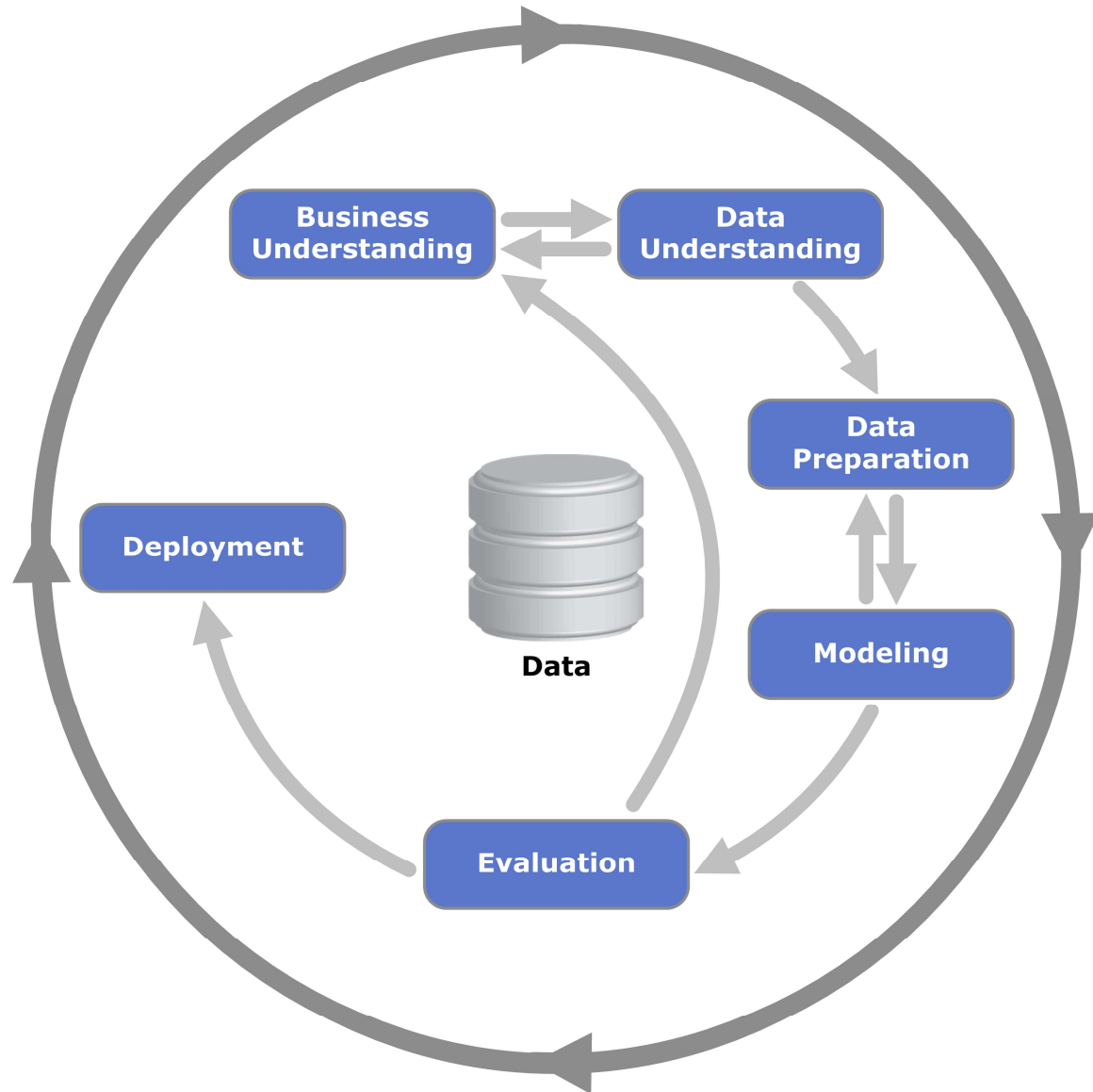


Clasificación de tipos de viajes Walmart



Metodología CRISP-DM

1. Comprensión del Negocio
2. Comprensión de los Datos
3. Preparación de los datos
4. Modelado
5. Evaluación
6. Implantación



Metodología CRISP-DM

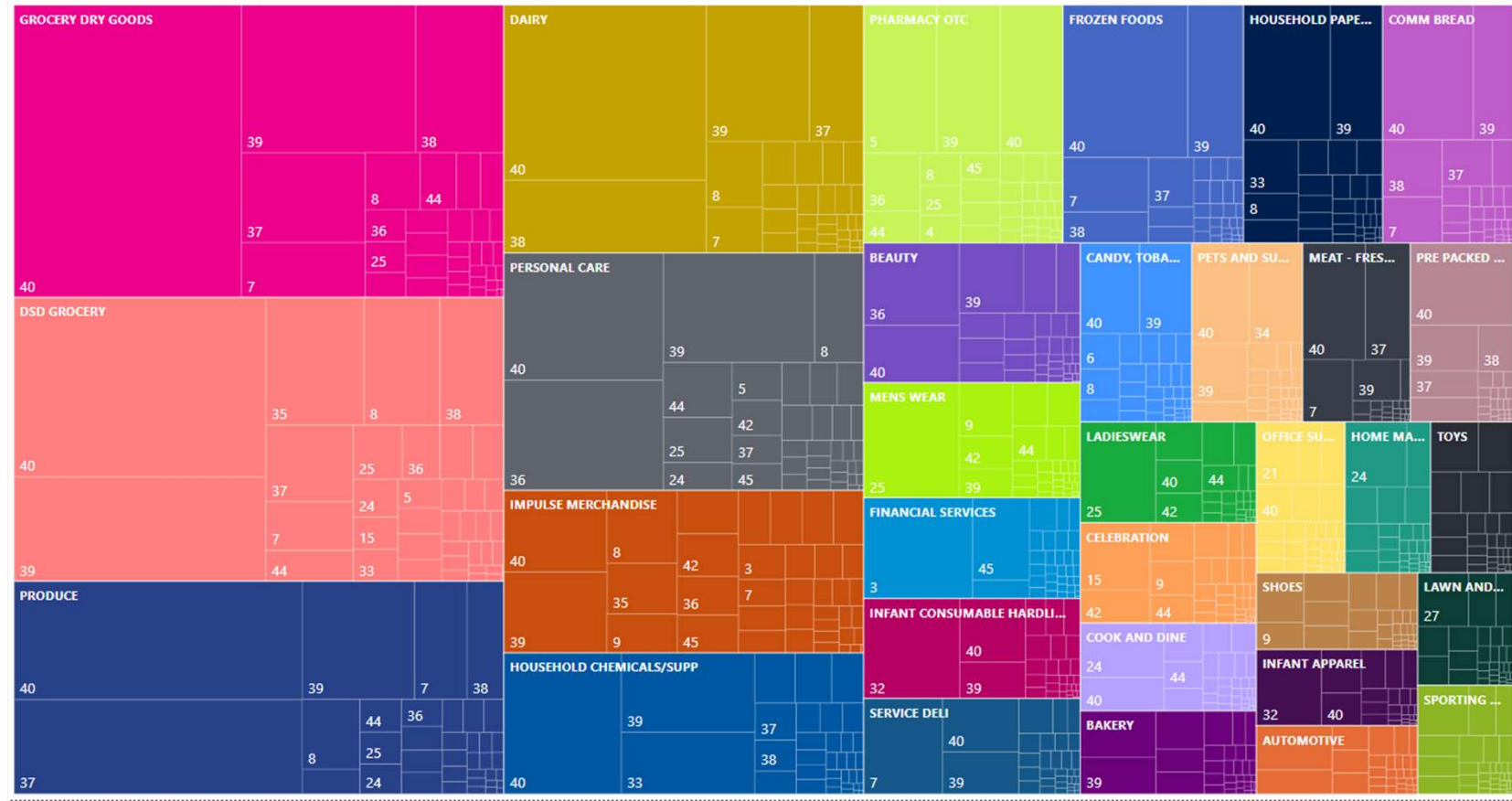


6. Implementación

Intentamos correr Flask pero por limitaciones de tiempo no se logró esa parte del proyecto.

2 COMPRENSIÓN DE LOS DATOS

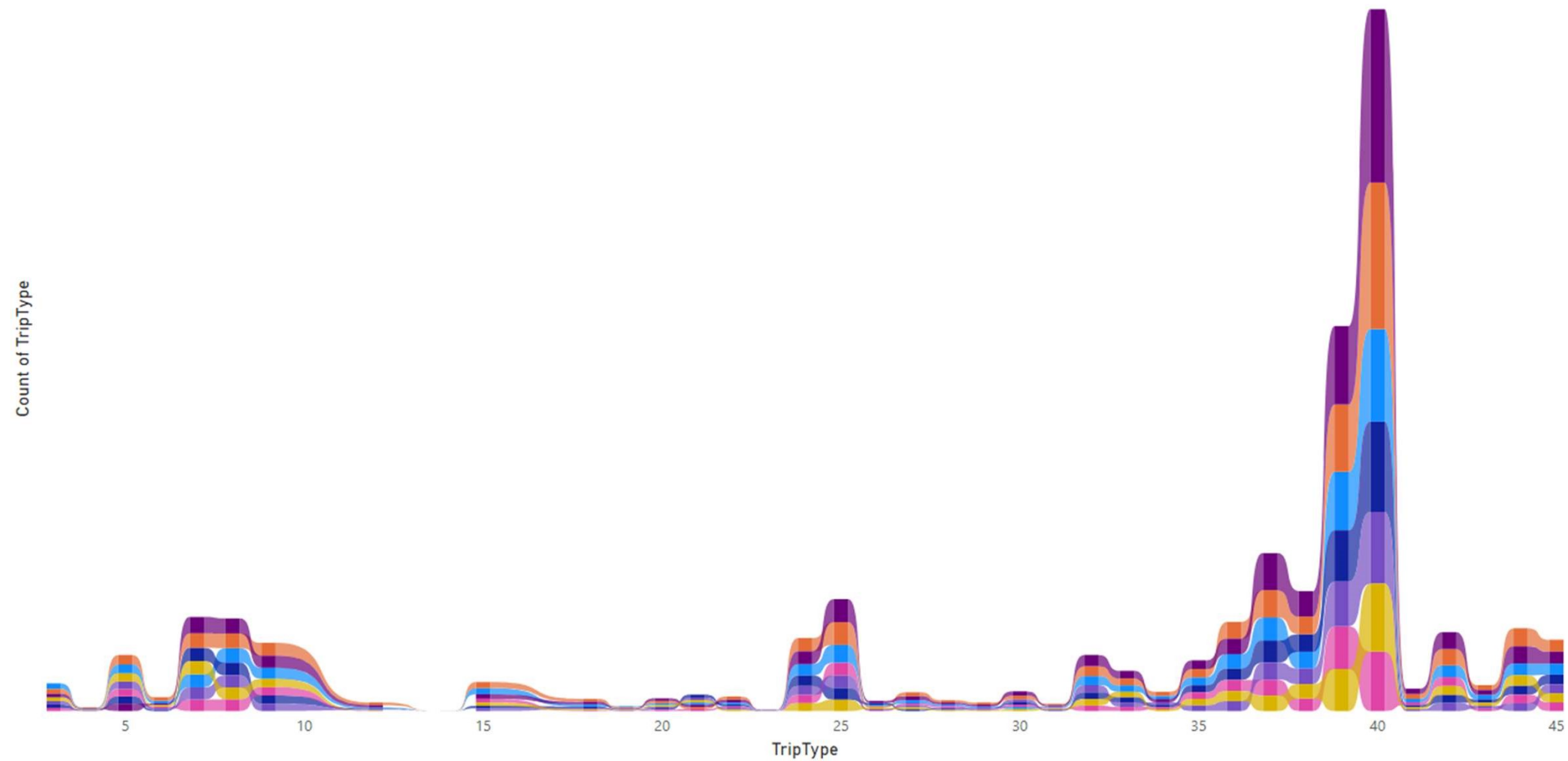
Count of VisitNumber by DepartmentDescription and TripType



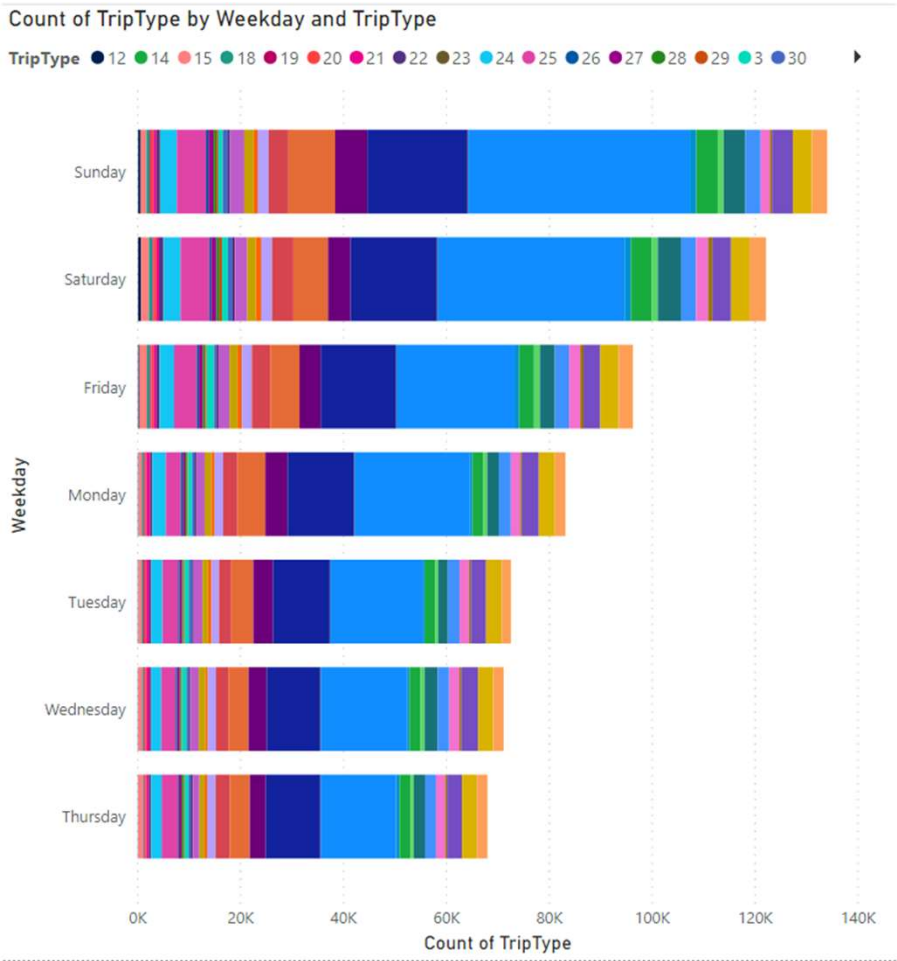
2.EDA

Count of TripType by TripType and Weekday

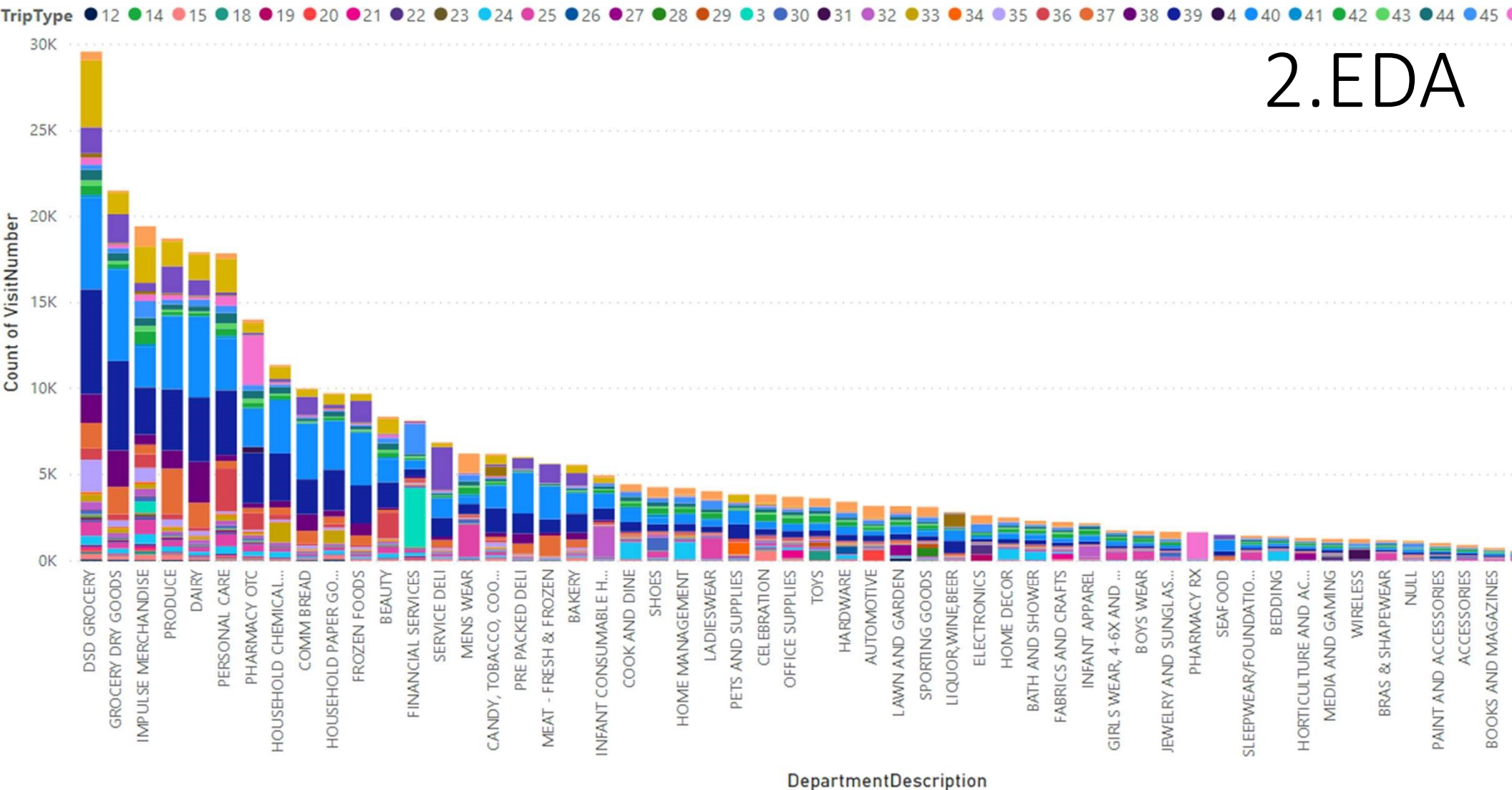
Weekday ● Friday ● Monday ● Saturday ● Sunday ● Thursday ● Tuesday ● Wednesday

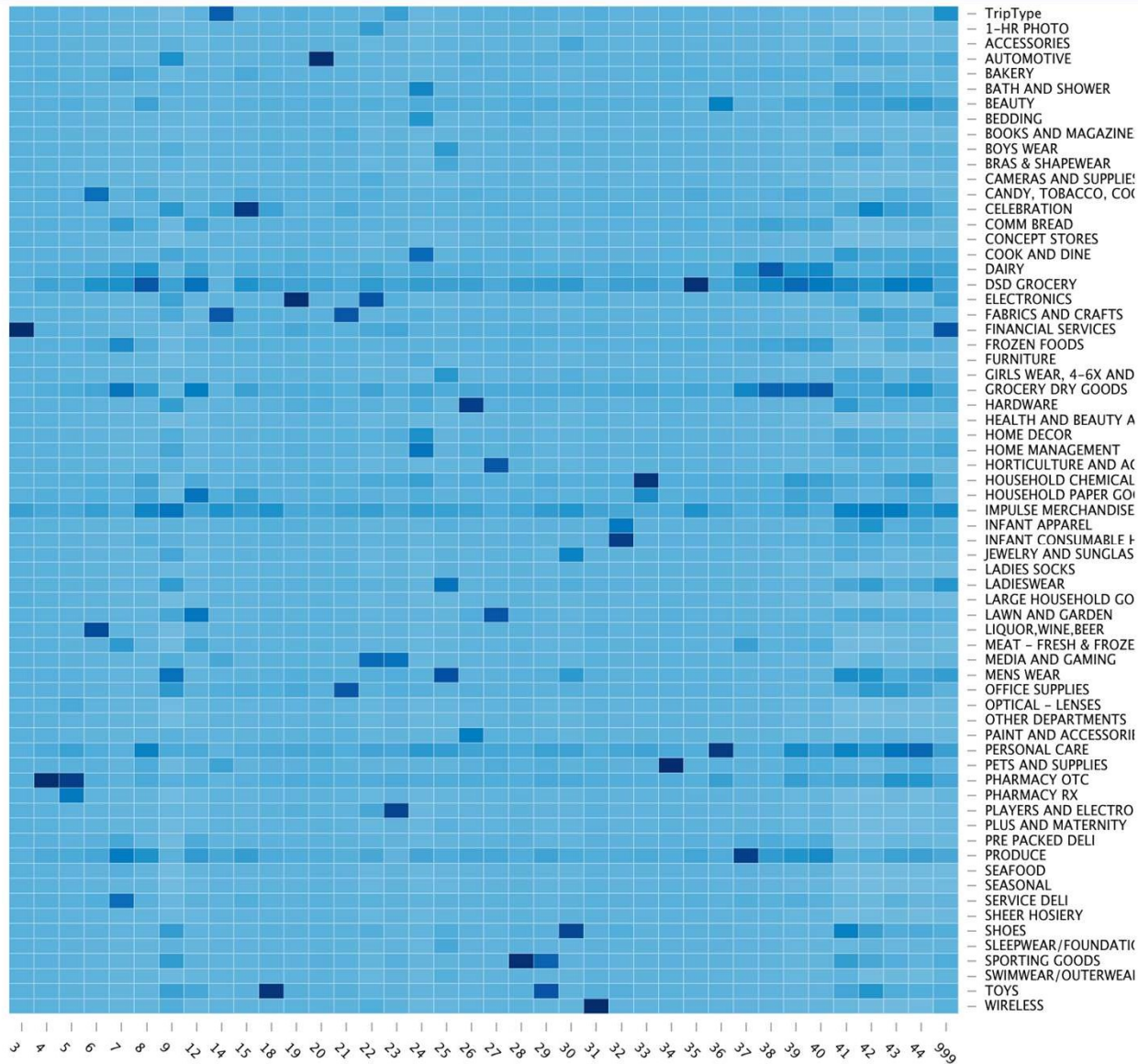


2.EDA



Count of VisitNumber by DepartmentDescription and TripType





3. PREPARACIÓN DE LOS DATOS

4. MODELADO

Microsoft Azure Machine Learning Studio (classic) Javier Vale

Training experiment Predictive experiment

Experiment created on 12/18/2019 Finished r

Search experiment items

- ▶ Saved Datasets
- ▶ Trained Models
- ▶ Data Format Conversions
- ▶ Data Input and Output
- ▶ Data Transformation
- ▶ Feature Selection
- ▶ Machine Learning
- ▶ OpenCV Library Modules
- ▶ Python Language Modules
- ▶ R Language Modules
- ▶ Statistical Functions
- ▶ Text Analytics
- ▶ Time Series
- ▶ Web Service

Flowchart illustrating the training experiment process:

```
graph TD; A[datos_sin_NAs.csv] --> B[Split Data]; A --> C[Multiclass Logistic Regression]; B --> C; B --> D[Train Model]; B --> E[Score Model]; C --> D; D --> E; E --> F[Evaluate Model]; F --> G[Convert to CSV];
```

The flowchart shows the following steps:

- Input data: `datos_sin_NAs.csv`
- Split Data (✓)
- Multiclass Logistic Regression (✓)
- Train Model (✓)
- Score Model (✓)
- Evaluate Model (✓)
- Convert to CSV (✓)

5. EVALUACIÓN Y RESULTADOS

9 submissions for Daniela Pinto Veizaga

Sort by Private Score

All Successful Selected

| Submission and Description | Private Score | Public Score | Use for Final Score |
|--|---------------|--------------|--------------------------|
| pre_GB_std.csv 6 hours ago by Daniela Pinto Veizaga add submission details | 1.32845 | 1.33002 | <input type="checkbox"/> |
| pre_GB.csv 6 hours ago by Daniela Pinto Veizaga XGboost | 1.40185 | 1.40690 | <input type="checkbox"/> |
| pre_logit.csv 8 hours ago by Daniela Pinto Veizaga third try | 1.42031 | 1.42349 | <input type="checkbox"/> |

6. CONCLUSIONES

1. Feature Engineering fue clave para resolver este problema, pues permitió reducir el error del modelo significativamente.
2. El problema más grande encontrado fue el tamaño del dataset, lo cual nos dificultó un poco hacer las transformaciones necesarias.
3. Usando herramientas como Azure ML Studio facilita y reduce el tiempo para entrenar y probar los modelos.
4. El mejor modelo fue gradient boosting después de aplicar feature engineering con un error de 1.32 lo cual nos dejó en el lugar 534

| | | | | | | |
|-----|-----|-------------|---|---------|---|----|
| 531 | ▼ 1 | Munch |  | 1.32293 | 9 | 4y |
| 532 | ▲ 1 | joaop |  | 1.32599 | 1 | 4y |
| 533 | ▲ 4 | Yilun Zhang |  | 1.32637 | 5 | 4y |
| 534 | ▲ 1 | NileshGupta |  | 1.33045 | 1 | 4y |

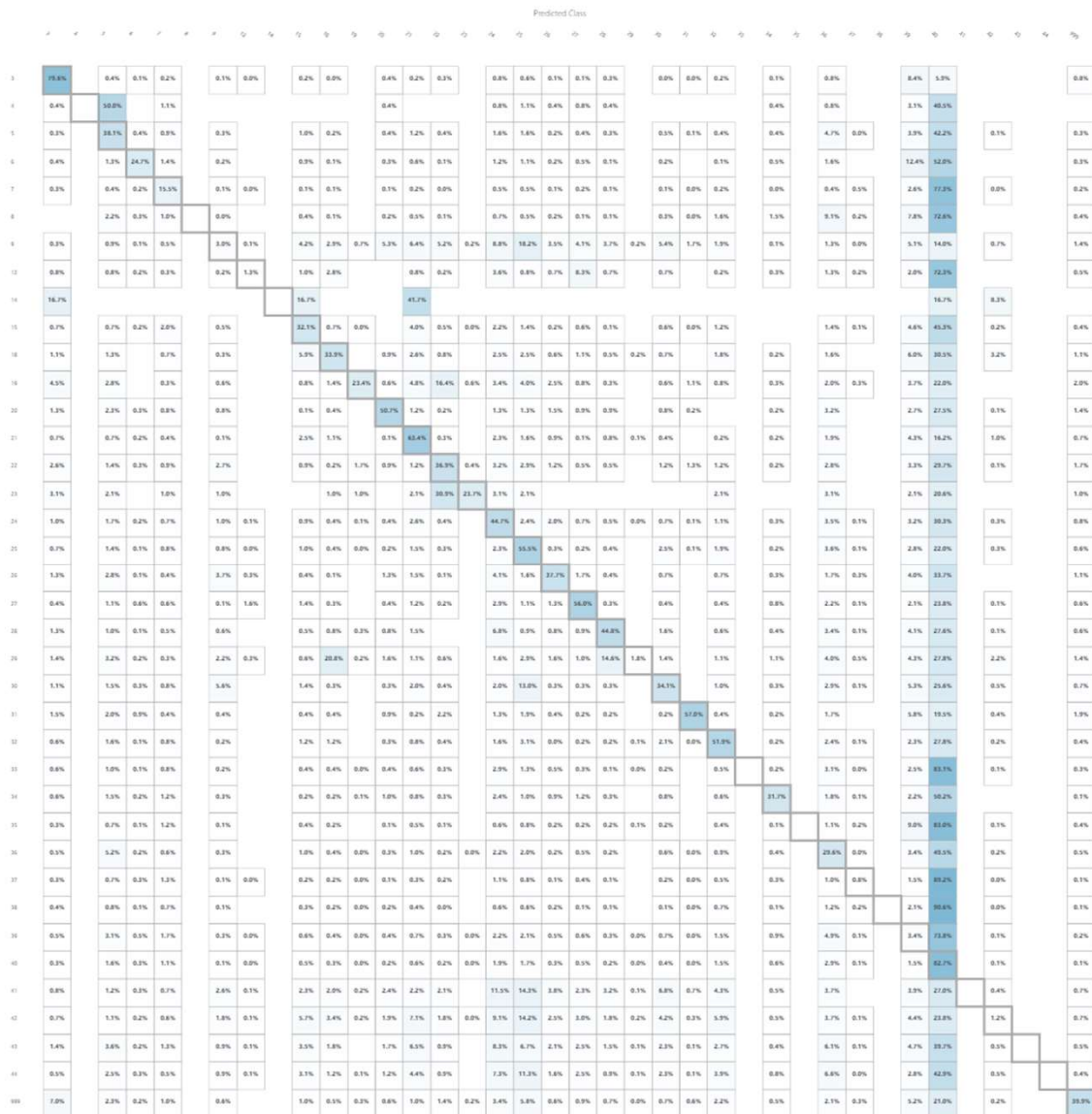
BONUS SLIDES - ANEXO

Hackeando Kaggle: resultados interesantes

7 submissions for [Javier Valencia Goujon](#) Sort by Private Score ▼

All **Successful** Selected

| Submission and Description | Private Score | Public Score | Use for Final Score |
|--|---------------|--------------|--------------------------|
| GB_std.csv a day ago by Javier Valencia Goujon prueba XGBOOST | 2.67421 | 2.70124 | <input type="checkbox"/> |
| logit_std.csv a day ago by Javier Valencia Goujon prueba modelo logistico | 3.47424 | 3.46926 | <input type="checkbox"/> |
| submission.csv a day ago by Javier Valencia Goujon Prueba con predicciones en cero mejor que el promedio de kaggle | 3.63758 | 3.63758 | <input type="checkbox"/> |
| logit_std.csv a day ago by Javier Valencia Goujon prueba dos modelo logistico | 12.39411 | 12.26618 | <input type="checkbox"/> |



Resultados Azure ML Studio

Metrics

| | |
|--------------------------|----------|
| Overall accuracy | 0.348963 |
| Average accuracy | 0.965735 |
| Micro-averaged precision | 0.348963 |
| Macro-averaged precision | NaN |
| Micro-averaged recall | 0.348963 |
| Macro-averaged recall | 0.262886 |