# Dimension 'Property'

## Introduction

#### **Hierarchies**

This dimension has one hierarchy: 'property type category' > 'property type'.

```
In [1]: # libraries
        import psycopg2 as pg
        import pandas as pd
        import numpy as np
        import seaborn as sns
        import matplotlib.pyplot as plt
        from scipy import stats
        from scipy.stats import norm
        import os
        import psycopg2.extras
        # set environment variable for psycog2 (for some systems)
        os.environ["PGGSSENCMODE"] = "disable"
        # get the stored passwords
        f=open("credentials.txt", "rt")
        pwd=f.readline().strip() #di password
        hpwd=f.readline().strip() #home password
        f.close()
```

## **Exploratory data analysis**

```
In [2]: def missing data(df, columns):
            df graph = df[columns]
            # missing data
            total = df graph.isnull().sum().sort values(ascending=False)
            percent = (df graph.isnull().sum()/df graph.isnull().count()).sort values(ascending=False)
            missing data = pd.concat([total, percent], axis=1, keys=['Count', 'Percentage (%)'])
            print(missing data)
            plt.figure(figsize=(10,5))
            ax = sns.barplot(x='index',y = 'Percentage (%)', data= missing data.reset index())
            ax.set xticklabels(ax.get xticklabels(), rotation=45)
            return ax
        #IN BAR CHARTS, SET PERCENTAGES ABOVE EACH BAR
        def set bar percentage(ax, df):
            for p in ax.patches:
                ax.annotate(f"{p.get_height() * 100 / df.shape[0]:.2f}%", (p.get_x() + p.get_width() / 2., p.get_height()),
                  ha='center', va='center', rotation=0, xytext=(0, 10),
                  textcoords='offset points')
        #IN BAR CHARTS, SET THE VALUE ABOVE EACH BAR
        def set bar amount(ax):
            for p in ax.patches:
                height = p.get height()
                ax.annotate("{0:.2f}%".format(p.get_height()*100), (p.get_x() + p.get_width() / 2., p.get_height()),
                    ha='center', va='center', rotation=0, xytext=(0, 10),
                    textcoords='offset points')
        # simple bar plot
        def simple plot(df,column):
            bx = sns.catplot(x=column,data=df, kind='count')
            (bx.set_axis_labels(column, "Count")
                .set titles("{col name} {col var}")
                .despine(left=True))
```

#### Missing data

```
In [3]: listings_al_file_path = '../data/listings_al.csv'
df_listings_al = pd.read_csv(listings_al_file_path)
```

C:\Users\joao\_\Anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3063: DtypeWarning: Columns (61,62) have mixed types.
Specify dtype option on import or set low\_memory=False.
interactivity=interactivity, compiler=compiler, result=result)

```
In [4]: columns = ['property type',
                   'room type',
                   'accommodates',
                   'bathrooms',
                   'bedrooms',
                   'beds',
                   'bed_type']
        ax = missing data(df listings al,columns)
        set bar amount(ax)
        ax.set xlabel('Columns')
        ax.set ylabel('Missing data (%)')
        plt.ylim(0,1)
        ax.spines['right'].set visible(False)
        ax.spines['top'].set visible(False)
        ax.set title('Property dimension missing values')
                       Count Percentage (%)
                          58
                                    0.003378
        beds
        bedrooms
                           8
                                    0.000466
```

```
        Count
        Percentage (%)

        beds
        58
        0.003378

        bedrooms
        8
        0.000466

        bathrooms
        5
        0.000291

        bed_type
        0
        0.000000

        accommodates
        0
        0.000000

        room_type
        0
        0.000000

        property_type
        0
        0.000000
```

Out[4]: Text(0.5, 1.0, 'Property dimension missing values')



The detected missing data will be removed, as it constitutes a small percentage of the overall information.

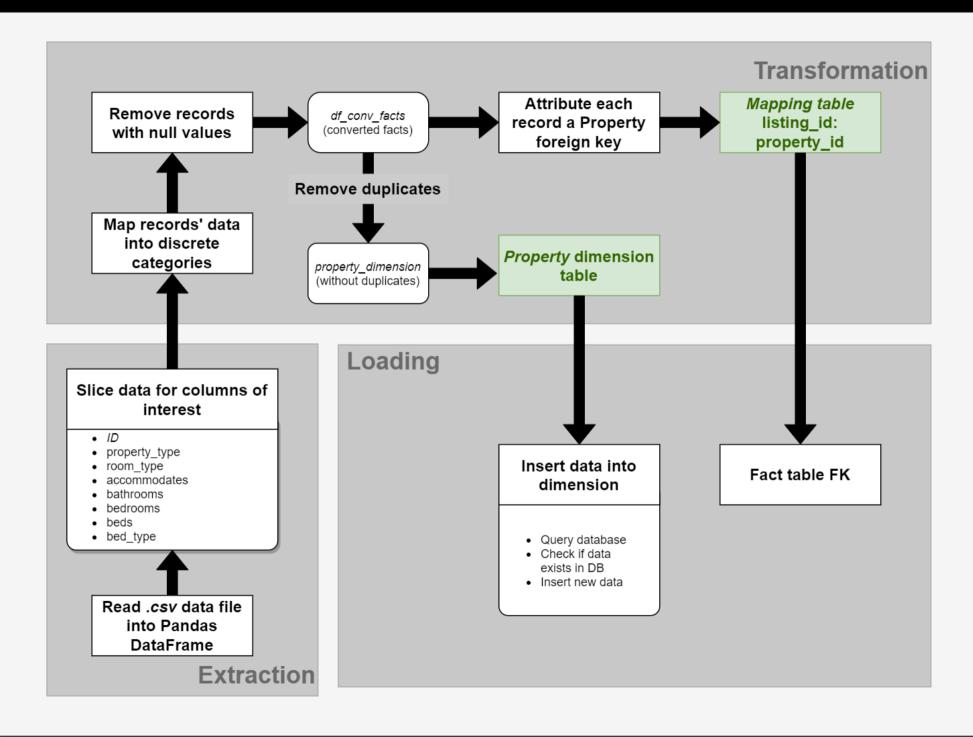
## **Preprocessing**

The defined 'Property' dimension table is as follows. Price per meter squared was previously a column in this table, but data was lacking to determine that attribute at the property level, and therefore was removed.

m property						
驔 property_id	integer					
property_type_category varchar(30)						
property_type	varchar(30)					
III room_type	varchar(30)					
accommodates	varchar(30)					
<b>■</b> bathrooms	varchar(30)					
<b>■</b> bedrooms	varchar(10)					
beds	varchar(10)					
■ bed_type	varchar(30)					

The ETL pipeline is described below:

## **ETL Process**



```
from db connection import dbconnection
In [6]: # test connection
        conn = psycopg2.connect(host = dbconnection.server host,database = dbconnection.dbname, user = dbconnection.dbusername,password =
        conn.close()
In [7]: def delete null rows(df,columns):
            """Removes all records with any empty cells from input DataFrame"""
            processed df = df[columns].copy()
            total rows = processed df.shape[0]
            delete rows = []
            if processed df.isnull().values.any(): # if there are any null values in DataFrame, process DataFrame
                for index, row in enumerate(processed df.itertuples(), start = 0):
                    if (pd.Series(row).isnull().values.any()): # if row has any null value
                        delete rows.append(index) # add row index to delete list
            processed_df.drop(df.index[delete_rows], inplace = True) # delete rows fr
            processed df = processed df.reset index().drop('index', axis = 1)
            print('DataFrame contains {} rows. Deleted {} rows ({}% of total rows)'.format(processed_df.shape[0], len(delete_rows), round(
            return processed df
        property columns = ['id',
In [8]:
                             'property_type',
                            'room type',
                             'accommodates',
                             'bathrooms',
                             'bedrooms',
                            'beds',
                             'bed type']
        df property = delete null rows(df listings al,property columns)
        DataFrame contains 17097 rows. Deleted 71 rows (0.41% of total rows)
```

In [5]: # Load connection settings

From df\_property, we create df\_conv\_fact, which will be the table containing a record for each fact ID converted to the dimension format. property\_dimension will be the *de facto* property dimension from which we can correspond each fact ID record to the appropriate dimension foreign key according to its set of attributes.

## 'property\_type\_category'

We start by creating 'property type category', a set of four categories from 'property type'.

```
In [9]: def get_property_type_category(value):
# change 'property_type' to four categories, exclude exceptions
    if value in ["Apartment", "Serviced apartment", "Aparthotel","Loft"]:
        return 'Apartment'
    if value in ["House","Townhouse","Villa","Dome house","Vacation home","Lighthouse","Casa particular (Cuba)","Tiny house","Farm
        return 'House'
    if value in ["Guesthouse","Guest suite"]:
        return 'Guesthouse'
    if value in ["Hostel","Bed and breakfast","Boutique hotel","Hotel"]:
        return 'Hotel/Hostel'
        return None # for the case of missing or invalid values
```

## 'room\_type'

This attribute is already defined in satisfactory categories.

### 'accommodates'

```
In [10]: def get_accommodates(value):
    # change 'property_type' to four categories, exclude exceptions
    if value in [1,2]: return 'Up to 2 guests'
    if value in [3,4]: return 'Up to 4 guests'
    if value in [5,6]: return 'Up to 6 guests'
    return 'Up to 7 guests or more' # no missing values in this column
```

## 'bathrooms'

```
In [11]: def get_bathrooms(value):
    # change 'bathrooms' to five categories, exclude exceptions
    if value < 0.5: return 'No bathrooms'
    if value < 1.5: return '1 bathroom'
    if value < 2.5: return '2 bathrooms'
    if value < 3.5: return '3 bathrooms'
    if value >= 3.5: return '4+ bathrooms'
    return None # for the case of missing or invalid values
```

#### 'bedrooms'

```
In [12]: def get_bedrooms(value):
    # change 'bedrooms' to five categories, exclude exceptions
    if value == 0: return 'T0'
    if value == 1: return 'T1'
    if value == 2: return 'T2'
    if value == 3: return 'T3'
    if value >= 4: return 'T4+'
    return None # for the case of missing or invalid values
```

#### 'beds'

```
In [13]: def get_beds(value):
    # change 'beds' to five categories, exclude exceptions
    if value == 0: return 'No beds'
    if value == 1: return '1 bed'
    if value == 2: return '2 beds'
    if value == 3: return '3 beds'
    if value >= 4: return '4+ beds'
    return None # for the case of missing or invalid values
```

```
In [14]: def convert facts property(df non null facts):
             """Performs preprocessing in facts to a Property dimension format"""
             dimension = {}
             dimension['ID'] = [value for value in df property['id']]
             dimension['property type category'] = [get property type category(value) for value in df property['property type']]
             dimension['property type'] = [value for value in df property['property type']]
             dimension['room type'] = [value for value in df property['room type']]
             dimension['accommodates'] = [get accommodates(value) for value in df property['accommodates']]
             dimension['bathrooms'] = [get bathrooms(value) for value in df property['bathrooms']]
             dimension['bedrooms'] = [get bedrooms(value) for value in df property['bedrooms']]
             dimension['beds'] = [get beds(value) for value in df property['beds']]
             dimension['bed type'] = [value for value in df property['bed type']]
             df conv facts = pd.DataFrame(dimension)
             df conv facts.index += 1
             df_conv_facts = delete_null_rows(df_conv_facts,df_conv_facts.columns) # removes any rows that during processing gained null va
             return df conv facts
```

DataFrame contains 16626 rows. Deleted 471 rows (2.75% of total rows)

### Out[15]:

	ID	property_type_category	property_type	room_type	accommodates	bathrooms	bedrooms	beds	bed_type
0	25659	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	1 bed	Real Bed
1	29248	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	2 beds	Real Bed
2	29396	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	1 bed	Real Bed
3	29720	Apartment	Apartment	Entire home/apt	Up to 7 guests or more	4+ bathrooms	T4+	4+ beds	Real Bed
4	27899698	Apartment	Apartment	Entire home/apt	Up to 4 guests	2 bathrooms	T2	2 beds	Real Bed
16621	41870065	Apartment	Apartment	Entire home/apt	Up to 2 guests	1 bathroom	T1	1 bed	Real Bed
16622	41879410	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	3 beds	Real Bed
16623	41882911	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	3 beds	Real Bed
16624	41879552	Apartment	Loft	Entire home/apt	Up to 2 guests	1 bathroom	T1	1 bed	Real Bed
16625	41882713	Hotel/Hostel	Boutique hotel	Private room	Up to 2 guests	1 bathroom	T1	2 beds	Real Bed

16626 rows × 9 columns

df\_conv\_facts now contains all facts of interest converted to the appropriate dimension format.

# **Creating the dimension table**

The create table command contains the integrity constraints essential for modelling the dimension.

In [18]: excuteSingleSQLstatement(create\_table, dbconnection.server\_host, dbconnection.dbname, dbconnection.dbusername, dbconnection.dbpass

The dimension table is created empty, to be populated later with incoming data.

## Adding new data to dimension table

```
In [19]: def create_property_dimension(df_conv_facts,columns):
    """Creates a Property dimension table from a converted facts DataFrame"""
    property_dimension = df_conv_facts[columns].drop_duplicates().copy()
    property_dimension = property_dimension.reset_index().drop('index', axis = 1)
    property_dimension.index += 1
    return property_dimension
```

### Out[20]:

	property_type_category	property_type	room_type	accommodates	bathrooms	bedrooms	beds	bed_type
1	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	1 bed	Real Bed
2	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	2 beds	Real Bed
3	Apartment	Apartment	Entire home/apt	Up to 7 guests or more	4+ bathrooms	T4+	4+ beds	Real Bed
4	Apartment	Apartment	Entire home/apt	Up to 4 guests	2 bathrooms	T2	2 beds	Real Bed
5	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	No beds	Real Bed
1020	House	Farm stay	Private room	Up to 6 guests	1 bathroom	T2	4+ beds	Real Bed
1021	Apartment	Apartment	Shared room	Up to 4 guests	1 bathroom	T1	No beds	Real Bed
1022	House	Casa particular (Cuba)	Entire home/apt	Up to 4 guests	1 bathroom	T2	4+ beds	Real Bed
1023	Hotel/Hostel	Bed and breakfast	Private room	Up to 6 guests	2 bathrooms	T1	3 beds	Real Bed
1024	Apartment	Serviced apartment	Entire home/apt	Up to 4 guests	3 bathrooms	T2	2 beds	Real Bed

1024 rows × 8 columns

property\_dimension\_new contains all the potential new records for the dimension. This DataFrame is compared to existing dimension data in the database to assess which new records need to be added, and which are already accounted for in the dimension table.

```
In [21]: # function to query table and convert it to pandas dataframe
         def guery table(conn, table name):
             """Returns DataFrame with queried database table"""
             sql = "select * from {};".format(table name)
             #return dataframe
             return pd.read sql query(sql, conn)
         # for this function to run, the dataframes must have the same columns, in the same order
         def get data to insert(df etl, df sql):
             """Returns data valid for insertion in dimension from a new ETL-processed DataFrame"""
             return df etl[~df etl.isin(df sql)].dropna(how = 'all') # checks which rows are not yet in the dimension
         # function for bulk insert
         def insert data(df, table name, conn):
             """Inserts selected data into dimension table in database"""
             df columns = list(df)
             columns = ",".join(df columns)
             values = "VALUES({})".format(",".join(["%s" for _ in df_columns]))
             insert stmt = "INSERT INTO {} ({}) {}".format(table name,columns,values)
             success = True
             try:
                 cursor = conn.cursor()
                 psycopg2.extras.execute batch(cursor, insert stmt, df.values)
                 conn.commit()
                 success = True
             except pg.DatabaseError as error:
                 success = False
                 print(error)
             finally:
                 if conn is not None:
                     conn.close()
             return success
```

## Out[22]:

property\_type\_category property\_type room\_type accommodates bathrooms bedrooms beds bed\_type

property\_id

With property\_dimension\_old , we can now confirm which new records to add to the 'Property' dimension. When first populating the database, property dimension old yields an empty DataFrame.

```
In [23]: # checks which rows from new data will be inserted into database dimension table
def get_data_to_insert(df1, df2, columns):
    """Checks if rows in df1 are already present in df2"""
    return df1[~df1[columns].apply(tuple,1).isin(df2[columns].apply(tuple,1))]
```

In [24]: # compares rows (excluding unique IDs) of new data and database dimension
dimension\_insert = get\_data\_to\_insert(property\_dimension\_new,property\_dimension\_old,property\_dimension\_new.columns[1:])
dimension\_insert.index += property\_dimension\_old.shape[0] # changes the ID (index) to follow the last row in the database table

dimension insert contains all records that will be added to the database dimension table.

```
In [25]: def insert data(df, table name, conn):
             # Method for bulk insert
             df columns = list(df)
             columns = ",".join(df columns)
             values = "VALUES({})".format(",".join(["%s" for _ in df_columns]))
             insert stmt = "INSERT INTO {} ({}) {}".format(table name,columns,values)
             success = True
             try:
                  cursor = conn.cursor()
                 pg.extras.execute batch(cursor, insert stmt, df.values)
                 conn.commit()
                 success = True
             except pg.DatabaseError as error:
                  success = False
                 print(error)
             finally:
                 if conn is not None:
                     conn.close()
             return success
```

Data inserted succefully

The valid records are inserted into the dimension.

#### Critical assessment

a)

The ETL procedure above does *not* account for this case, as it would require an added level of complexity to the pipeline, but it was also not seen as having high priority, since, for the most part, the characteristics that describe a property in the dimension table are not frequently changed.

If missing data is very prevalent, the ETL pipeline will eliminate many records.

c)

This pipeline achieves good **reusability**. It was designed in a way that ensures that multiple runs of this code are not redundant, that is, only new data that truly needs to be added is added to the warehouse.

The dimension grows as new combinations of property attributes are found in the incoming facts. It could be considered a Type I or Type II slow-changing dimension. Maintaining history by implementing a Type II methodology could be benefitial for queries such as comparing bookings before and after changes in a property (adding beds, building bedrooms/bathrooms).

The determination of which records to use to expand the dimension is based on comparing each processed fact row with the existing dimension rowsm excluding IDs. In this way, we can ensure new types of property are always accounted for in the dimension. However, IDs should be verified, since duplicate listing IDs in the facts table can be indicative of an update to the listing's information.

## **Attributing dimension keys to facts**

Having our dimension DataFrame property\_dimension ready, we can now use our converted listings table, df\_conv\_facts (which contains 'Property' dimension attributes for all valid facts), to assign each fact its corresponding 'Property' dimension foreign key.

## Out[27]:

	property_id	property_type_category	property_type	room_type	accommodates	bathrooms	bedrooms	beds	bed_type
0	1	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	1 bed	Real Bed
1	2	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	2 beds	Real Bed
2	3	Apartment	Apartment	Entire home/apt	Up to 7 guests or more	4+ bathrooms	T4+	4+ beds	Real Bed
3	4	Apartment	Apartment	Entire home/apt	Up to 4 guests	2 bathrooms	T2	2 beds	Real Bed
4	5	Apartment	Apartment	Entire home/apt	Up to 4 guests	1 bathroom	T1	No beds	Real Bed
1019	1020	House	Farm stay	Private room	Up to 6 guests	1 bathroom	T2	4+ beds	Real Bed
1020	1021	Apartment	Apartment	Shared room	Up to 4 guests	1 bathroom	T1	No beds	Real Bed
1021	1022	House	Casa particular (Cuba)	Entire home/apt	Up to 4 guests	1 bathroom	T2	4+ beds	Real Bed
1022	1023	Hotel/Hostel	Bed and breakfast	Private room	Up to 6 guests	2 bathrooms	T1	3 beds	Real Bed
1023	1024	Apartment	Serviced apartment	Entire home/apt	Up to 4 guests	3 bathrooms	T2	2 beds	Real Bed

1024 rows × 9 columns

```
In [28]: def key_mapping(df1, df2, pk1, pk2):
    # This function merges two dataframes, and creates a map linking their keys
    df_merged = df1.merge(df2, how='outer')
    df_map = pd.DataFrame()
    df_map[pk1]= df_merged[pk1]
    df_map[pk2]= df_merged[pk2]

df_map = delete_null_rows(df_map,df_map.columns)
    return df_map
```

```
In [29]: # match fact IDs with FKs in dimension
df_mapping = key_mapping(df_conv_facts, property_dimension_updated, 'ID', 'property_id')
df_mapping.to_csv('../processed_dt/df_listings_property.csv')
```

DataFrame contains 16626 rows. Deleted 0 rows (0.0% of total rows)

This table can then be merged with all other corresponding tables for the remaining dimensions to produce each fact record in the facts table. It will be used both in the *Listings* and *Availability* facts tables.