





```
In [43]: df2 = df2.reset_index(drop = True)

Out [44]:
```

	diag_cols
0	['DIAG_1', 'DIAG_2', 'DIAG_3']

```
Out [45]:
```

	ord('B') , ord('V')
0	(65, 86)

```
In [46]: for i in diag_cols:
    for j in range(len(df2[i])):
        if str(df2.loc[j, i])[3] == '250':
            df2.loc[j, i] = 'Diabetes'
        elif ord(str(df2.loc[j, i])[0]) in range(69, 87):
            df2.loc[j, i] = 'External causes of injury'
        elif df2.loc[j, i] == '?':
            df2.loc[j, i] = 'Missing'
        else:
            x = float(df2.loc[j, i])
            if x in range(130, 460) or x == 785:
                df2.loc[j, i] = 'Circulatory'
            elif x in range(460, 520) or x == 786:
                df2.loc[j, i] = 'Respiratory'
            elif x in range(520, 580) or x == 787:
                df2.loc[j, i] = 'Digestive'
            elif x in range(800, 1000):
                df2.loc[j, i] = 'Injury and Poisoning'
            elif x in range(710, 740):
                df2.loc[j, i] = 'Musculoskeletal System and Connective Tissue'
            elif x in range(150, 630) or x == 788:
                df2.loc[j, i] = 'Genitourinary'
            elif x in range(140, 230):
                df2.loc[j, i] = 'Neoplasms'
            elif x in [780, 781, 784] or x in range(790, 800):
                df2.loc[j, i] = 'Contagious'
            elif x in range(240, 280):
                df2.loc[j, i] = 'Endocrine, Nutritional, Metabolic, Immunity'
            elif x in range(660, 710) or x == 782:
                df2.loc[j, i] = 'Skin and Subcutaneous Tissue'
            elif x in range(1, 140):
                df2.loc[j, i] = 'Infectious and Parasitic'
            elif x in range(290, 320):
                df2.loc[j, i] = 'Mental Disorders'
            elif x in range(200, 290):
                df2.loc[j, i] = 'Blood and Blood-Forming Organs'
            elif x in range(320, 360):
                df2.loc[j, i] = 'Nervous'
            elif x in range(630, 680):
                df2.loc[j, i] = 'Pregnancy, Childbirth'
            elif x in range(360, 390):
                df2.loc[j, i] = 'Sense Organs'
            else:
                df2.loc[j, i] = 'Congenital Anomalies'
```

```
In [47]: for i in range(len(df2['DIAG_2'])):
    if df2.loc[i, 'DIAG_2'] == 'Missing':
        df2.loc[i, 'DIAG_2'] = 'Missing'
    elif df2.loc[i, 'DIAG_2'] == 'Not Required':
        df2.loc[i, 'DIAG_2'] = 'Not Required'
```

```
In [48]: for i in range(len(df2['DIAG_3'])):
    if df2.loc[i, 'DIAG_3'] == 'Missing':
        df2.loc[i, 'DIAG_3'] = 'Missing'
    elif df2.loc[i, 'DIAG_3'] == 'Not Required':
        df2.loc[i, 'DIAG_3'] = 'Not Required'
```

```
In [49]: df3 = pd.DataFrame(df2['DIAG_1'].value_counts())
df4 = df3.merge(pd.DataFrame(df2['DIAG_2'].value_counts()), how = 'outer', left_index = True, right_index = True)
df4.merge(pd.DataFrame(df2['DIAG_3'].value_counts()), how = 'outer', left_index = True, right_index = True)
```

DIAG_1					DIAG_2	DIAG_3
Blood and Blood-Forming Organs					282	701
Circulatory					7623	8110
Congenital Anomalies					210	146
Diabetes					2183	2918
Digestive					2244	1090
Endocrine, Nutritional, Metabolic, Immunity					694	2120
External causes of injury					243	660
Genitourinary					1191	1804
Infectious and Parasitic					557	472
Injury and Poisoning					1801	610
Mental Disorders					539	662
Musculoskeletal System and Connective Tissue					1613	553
Neoplasms					825	646
Nervous					240	323
Other Symptoms					680	638
Pregnancy, Childbirth					229	149
Respiratory					3441	2634
Sense Organs					60	40
Skin and Subcutaneous Tissue					682	1061
					619	

```
In [50]: len(df2[df2['DIAG_1'] == 'Missing'])
```

```
Out [50]: 0
```

20. **MAX\_GLU\_SERUM**: Indique la plage du résultat ou si le test n'a pas été effectué.

Valeurs: > 200", > 300", "normal" et "None" si non mesuré

```
In [52]: df2.MAX_GLU_SERUM.value_counts()
```

```
Out [52]: None    25355
Norm      2
Name: MAX_GLU_SERUM, dtype: int64
```

21. **A1CRESULT**: Indique la plage du résultat ou si le test n'a pas été effectué.

```
In [53]: df2.A1CRESULT.value_counts()
```

```
Out [53]: None    21217
Norm    1914
Norm    1122
Norm    1084
Name: A1CRESULT, dtype: int64
```

Valeurs: > 8" si le résultat est supérieur à 8%

»> 7" si le résultat est supérieur à 7% mais inférieur à 8%

«Normal» si le résultat est inférieur à 7%

«Aucun» s'il n'est pas mesuré.

22. **caractéristiques pour les médicaments**.

Pour les noms génériques: METFORMIN, REPAGLINIDE, NATEGLINIDE, CHLORPROPAMIDE, GLIMEPRIDE, ACETOHEXAMIDE, GLUPIZIDE, GLYBURIDE, TOLAZAMIDE, PIOGLITAZONE, ROSIGLITAZONE, ACARBOSE, MIGLITOL, TROGLITAZONE, TOLAZAMIDE, EXAMIDE, CITOGLOPTON, INSULIN, GLYBURIDE, METFORMIN, GLUPIZIDE, METFORMIN, GLIMEPRIDE, PIOGLITAZONE, METFORMIN, ROSIGLITAZONE, et METFORMIN, PIOGLITAZONE

La fonction indique si le médicament a été prescrit ou si il y a eu un changement de posologie.

Valeurs: «up» si la posologie a été augmentée pendant la rencontre

«Down» si le dosage a été diminué

«Stable» si la posologie n'a pas changé

«Non» si le médicament n'a pas été prescrit

```
In [54]: meds = ['METFORMIN', 'REPAGLINIDE', 'NATEGLINIDE', 'CHLORPROPAMIDE',
             'GLIMEPRIDE', 'ACETOHEXAMIDE', 'GLIPIZIDE', 'GLYBURIDE', 'TOLBUTAMIDE',
             'PIOGLITAZONE', 'ROSIGLITAZONE', 'ACARBOSE', 'MIGLITOL', 'TROGLITAZONE',
             'TOLAZAMIDE', 'EXAMIDE', 'CITOGLOPTON', 'INSULIN',
             'GLYBURIDE', 'METFORMIN', 'GLIMEPRIDE', 'PIOGLITAZONE',
             'GLIMEPRIDE', 'PIOGLITAZONE', 'METFORMIN', 'ROSIGLITAZONE',
             'METFORMIN', 'PIOGLITAZONE']

for i in meds:
    print(i, df2[i].unique())

METFORMIN ['No' 'Steady' 'Up' 'Down']
REPAGLINIDE ['No' 'Steady' 'Up' 'Down']
NATEGLINIDE ['No' 'Steady' 'Up' 'Down']
CHLORPROPAMIDE ['No' 'Steady' 'Up' 'Down']
GLIMEPRIDE ['No' 'Steady' 'Up' 'Down']
ACETOHEXAMIDE ['No']
GLIPIZIDE ['No' 'Steady' 'Down' 'Up']
GLYBURIDE ['No' 'Steady' 'Up' 'Down']
TOLBUTAMIDE ['No']
PIOGLITAZONE ['No' 'Up' 'Steady' 'Down']
ROSIGLITAZONE ['No' 'Steady' 'Up' 'Down']
ACARBOSE ['No' 'Steady' 'Up' 'Down']
MIGLITOL ['No' 'Steady' 'Up' 'Down']
TROGLITAZONE ['No']
TOLAZAMIDE ['No' 'Steady']
EXAMIDE ['No']
CITOGLOPTON ['No']
INSULIN ['Steady' 'No' 'Up', 'Down']
GLYBURIDE, METFORMIN ['No' 'Steady' 'Up' 'Down']
GLIPIZIDE, METFORMIN ['No' 'Steady']
GLIMEPRIDE, PIOGLITAZONE ['No']
METFORMIN, ROSIGLITAZONE ['No']
METFORMIN, PIOGLITAZONE ['No' 'Steady']
```

Nous pouvons supprimer examide, citoglopton et metformin-rosiglitazone car ils ne sont pas prescrit à aucun patient.

```
In [55]: df2.drop(['EXAMIDE', 'CITOGLOPTON', 'METFORMIN_ROSIGLITAZONE'], axis = 1, inplace = True)
```

```
In [56]: meds = ['METFORMIN', 'REPAGLINIDE', 'NATEGLINIDE', 'CHLORPROPAMIDE',
             'GLIMEPRIDE', 'ACETOHEXAMIDE', 'GLIPIZIDE', 'GLYBURIDE', 'TOLBUTAMIDE',
             'PIOGLITAZONE', 'ROSIGLITAZONE', 'ACARBOSE', 'MIGLITOL', 'TROGLITAZONE',
             'TOLAZAMIDE', 'INSULIN', 'GLYBURIDE', 'METFORMIN', 'GLIMEPRIDE', 'PIOGLITAZONE',
             'GLIMEPRIDE', 'PIOGLITAZONE', 'METFORMIN', 'PIOGLITAZONE']
```

```
In [57]: # Les médicaments peuvent être codés comme ci-dessous
for i in meds:
    df2[i] = df2[i].replace({'No': -2,
                             'Down': -1,
                             'Steady': 0,
                             'Up': 1})
```

```
In [58]: df2[i] = df2[i].astype('int64')
```

23. **CHANGE**: Indique s'il y a eu un changement dans les médicaments diabétiques (soit la posologie, soit le nom générique)

```
In [59]: df2.CHANGE.value_counts()
```

```
Out [59]: Ch    13823
No     11514
Name: CHANGE, dtype: int64
```

```
In [60]: df2['CHANGE'] = df2['CHANGE'].replace({'No': 0, 'Ch': 1})
```

24. **DIABETESMED**: Indique si un médicament diabétique a été prescrit.

```
In [61]: df2.DIABETESMED.value_counts()
```

```
Out [61]: Yes    20793
No     4547
Name: DIABETESMED, dtype: int64
```

```
In [62]: df2['DIABETESMED'] = df2['DIABETESMED'].replace({'Yes': 1, 'No': 0})
```

25. **READMITTED**: Outcome variable

Le résultat que nous examinons est de savoir si le patient est réadmis à l'hôpital dans les 30 jours ou non.

```
In [63]: df2.READMITTED.value_counts()
```

```
Out [63]: No    14044
>30    8539
<30    2754
Name: READMITTED, dtype: int64
```

Le résultat a en fait des catégories <30, >30 et Aucune réadmission, notre problème, nous allons le convertir en 2 catégories, c'est-à-dire pas de réadmission ou de réadmission dans 30 jours

```
In [64]: df2['READMITTED'] = df2['READMITTED'].replace({'>30', 'No'}, {'<30': 1})
df2['READMITTED'] = df2['READMITTED'].replace({'No': 0, '<30': 1}).astype(int)
```

```
In [65]: df2.reset_index(drop = True, inplace = True)
```

(ADDITIONAL) Création de nouvelles fonctionnalités et suppression des fonctionnalités redondantes

1. Visites de l'année précédente:

NUMBER\_OUTPATIENT, NUMBER\_EMERGENCY, NUMBER\_INPATIENT ont été combinés en les additionnant pour obtenir le nombre total de visites du patient au cours de l'année précédente, puis en supprimant les fonctionnalités de combinaison.

```
In [66]: df2['preceding_year_visits']=df2['NUMBER_OUTPATIENT']+df2['NUMBER_EMERGENCY']+df2['NUMBER_INPATIENT']
```

```
In [67]: # Dropping 'number_outpatient', 'number_emergency', 'number_inpatient'
df2 = df2.drop(columns=['NUMBER_OUTPATIENT', 'NUMBER_EMERGENCY', 'NUMBER_INPATIENT'])
```

2. Nombre de changements de médicaments

```
In [68]: df2['number_changes'] = np.nan
for i in range(len(df2)):
    n = 0
    for j in meds:
        if df2.loc[i, j] == -1 or df2.loc[i, j] == 1:
            n += 1
    df2.loc[i, 'number_changes'] = n
```

```
In [69]: df2['number_changes'].value_counts()
```

```
Out [69]: 0.0    16388
1.0     8417
2.0     485
3.0      44
4.0       0
Name: number_changes, dtype: int64
```

3. Nombre de médicaments diabétiques

```
In [70]: df2['number_diab_meds'] = np.nan
for i in range(len(df2)):
    n = 0
    for j in meds:
        if df2.loc[i, j] != -2:
            n += 1
    df2.loc[i, 'number_diab_meds'] = n
df2.number_diab_meds=df2.number_diab_meds.astype('int64')
df2['number_diab_meds'].value_counts()
```

```
Out [70]: 1    11618
2     6017
0     4547
3     2597
4     535
5       22
Name: number_diab_meds, dtype: int64
```

4. Insulin treatment

```
In [71]: df2['insulin_treatment'] = np.nan

In [72]: for i in range(len(df2)):
    if df2.loc[i, 'DIAG_1'] == '2' and df2.loc[i, 'number_diab_meds'] == 1:
        df2.loc[i, 'insulin_treatment'] = 'insulin_only'
    elif df2.loc[i, 'INSULIN'] != -2 and df2.loc[i, 'number_diab_meds'] > 1:
        df2.loc[i, 'insulin_treatment'] = 'insulin_combo'
    elif df2.loc[i, 'INSULIN'] == -2 and df2.loc[i, 'number_diab_meds'] == 0:
        df2.loc[i, 'insulin_treatment'] = 'no_meds'
    else:
        df2.loc[i, 'insulin_treatment'] = 'other_meds'
df2['insulin_treatment'].value_counts()
```

```
Out [72]: insulin_only    8128
insulin_combo    7511
other_meds    4547
no_meds    0
Name: insulin_treatment, dtype: int64
```

number\_diab\_meds a été créé pour faire le traitement de l'insuline de la colonne, par conséquent nous les déposons

```
In [73]: df2=df2.drop('number_diab_meds',1)
df2.reset_index(drop = True, inplace = True)
```

Tests statistiques

```
In [74]: # Caractéristiques numériques
num_cols = ['AGE_INT', 'TIME_IN_HOSPITAL', 'NUM_LAB_PROCEDURES', 'NUM_PROCEDURES', 'NUM_MEDICATIONS',
            'NUMBER_DIAGNOSES', 'preceding_year_visits', 'number_changes']
```

```
In [75]: # Caractéristiques catégoriques
cat_cols = ['DIAG_1', 'DIAG_2', 'DIAG_3', 'MAX_GLU_SERUM', 'DISCHARGE_DISPOSITION', 'ADMISSION_SOURCE',
            'ADMISSION_TYPE', 'DIAG_1', 'DIAG_2', 'DIAG_3', 'MAX_GLU_SERUM', 'A1CRESULT', 'METFORMIN', 'REPAGLINIDE',
            'NATEGLINIDE', 'CHLORPROPAMIDE', 'GLIMEPRIDE', 'ACETOHEXAMIDE', 'GLIPIZIDE', 'GLYBURIDE',
            'TOLBUTAMIDE', 'PIOGLITAZONE', 'ROSIGLITAZONE', 'ACARBOSE', 'MIGLITOL', 'TROGLITAZONE',
            'TOLAZAMIDE', 'INSULIN', 'GLYBURIDE', 'METFORMIN', 'GLIMEPRIDE', 'PIOGLITAZONE',
            'GLIMEPRIDE', 'PIOGLITAZONE', 'METFORMIN', 'ROSIGLITAZONE', 'METFORMIN', 'PIOGLITAZONE',
            'CHANGE', 'DIABETESMED', 'insulin_treatment']
```

```
In [76]: df2.shape
```

```
Out [76]: (25337, 43)
```

```
In [77]: len(num_cols)+len(cat_cols)
```

```
Out [77]: 42
```

Chi-2 pour les catégories vs significatives et one way anova pour les catégories vs numériques

Pour voir si les fonctionnalités sont significatives par rapport à la variable cible (readmitted)

```
In [78]: #Statistical Tests (Chi Square and Anova)
p_val = []
sig = []
for i in df2.columns:
    if i in num_cols:
        stat, p = chi2_contingency(df2[df2['READMITTED'] == 0][i], df2[df2['READMITTED'] == 1][i])
    else:
        ct = pd.crosstab([df2[i], df2['READMITTED']])
        stat, p, dof, exp = stats.chi2_contingency(ct)
        p_val.append(p)
        if p < 0.05:
            sig.append('Significant')
        else:
            sig.append('NotSignificant')
stats_df = pd.DataFrame({'columns': df2.columns, 'p_value': p_val, 'significance': sig})
stats_df.sort_values(by='p_value', ascending = True)
```

```
Out [78]:
```

	columns	p_value	significance
39	READMITTED	0.00000e+00	
40	preceding_year_visits	9.38727e-128	Significant
4	DISCHARGE_DISPOSITION	1.277504e-29	Significant
32	INSULIN	2.38248e-22	Significant
14	NUMBER_DIAGNOSES	4.71365e-22	Significant
42	insulin_treatment	2.86496e-19	Significant
7	MEDICAL_SPECIALTY	3.29907e-17	Significant
41	number_changes	6.35779e-17	Significant
38	DIABETESMED	6.02185e-14	Significant
37	CHANGE	3.972180e-13	Significant
11	DIAG_1	5.83030e-13	Significant
13	DIAG_3	5.56185e-11	Significant
6	TIME_IN_HOSPITAL	3.56017e-08	Significant
10	NUM_MEDICATIONS	3.950002e-06	Significant
12	DIAG_2	3.845994e-05	Significant
5	ADMISSION_SOURCE	1.36742e-03	Significant
17	METFORMIN	1.004162e-02	Significant
29	MIGLITOL	2.402984e-02	Significant
16	A1CRESULT	2.85536e-02	Significant
26	PIOGLITAZONE	2.85536e-02	Significant
8	NUM_LAB_PROCEDURES	5.69281e-02	Significant
2	AGE_INT	6.031671e-02	Insigificant
18	REPAGLINIDE	8.278125e-02	Insigificant
24	GLYBURIDE	1.020866e-01	Insigificant
36	METFORMIN,PIOGLITAZONE	2.086789e-01	Insigificant
34	GLUPIZIDE,METFORMIN	2.086789e-01	Insigificant
23	GLUPIZIDE	5.06974e-01	Insigificant
21	GLIMEPRIDE	5.10105e-01	Insigificant
15	MAX_GLU_SERUM	5.208370e-01	Insigificant
28	ACARBOSE	5.234767e-01	Insigificant
19	NATEGLINIDE	5.251132e-01	Insigificant
9	NUM_PROCEDURES	5.700998e-01	Insigificant
3	ADMISSION_TYPE	6.691503e-01	Insigificant
20	CHLORPROPAMIDE	7.680758e-01	Insigificant
31	TOLAZAMIDE	8.417851e-01	Insigificant
1	GENDER	8.508377e-01	Insigificant
0	RACE	8.786289e-01	Insigificant
27	ROSIGLITAZONE	8.986886e-01	Insigificant
33	GLYBURIDE,METFORMIN	9.191034e-01	Insigificant
25	TOLBUTAMIDE	1.000000e+00	Insigificant
30	TROGLITAZONE	1.000000e+00	Insigificant
22	ACETOHEXAMIDE	1.000000e+00	Insigificant
35	GLIMEPRIDE,PIOGLITAZONE	1.000000e+00	Insigificant

```
In [79]: sig_cols = stats_df[stats_df['significance'] == 'Significant']['columns'].reset_index(drop = True)
insig_cols = stats_df[stats_df['significance'] == 'Insigificant']['columns'].reset_index(drop = True)
anova_cols = stats_df[stats_df['significance'] == 'Insigificant']
print(sig_cols)
```

Statistically significant features are :

```
0    DISCHARGE_DISPOSITION
1    ADMISSION_SOURCE
2    TIME_IN_HOSPITAL
3    MEDICAL_SPECIALTY
4    NUM_MEDICATIONS
5    DIAG_1
6    DIAG_2
7    DIAG_3
8    NUMBER_DIAGNOSES
9    A1CRESULT
10   METFORMIN
11   PIOGLITAZONE
12   MIGLITOL
13   INSULIN
14   CHANGE
15   DIABETESMED
16   READMITTED
17   preceding_year_visits
18   number_changes
19   insulin_treatment
Name: columns, dtype: object
```

Nous pouvons supprimer toutes ces variables

```
In [80]: df2 = df2.drop(['NATEGLINIDE', 'CHLORPROPAMIDE', 'ACETOHEXAMIDE', 'TOLBUTAMIDE',
                       'ACARBOSE', 'MIGLITOL', 'TROGLITAZONE', 'TOLAZAMIDE',
                       'GLYBURIDE,METFORMIN', 'GLIMEPRIDE,METFORMIN',
                       'GLIMEPRIDE,PIOGLITAZONE', 'METFORMIN,PIOGLITAZONE'], axis=1)
```

CHANGE et number\_change décrivent les mêmes informations et les deux sont significatives, nous en abandonnons une basée sur le fait que le nombre de changements est un peu plus significatif que le changement.

```
In [81]: df2 = df2.drop(['CHANGE'], axis=1)
```

```
In [82]: df2['ADMISSION_TYPE'].value_counts()
```

```
Out [82]: Emergency    17387
Elective      7913
Not Available      3
Name: ADMISSION_TYPE, dtype: int64
```

```
In [83]: df2['ADMISSION_SOURCE'].value_counts()
```

```
Out [83]: Emergency    13287
Referral      9637
Transferred from another health care facility    2413
Name: ADMISSION_SOURCE, dtype: int64
```

ADMISSION\_SOURCE et ADMISSION\_TYPE sont similaires car les détails de ADMISSION\_TYPE sont présents dans l'identifiant de la source d'admission. Nous pouvons donc supprimer ADMISSION\_TYPE et simplement conserver ADMISSION\_SOURCE

```
In [84]: df2 = df2.drop(['ADMISSION_TYPE'], axis=1)
```

```
In [85]: df2.shape
```

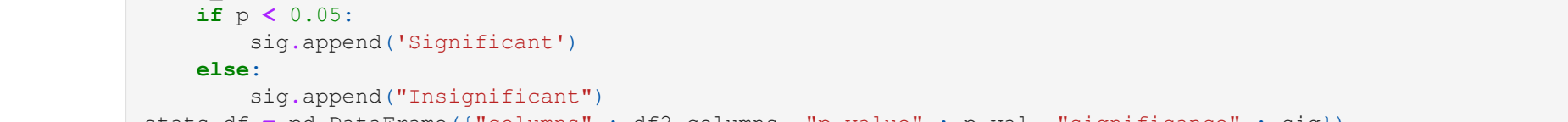
```
Out [85]: (25337, 29)
```

```
In [86]: df2.columns
```

```
Out [86]: Index(['RACE', 'GENDER', 'AGE_INT', 'DISCHARGE_DISPOSITION',
        'ADMISSION_SOURCE', 'TIME_IN_HOSPITAL', 'MEDICAL_SPECIALTY',
        'NUM_LAB_PROCEDURES', 'NUM_PROCEDURES', 'NUM_MEDICATIONS', 'DIAG_1',
        'DIAG_2', 'DIAG_3', 'NUMBER_DIAGNOSES', 'MAX_GLU_SERUM', 'A1CRESULT',
        'METFORMIN', 'REPAGLINIDE', 'GLIMEPRIDE', 'GLIPIZIDE', 'GLYBURIDE',
        'TOLBUTAMIDE', 'PIOGLITAZONE', 'INSULIN', 'DIABETESMED', 'READMITTED',
        'preceding_year_visits', 'number_changes', 'insulin_treatment'],
        dtype='object')
```

Multi-Collinearity

```
In [87]: plt.figure(figsize=(6,5))
sns.pairplot(df2[num_cols])
sns.heatmap(corr=df2[num_cols].corr, cmap='YlGnBu')
b=plt.ylim()
plt.ylim(b*0.5, b*0.5)
plt.show()
```



```
In [88]: from statsmodels.stats.outliers_influence import variance_inflation_factor
vif=pd.DataFrame()
vif['VIF']=variance_inflation_factor(df2[num_cols].values,1)
for i in range(df2[num_cols].shape[1]):
    vif['feature'+str(df2[num_cols].columns[i])] = variance_inflation_factor(df2[num_cols].values,i+1)
vif=vif.sort_values('VIF',ascending=False)
vif
```

```
Out [88]:
```

	VIF	feature
5	14.765046	NUMBER_DIAGNOSES
0	12.129846	AGE_INT
6	6.855368	NUM_MEDICATIONS
2	5.869225	NUM_LAB_PROCEDURES
1	4.273706	TIME_IN_HOSPITAL
7	1.688623	number_changes