/Null values in test'.format(test.isnull().values.sum()))
```

We have 0 NaN/Null values in train
We have 0 NaN/Null values in test

3. Check for data imbalance

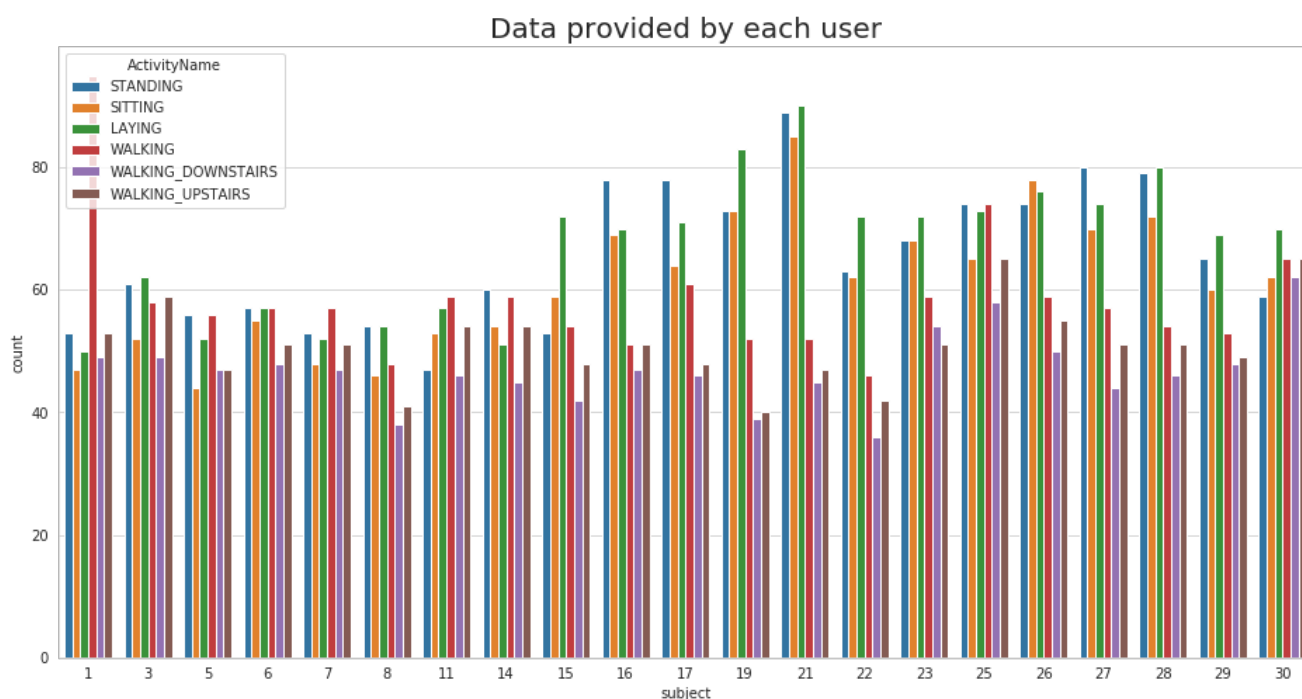
In [8]:

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
sns.set_style('whitegrid')
plt.rcParams['font.family'] = 'Dejavu Sans'
```

In [9]:

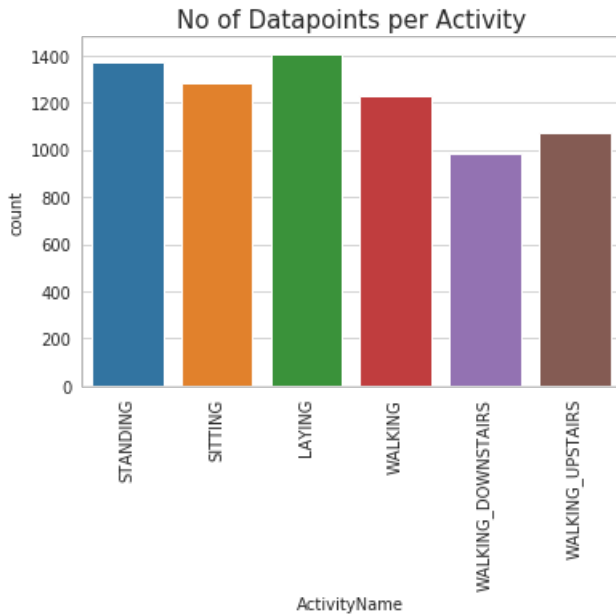
```
plt.figure(figsize=(16,8))
plt.title('Data provided by each user', fontsize=20)
sns.countplot(x='subject',hue='ActivityName', data = train)
plt.show()
```



We have got almost same number of reading from all the subjects

In [10]:

```
plt.title('No of Datapoints per Activity', fontsize=15)
sns.countplot(train.ActivityName)
plt.xticks(rotation=90)
plt.show()
```



Observation

Our data is well balanced (almost)

4. Changing feature names

In [11]:

```
columns = train.columns

# Removing '()' from column names
columns = columns.str.replace('()', '')
columns = columns.str.replace('[-]', '')
columns = columns.str.replace('[,]', '')

train.columns = columns
test.columns = columns

test.columns
```

Out[11]:

```
Index(['tBodyAccmeanX', 'tBodyAccmeanY', 'tBodyAccmeanZ', 'tBodyAccstdX',
      'tBodyAccstdY', 'tBodyAccstdZ', 'tBodyAccmadX', 'tBodyAccmadY',
      'tBodyAccmadZ', 'tBodyAccmaxX',
      ...,
      'angletBodyAccMeangravity', 'angletBodyAccJerkMeangravityMean',
      'angletBodyGyroMeangravityMean', 'angletBodyGyroJerkMeangravityMean',
      'angleXgravityMean', 'angleYgravityMean', 'angleZgravityMean',
      'subject', 'Activity', 'ActivityName'],
      dtype='object', length=564)
```

5. Save this dataframe in a csv files

In [12]:

```
train.to_csv('UCI_HAR_Dataset/csv_files/train.csv', index=False)
test.to_csv('UCI_HAR_Dataset/csv_files/test.csv', index=False)
```

Exploratory Data Analysis

"Without domain knowledge EDA has no meaning, without EDA a problem has no soul."

1. Featuring Engineering from Domain Knowledge

- **Static and Dynamic Activities**

- In static activities (sit, stand, lie down) motion information will not be very useful.
- In the dynamic activities (Walking, WalkingUpstairs, WalkingDownstairs) motion info will be significant.

2. Stationary and Moving activities are completely different

In [13]:

```
sns.set_palette("Set1", desat=0.80)
facetgrid = sns.FacetGrid(train, hue='ActivityName', size=6, aspect=2)
facetgrid.map(sns.distplot, 'tBodyAccMagmean', hist=False)\
    .add_legend()
plt.annotate("Stationary Activities", xy=(-0.956, 17), xytext=(-0.9, 23), size=20,\
    va='center', ha='left',\
    arrowprops=dict(arrowstyle="simple", connectionstyle="arc3,rad=0.1"))

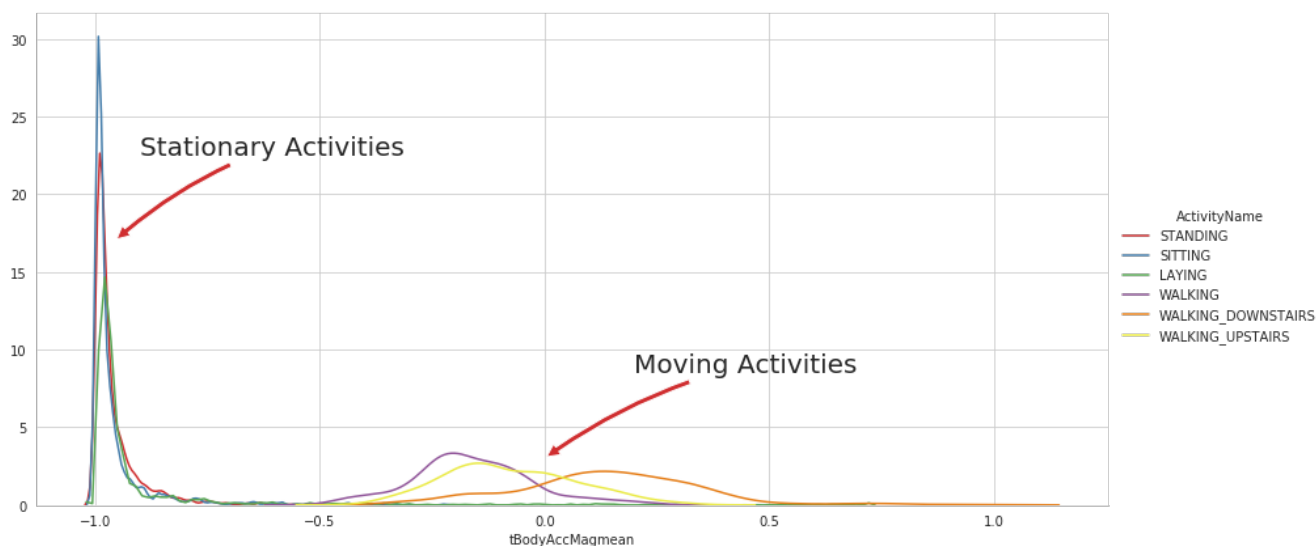
plt.annotate("Moving Activities", xy=(0, 3), xytext=(0.2, 9), size=20,\
    va='center', ha='left',\
    arrowprops=dict(arrowstyle="simple", connectionstyle="arc3,rad=0.1"))
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\seaborn\axisgrid.py:230: UserWarning: The `size` paramter has been renamed to `height`; please update your code.

warnings.warn(msg, UserWarning)

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval



In [14]:

```
# for plotting purposes taking datapoints of each activity to a different dataframe
df1 = train[train['Activity']==1]
df2 = train[train['Activity']==2]
df3 = train[train['Activity']==3]
df4 = train[train['Activity']==4]
df5 = train[train['Activity']==5]
df6 = train[train['Activity']==6]

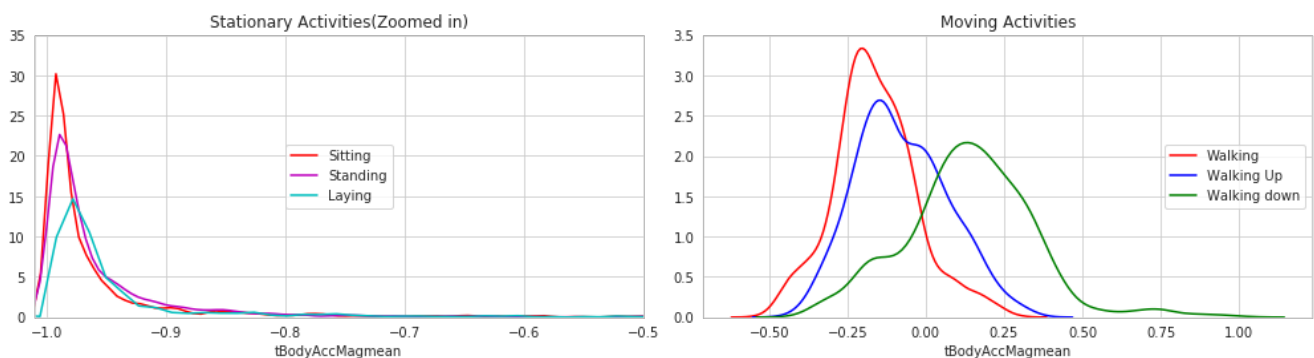
plt.figure(figsize=(14,7))
plt.subplot(2,2,1)
plt.title('Stationary Activities(Zoomed in)')
sns.distplot(df4['tBodyAccMagmean'],color = 'r',hist = False, label = 'Sitting')
sns.distplot(df5['tBodyAccMagmean'],color = 'm',hist = False,label = 'Standing')
sns.distplot(df6['tBodyAccMagmean'],color = 'c',hist = False, label = 'Laying')
plt.axis([-1.01, -0.5, 0, 35])
plt.legend(loc='center')

plt.subplot(2,2,2)
plt.title('Moving Activities')
sns.distplot(df1['tBodyAccMagmean'],color = 'red',hist = False, label = 'Walking')
sns.distplot(df2['tBodyAccMagmean'],color = 'blue',hist = False,label = 'Walking Up')
sns.distplot(df3['tBodyAccMagmean'],color = 'green',hist = False, label = 'Walking down')
plt.legend(loc='center right')

plt.tight_layout()
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

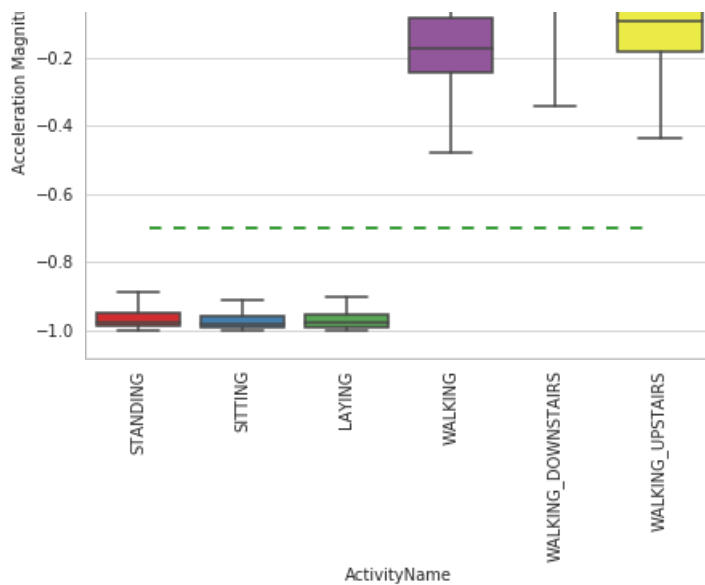


3. Magnitude of an acceleration can saperate it well

In [15]:

```
plt.figure(figsize=(7,7))
sns.boxplot(x='ActivityName', y='tBodyAccMagmean',data=train, showfliers=False, saturation=1)
plt.ylabel('Acceleration Magnitude mean')
plt.axhline(y=-0.7, xmin=0.1, xmax=0.9,dashes=(5,5), c='g')
plt.axhline(y=-0.05, xmin=0.4, dashes=(5,5), c='m')
plt.xticks(rotation=90)
plt.show()
```



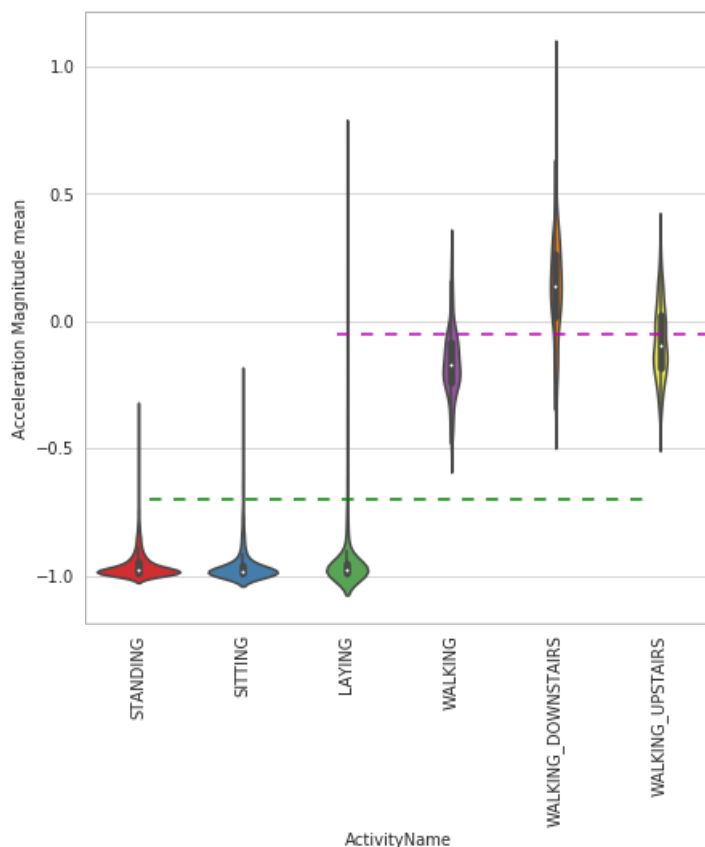


In [23]:

```
plt.figure(figsize=(7,7))
sns.violinplot(x='ActivityName', y='tBodyAccMagmean',data=train, showfliers=False, saturation=1)
plt.ylabel('Acceleration Magnitude mean')
plt.axhline(y=-0.7, xmin=0.1, xmax=0.9,dashes=(5,5), c='g')
plt.axhline(y=-0.05, xmin=0.4, dashes=(5,5), c='m')
plt.xticks(rotation=90)
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



Observations:

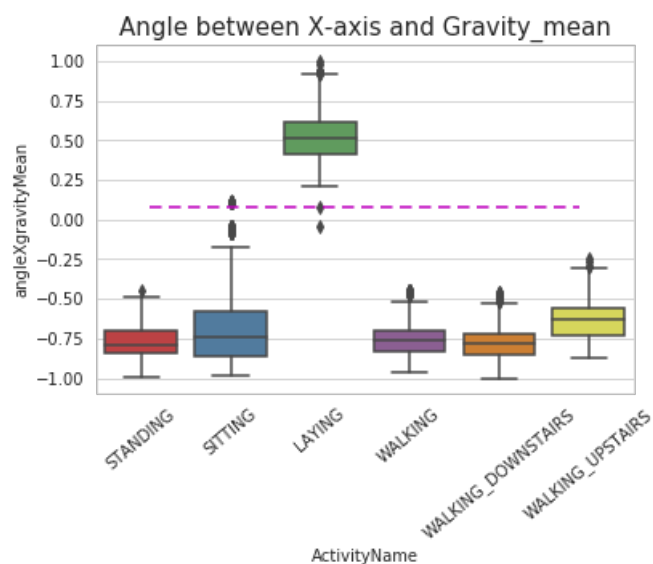
1. The mean acceleration magnitude for walking is significantly higher than for standing, sitting, and laying.

- If tAccMean is < -0.8 then the Activities are either Standing or Sitting or Laying.
- If tAccMean is > -0.6 then the Activities are either Walking or WalkingDownstairs or WalkingUpstairs.
- If tAccMean > 0.0 then the Activity is WalkingDownstairs.
- We can classify 75% the Activity labels with some errors.

4. Position of GravityAccelerationComponents also matters

In [16]:

```
sns.boxplot(x='ActivityName', y='angleXgravityMean', data=train)
plt.axhline(y=0.08, xmin=0.1, xmax=0.9, c='m', dashes=(5,3))
plt.title('Angle between X-axis and Gravity_mean', fontsize=15)
plt.xticks(rotation = 40)
plt.show()
```

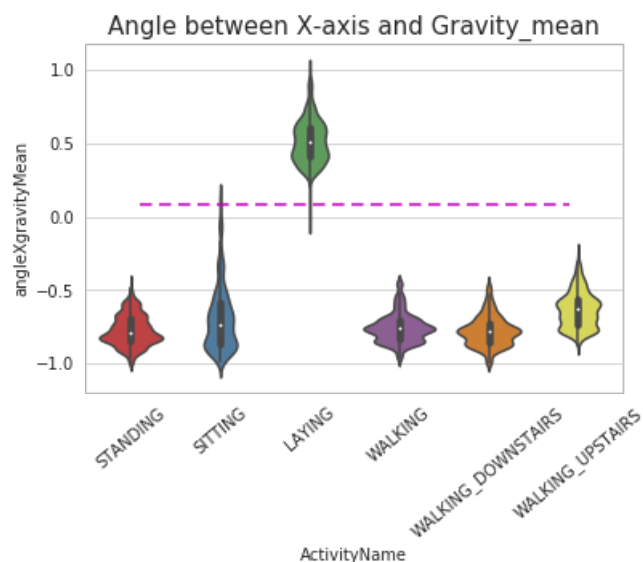


In [24]:

```
sns.violinplot(x='ActivityName', y='angleXgravityMean', data=train)
plt.axhline(y=0.08, xmin=0.1, xmax=0.9, c='m', dashes=(5,3))
plt.title('Angle between X-axis and Gravity_mean', fontsize=15)
plt.xticks(rotation = 40)
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

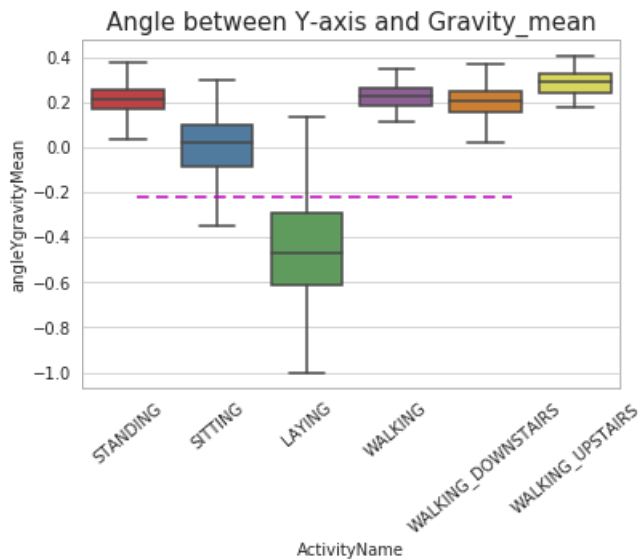


Observations:

- If $\text{angleXgravityMean} > 0$ then Activity is Laying.
- We can classify all datapoints belonging to Laying activity with just a single if else statement.

In [17]:

```
sns.boxplot(x='ActivityName', y='angleYgravityMean', data = train, showfliers=False)
plt.title('Angle between Y-axis and Gravity_mean', fontsize=15)
plt.xticks(rotation = 40)
plt.axhline(y=-0.22, xmin=0.1, xmax=0.8, dashes=(5,3), c='m')
plt.show()
```

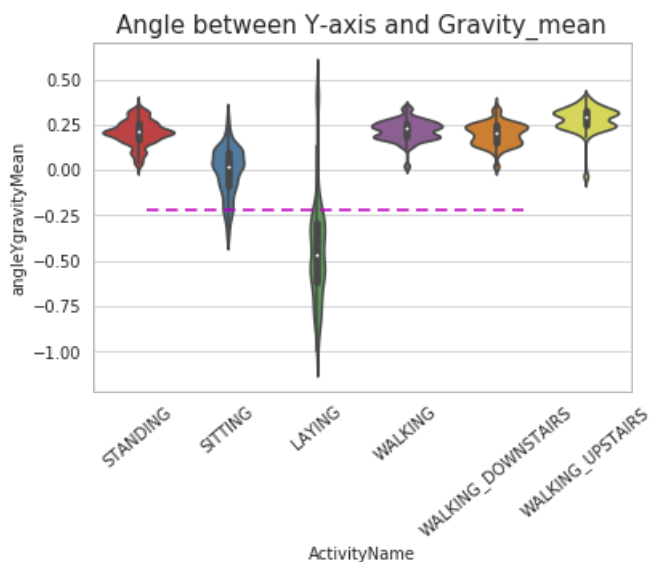


In [25]:

```
sns.violinplot(x='ActivityName', y='angleYgravityMean', data = train, showfliers=False)
plt.title('Angle between Y-axis and Gravity_mean', fontsize=15)
plt.xticks(rotation = 40)
plt.axhline(y=-0.22, xmin=0.1, xmax=0.8, dashes=(5,3), c='m')
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

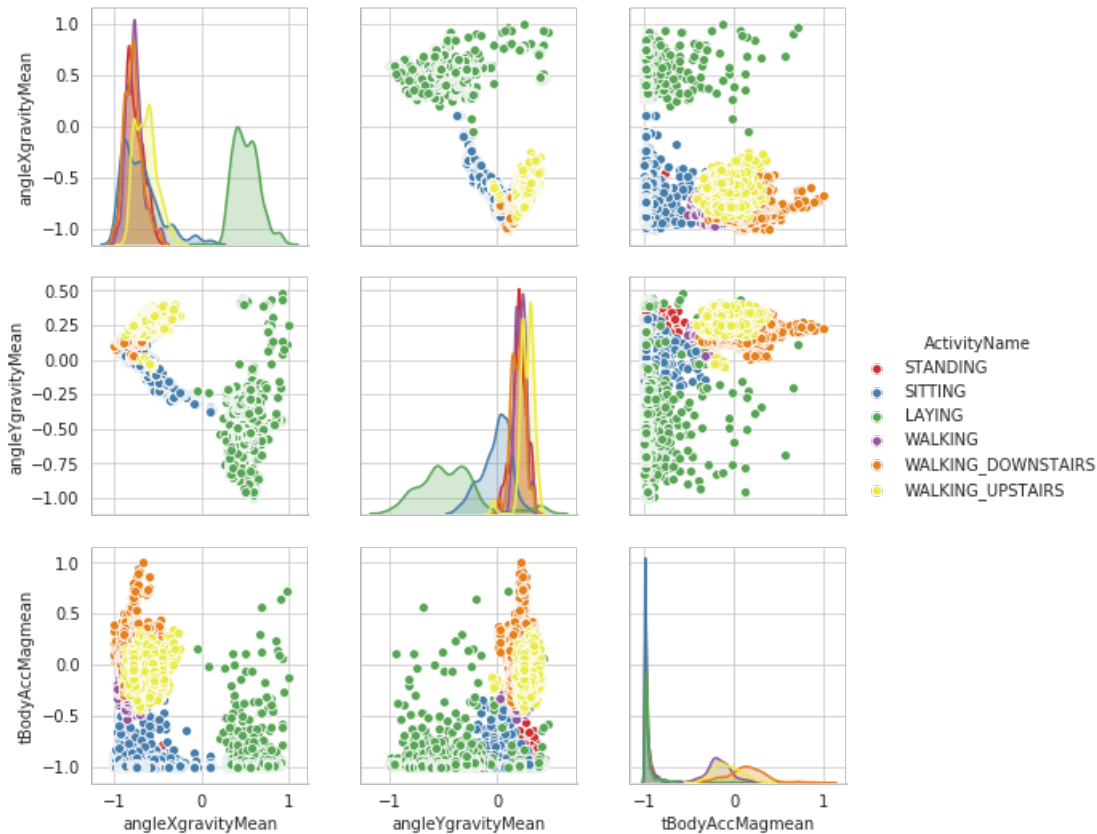


In [22]:

```
n = train.shape[0]
sns.pairplot(train[['angleXgravityMean', 'angleYgravityMean', 'tBodyAccMagmean', 'ActivityName']][0:n], hue='ActivityName', vars=['angleXgravityMean', 'angleYgravityMean', 'tBodyAccMagmean'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



In [26]:

```
train.columns
```

Out[26]:

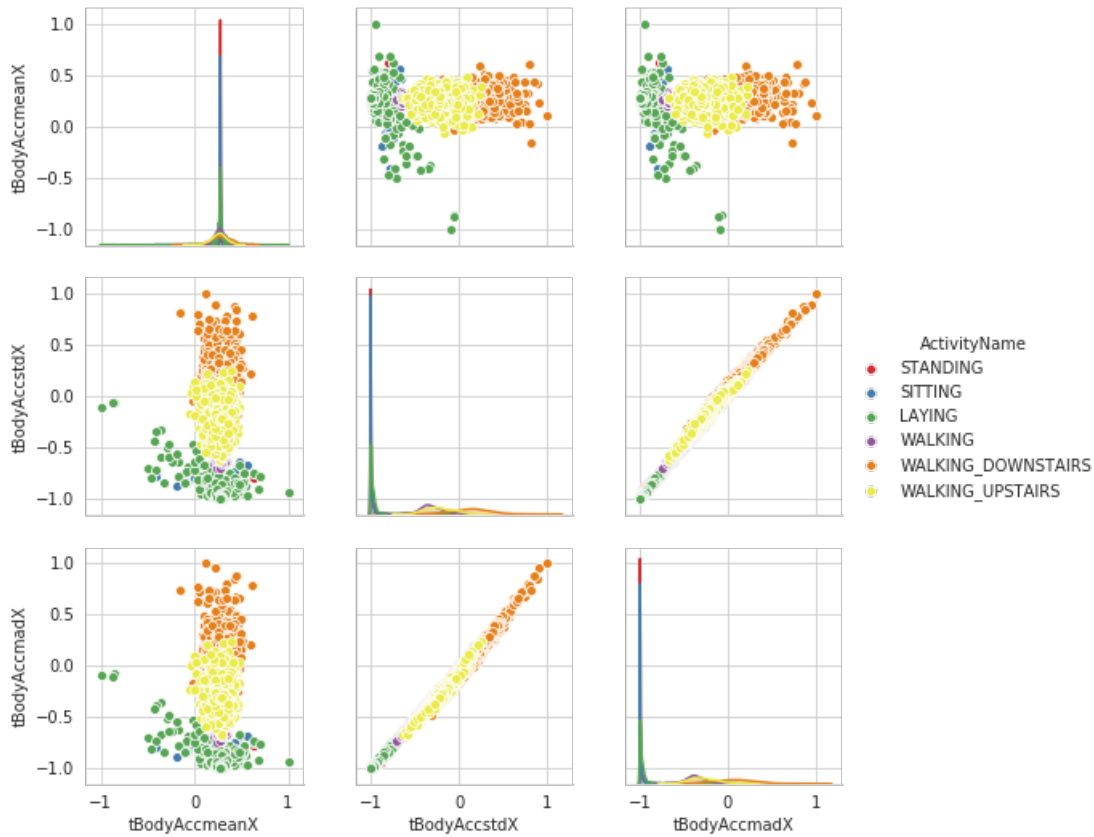
```
Index(['tBodyAccmeanX', 'tBodyAccmeanY', 'tBodyAccmeanZ', 'tBodyAccstdX',
      'tBodyAccstdY', 'tBodyAccstdZ', 'tBodyAccmadX', 'tBodyAccmadY',
      'tBodyAccmadZ', 'tBodyAccmaxX',
      ...,
      'angletBodyAccMeangravity', 'angletBodyAccJerkMeangravityMean',
      'angletBodyGyroMeangravityMean', 'angletBodyGyroJerkMeangravityMean',
      'angleXgravityMean', 'angleYgravityMean', 'angleZgravityMean',
      'subject', 'Activity', 'ActivityName'],
      dtype='object', length=564)
```

In [27]:

```
n = train.shape[0]
sns.pairplot(train[['tBodyAccmeanX', 'tBodyAccstdX', 'tBodyAccmadX', 'ActivityName']][0:n], hue='ActivityName', vars=['tBodyAccmeanX', 'tBodyAccstdX', 'tBodyAccmadX'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

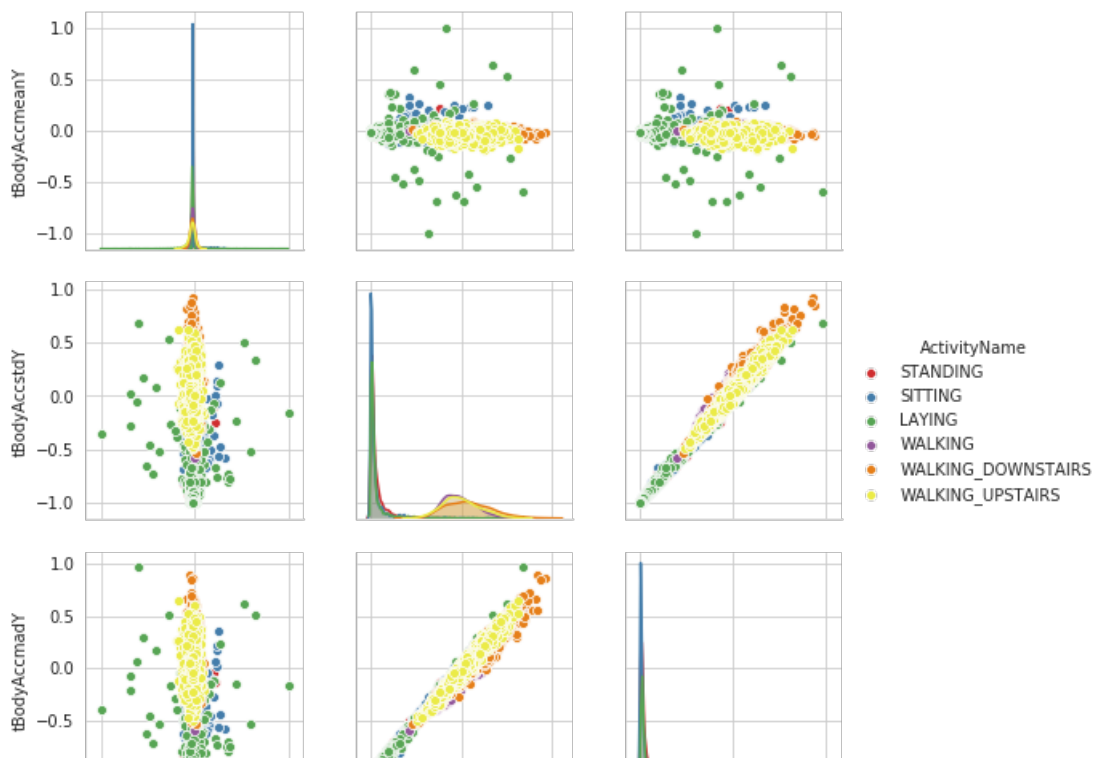


In [28]:

```
n = train.shape[0]
sns.pairplot(train[['tBodyAccmeanY', 'tBodyAccstdY', 'tBodyAccmadY', 'ActivityName']][0:n], hue='ActivityName', vars=['tBodyAccmeanY', 'tBodyAccstdY', 'tBodyAccmadY'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



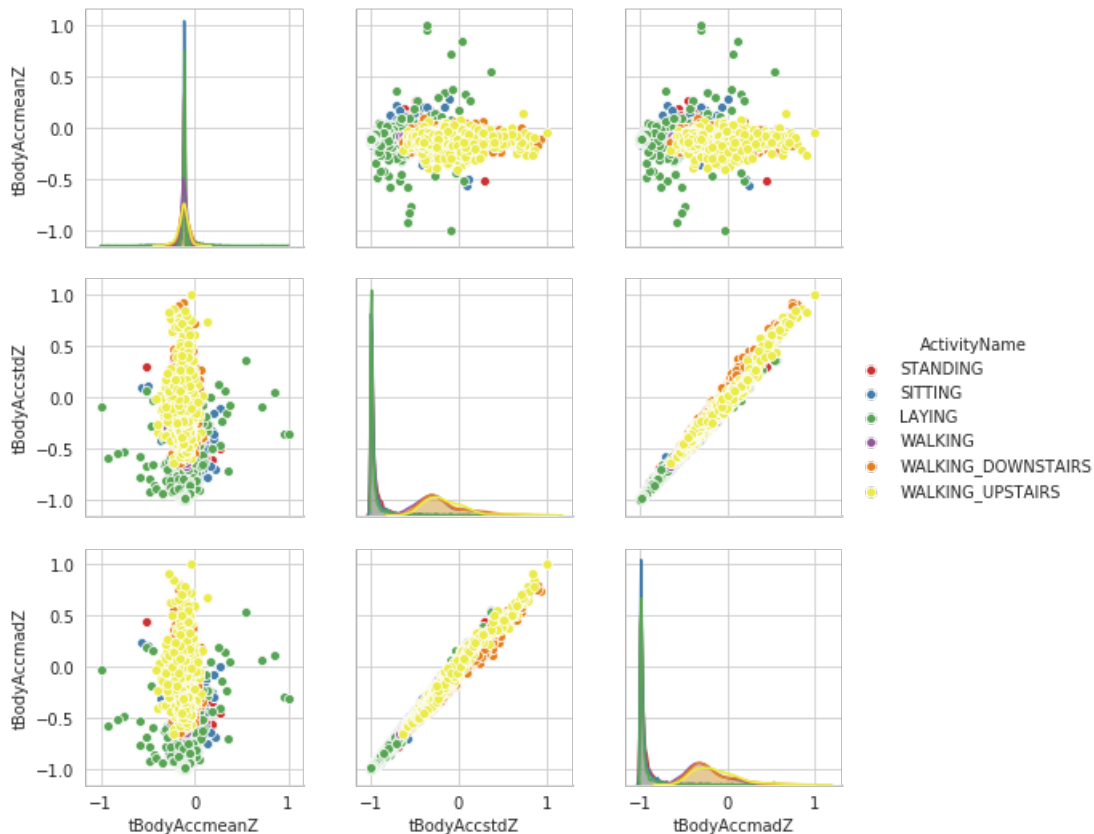


In [30]:

```
n = train.shape[0]
sns.pairplot(train[['tBodyAccmeanZ', 'tBodyAccstdZ', 'tBodyAccmadZ', 'ActivityName']][0:n], hue='ActivityName', vars=['tBodyAccmeanZ', 'tBodyAccstdZ', 'tBodyAccmadZ'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```

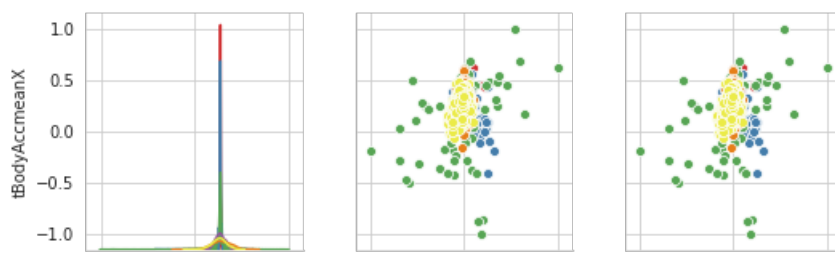


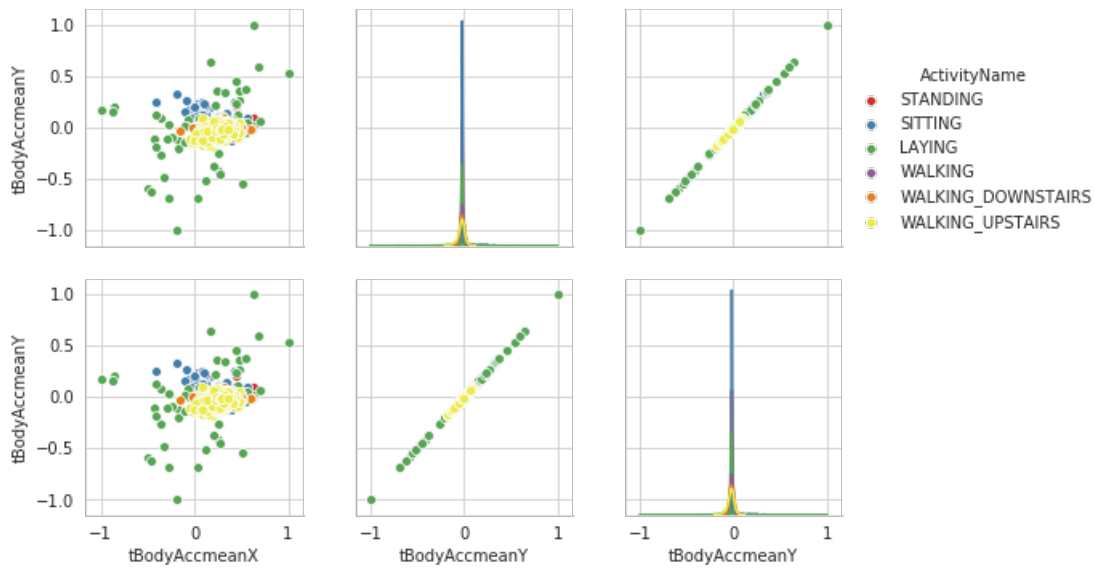
In [31]:

```
n = train.shape[0]
sns.pairplot(train[['tBodyAccmeanX', 'tBodyAccmeanY', 'tBodyAccmeanZ', 'ActivityName']][0:n], hue='ActivityName', vars=['tBodyAccmeanX', 'tBodyAccmeanY', 'tBodyAccmeanZ'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



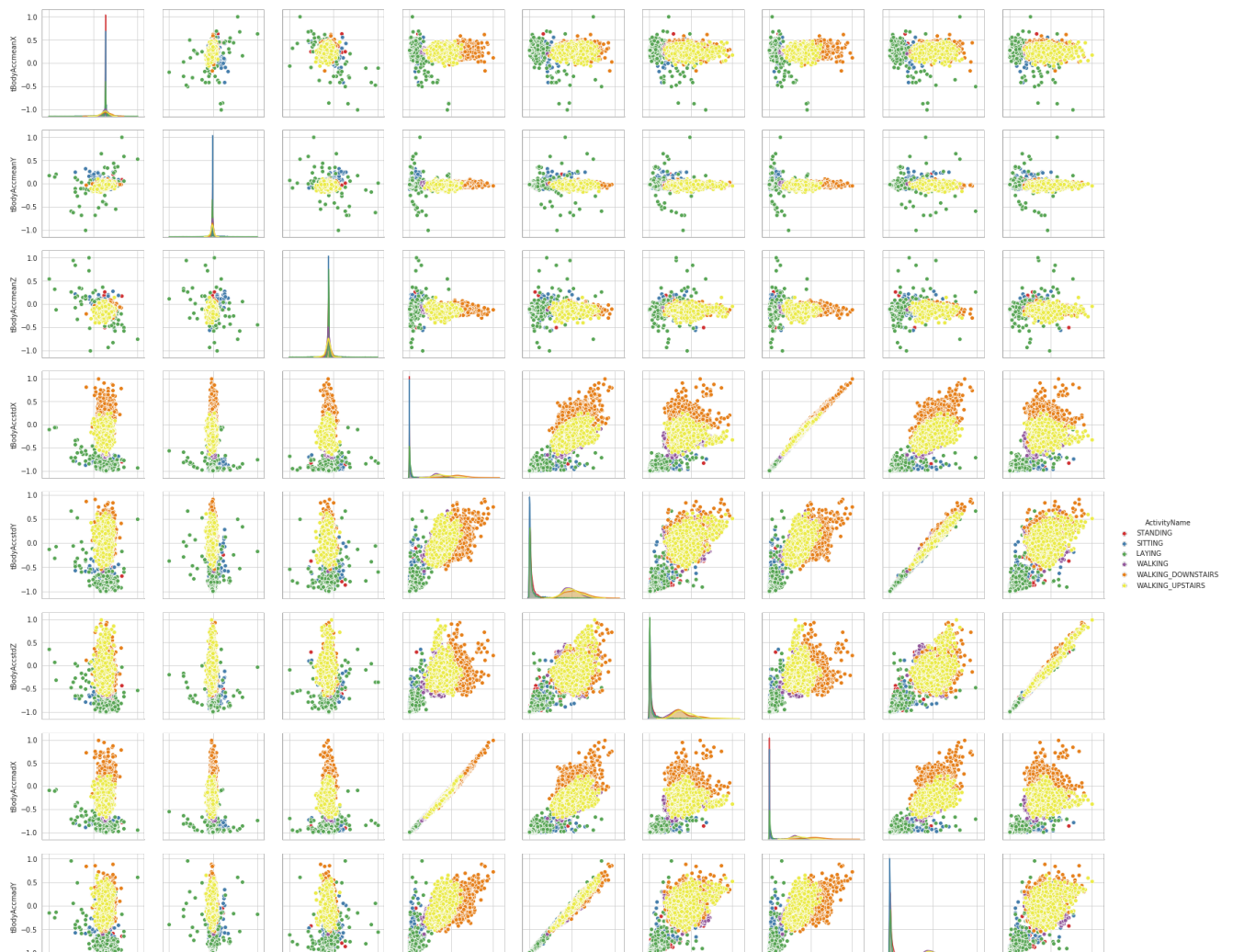


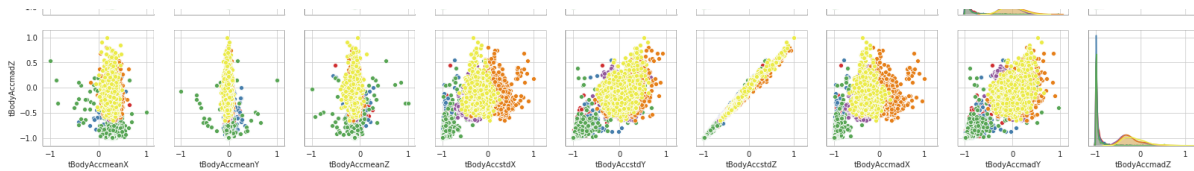
In [34]:

```
n = train.shape[0]
sns.pairplot(train[['tBodyAccmeanX', 'tBodyAccmeanY', 'tBodyAccmeanZ', 'tBodyAccstdX', 'tBodyAccstdY',
'tBodyAccstdZ', 'tBodyAccmadX', 'tBodyAccmadY', 'tBodyAccmadZ', 'ActivityName']][0:n], hue='ActivityName',
vars=['tBodyAccmeanX', 'tBodyAccmeanY', 'tBodyAccmeanZ', 'tBodyAccstdX', 'tBodyAccstdY',
'tBodyAccstdZ', 'tBodyAccmadX', 'tBodyAccmadY', 'tBodyAccmadZ'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



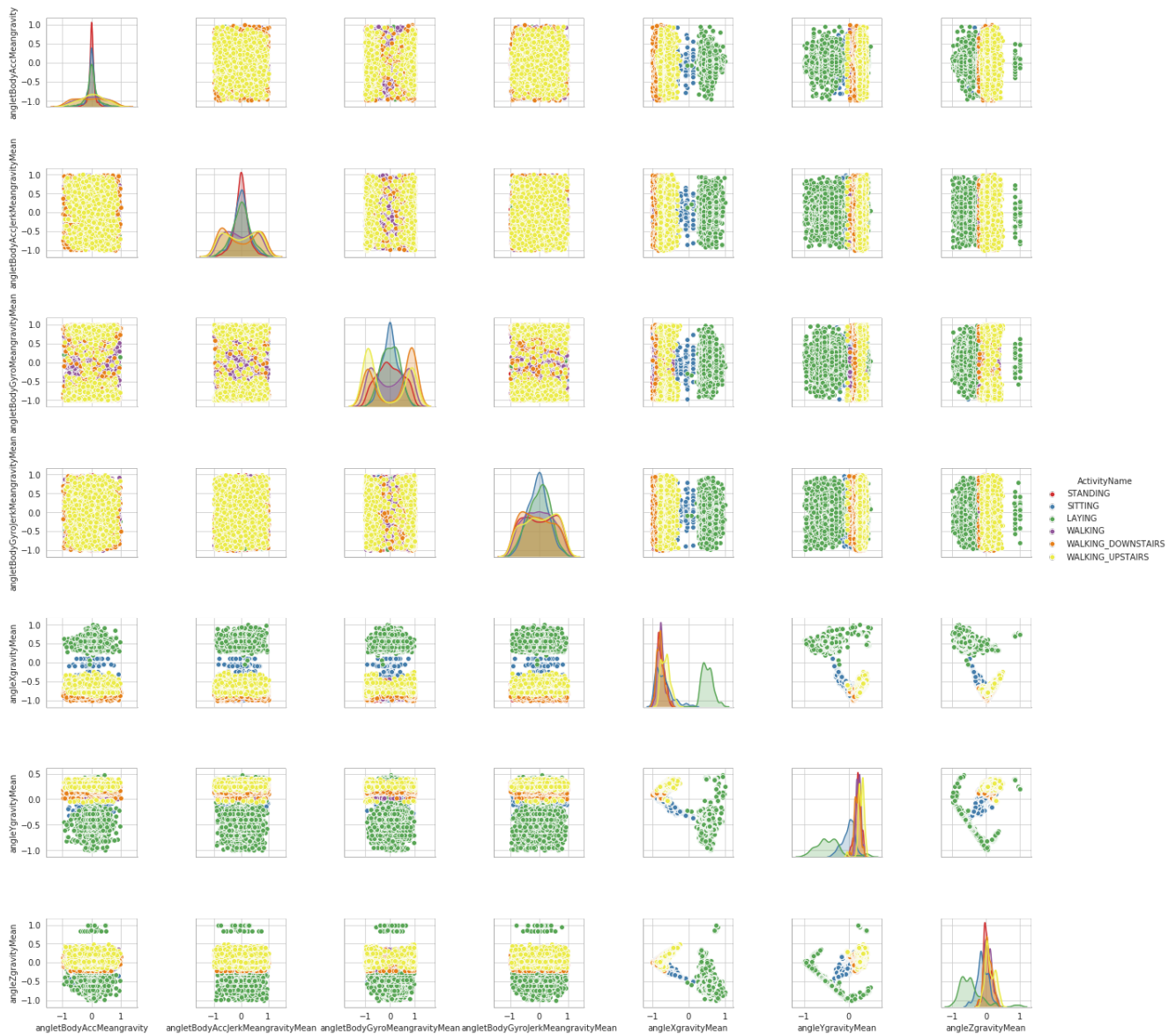


In [35]:

```
n = train.shape[0]
sns.pairplot(train[['angletBodyAccMeangravity', 'angletBodyAccJerkMeangravityMean',
                    'angletBodyGyroMeangravityMean', 'angletBodyGyroJerkMeangravityMean',
                    'angleXgravityMean', 'angleYgravityMean', 'angleZgravityMean', 'ActivityName']] [0:n], hue='ActivityName',
              vars=['angletBodyAccMeangravity', 'angletBodyAccJerkMeangravityMean',
                    'angletBodyGyroMeangravityMean', 'angletBodyGyroJerkMeangravityMean',
                    'angleXgravityMean', 'angleYgravityMean', 'angleZgravityMean'])
plt.show()
```

C:\Users\krush\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
return np.add.reduce(sorted[indexer] * weights, axis=axis) / sumval
```



Apply t-sne on the data

In [36]:

```
import numpy as np
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt
import seaborn as sns
```

In [37]:

```
# performs t-sne with different perplexity values and their repective plots..

def perform_tsne(X_data, y_data, perplexities, n_iter=1000, img_name_prefix='t-sne'):

    for index,perplexity in enumerate(perplexities):
        # perform t-sne
        print('\nperforming tsne with perplexity {} and with {} iterations at max'.format(perplexity, n_iter))
        X_reduced = TSNE(verbose=2, perplexity=perplexity).fit_transform(X_data)
        print('Done..')

        # prepare the data for seaborn
        print('Creating plot for this t-sne visualization..')
        df = pd.DataFrame({'x':X_reduced[:,0], 'y':X_reduced[:,1], 'label':y_data})

        # draw the plot in appropriate place in the grid
        sns.lmplot(data=df, x='x', y='y', hue='label', fit_reg=False, size=8,\
                    palette="Set1",markers=['^','v','s','o', '1','2'])
        plt.title("perplexity : {} and max_iter : {}".format(perplexity, n_iter))
        img_name = img_name_prefix + '_perp_{}_iter_{}.png'.format(perplexity, n_iter)
        print('saving this plot as image in present working directory...')
        plt.savefig(img_name)
        plt.show()
        print('Done')
```

In [38]:

```
X_pre_tsne = train.drop(['subject', 'Activity','ActivityName'], axis=1)
y_pre_tsne = train['ActivityName']
perform_tsne(X_data = X_pre_tsne,y_data=y_pre_tsne, perplexities =[2,5,10,20,50])
```

```
performing tsne with perplexity 2 and with 1000 iterations at max
[t-SNE] Computing 7 nearest neighbors...
[t-SNE] Indexed 7352 samples in 0.617s...
[t-SNE] Computed neighbors for 7352 samples in 62.024s...
[t-SNE] Computed conditional probabilities for sample 1000 / 7352
[t-SNE] Computed conditional probabilities for sample 2000 / 7352
[t-SNE] Computed conditional probabilities for sample 3000 / 7352
[t-SNE] Computed conditional probabilities for sample 4000 / 7352
[t-SNE] Computed conditional probabilities for sample 5000 / 7352
[t-SNE] Computed conditional probabilities for sample 6000 / 7352
[t-SNE] Computed conditional probabilities for sample 7000 / 7352
[t-SNE] Computed conditional probabilities for sample 7352 / 7352
[t-SNE] Mean sigma: 0.635855
[t-SNE] Computed conditional probabilities in 0.169s
[t-SNE] Iteration 50: error = 124.7786331, gradient norm = 0.0249788 (50 iterations in 9.082s)
[t-SNE] Iteration 100: error = 107.0525131, gradient norm = 0.0274130 (50 iterations in 5.152s)
[t-SNE] Iteration 150: error = 100.7961349, gradient norm = 0.0203417 (50 iterations in 4.048s)
[t-SNE] Iteration 200: error = 97.4336472, gradient norm = 0.0177385 (50 iterations in 3.876s)
[t-SNE] Iteration 250: error = 95.1421814, gradient norm = 0.0137455 (50 iterations in 3.791s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 95.142181
[t-SNE] Iteration 300: error = 4.1165681, gradient norm = 0.0015623 (50 iterations in 3.666s)
[t-SNE] Iteration 350: error = 3.2067621, gradient norm = 0.0010015 (50 iterations in 3.752s)
[t-SNE] Iteration 400: error = 2.7779393, gradient norm = 0.0007213 (50 iterations in 3.865s)
[t-SNE] Iteration 450: error = 2.5142267, gradient norm = 0.0005690 (50 iterations in 3.589s)
[t-SNE] Iteration 500: error = 2.3316514, gradient norm = 0.0004754 (50 iterations in 3.631s)
[t-SNE] Iteration 550: error = 2.1938930, gradient norm = 0.0004095 (50 iterations in 3.638s)
[t-SNE] Iteration 600: error = 2.0844235, gradient norm = 0.0003672 (50 iterations in 3.573s)
[t-SNE] Iteration 650: error = 1.9945434, gradient norm = 0.0003335 (50 iterations in 3.804s)
[t-SNE] Iteration 700: error = 1.9190723, gradient norm = 0.0003025 (50 iterations in 3.549s)
[t-SNE] Iteration 750: error = 1.8541486, gradient norm = 0.0002795 (50 iterations in 3.622s)
[t-SNE] Iteration 800: error = 1.7977734, gradient norm = 0.0002598 (50 iterations in 3.589s)
[t-SNE] Iteration 850: error = 1.7480551, gradient norm = 0.0002391 (50 iterations in 3.766s)
```

```
[t-SNE] Iteration 900: error = 1.7100001, gradient norm = 0.0002259 (50 iterations in 3.700s)
[t-SNE] Iteration 950: error = 1.6635556, gradient norm = 0.0002130 (50 iterations in 4.173s)
[t-SNE] Iteration 1000: error = 1.6273054, gradient norm = 0.0002003 (50 iterations in 3.781s)
[t-SNE] KL divergence after 1000 iterations: 1.627305
Done..
Creating plot for this t-sne visualization..
```

C:\Users\krush\Anaconda3\lib\site-packages\seaborn\regression.py:546: UserWarning: The `size` parameter has been renamed to `height`; please update your code.
warnings.warn(msg, UserWarning)

saving this plot as image in present working directory...



Done

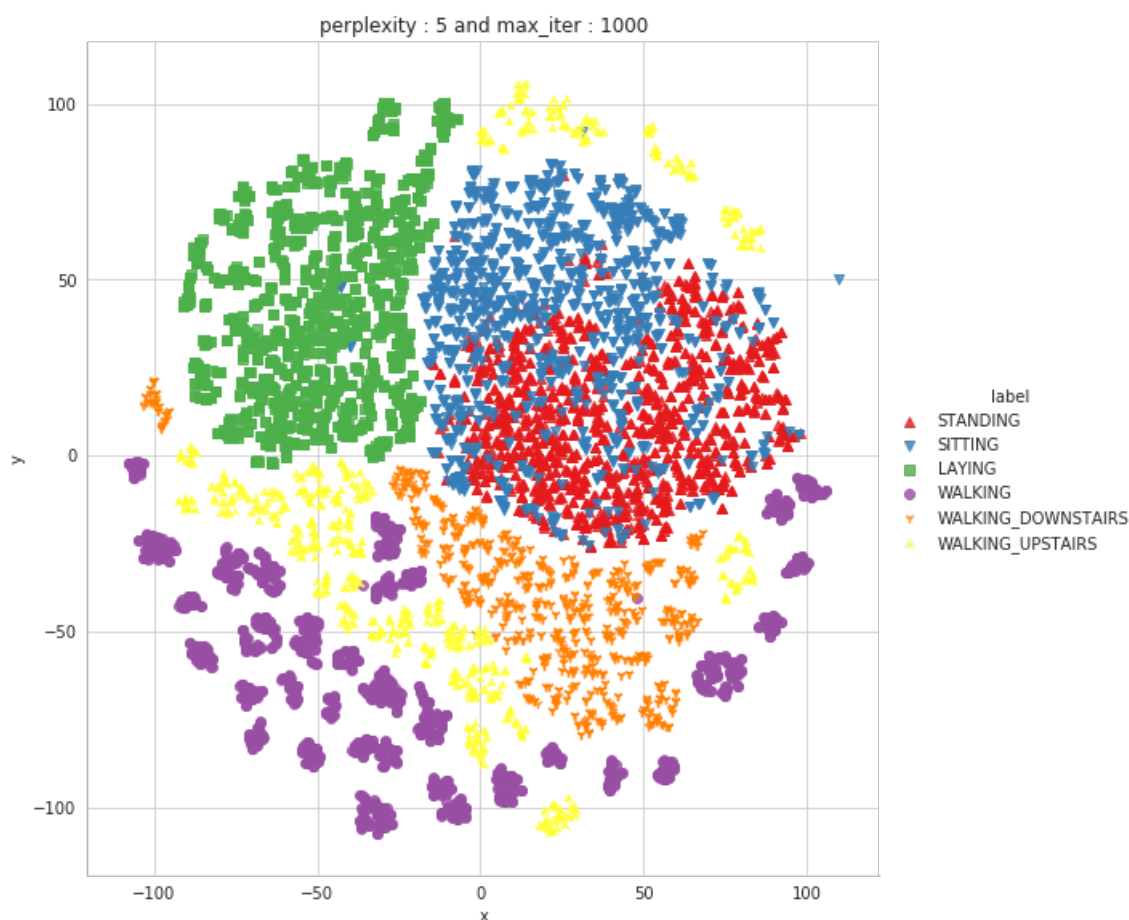
```
performing tsne with perplexity 5 and with 1000 iterations at max
[t-SNE] Computing 16 nearest neighbors...
[t-SNE] Indexed 7352 samples in 0.331s...
[t-SNE] Computed neighbors for 7352 samples in 63.594s...
[t-SNE] Computed conditional probabilities for sample 1000 / 7352
[t-SNE] Computed conditional probabilities for sample 2000 / 7352
[t-SNE] Computed conditional probabilities for sample 3000 / 7352
[t-SNE] Computed conditional probabilities for sample 4000 / 7352
[t-SNE] Computed conditional probabilities for sample 5000 / 7352
[t-SNE] Computed conditional probabilities for sample 6000 / 7352
[t-SNE] Computed conditional probabilities for sample 7000 / 7352
[t-SNE] Computed conditional probabilities for sample 7352 / 7352
[t-SNE] Mean sigma: 0.961265
[t-SNE] Computed conditional probabilities in 0.081s
[t-SNE] Iteration 50: error = 114.1141434, gradient norm = 0.0201779 (50 iterations in 7.367s)
[t-SNE] Iteration 100: error = 97.4559097, gradient norm = 0.0147639 (50 iterations in 5.211s)
[t-SNE] Iteration 150: error = 93.1413498, gradient norm = 0.0095241 (50 iterations in 4.341s)
[t-SNE] Iteration 200: error = 91.1937943, gradient norm = 0.0080519 (50 iterations in 4.635s)
[t-SNE] Iteration 250: error = 90.0317841, gradient norm = 0.0050512 (50 iterations in 4.003s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 90.031784
[t-SNE] Iteration 300: error = 3.5684376, gradient norm = 0.0014608 (50 iterations in 3.750s)
[t-SNE] Iteration 350: error = 2.8115726, gradient norm = 0.0007456 (50 iterations in 3.650s)
[t-SNE] Iteration 400: error = 2.4316602, gradient norm = 0.0005254 (50 iterations in 3.624s)
[t-SNE] Iteration 450: error = 2.2146931, gradient norm = 0.0004030 (50 iterations in 3.591s)
[t-SNE] Iteration 500: error = 2.0666627, gradient norm = 0.0003050 (50 iterations in 3.500s)
[t-SNE] Iteration 550: error = 1.9666667, gradient norm = 0.0002222 (50 iterations in 3.400s)
[t-SNE] Iteration 600: error = 1.8666667, gradient norm = 0.0001444 (50 iterations in 3.300s)
[t-SNE] Iteration 650: error = 1.7666667, gradient norm = 0.0000666 (50 iterations in 3.200s)
[t-SNE] Iteration 700: error = 1.6666667, gradient norm = 0.0000000 (50 iterations in 3.100s)
[t-SNE] Iteration 750: error = 1.5666667, gradient norm = 0.0000000 (50 iterations in 3.000s)
[t-SNE] Iteration 800: error = 1.4666667, gradient norm = 0.0000000 (50 iterations in 2.900s)
[t-SNE] Iteration 850: error = 1.3666667, gradient norm = 0.0000000 (50 iterations in 2.800s)
[t-SNE] Iteration 900: error = 1.2666667, gradient norm = 0.0000000 (50 iterations in 2.700s)
[t-SNE] Iteration 950: error = 1.1666667, gradient norm = 0.0000000 (50 iterations in 2.600s)
[t-SNE] Iteration 1000: error = 1.0666667, gradient norm = 0.0000000 (50 iterations in 2.500s)
[t-SNE] KL divergence after 1000 iterations: 1.0666667
```

```
[t-SNE] Iteration 500: error = 2.0698187, gradient norm = 0.0003310 (50 iterations in 3.528s)
[t-SNE] Iteration 550: error = 1.9649231, gradient norm = 0.0002840 (50 iterations in 3.471s)
[t-SNE] Iteration 600: error = 1.8837612, gradient norm = 0.0002473 (50 iterations in 3.871s)
[t-SNE] Iteration 650: error = 1.8185452, gradient norm = 0.0002170 (50 iterations in 5.314s)
[t-SNE] Iteration 700: error = 1.7646865, gradient norm = 0.0001979 (50 iterations in 3.994s)
[t-SNE] Iteration 750: error = 1.7192082, gradient norm = 0.0001816 (50 iterations in 3.857s)
[t-SNE] Iteration 800: error = 1.6805396, gradient norm = 0.0001645 (50 iterations in 3.801s)
[t-SNE] Iteration 850: error = 1.6468337, gradient norm = 0.0001534 (50 iterations in 4.335s)
[t-SNE] Iteration 900: error = 1.6173676, gradient norm = 0.0001425 (50 iterations in 3.774s)
[t-SNE] Iteration 950: error = 1.5910425, gradient norm = 0.0001345 (50 iterations in 3.749s)
[t-SNE] Iteration 1000: error = 1.5675385, gradient norm = 0.0001257 (50 iterations in 3.951s)
[t-SNE] KL divergence after 1000 iterations: 1.567538
Done..
```

Creating plot for this t-sne visualization..

```
C:\Users\krush\Anaconda3\lib\site-packages\seaborn\regression.py:546: UserWarning: The `size` parameter has been renamed to `height`; please update your code.
  warnings.warn(msg, UserWarning)
```

saving this plot as image in present working directory...



Done

performing tsne with perplexity 10 and with 1000 iterations at max

```
[t-SNE] Computing 31 nearest neighbors...
[t-SNE] Indexed 7352 samples in 0.359s...
[t-SNE] Computed neighbors for 7352 samples in 63.623s...
[t-SNE] Computed conditional probabilities for sample 1000 / 7352
[t-SNE] Computed conditional probabilities for sample 2000 / 7352
[t-SNE] Computed conditional probabilities for sample 3000 / 7352
[t-SNE] Computed conditional probabilities for sample 4000 / 7352
[t-SNE] Computed conditional probabilities for sample 5000 / 7352
[t-SNE] Computed conditional probabilities for sample 6000 / 7352
[t-SNE] Computed conditional probabilities for sample 7000 / 7352
[t-SNE] Computed conditional probabilities for sample 7352 / 7352
[t-SNE] Mean sigma: 1.133828
[t-SNE] Computed conditional probabilities in 0.154s
[t-SNE] Iteration 50: error = 105.8459625, gradient norm = 0.0166295 (50 iterations in 7.573s)
[t-SNE] Iteration 100: error = 90.3692703, gradient norm = 0.0102542 (50 iterations in 4.903s)
[t-SNE] Iteration 150: error = 87.2760925, gradient norm = 0.0054253 (50 iterations in 4.190s)
```

```

[t-SNE] Iteration 200: error = 86.0467682, gradient norm = 0.0035620 (50 iterations in 4.013s)
[t-SNE] Iteration 250: error = 85.3485184, gradient norm = 0.0062483 (50 iterations in 4.041s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 85.348518
[t-SNE] Iteration 300: error = 3.1337800, gradient norm = 0.0013914 (50 iterations in 3.846s)
[t-SNE] Iteration 350: error = 2.4893336, gradient norm = 0.0006505 (50 iterations in 3.776s)
[t-SNE] Iteration 400: error = 2.1694570, gradient norm = 0.0004236 (50 iterations in 3.684s)
[t-SNE] Iteration 450: error = 1.9849304, gradient norm = 0.0003134 (50 iterations in 3.950s)
[t-SNE] Iteration 500: error = 1.8666941, gradient norm = 0.0002546 (50 iterations in 3.756s)
[t-SNE] Iteration 550: error = 1.7835799, gradient norm = 0.0002102 (50 iterations in 4.184s)
[t-SNE] Iteration 600: error = 1.7211800, gradient norm = 0.0001810 (50 iterations in 4.190s)
[t-SNE] Iteration 650: error = 1.6720921, gradient norm = 0.0001592 (50 iterations in 3.721s)
[t-SNE] Iteration 700: error = 1.6323837, gradient norm = 0.0001456 (50 iterations in 3.774s)
[t-SNE] Iteration 750: error = 1.5994935, gradient norm = 0.0001287 (50 iterations in 4.097s)
[t-SNE] Iteration 800: error = 1.5716071, gradient norm = 0.0001186 (50 iterations in 3.882s)
[t-SNE] Iteration 850: error = 1.5481973, gradient norm = 0.0001116 (50 iterations in 3.979s)
[t-SNE] Iteration 900: error = 1.5282161, gradient norm = 0.0001034 (50 iterations in 3.836s)
[t-SNE] Iteration 950: error = 1.5108240, gradient norm = 0.0000982 (50 iterations in 3.859s)
[t-SNE] Iteration 1000: error = 1.4956614, gradient norm = 0.0000902 (50 iterations in 3.970s)
[t-SNE] KL divergence after 1000 iterations: 1.495661

```

Done..

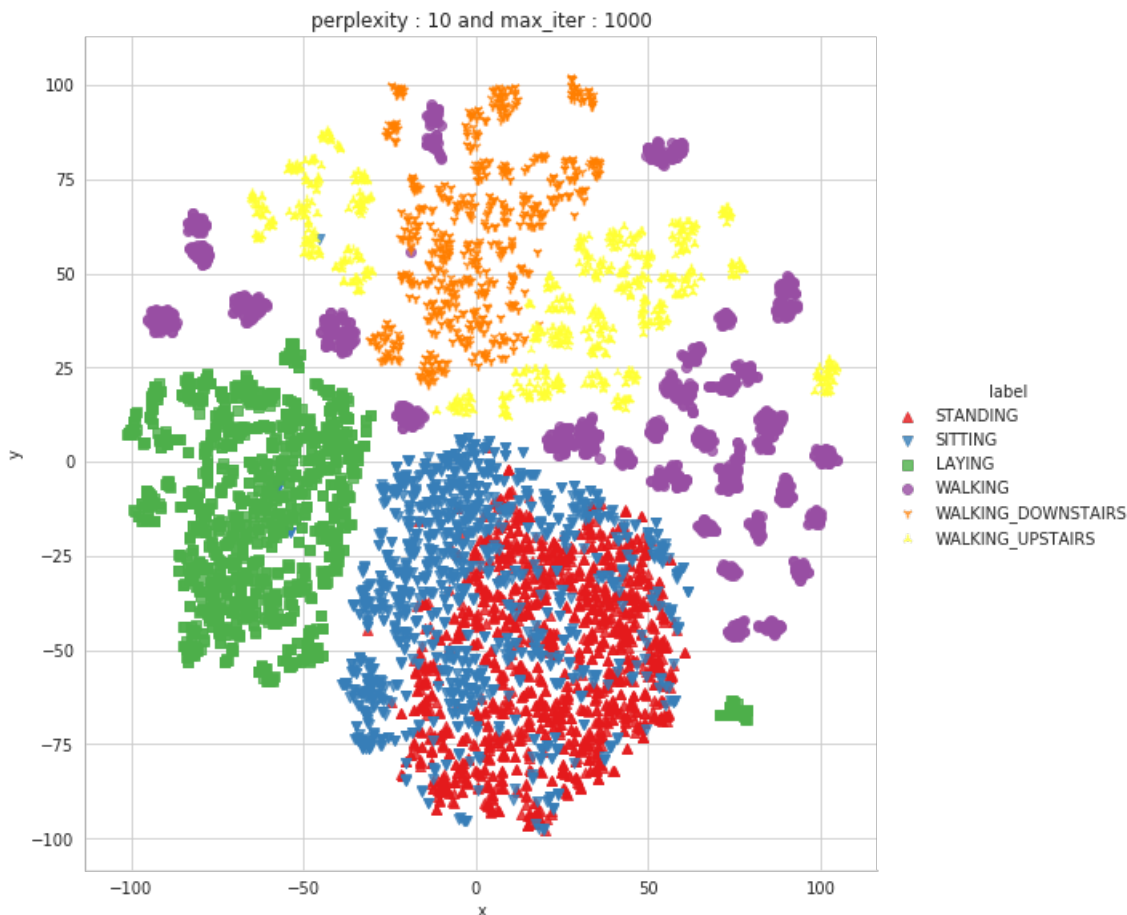
Creating plot for this t-sne visualization..

```

C:\Users\krush\Anaconda3\lib\site-packages\seaborn\regression.py:546: UserWarning: The `size` parameter has been renamed to `height`; please update your code.
  warnings.warn(msg, UserWarning)

```

saving this plot as image in present working directory...



Done

```

performing tsne with perplexity 20 and with 1000 iterations at max
[t-SNE] Computing 61 nearest neighbors...
[t-SNE] Indexed 7352 samples in 0.351s...
[t-SNE] Computed neighbors for 7352 samples in 63.778s...
[t-SNE] Computed conditional probabilities for sample 1000 / 7352
[t-SNE] Computed conditional probabilities for sample 2000 / 7352
[t-SNE] Computed conditional probabilities for sample 3000 / 7352
[t-SNE] Computed conditional probabilities for sample 4000 / 7352
[t-SNE] Computed conditional probabilities for sample 5000 / 7352
[t-SNE] Computed conditional probabilities for sample 6000 / 7352

```

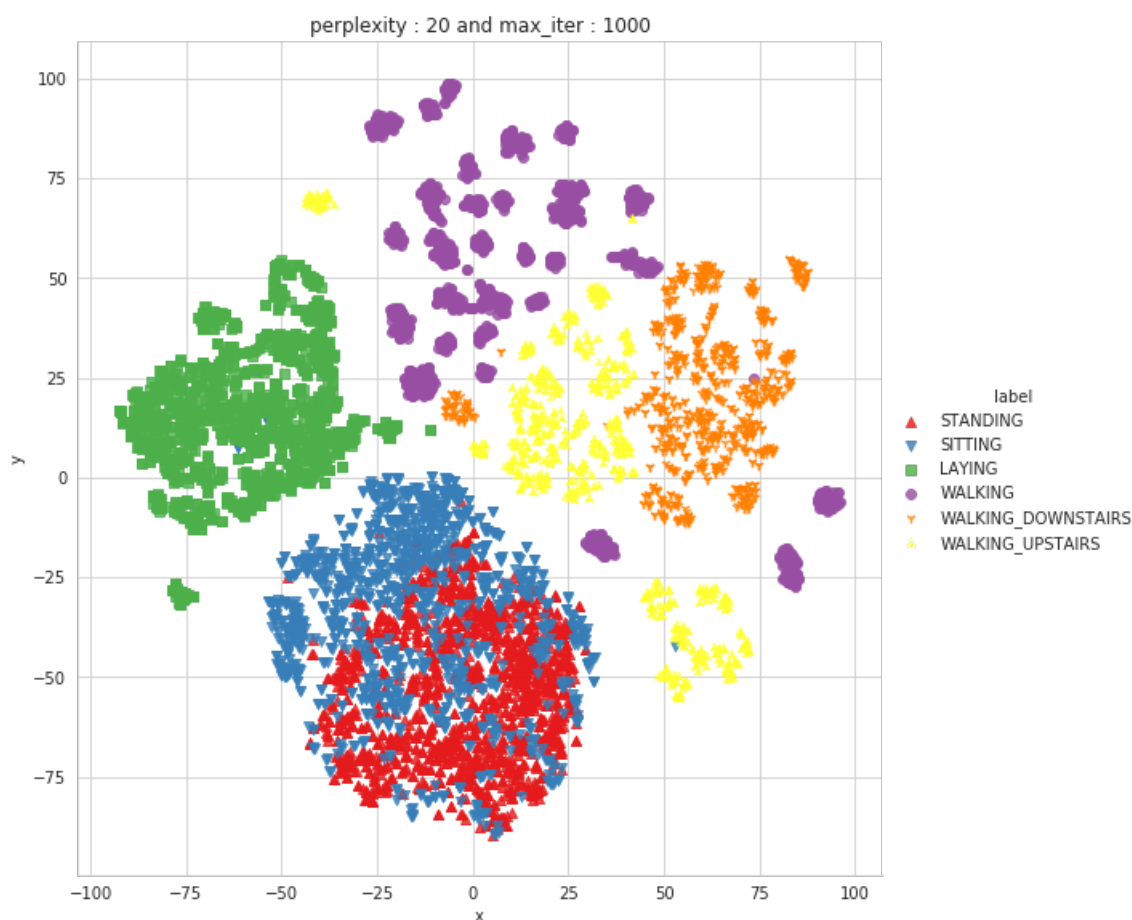
```

[t-SNE] Computed conditional probabilities for sample 6000 / 7352
[t-SNE] Computed conditional probabilities for sample 7000 / 7352
[t-SNE] Computed conditional probabilities for sample 7352 / 7352
[t-SNE] Mean sigma: 1.274335
[t-SNE] Computed conditional probabilities in 0.261s
[t-SNE] Iteration 50: error = 97.5965195, gradient norm = 0.0193359 (50 iterations in 6.702s)
[t-SNE] Iteration 100: error = 84.0039368, gradient norm = 0.0063700 (50 iterations in 4.966s)
[t-SNE] Iteration 150: error = 81.9831009, gradient norm = 0.0037166 (50 iterations in 4.459s)
[t-SNE] Iteration 200: error = 81.2175064, gradient norm = 0.0026628 (50 iterations in 4.396s)
[t-SNE] Iteration 250: error = 80.8371353, gradient norm = 0.0018526 (50 iterations in 4.322s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 80.837135
[t-SNE] Iteration 300: error = 2.7057853, gradient norm = 0.0013167 (50 iterations in 4.140s)
[t-SNE] Iteration 350: error = 2.1680186, gradient norm = 0.0005783 (50 iterations in 4.150s)
[t-SNE] Iteration 400: error = 1.9174963, gradient norm = 0.0003500 (50 iterations in 4.176s)
[t-SNE] Iteration 450: error = 1.7707605, gradient norm = 0.0002497 (50 iterations in 4.251s)
[t-SNE] Iteration 500: error = 1.6766721, gradient norm = 0.0001927 (50 iterations in 4.540s)
[t-SNE] Iteration 550: error = 1.6127187, gradient norm = 0.0001565 (50 iterations in 4.126s)
[t-SNE] Iteration 600: error = 1.5663614, gradient norm = 0.0001345 (50 iterations in 4.480s)
[t-SNE] Iteration 650: error = 1.5310793, gradient norm = 0.0001198 (50 iterations in 4.184s)
[t-SNE] Iteration 700: error = 1.5032765, gradient norm = 0.0001069 (50 iterations in 4.175s)
[t-SNE] Iteration 750: error = 1.4817042, gradient norm = 0.0000987 (50 iterations in 4.197s)
[t-SNE] Iteration 800: error = 1.4646833, gradient norm = 0.0000916 (50 iterations in 4.397s)
[t-SNE] Iteration 850: error = 1.4509033, gradient norm = 0.0000867 (50 iterations in 4.554s)
[t-SNE] Iteration 900: error = 1.4396840, gradient norm = 0.0000820 (50 iterations in 5.323s)
[t-SNE] Iteration 950: error = 1.4304085, gradient norm = 0.0000767 (50 iterations in 4.682s)
[t-SNE] Iteration 1000: error = 1.4219664, gradient norm = 0.0000729 (50 iterations in 4.141s)
[t-SNE] KL divergence after 1000 iterations: 1.421966
Done..
Creating plot for this t-sne visualization..

```

C:\Users\krush\Anaconda3\lib\site-packages\seaborn\regression.py:546: UserWarning: The `size` parameter has been renamed to `height`; please update your code.
warnings.warn(msg, UserWarning)

saving this plot as image in present working directory...



Done

performing tsne with perplexity 50 and with 1000 iterations at max
[t-SNE] Computing 151 nearest neighbors...

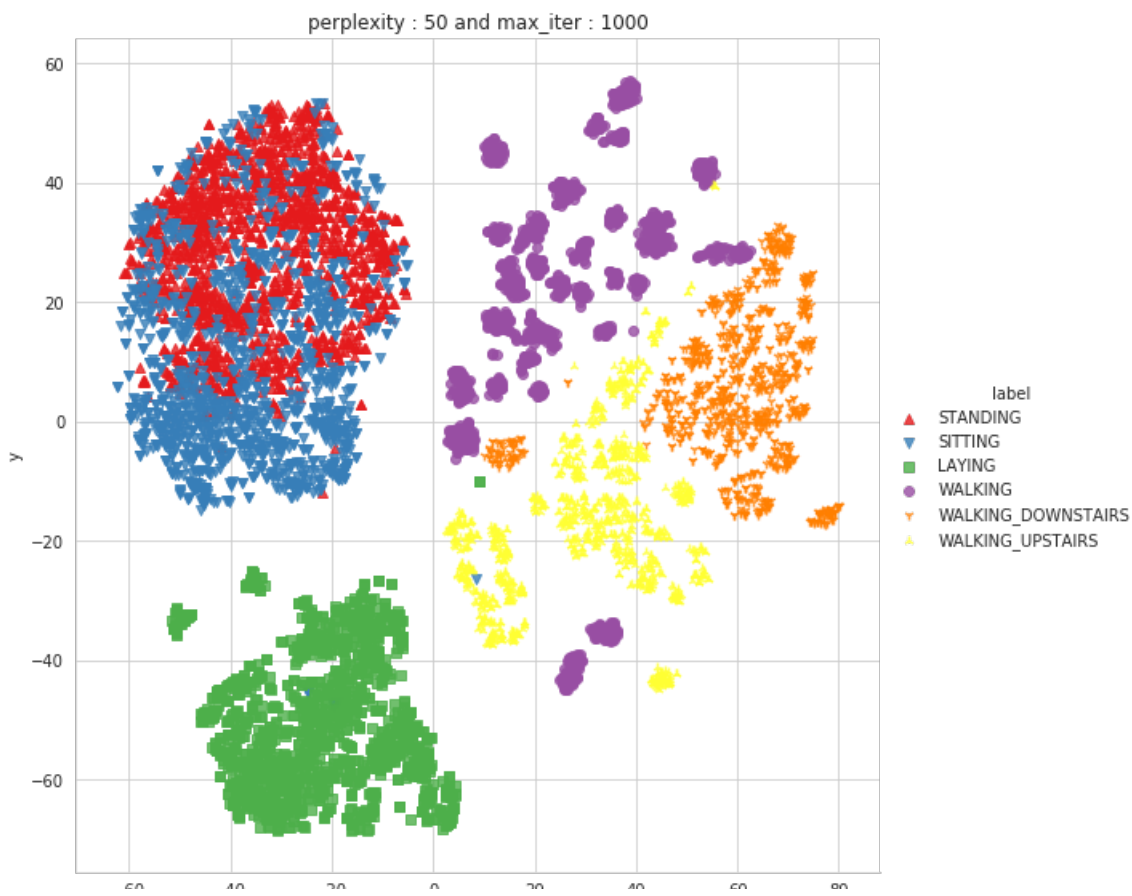

```

[t-SNE] Indexed 7352 samples in 0.396s...
[t-SNE] Computed neighbors for 7352 samples in 66.992s...
[t-SNE] Computed conditional probabilities for sample 1000 / 7352
[t-SNE] Computed conditional probabilities for sample 2000 / 7352
[t-SNE] Computed conditional probabilities for sample 3000 / 7352
[t-SNE] Computed conditional probabilities for sample 4000 / 7352
[t-SNE] Computed conditional probabilities for sample 5000 / 7352
[t-SNE] Computed conditional probabilities for sample 6000 / 7352
[t-SNE] Computed conditional probabilities for sample 7000 / 7352
[t-SNE] Computed conditional probabilities for sample 7352 / 7352
[t-SNE] Mean sigma: 1.437672
[t-SNE] Computed conditional probabilities in 0.720s
[t-SNE] Iteration 50: error = 86.2460175, gradient norm = 0.0219347 (50 iterations in 9.061s)
[t-SNE] Iteration 100: error = 75.6590500, gradient norm = 0.0042137 (50 iterations in 7.687s)
[t-SNE] Iteration 150: error = 74.5999908, gradient norm = 0.0025881 (50 iterations in 6.858s)
[t-SNE] Iteration 200: error = 74.2365189, gradient norm = 0.0014762 (50 iterations in 6.931s)
[t-SNE] Iteration 250: error = 74.0559464, gradient norm = 0.0013726 (50 iterations in 7.091s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 74.055946
[t-SNE] Iteration 300: error = 2.1624236, gradient norm = 0.0012075 (50 iterations in 6.177s)
[t-SNE] Iteration 350: error = 1.7617375, gradient norm = 0.0004874 (50 iterations in 5.295s)
[t-SNE] Iteration 400: error = 1.5920976, gradient norm = 0.0002845 (50 iterations in 5.357s)
[t-SNE] Iteration 450: error = 1.4977177, gradient norm = 0.0001921 (50 iterations in 5.368s)
[t-SNE] Iteration 500: error = 1.4376321, gradient norm = 0.0001428 (50 iterations in 5.418s)
[t-SNE] Iteration 550: error = 1.3960398, gradient norm = 0.0001147 (50 iterations in 5.445s)
[t-SNE] Iteration 600: error = 1.3667738, gradient norm = 0.0000951 (50 iterations in 5.440s)
[t-SNE] Iteration 650: error = 1.3455197, gradient norm = 0.0000832 (50 iterations in 5.192s)
[t-SNE] Iteration 700: error = 1.3300080, gradient norm = 0.0000774 (50 iterations in 5.419s)
[t-SNE] Iteration 750: error = 1.3182805, gradient norm = 0.0000708 (50 iterations in 5.449s)
[t-SNE] Iteration 800: error = 1.3093848, gradient norm = 0.0000645 (50 iterations in 5.518s)
[t-SNE] Iteration 850: error = 1.3020622, gradient norm = 0.0000593 (50 iterations in 5.435s)
[t-SNE] Iteration 900: error = 1.2960280, gradient norm = 0.0000570 (50 iterations in 5.595s)
[t-SNE] Iteration 950: error = 1.2910209, gradient norm = 0.0000556 (50 iterations in 5.371s)
[t-SNE] Iteration 1000: error = 1.2867339, gradient norm = 0.0000528 (50 iterations in 5.361s)
[t-SNE] KL divergence after 1000 iterations: 1.286734
Done..
Creating plot for this t-sne visualization..

```

C:\Users\krush\Anaconda3\lib\site-packages\seaborn\regression.py:546: UserWarning: The `size` parameter has been renamed to `height`; please update your code.
 warnings.warn(msg, UserWarning)

saving this plot as image in present working directory...



-60 -40 -20 0 20 40 60 80

x

Done