

Proyecto Final Inteligencia Artificial

Flores González, Jesús Eduardo.*

* Universidad Nacional Autónoma de México, Facultad de Ingeniería,
e-mail: jef007.flores@gmail.com

Resumen: El presente documento muestra la documentación y explicación del proyecto analizador de datos para la materia de Inteligencia Artificial impartida en la Facultad de Ingeniería de la UNAM, en el se abordan los objetivos del proyecto y se describe sobre el uso de las tecnologías implementadas junto con el desarrollo de esta software para computadora. El como surge el diseño de la interfaz planteada así como una explicación de los aspectos mas importantes del código.

Palabras clave: Inteligencia artificial, algoritmo, python, apriori, clustering, correlacional, métricas de similitud, regresión logística, interfaz, pandas, matplotlib, software, PyQt5.

1. OBJETIVO

Diseñar software que por medio de una interfaz se puedan analizar datos donde el usuario proporcione un archivo con dicha información y estos puedan ser analizados mediante los algoritmos de inteligencia artificial vistos en clase poniendo en práctica los conocimientos adquiridos en la prácticas así como conocimientos que se han adquirido a lo largo de la carrera de Ingeniería en Computación.

2. REQUERIMIENTOS

- Carga de archivos proporcionados por el usuario en formatos .txt, .csv, y .xls.
- Visualización de los datos en la interfaz.
- Poder seleccionar el algoritmo que se desea emplear para poder analizar los datos.
- Mostrar las salidas que puedan generar los diversos algoritmos (texto y gráficos en caso de ser necesario)
- Poder seleccionar parámetros de entrada para los algoritmos con el fin de tener un análisis personalizado al gusto del usuario

3. HERRAMIENTAS DE DESARROLLO

Lenguaje de alto nivel Python en su versión 3.8.7 el cual es un lenguaje interpretado el cual se trata de un lenguaje de programación multiparadigma, ya que soporta orientación a objetos, programación imperativa y, en menor medida, programación funcional.

Otra herramienta utilizada es PyQt5 junto con Qt Designer el cual es un binding de la biblioteca gráfica Qt para el lenguaje de programación Python. La biblioteca está desarrollada por la firma británica Riverbank Computing y está disponible para Windows, GNU/Linux y Mac OS X bajo diferentes licencias.

4. REQUISITOS DE INSTALACIÓN

- Python 3.x
- Modulos de Python
 - Sys
 - Pandas
 - Numpy

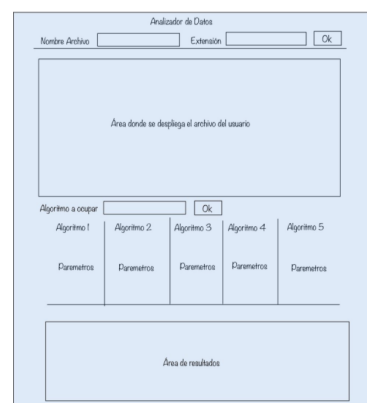
- Math
- Matplotlib
- Seaborn
- Scipy
- Axes3D
- Kmeans
- PyQt5
- Apyori
- Knead
- Sistema Operativo
 - Windows 7, 8, 8.1 y 10
 - Ubuntu 14.04, 16.04 y 16.10 (y derivados)
 - Mac OS X 10.7 (Lion) o superior

Se recomienda de un mínimo de 6Gb RAM y un procesador de 4 núcleos (QuadCore) ya que es el equipo donde se desarrollo y no se presentaban problemas de velocidad para ejecutar el software.

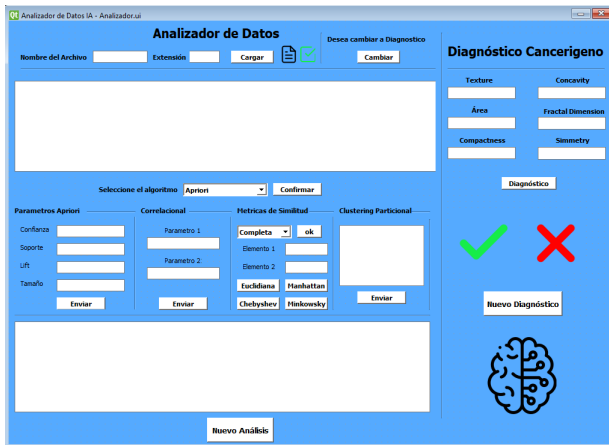
5. SOBRE LA INTERFAZ

Para la aplicación se diseño una ventana en la cual se pudiera cargar una archivo especificando su nombre y extensión, además que en esta ventana se pudiera seleccionar una algoritmo de trabajo con el cual analizar los datos previamente cargados.

A continuación se muestra un borrador de una versión temprana del desarrollo.



Posteriormente conforme el desarrollo fue avanzando y se tenían todos los algoritmos a implementar la interfaz final por la cual se optó fue la siguiente.

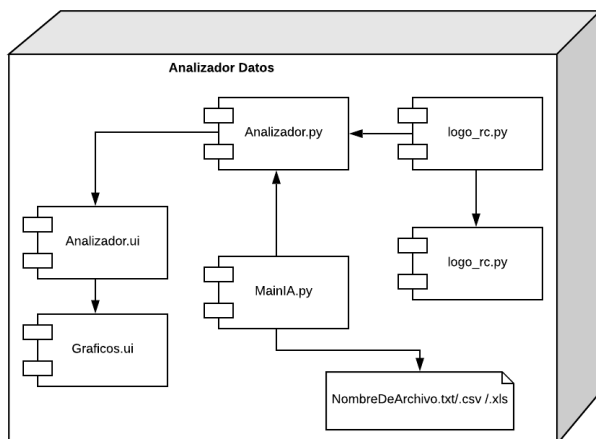


En dicha interfaz podemos observar elementos como:

- Cajas de texto: Las cuales están encargadas de recibir los datos introducidos por el usuario.
- Botones: Los cuales generar eventos al presionarse dichos eventos mandan a llamar a métodos implementados en el código dichos métodos van desde selección de algoritmos, el control de la interfaz y la ejecución de los algoritmos.
- Áreas de texto: Las cuales sirven para mostrar los resultados de los algoritmos como matrices, reglas de asociación, cálculos de distancias entre otras.
- Etiquetas: Se emplean para mostrar texto que no cambia es decir el que muestra la petición de datos, títulos así como imágenes que pueden encontrarse activadas o desactivadas dependiendo de la situación
- Gráficos: Estos elementos permiten que los datos de salida sean mas claros de visualizar aparte de las matrices que se muestran en el áreas de texto.

6. DIAGRAMA COMPONENTES

Para el desarrollo del proyecto se creó el siguiente diagrama de componentes UML el cual muestra la interacción de los archivos de código entre sí junto con la entrada en este caso un archivo el cual como se especificó en los requerimientos puede contar con extensiones .txt, .csv y .xls.



7. FUNCIONAMIENTO DEL CÓDIGO

7.1 Lectura de archivos

Para poder implementar esto un método dentro de la clase `AnalizadorDataApp` se encarga de obtener el nombre del archivo y extensión cada uno con una caja de texto concatenándolos para así saber que tipo de archivo es el que leerá y entra en un serie de condiciones donde se detecta si es .txt, .csv o .xls.. Una vez dentro de esto se realiza la lectura y se almacena y se muestra en el área de texto 1 para visualización del usuario.

```
1 #Obtencion del nombre del archivo
2 NombreArchivo = self.LE_NombreArchivo.text()
3 #Obtencion de la extensión del archivo
4 Extension = self.LE_Extension.text()
5 Extension.lower
6 global Data
7 if Extension == ".txt" :
8     File = NombreArchivo + Extension
9     Data = pd.read_table(File)
10    DataStr = str(Data)
11    self.txt_Texto.setText(DataStr)
```

7.2 Selección de algoritmo

Este método se encarga de detectar el algoritmo con el que se desea trabajar, aparte de detectar esto se encarga de activar y desactivar elementos de la interfaz que no se ocuparan en ese momento así como de empezar a crear algunas paso previos a la ejecución del algoritmo seleccionado. Es importante mencionar que se han omitido las líneas de activación y desactivación sin embargo pueden observar en el repositorio donde se aloja todo el código fuente del proyecto.

```
def Algoritmos(self): #Función que controla la
    QComboBox para la selección de los distintos
    algoritmos implementados
2     AlgoritmoSeleccionado = str(self.
        cb_Algoritmo.currentText())
3     print(AlgoritmoSeleccionado)
4
5     if AlgoritmoSeleccionado == "Apriori":
6         #Se habilitan todas las entradas (LINEAS
            OMITIDAS)
7         global x,y
8         x = len(Data.index)
9         y = len(Data.columns)
10        print(x,y)
11        # Se determinan las transacciones
            realizadas
12        global Transacciones
13        for i in range(0, int(x)):
14            Transacciones.append([str(Data.
                values[i, j]) for j in range(0,
                    int(y))])
15        print(Transacciones) #Para ver las
            transacciones ya que tardan en
            consola
16        if AlgoritmoSeleccionado == "Correlacional":
17            #Se habilitan todas las entradas (LINEAS
                OMITIDAS)
18            Matriz = Data.corr(method='pearson')
19            self.txt_Texto_2.insertPlainText(str(
                Matriz) + "\n")
20            plt.matshow(Matriz)
```

```

21 plt.show()
22 if AlgoritmoSeleccionado == "Metricas de
    Similitud":
23     #Se habilitan todas las entradas (LINEAS
        OMITIDAS)
24     self.pB_ok.clicked.connect(self.
        OptMetricas)
25
26 if AlgoritmoSeleccionado == "Clustering
    Particional":
27     #Se habilitan todas las entradas (
        LINEAS OMITIDAS)
28     #Obtencion del nombre del archivo
29     NombreArchivo = self.LE_NombreArchivo.
        text()
30     #Obtencion de la extensi n del archivo
31     Extension = self.LE_Extension.text()
32     Extension.lower
33     global DataCluster
34     if Extension == ".csv":
35         File = NombreArchivo + Extension
36         DataCluster = pd.read_csv(File)
37         DataStr = str(Data)
38         self.txt_Texto.setText(DataStr)
39         self.label_19.setHidden(False) #Se
            activan iconos de que se ley
            correctamente el archivo
40         #Se habilitan todas las entradas (
            LINEAS OMITIDAS)
41         for column in DataCluster:
42             self.LW_SelectVar.addItem(column)
43         if Extension == ".xls":
44             File = NombreArchivo + Extension
45             DataCluster = pd.read_excel(File)
46             DataStr = str(Data)
47             self.txt_Texto.setText(DataStr)
48             #Se habilitan todas las entradas (
                LINEAS OMITIDAS)

```

7.3 Algoritmo Apriori

Este algoritmo es utilizado en Inteligencia Artificial y Minería de Datos sobre bases de datos transaccionales, que permite encontrar de forma eficiente conjuntos de ítems frecuentes", los cuales sirven de base para generar reglas de asociación. En este caso no se trabajo con una conexión una base de datos, esto para que un usuario pudiera cargar los archivos desde una hoja de calculo o un documento plano el cual puede generar desde una base de datos. Una vez que se carga el archivo el algoritmo todavía no entra en ejecución pues necesita de los siguientes parámetros para su funcionamiento:

- Soporte
- Confianza
- Lift
- Tamaño

Una vez introducidos los valores el algoritmo entra en acción mostrando el numero total de reglas generadas junto con cada una de la reglas juntos con los valores calculados para cada una de estas. (Esto puede visualizarse mejor en el funcionamiento del vídeo).

```

def LeerParametrosApriori(self): #Funcion de lectura
    de los parametros necesarios para poder
    ejecutar el algoritmo apriorio

```

```

2     #Se obtienen los parametros de Apriori
        provenientes de las cajas de texto y sea
        castean a flotantes
3     #En caso de generar un error se muestra un
        error
4     try:
5         global Confianza, Soporte, Lift, Tamano
6         Confianza = float(self.LE_Confianza.text
            ())
7         Soporte = float(self.LE_Soporte.text())
8         Lift = float(self.LE_lift.text())
9         Tamano = float(self.LE_Tamano.text())
10    except:
11        QMessageBox.critical(self, "ERROR", "No
            fue posible realizar ningun calculo,
            verifique sus datos.")
12
13    #Se generan todas la posibles reglar con
        los datos introducidos por el usuario
14    Reglas = apriori(Transacciones, min_support=
        Soporte, min_confidence=Confianza,
        min_lift=Lift, min_length=Tamano)
15    Resultados = list(Reglas)
16
17    #Se muestra el numero de reglas obtenidas en
        el area de texto junto con cada una de
        las reglas obtenidas
18    self.txt_Texto_2.insertPlainText("El total
        de reglas son:" + str(len(Resultados)) +
        "\n")
19    for item in Resultados:
20        pair = item[0]
21        items = [x for x in pair]
22        self.txt_Texto_2.insertPlainText("Regla:
            " + items[0] + " -> " + items[1] +
            "\n")
23        self.txt_Texto_2.insertPlainText("
            Soporte: " + str(item[1]) + "\n")
24        self.txt_Texto_2.insertPlainText("
            Confianza: " + str(item[2][0][2]) +
            "\n")
25        self.txt_Texto_2.insertPlainText("Lift:
            " + str(item[2][0][3]) + "\n")
26        self.txt_Texto_2.insertPlainText("
            =====
            ====="
27        + "\n")

```

7.4 Algoritmo Correccional

Este algoritmo se encarga de generar una matriz de correlaciones la cual aporta información sobre la relación entre pares de variables obteniendo un subconjunto de variables representativas que no tengan dependencia entre si. Una vez esto los datos se pueden graficar en pares obteniendo una gráfica de correlaciones de la cual se debe hacer una inspección visual para interpretar si los datos tiene correlación junto con el coeficiente de correlación calculado para esos datos que se encuentra en la matriz de correlaciones.

La gráfica de correlaciones se genera después de leer el archivo y esta también ser muestra en un gráfico de colores el cual sirve de apoyo para el usuario.

```

1     Matriz = Data.corr(method='pearson')
2     self.txt_Texto_2.insertPlainText(str(Matriz) + "
        \n")
3     plt.matshow(Matriz)

```

```
4 plt.show()
```

Una vez que se tiene esto en la interfaz se habilita el área de parámetros estos son de los que se obtendrá el gráfico de correlación para inspección del usuario. En caso de que se introduzca de forma errónea el nombre del parámetro o que este no exista en los datos se le notifica un error al usuario.

```
def LeerParametrosCorrelacional(self): #Funci n de
    lectura de parmetros para el m todo de
    correlaciones donde de obtiene la gr fica de
    dichos datos
2     Parametro1 = self.LE_Parametro.text()
3     Parametro2 = self.LE_Parametro_2.text()
4     try:
5         plt.close() #Limpiar pantalla grafico
6         plt.plot(Data[Parametro1], Data[
            Parametro2], 'r+')
7         plt.ylabel(Parametro2)
8         plt.xlabel(Parametro1)
9         plt.show()
10    except:
11        QMessageBox.critical(self, "ERROR", "El
            par metro no existe, no es posible
            graficar")
```

7.5 Algoritmo Métricas de Similitud

Este algoritmo se encarga de calcular la distancia que existen entre datos junto con una matriz de similitud es decir el calculo de todas las distancias. Para esto es necesario generar matrices para cálculos y una donde se almacenara la matriz de similitudes esto es un proceso que se repite para los cuatro tipos de métricas disponibles en este software.

```
1DataMetricas = 0
2DataMetricas = Data
3DataMetricas.drop(0)
4x = len(DataMetricas.index)
5y = len(DataMetricas.columns)
6
7DataCalc = {}
8for i in range(0,y):
9    Lista = []
10   for j in range(0,y):
11       Lista.append('--')
12       DataCalc[i] = Lista
13       MatrizSimilitudes = pd.DataFrame(data =
           DataCalc)
14       for i in range(0, y):
15           for j in range(0, i+1):
16               E1 = DataMetricas.iloc[i]
17               E2 = DataMetricas.iloc[j]
18               Temp1 = []
19               Temp2 = []
20               for k in range(0, len(E1)):
21                   try:
22                       Temp1.append(float(E1[k]))
23                   except ValueError:
24                       pass
25                   for k in range(0, len(E2)):
26                       try:
27                           Temp2.append(float(E2[k]
                               ))
28                       except ValueError:
29                           pass
```

```
30         #Linea de calculo segun la metrica
           seleccionada
31         MatrizSimilitudes.iloc[i, j] = dist
```

La linea 350 cambia para cada una de las métricas

■ Euclidiana

```
1     dist = "{0:.05f}".format(math.sqrt(sum((
        Temp1-Temp2)**2 for Temp1, Temp2 in zip
        (Temp1, Temp2))))
```

■ Manhattan

```
1     dist = "{0:.5f}".format(distance.cityblock(
        Temp1, Temp2))
```

■ Minkowsky

```
1     dist = "{0:.5f}".format(distance.minkowski(
        Temp1, Temp2))
```

■ Chebyshev

```
1     dist = "{0:.5f}".format(distance.chebyshev(
        Temp1, Temp2))
```

7.6 Algoritmo Clustering Particional

Para la implementación de dicho algoritmo utilizado en la inteligencia artificial aplicada, consiste en la segmentación de grupos de elementos de manera automática que son unidos por características comunes que estos comparten. El método que realiza el funcionamiento del algoritmo genera el gráfico del método del código así como el de los clusters junto con la selección de variables y muestra en el área de texto resultados como los clusters obtenidos.

```
def ClusteringParticional(self):
2     VarSelec = self.LW_SelectVar.selectedItems()
3     select = []
4     row =[]
5     for x in range(len(VarSelec)):
6         select.append(self.LW_SelectVar.
            selectedItems()[x].text())
7     print(select)
8
9     Matriz = DataCluster.corr(method='pearson')
10    VariablesModelo = np.array(DataCluster[
        select])
11    print(VariablesModelo)
12
13    SSE = []
14    for i in range(2, 16):
15        km = KMeans(n_clusters=i, random_state
            =0)
16        km.fit(VariablesModelo)
17        SSE.append(km.inertia_)
18
19    #Se grafica SSE en funci n de k
20    plt.figure(figsize=(7, 3))
21    plt.plot(range(2, 16), SSE, marker='o')
22    plt.xlabel('Cantidad de clusters *k*')
23    plt.ylabel('SSE')
```

```

24 plt.title('Elbow Method')
25 plt.show()
26
27 #Se crean los clusters
28 #random_state se utiliza para inicializar el
    generador interno de numeros aleatorios
    (mismo resultado)
29 MParticional = KMeans(n_clusters=5,
    random_state=0).fit(VariablesModelo)
30 MParticional.predict(VariablesModelo)
31
32 DataCluster['clusterP'] = MParticional.
    labels_
33 self.txt_Texto_2.insertPlainText("Los
    Clusters son:\n")
34 self.txt_Texto_2.insertPlainText(str(
    DataCluster.groupby(['clusterP'])['
    clusterP'].count()) + "\n")
35 self.txt_Texto_2.insertPlainText("
    =====
36 =====" + "\n")
37
38 CentroidesP = MParticional.cluster_centers_
39
40 self.txt_Texto_2.insertPlainText("La matriz
    de centroides es:\n")
41 self.txt_Texto_2.insertPlainText(str(pd.
    DataFrame(CentroidesP.round(4))) + "\n")
42
43 # Grqfica de los elementos y los centros de
    los clusters
44 plt.rcParams['figure.figsize'] = (7, 3)
45 plt.style.use('ggplot')
46 colores=['red', 'blue', 'cyan', 'green', '
    yellow']
47 asignar=[]
48 for row in MParticional.labels_:
49     asignar.append(colores[row])
50
51 fig = plt.figure()
52 ax = Axes3D(fig)
53 ax.scatter(VariablesModelo[:, 0],
    VariablesModelo[:, 1], VariablesModelo
   [:, 2], marker='o', c=asignar, s=60)
54 ax.scatter(CentroidesP[:, 0], CentroidesP[:,
    1], CentroidesP[:, 2], marker='*', c=
    colores, s=1000)
55 plt.show()

```

7.7 Algoritmo Regresión Logística

El uso de este algoritmo dentro de la interfaz es para la predicción de si un paciente puede tener un tumor maligno o benigno para esto el modelo se entrena a parte con una cantidad de datos muy grande y una vez que se tiene el modelo de regresión logística este se implanta en la interfaz para poder trabajar directamente con el es decir el modelo no se genera en el programa esto por cuestiones de optimizaciones y ya que solamente se requiere como herramienta de diagnostico. El método encargado en el código para realizar esto es:

```

def LeerParametrosDiagnostico(self): #Funcion que
    implementa la lectura y ejecucion del
    diagnostico con los datos propocionados
2     global Texture, Area, Compactness, Concavity
    , FractalDim, Symmetry
3     try:

```

```

4     Texture = float(self.LE_Texture.text())
5     Area = float(self.LE_Area.text())
6     Compactness = float(self.LE_Compactness.
    text())
7     Concavity = float(self.LE_Concavity.text
    ())
8     FractalDim = float(self.LE_FractalDim.
    text())
9     Symmetry = float(self.LE_Symmetry.text()
    )
10    except:
11        QMessageBox.critical(self, "ERROR", "
    Verifique los datos introducidos")
12
13    ModeloDeClasificacion = 11.72-0.19*Texture
    -0.01*Area-2.27*Compactness-3.08*
    Concavity-0.88*Symmetry-0.21*FractalDim
14    Probabilidad = 1/(1+np.exp(
    ModeloDeClasificacion))
15
16    if(Probabilidad < 0.5):
17        #Descativa las cajas y etiquetas donde
    se introducen los datos
18        self.label_13.setHidden(False)
19        self.lb_Diagnostico.setText('Tumor
    benigno')
20    else:
21        #Descativa las cajas y etiquetas donde
    se introducen los datos
22        self.label_12.setHidden(False)
23        self.lb_Diagnostico.setText('Tumor
    Maligno')

```

7.8 Sobre otros apartado dentro del código

El resto de métodos implementado en el código son para el control de la interfaz, aspectos, orientación de elementos y manejo de errores que puedan surgir.

8. CONCLUSIONES DEL PROYECTO

Considero de gran utilidad la herramienta desarrollado durante el curso ya que sirve como un visualizador de datos mediante algoritmos de inteligencia artificial es basten te funcional, sin embargo me gustaría poder optimizar el código haciendo que los algoritmos se ejecuten cuando se tiene la opción en equipo instalado por medio de un GPU a pesar de que en el equipo donde se desarrollo no se presenta ro problemas es cierto que cuando la cantidad de datos es mayor los tiempos de ejecución de los algoritmos penden elevarse. Aparte de esto seria conveniente añadir un barra de progreso cuando se esta ejecutando los algoritmos y estos tarden esto para que el usuario sepa que se están calculando los datos y no crea que la aplicación se ha trabado o genere esa sensación.

9. BIBLIOGRAFÍA

Python Org., (2020) *Python 3.8.7 documentation*. [Online]. Disponible en <https://docs.python.org/3.8/>

Python Org., (2021) *PyQt5 Reference Guide*. [Online]. Disponible en <http://www.latex-project.org/>