In [18]:

```
from pickletools import uint8
from turtle import color
import numpy as np
import matplotlib.pyplot as plt
import os
import sys
from MNISTManager import MNISTManager
from operator import itemgetter
from terminaltables import AsciiTable
```

Import each matrices in a dictionnary like this :

- trainingImages: 60000 training 28*28 images
- trainingLabels: 60000 labels for each training images
- testingImages: 10000 testing 28*28 images
- testingLabels: 10000 labels for each testing images

In [51]:

```
from MNISTManager import MNISTManager

matrices = MNISTManager(
   os.path.abspath('') + "/datasets/emnist-balanced-train-images-idx3-ubyte.gz",
   os.path.abspath('') + "/datasets/emnist-balanced-train-labels-idx1-ubyte.gz",
   os.path.abspath('') + "/datasets/emnist-balanced-test-images-idx3-ubyte.gz",
   os.path.abspath('') + "/datasets/emnist-balanced-test-labels-idx1-ubyte.gz",
).getMatrices()
```

Define utilitaries functions

```
In [52]:
```

```
def getEuclidianDistance(x, y): return np.sqrt(np.sum(x - y)**2) #get the euclidian dist
ance between 2 matrices
def getKNearestNeighbors(k, sortedDistances): #given a value of k and a sorteed distances
list, returns the k neirest neighbors in an new list
    kDistances = sortedDistances[:k]
   return [x[1] for x in kDistances]
def getAccuraryPercentage(predictedLabels, realLabels): #get the acccuracy of the given 1
abels list compared to the training labels
   accuracy = 0
    for i in range(len(predictedLabels)):
        if realLabels[i] == predictedLabels[i]:
            accuracy += 1
    return (accuracy / len(predictedLabels) * 100)
def printProgression(x):
    print(f'' \setminus \{x\} \%'', end='\r')
    sys.stdout.flush()
```

Define the getKNearestCorrespondingLabels function

```
def getKNearestCorrespondingLabels(k, imageDataset, labelDataset, kValuesData=None): #kV
aluesData = tuple(kIndex, kValuesLength)
   correspondingLabels = []
   batchSize = len(imageDataset)
   for i in range(batchSize):
       distances = []
       printProgression(
          f"\t{ round( ((kValuesData[0] / kValuesData[1] * 100) + ( (i / batchSize)
) / kValuesData[1] * 100 ), 5 ) }%"
          if kValuesData != None
          else f"\t{round( (i / batchSize) * 100, 5 )}%",
       for j in range(len(imageDataset)):
          labelDataset[j] ) )
       distances.sort(key = lambda \times x : x[0])
       kNearestNeighbors = getKNearestNeighbors(k, distances)
       correspondingLabels.append( max( kNearestNeighbors, key=kNearestNeighbors.count
) ) #for each image, we add the most propable neighbor for the current value of k
   return correspondingLabels
```

KNN algorithm to get the best value of K

```
In [54]:
```

```
def train(batchSize = 600): #returns K
   kValues = [(3, 0), (5, 0), (7, 0), (9, 0)] #[(kNeigbors, accuracyPercentage)] possib
le k values
   for kIndex in range(len(kValues)):
       correspondingLabels = getKNearestCorrespondingLabels(
          kValues[kIndex][0], matrices['trainingImages'][:batchSize],
          matrices['trainingLabels'][:batchSize],
          (kIndex, len(kValues))
       kValues[kIndex] = (kValues[kIndex][0], getAccuraryPercentage(correspondingLabels
, matrices['trainingLabels'][:batchSize]))
   #print a basic table of k stats
   print(kValues)
   kValuesTable = AsciiTable(
          ["K"] + [x[0]  for x in kValues],
          ["Accuracy"] + [str(round(x[1], 3)) + "%" for x in kValues]
   print(kValuesTable.table)
   return max(kValues, key=itemgetter(1))[0] # we keep only the best accuracy of each k
values
k = train(600)
print(f"K value is : {k}")
[(3, 94.333333333333), (5, 81.666666666666), (7, 62.333333333333), (9, 48.16666666
66667)]
+----+
     | 3 | 5 | 7 | 9
| K
+----+
| Accuracy | 94.333% | 81.667% | 62.333% | 48.167% |
+----+
```

We store the predicted values using the testing set and the k value

```
In [55]:
```

K value is: 3

```
from numpy import append
```

```
testingBatchSize = 500

def getPredictedValues(k, batchSize = 100): #we set a limit to limit time process
    predictedValues = []
    for i in range(batchSize):
        distances = []
        printProgression( round( (i / batchSize) * 100, 5 ))
        for j in range(batchSize):
            distances.append( ( getEuclidianDistance(matrices['testingImages'][i], matrices['testingImages'][j]) , matrices['testingLabels'][j]) )
            distances.sort(key = lambda x : x[0]) #sort by distance
            kNearestNeighbors = getKNearestNeighbors(k, distances)
            predictedValues.append(max( kNearestNeighbors, key=kNearestNeighbors.count ))
    return predictedValues

predictedValues = getPredictedValues(k, testingBatchSize)
```

99.8%

Plot the confusion matrix

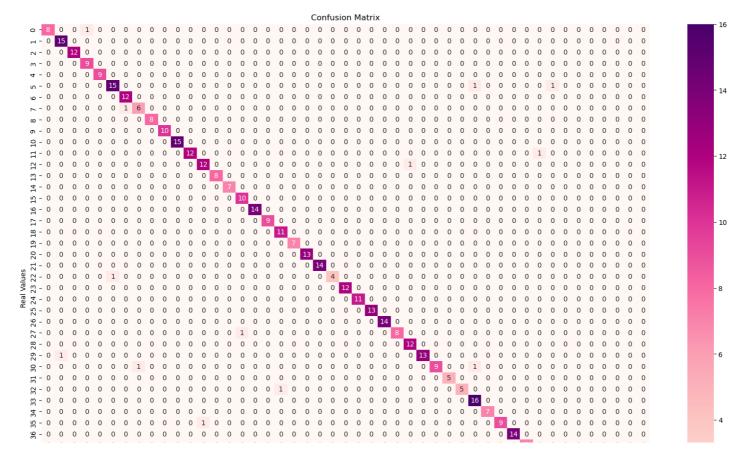
In [56]:

```
from sklearn import metrics
import seaborn as sn

def plotConfusionMatrix(realValues, predictedValues):
    cm = metrics.confusion_matrix(realValues, predictedValues)
    plt.figure(figsize=(20, 14))
    plt.title("Confusion Matrix")
    snHeatmap = sn.heatmap(cm, annot=True, cmap='RdPu')
    snHeatmap.set(xlabel='Predicted Values', ylabel='Real Values')

plotConfusionMatrix(matrices['testingLabels'][:testingBatchSize].tolist(), predictedValue s)
    print(f"Accuracy : {getAccuraryPercentage(predictedValues, matrices['testingLabels'][:testingBatchSize])}%")
```

Accuracy: 96.2%



<u>۾</u> -	0	0	0	0	0	0	0				0																																			
99 -	0	0	0	0	0	0	0																																						0	
우 -	0	0	0	0	0	0	0																																						0	
- 41	0	0	0	0	0	0	0																																				0	0	0	0
- 45	1	0	0	0	0	0	0				0																																6	0	0	0
43	0	0	0	0	0	0	0				0																																0	10	0	0
4 -	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	15	0
- 42	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	6
- 46	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	ó	i	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
																						Pre	edict	ted	Valu	es																				

In []: