

Trabajo Práctico 2: Machine Learning

[75.06 / 95.58] Organización de Datos
Segundo cuatrimestre de 2018

Grupo Datatouille

Alumno	Padrón	Mail
Bojman, Camila	101055	camiboj@gmail.com
del Mazo, Federico	100029	delmazofederico@gmail.com
Hortas, Cecilia	100687	ceci.hortas@gmail.com
Souto, Rodrigo	97649	rnsoutob@gmail.com

<https://github.com/FdelMazo/7506-Datos/>

<https://kaggle.com/datatouille2018/>

Curso 01

- Argerich, Luis Argerich
- Golmar, Natalia
- Martinelli, Damina Ariel
- Ramos Mejia, Martín Gabriel

Contents

1	Introducción	1
2	Investigación	1
3	Feature engineering	2
3.1	Features básicos	2
3.2	Suma total de eventos	2
3.3	Cantidad de eventos por mes	3
4	Algoritmos utilizados	4
5	Desarrollo	4
6	Resultados obtenidos	4
7	Conclusiones	4

1 Introducción

Se propone exponer en el siguiente informe el desarrollo del Trabajo Práctico para predecir la probabilidad de que un usuario de Trocafone realice una conversión en un período determinado de tiempo. Con dicho propósito se utilizan como base dos archivos csv brindados por la empresa. En un primer lugar el archivo `events_up_to_01062018.csv` que contiene la información de los eventos realizados por un conjunto de usuarios hasta el 31 de mayo de 2018. En un segundo lugar el archivo `labels_training_set.csv` que determina si un conjunto de usuarios realizó o no una conversión desde el 01/06/2018 hasta el 15/06/2018.

Se propone como objetivo desarrollar una métrica a fin de evaluar los distintos resultados obtenidos. Se presentan en el informe las distintas decisiones adoptadas para elegir un resultado por sobre otro. Así mismo, se propuso utilizar distintos algoritmos de la librería `sklearn` para elegir el que modelice mejor.

Por otro lado se crearon distintos csv a fin de poder desarrollar de mejor manera el proceso de feature engineering y obtener un pre-procesamiento de los datos que permita encontrar los resultados más altos.

2 Investigación

Como se mencionó en la *Introducción*, se realizó un proceso de exploración de los datos para crear nuevos atributos y extraer nuevos features. De esta manera se crearon distintos csv con información extraída tanto de internet como de los propios datos para luego concatenar al set de datos final. Se procede a dar una breve explicación de la funcionalidad de cada uno de ellos. Muchos de ellos se encuentran detallados en el 'Notebook Anexo' del TP1¹.

- `brands.csv`

Se agregó una columna al dataframe que detalla qué marca está involucrada en el evento del usuario.

- `os.csv`

Se agregó una columna al dataframe que detalla qué sistema operativo está involucrada en el evento del usuario.

- `browsers.csv`

Se agregó una columna al dataframe que detalla qué explorador de internet se accede al sitio.

- `sessions.csv`

Se agregó el concepto de sesión, que se define como la agrupación de una serie de eventos por usuario, los cuales están todos con menos de 30 minutos de inactividad entre el actual y el anterior. Esto fue fijado con un criterio arbitrario a fin de poder discretizar el tiempo y definir este concepto.

- `prices.csv`

¹<https://fdelmazo.github.io/7506-Datos/TP1/TP1.html>

Se agregó una columna al dataframe que indica el precio del producto involucrado en el evento del usuario. Para ello se extrajeron los precios de la página de Trocafone² considerando el sku, el modelo, el color, la capacidad de almacenamiento y la condición.

3 Feature engineering

A partir del nuevo dataframe obtenido con la unión de todos los csv descritos en el inciso anterior se procedió a la búsqueda de features. En esta etapa del desarrollo del Trabajo Práctico se buscó explotar las distintas ideas y después con un proceso de selección que será explicado más adelante elegir los features pertinentes y más útiles al modelo.

3.1 Features básicos

Se detallan los features generales considerados como pertinentes al modelo.

- `is_viewed_product`: el usuario vió un producto
- `is_checkout`: el usuario llegó a checkout con un producto
- `is_conversion`: el usuario compró un producto
- `session_checkout_first`: el usuario en su primera sesión realizó un checkout
- `session_conversion_first`: el usuario en su primera sesión realizó una conversión
- `session_ad_first`: el usuario en su primera sesión llegó con una campaña publicitaria
- `session_ad_checkout_event`: el usuario en su primera sesión llegó con una campaña publicitaria e hizo checkout
- `session_ad_conversion_event`: el usuario en su primera sesión llegó con una campaña publicitaria y compró el producto

3.2 Suma total de eventos

A los features agregados como *features básicas* se le calcula el total por usuario y se obtienen el siguiente listado de features:

- `total_viewed_products`: cantidad de productos que vio el usuario en el período de tiempo determinado.
- `total_checkouts`: cantidad de veces que el usuario hizo checkout en el período de tiempo determinado.
- `total_conversions`: cantidad de compras que realizó el usuario en el período de tiempo determinado.

²<https://www.trocafone.com/>

- `total_events`: cantidad de eventos totales que el usuario hizo en el período de tiempo determinado.
- `total_sessions`: cantidad total de sesiones del usuario
- `total_session_checkout`: cantidad total de sesiones donde el usuario hizo checkout
- `total_session_conversion`: cantidad total de sesiones donde el usuario convirtió.
- `total_events_ad_session`: cantidad total de sesiones donde el usuario ingresó por una campaña publicitaria.
- `total_ad_sessions`: cantidad total de sesiones donde el usuario ingresó por primera vez por una campaña publicitaria.

A partir de estos features se deducen los siguientes:

- `avg_events_per_session`: porcentaje de cantidad total de eventos sobre cantidad de sesiones
- `avg_events_per_ad_session`: porcentaje de cantidad total de eventos donde el usuario ingresó por una campaña publicitaria sobre cantidad total de sesiones donde el usuario ingresó por una campaña publicitaria
- `percentage_session_ad`: porcentaje de cantidad total de sesiones donde el usuario ingresó por primera vez por una campaña publicitaria sobre el total de sesiones
- `percentage_session_conversion`: porcentaje de cantidad total de sesiones donde el usuario ingresó por primera vez y compró sobre la cantidad total de sesiones

3.3 Cantidad de eventos por mes

Se agregan una serie de features relacionados a la cantidad de eventos y sesiones por mes que se consideraron pertinentes al modelo.

- `total_viewed_products_month`: cantidad de productos vistos por mes por usuario
- `total_checkouts_month`: cantidad de productos que llegaron a checkout por mes por usuario
- `total_conversions_month`: cantidad de productos que llegaron a ser comprados por mes por usuario
- `total_events_month`: cantidad de eventos por mes por usuario
- `total_sessions_month`: cantidad total de sesiones por mes
- `total_session_checkouts_month`: cantidad total de sesiones donde el usuario hace checkout por mes

- `total_session_conversions_month_`: cantidad total de sesiones donde el usuario compra un producto por mes
- `total_events_ad_session_month_`: cantidad total de sesiones donde el usuario ingresa a la página por una campaña publicitaria por mes
- `total_ad_sessions_month_`: cantidad total de sesiones donde el usuario ingresa a la página por primera vez por una campaña publicitaria por mes

3.4 Eventos sin contar mayo

3.5 Eventos en última semana

4 Algoritmos utilizados

5 Desarrollo

6 Resultados obtenidos

7 Conclusiones