



CIENCIA DE DATOS EN PYTHON

PROYECTO FINAL

BALMORE hERNÁNDEZ BERNAL
ANDREA hERNÁNDEZ MARROQUÍN

Estructura de Nuestro Proyecto

Consta de 4 Archivos

1. Diseño Transaccional y carga datos.ipynb
2. Diseño Dimensional.ipynb
3. ETL dimensional con inserts a el DW.ipynb
4. Respuestas Preguntas del Negocio.ipynb

1. Diseño Transaccional y carga datos.ipynb

```
1 aws_rds_conn = boto3.client('rds', aws_access_key_id = config.get('IAM','ACCESS_KEY'),
2                               aws_secret_access_key = config.get('IAM','SECRET_KEY'),
3                               region_name = 'us-east-2')
```

```
1 rds_instances_ids = []
2
3 aws_response = aws_rds_conn.describe_db_instances() # retorna un objeto [diccionario] iterable de las instancias
4
5 for response in aws_response['DBInstances']:
6     rds_instances_ids.append(response['DBInstanceIdentifier'])
7
8 print(f'Instancias disponibles: {rds_instances_ids}')
```

Instancias disponibles: ['dbdim', 'dbtienda', 'sakila-db-pg-v']

```
1 try:
2     response = aws_rds_conn.create_db_instance(
3         DBInstanceIdentifier=config.get('DBTIENDA','DB_INSTANCE_ID'),
4         DBName=config.get('DBTIENDA','DB_NAME'),
5         MasterUsername=config.get('DBTIENDA','DB_USERNAME'),
6         MasterUserPassword=config.get('DBTIENDA','DB_PASSWORD'),
7         Port=int(config.get('DBTIENDA','DB_PORT')),
8         DBInstanceClass='db.t3.micro',
9         Engine=config.get('DBTIENDA','DB_ENGINE'),
10        PubliclyAccessible=True,
11        AllocatedStorage=20,
12        VpcSecurityGroupIds=[config.get('VPC','SECURITY_GROUP')],
13    )
```

Access point

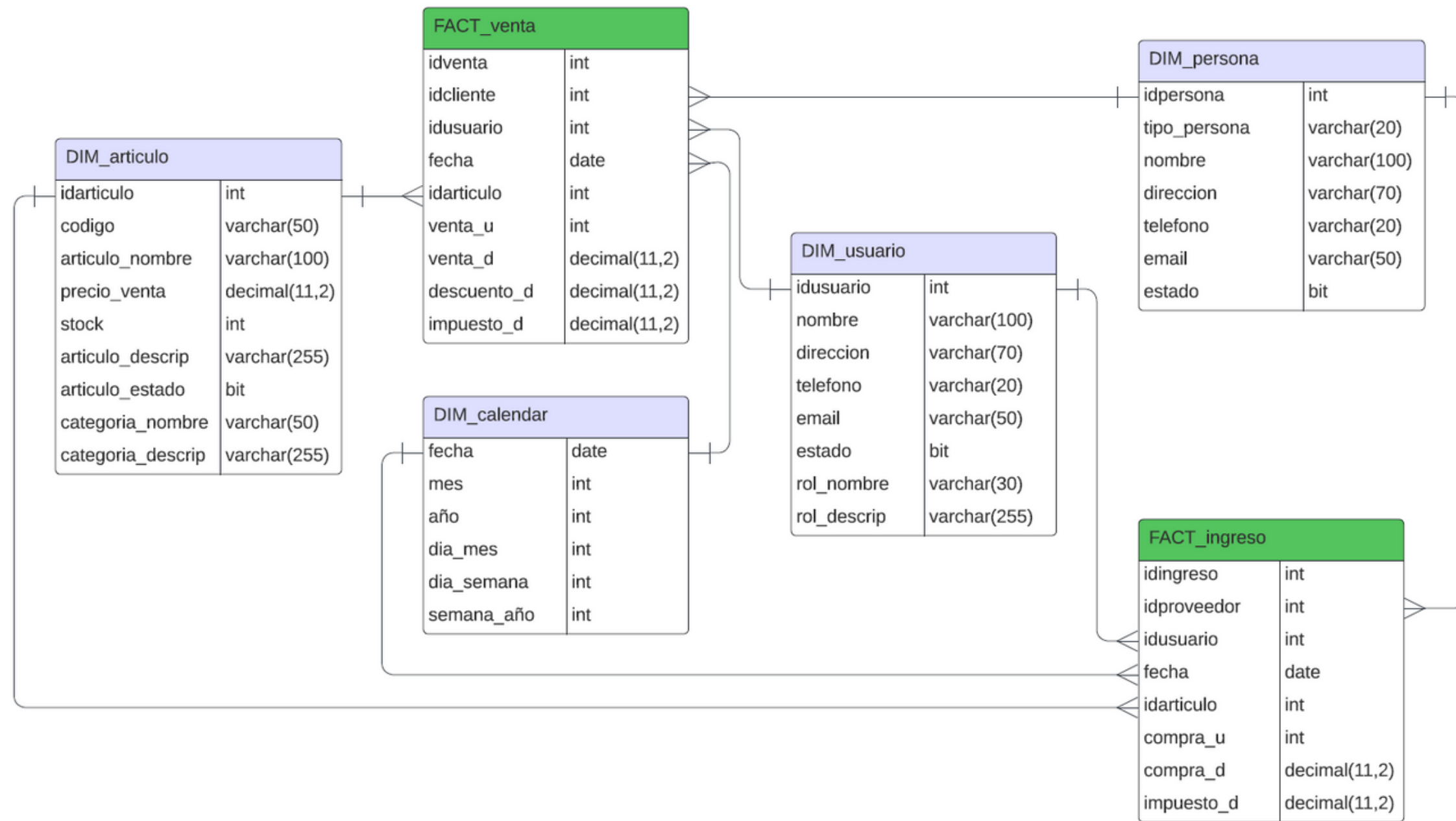
```
1 try:
2     instance = aws_rds_conn.describe_db_instances(DBInstanceIdentifier=config.get('DBTIENDA','DB_INSTANCE_ID'))
3     RDS_HOSTNAME = instance.get('DBInstances')[0].get('Endpoint').get('Address')
4     print(RDS_HOSTNAME)
5 except Exception as ex:
6     print('Error!!!!',ex)
```

dbtienda.cp6geq8ycm59.us-east-2.rds.amazonaws.com

Conexion a BD transaccional en mysql y creacion de tablas

```
1 import ddl_transacdb # py donde se encuentra el ddl
2 #ddl_transdb.ddl
```

```
1 # TABLA ARTICULO
2 for _ in range(100): #Insertamos 100 datos
3     idcategoria = random.randint(1, 6) # Le colocamos hasta el 5 porque solo colocamos 5 categorias
4     codigo = faker.uuid4()[:6] # Generamos un codigo uuid
5     nombre_articulo = faker.word() # Le colocamos nombre
6     precio_venta = round(random.uniform(100, 1000), 2) # Le colocamos un precio entre 100 y mil que tenga 2 decimales como maximo
7     stock = random.randint(1, 100)
8     descripcion_articulo = faker.text() # Le colocamos la descripcion del articulo
9     imagen = faker.file_name(category='image', extension='jpg') # Agregamos una imagen con extension jpg
10    estado_articulo = random.choice([0, 1]) # Estado aleatorio (0 o 1)
11
12    # Insertamos los datos a la base de datos
13    sql_articulo = "INSERT INTO articulo (idcategoria, codigo, nombre, precio_venta, stock, descripcion, imagen, estado) VALUES (%s, %s, %s, %s, %s, %s, %s"
14    val_articulo = (idcategoria, codigo, nombre_articulo, precio_venta, stock, descripcion_articulo, imagen, estado_articulo)
15
16    # Ejecutar la consulta para la tabla "articulo"
17    cursor.execute(sql_articulo, val_articulo)
```



2. Diseño Dimensional.ipynb

Access point dimensional

```
1 try:
2     instance = aws_rds_conn.describe_db_instances(DBInstanceIdentifier=config.get('DIM','DB_INSTANCE_ID'))
3     RDS_HOST_DBDIM = instance.get('DBInstances')[0].get('Endpoint').get('Address')
4     print(RDS_HOST_DBDIM)
5 except Exception as ex:
6     print('Error!!!!',ex)
```

dbdim.cp6geq8ycm59.us-east-2.rds.amazonaws.com

Conexion a DB y creacion de tablas

```
1 import ddl_dimdb # py donde se encuentra el ddl
```

```
1 try:
2     db_pg_conn = psycopg2.connect(
3         database=config.get('DIM','DB_NAME'),
4         user=config.get('DIM','DB_USERNAME'),
5         password=config.get('DIM','DB_PASSWORD'),
6         port=config.get('DIM','DB_PORT'),
7         host=RDS_HOSTNAME
8     )
9     cursor = db_pg_conn.cursor()
10    cursor.execute(ddl_dimdb.ddl)
11    db_pg_conn.commit()
12 except Exception as ex:
```


3. ETL dimensional con inserts a el DW.ipynb

Transformar tablas dimensionales e insercion de datos

Dimension Persona

```
1 #Seleccionamos los campos que necesitamos para la dimension
2 dim_persona = df_persona.loc[:, ["idpersona","tipo_persona","nombre","direccion","telefono","email"]]
3 #Le cambiamos los datos como aparecen en la tabla dimensional
4 dim_persona = dim_persona.rename(columns={'nombre': 'nombre_p', 'direccion': 'direccion_p', 'telefono': 'telefono_p', 'email': 'email_p'})
5 #insertamos datos a dim_persona.
6 dim_persona
7 #dim_persona = dim_persona.to_sql('dim_persona', postgres_driver, index=False, if_exists='append')
```

	idpersona	tipo_persona	nombre_p	direccion_p	telefono_p	email_p
0	1	PERSONA	Ashley Long	731 Hall Rest Apt. 345\nNew Stephanie, TX 45230	+1-612-213-0569x885	heatherbass@example.org
1	2	PERSONA	Joseph Floyd	827 Harris Squares\nSouth Michael, OH 15445	711.625.8479	ydavis@example.net
2	3	EMPRESA	Matthew Brown	Unit 5986 Box 8060\nDPO AE 31584	293.723.3014	taraholmes@example.org
3	4	EMPRESA	Scott Mayer	824 Medina Avenue Suite 336\nPort Terri, NH 27174	001-763-708-7537	jimenezalicia@example.net
4	5	EMPRESA	Cassandra Torres	32549 Mendoza Extension Apt. 991\nEast Timothy...	5579496223	melissapeters@example.com
5	6	EMPRESA	Michael Cox	01919 Richard Common Suite 092\nSouth Anna, OR...	412.396.3173x3157	donnakelley@example.org
6	7	EMPRESA	Kathy Perez	598 Michael Forges Apt. 030\nWest Christine, N...	712.735.8251	wbrown@example.org
7	8	PERSONA	Margaret Goodwin	PSC 8756, Box 2910\nAPO AE 34076	001-856-756-5240x009	wilsonjesse@example.net
8	9	EMPRESA	Jennifer Daugherty	9702 Ingram Curve\nLake Jeffery, CT 95522	(647)397-5018x3727	johnbrown@example.net
9	10	PERSONA	Samantha Ramirez	823 Brown Fords\nCollinsberg, MH 62101	+1-363-212-6561x8889	mark54@example.com

Dimension Usuario

```
1 dim_usuario['estado_u'].dtypes
```

Python

```
dtype('int64')
```

```
1 #Seleccionamos los campos que necesitamos para la dimension
2 dim_usuario = df_usuario.loc[:, ["idusuario","nombre","direccion","telefono","email","idrol","estado"]]
3
4 #Le cambiamos el nombre a los campos que tienen el nombre diferente
5 dim_usuario = dim_usuario.rename(columns={'nombre': 'nombre_u', 'direccion' : 'direccion_u', 'telefono': 'telefono_u', 'email': 'email_u', 'estado': 'estado_u'})
6
7 #Seleccionamos los campos que necesitamos para la dimension
8 dim_rol = df_rol.loc[:, ["idrol","nombre","descripcion"]]
9 #Le cambiamos el nombre a los campos que tienen el nombre diferente
10 dim_rol = dim_rol.rename(columns={'nombre': 'rol_nombre', 'descripcion' : 'rol_descrip'})
11
12 join_usuario_rol = pd.merge(dim_usuario, dim_rol, left_on='idrol', right_on='idrol', how='inner').drop_duplicates() # Eliminamos Duplicados
13 #Ordenamos los campos
14 dim_usuario = join_usuario_rol.loc[:, ["idusuario","nombre_u","direccion_u","telefono_u","email_u","rol_nombre","rol_descrip"]]
15 #dim_usuario['estado_u'] = dim_usuario['estado_u'].astype('bool')
16 dim_usuario
```

4. Respuestas Preguntas del Negocio.ipynb

¿Cuál es la persona que más ha comprado?

Samantha Ramirez

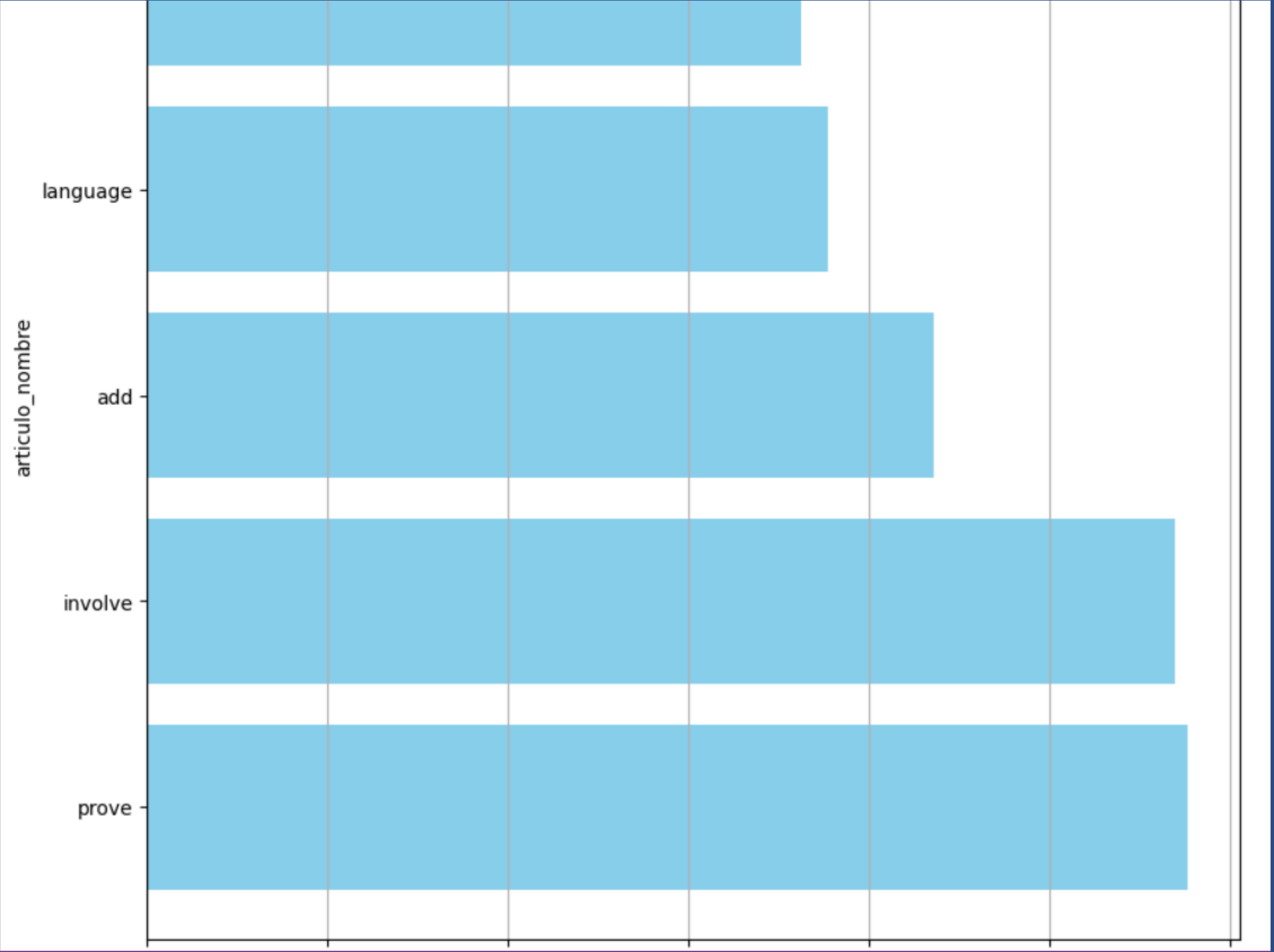
```
1 join_FACT_venta_dim_persona = pd.merge(df_FACT_venta, df_dim_persona, left_on='idcliente', right_on='idpersona', how='inner').drop_duplicates() # Eliminamos duplicados
2 pregunta_2= join_FACT_venta_dim_persona.groupby('nombre_p')['venta_d'].sum().reset_index()
3 # Ordenar de mayor a menor
4 pregunta_2 = pregunta_2.sort_values(by='venta_d', ascending=False)
5 pregunta_2 = pregunta_2.head(10)
```

Python

```
1 display(pregunta_2)
```

Python

	nombre_p	venta_d
8	Samantha Ramirez	190396.30
2	Jennifer Daugherty	182366.90
1	Cassandra Torres	178075.10
4	Kathy Perez	168078.03
7	Michael Cox	140611.37
5	Margaret Goodwin	125581.82
3	Joseph Floyd	120635.09
9	Scott Mayer	109964.56
6	Matthew Brown	109882.56



¿Qué categoría tiene más ventas?

Ropa

```
1 pregunta_3= join_FACT_venta_dim_articulo.groupby('categoria_nombre')['venta_d'].sum().reset_index()
2 # Ordenar de mayor a menor
3 pregunta_3 = pregunta_3.sort_values(by='venta_d', ascending=False)
4 pregunta_3 = pregunta_3.head(10)
```

[27]

```
1 display(pregunta_3)
```

[28]

...

	categoria_nombre	venta_d
5	ROPA	335650.56
3	HOGAR	273236.57
0	ALIMENTOS	246790.17
4	JUGUETES	243904.96
2	DEPORTES	175287.38
1	CALZADO	112887.24

¿Cual es el rol de usuario que mas transacciones ha hecho?

Jefe Tienda

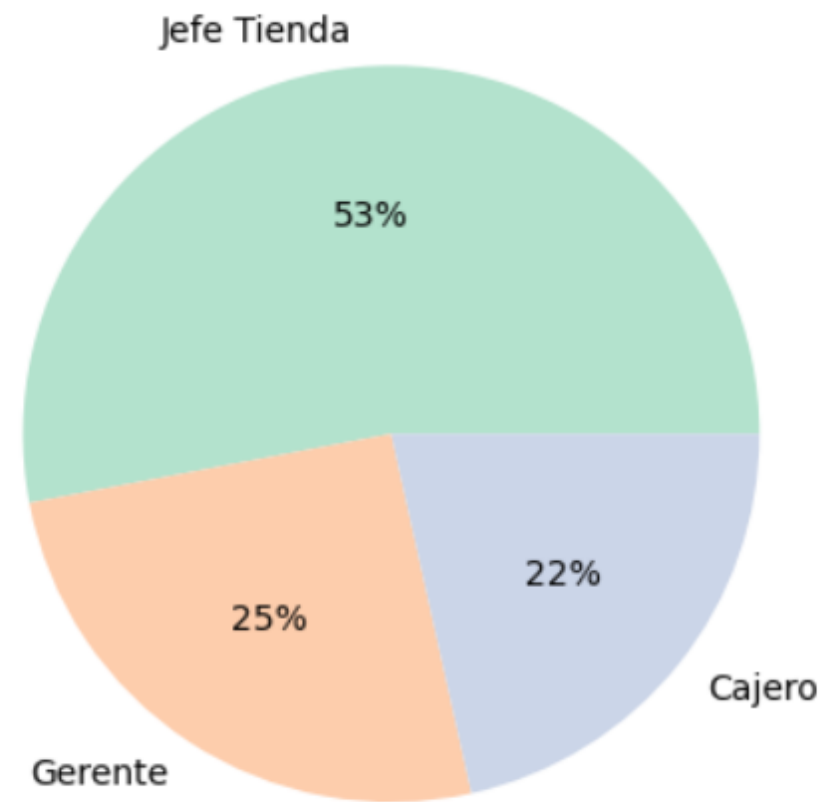
```
1 # sumar los registros unicos de venta e ingresos de cada usuario y luego agregar el rol del usuario a la tabla
2 pregunta_4 = (df_FACT_ingreso.groupby(by='idusuario')['idingreso'].nunique() + \
3               df_FACT_venta.groupby(by='idusuario')['idventa'].nunique()).reset_index().\
4               rename({0:'Transacciones'}, axis=1).\
5               merge(df_dim_usuario[['idusuario','rol_nombre']], on='idusuario', how='left')
6 # agrupar por roles y ordenar de mayor a menor
7 pregunta_4 = pregunta_4.groupby(by='rol_nombre', as_index=False)[['Transacciones']].sum().\
8               sort_values('Transacciones', ascending=False)
9 pregunta_4
```

	rol_nombre	Transacciones
2	Jefe Tienda	106
1	Gerente	51
0	Cajero	43

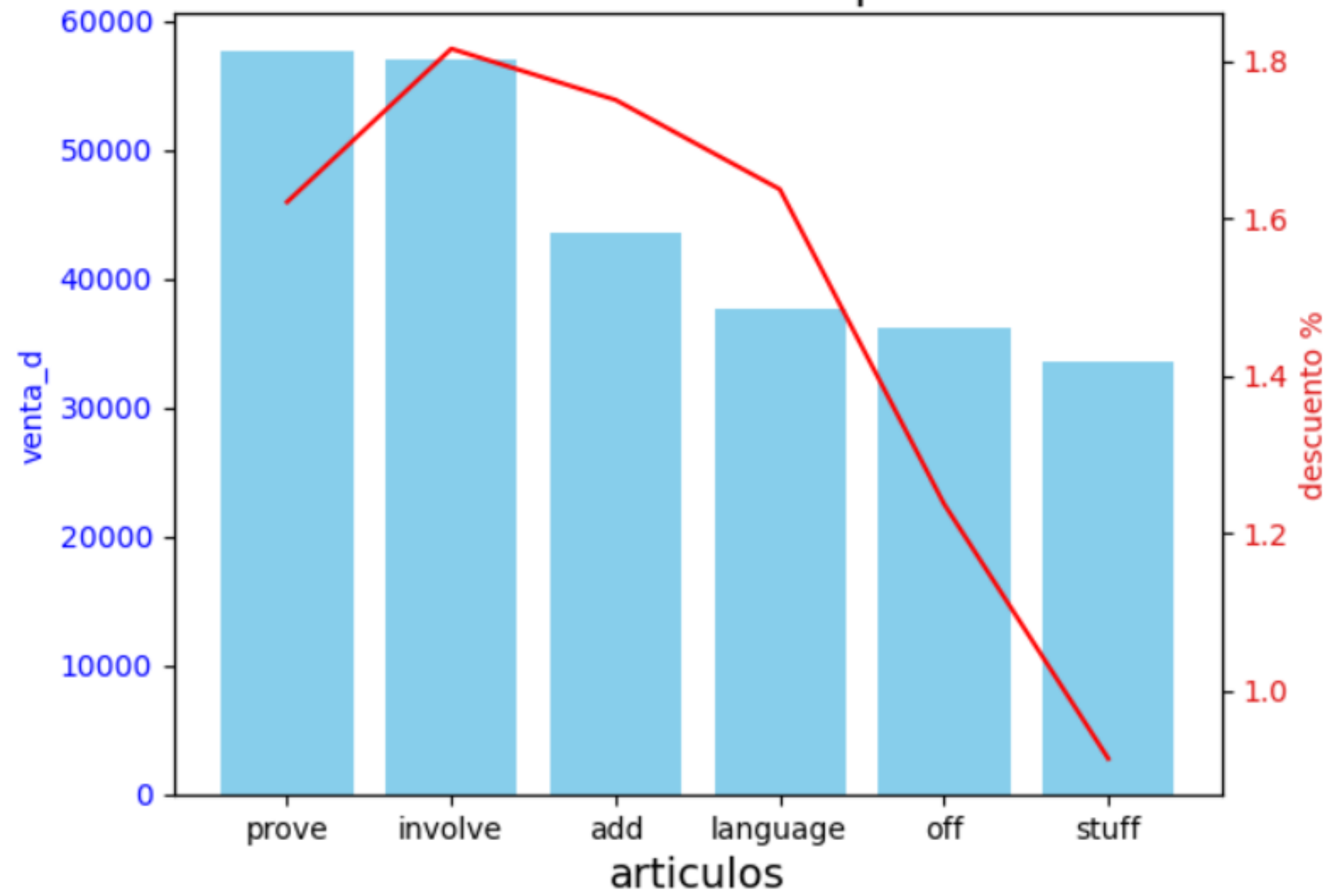
```
1 palette_color = sns.color_palette('Pastel2')
2 plt.pie(pregunta_4['Transacciones'], labels=pregunta_4['rol_nombre'], colors=palette_color, autopct='%0f%%')
3 plt.title('Distribución de las transacciones por rol de usuario', size=16)
```

```
Text(0.5, 1.0, 'Distribución de las transacciones por rol de usuario')
```

Distribución de las transacciones por rol de usuario



Articulos mas vendidos con su respectivo descuento %



Registros realizados por mes

