

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

Índice

1	Introducción	1
2	Organización del Trabajo	1
3	Features	2
3.1	Investigación previa	2
3.2	Creación de dataframes	3
4	Feature engineering	3
4.0.1	Features básicos	3
4.0.2	Suma total de eventos	4
4.0.3	Cantidad de eventos por mes	5
4.0.4	Eventos sin contar mayo	5
4.0.5	Eventos en última semana	5
4.0.6	Distribución mensual de las conversiones	5
4.0.7	Informacion de los últimos eventos registrados por usuario	5
4.0.8	Precios de la ultima conversion realizada por el usuario	6
4.0.9	Porcentaje de la actividad de la ultima semana	6
4.0.10	Porcentaje de la actividad del ultimo mes	6
4.0.11	Días entre el último checkout y última actividad	7
4.0.12	Estados de celulares	7
4.0.13	Varianza logarítmica de productos vistos	7
4.0.14	¿El usuario compró más de la media?	7
4.0.15	¿Cuántas veces vio el último modelo que compró?	7
4.0.16	¿Cuantas veces vio la última marca que compró?	7
4.0.17	Comportamiento en sesiones de las últimas semanas	7
4.1	Feature selection	8
5	Algoritmos utilizados	8
5.1	Parameter tuning	8
6	Desarrollo	8
6.1	Submission framework	8
6.2	Encontrando el mejor submit	8
7	Resultados obtenidos	9
8	Conclusiones	9

1. Introducción

El objetivo principal del trabajo es el de predecir la probabilidad de que un usuario de la empresa Trocafone realice una compra de un dispositivo celular (conversión).

La realización del trabajo se hace con algoritmos de Machine Learning, una disciplina que busca poder generar clasificaciones en base a un entrenamiento sobre información pasada, seguida de una validación de las predicciones generadas. En el trabajo se prueban distintos algoritmos, los cuales todos en distinta manera hacen uso de los datos (en particular, de sus atributos). Es por esto que es muy importante saber que datos usar, y buscar como codificarlos de tal forma que mejor se aprovechen.

Con dicho propósito se utilizan como base dos sets de datos 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 desde el 1ro de enero hasta el 31 de mayo de 2018 y servirá como entrenamiento de los algoritmos, y en un segundo lugar el archivo `labels_training_set.csv` que determina si el usuario realizó o no una conversión desde el primero hasta el quince de junio el cual servirá de validación.

2. Organización del Trabajo

Para un desarrollo más cómodo, se modularizó el trabajo en distintos notebooks de Jupyter, y luego se importan entre sí ¹, intentando emular la programación estructurada. Si bien esta forma de trabajar sirve bastante para no estar repitiendo código, el dividir en distintos notebooks que dependan unos de otros sube el acoplamiento del proyecto en su totalidad.

Los notebooks, en orden de lectura y corrida son:

1. Investigación Previa

<https://fdelmazo.github.io/7506-Datos/TP2/investigacion.html>

En este notebook se presentan las distintas exploraciones que se hicieron sobre el dataset del trabajo previo y del actual, en búsqueda de ideas útiles para el desarrollo de este.

2. Creación de dataframes

https://fdelmazo.github.io/7506-Datos/TP2/new_dataframes.html

En este notebook se crean los distintos dataframes necesarios para poder extraer atributos de los usuarios.

3. Feature engineering

https://fdelmazo.github.io/7506-Datos/TP2/feature_engineering.html En este notebook se agregan todos los features que se consideran que pueden ser pertinentes para el modelo. Este notebook es el que genera el `user-features.csv` del cual entrenan los modelos.

4. Feature selection

https://fdelmazo.github.io/7506-Datos/TP2/feature_selection.html

¹Haciendo uso del magnífico `nbimporter`

En este notebook se utilizan distintas formas de seleccionar los features con el objetivo de eliminar ruido y encontrar la mejor combinación de atributos a usar a la hora de entrenar modelos.

5. Parameter tuning

https://fdelmazo.github.io/7506-Datos/TP2/parameter_tuning.html

En este notebook se hacen diversas pruebas sobre cada algoritmo hasta encontrar los hiper parametros óptimos de cada uno.

6. Submission framework

https://fdelmazo.github.io/7506-Datos/TP2/submission_framework.html

La principal idea de este notebook es definir una serie de pasos para armar las postulaciones de predicciones del trabajo práctico.

7. Notebook principal

<https://fdelmazo.github.io/7506-Datos/TP2/TP2.html>

Finalmente, haciendo uso del framework previamente definido, se encuentra la combinación óptima de atributos y algoritmos para hacer una postulación de predicciones.

3. Features

3.1. Investigación previa

El primer paso del trabajo consistió en realizar una investigación sobre lo ya hecho en el trabajo anterior. El TP1 ² es un análisis exploratorio de datos de la empresa. Si bien no son exactamente los mismos datos que los trabajados acá, son de la misma índole, y la exploración de ellos dan a luz a patrones en los usuarios del sitio.

Esta investigación se compone de dos partes, una técnica y otra teórica.

Por el lado técnico, viendo que se uso otro set de datos para el TP1, se buscó alguna forma de integrar los datos anteriores con los nuevos (por ejemplo, buscar si hay usuarios compartidos entre los dos sets, o si hay compras para registrar), para poder usar una base de datos más grande tanto para el entrenamiento como la validación de las predicciones.

Luego de una serie de pasos, búsquedas y validaciones, se recopiló la siguiente información:

1. No se repiten usuarios en los datasets.
2. En el primer dataset (TP1) hay 27624 usuarios de los cuales 13967 tuvieron actividad en junio. Entre el 1 y el 15 (inclusive) de junio 82 usuarios compraron productos.
3. En el segundo dataset hay 19414 usuarios de los cuales 980 compraron en Junio.

Por lo tanto se concluyó que hacer un merge de los datos del TP1 con los del TP2 presentaría un *skewness* en el set de datos, por la despreciabilidad de estos.

²<https://fdelmazo.github.io/7506-Datos/TP1/TP1.html>

Por otro lado, en un marco teórico, se vió el análisis hecho en búsqueda de que patrones, ideas y conceptos pueden ser aplicados en este trabajo. En particular, se buscan atributos escondidos en el set original que puedan ser codificados de tal forma que luego los algoritmos de Machine Learning puedan utilizar a su favor. Los atributos encontrados son especificados en la sección de Feature Engineering.

TO-DO TSNE de Souto aca

3.2. Creación de dataframes

Este es mayoritariamente un re-trabajo sobre lo hecho para el TP1, en el Notebook Anexo³.

Como parte del feature engineering, se crean dataframes nuevos con información de los productos del sitio y de como se accede a este. Los dataframes generados son:

- **brands.csv**: Lista las marcas de los dispositivos de cada evento.
- **os.csv**: Lista los sistemas operativos desde los cuales se accedió al sitio.
- **browsers.csv**: Lista los exploradores desde los cuales se accedió 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.
- **prices.csv**: Lista los precios de los dispositivos de cada evento. Para lograr esto se hizo un *web-scraping* de la página de Trocafone de la cual se extrajo para cada conjunto de modelo, capacidad, color y condición el precio del dispositivo.

4. Feature engineering

Con todos los dataframes generados previamente y lo investigado del previo trabajo, se busca todo tipo de atributos de los usuarios, para que luego puedan ser seleccionados y aprovechados por los algoritmos a aplicar.

4.0.1. 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

³<https://fdelmazo.github.io/7506-Datos/TP1/anexo.html>

- **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

4.0.2. Suma total de eventos

A los features agregados como *features básicos* 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.
- **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

4.0.3. 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

4.0.4. Eventos sin contar mayo

4.0.5. Eventos en última semana

4.0.6. Distribución mensual de las conversiones

Se agrega en cuántos meses el usuario compró suponiendo que dicha distribución denota si el usuario es un comprador habitual o sólo compró alguna vez aisladamente. El feature se llama `amount_of_months_that_have_bought`.

4.0.7. Información de los últimos eventos registrados por usuario

Se busca extraer información de los días que transcurrieron hasta el último evento de un usuario. De esta manera se espera que el modelo aprenda un factor importante para la predicción. Por ejemplo, si un usuario vio un producto hace muchos días es muy probable que no lo compre pero si hizo checkout hace 1 día es probable que en un futuro cercano compre.

- `days_to_last_event`: cantidad de días hasta el último evento
- `days_to_last_checkout`: cantidad de días hasta el último checkout. Si el usuario no hizo checkout se considera un número mayor a la cantidad de días del período de tiempo comprendido.

- **days_to_last_conversion**: cantidad de días hasta la última compra del usuario. Si el usuario nunca compró se considera un número mayor a la cantidad de días del período de tiempo comprendido.
- **days_to_last_viewed_product**: cantidad de días hasta el último día que el usuario vio un producto. Si el usuario nunca vio un producto se considera un número mayor a la cantidad de días del período de tiempo comprendido.

En paralelo con estos features se consideran los días de la semana, del mes, del año y la semana del año donde ocurren estos últimos eventos.

4.0.8. Precios de la ultima conversion realizada por el usuario

Se consideró que podría considerarse el precio de la última conversión del usuario como un feature pero a la hora de la selección reflejó una importancia muy baja. Por lo tanto consideramos impertinente la descripción de la idea que habíamos pensado desarrollar.

4.0.9. Porcentaje de la actividad de la ultima semana

Aquí la idea pensada era reflejar la cantidad de eventos del usuario de la última semana sobre el total. Si el usuario ingresó muchas veces a la página en la última semana de mayo es muy probable que compre en la primera semana de junio. De la misma manera, si el usuario compró la última semana de mayo es probable que no compre por las siguientes dos.

Por lo tanto se pensaron los siguientes features:

- **percentage_last_week_activity**: porcentaje de la cantidad de eventos de esa semana sobre el total de eventos
- **percentage_last_week_conversions**: porcentaje de la cantidad de compras de esa semana sobre el total de eventos
- **percentage_last_week_checkouts**: porcentaje de la cantidad de checkouts de esa semana sobre el total de eventos
- **percentage_last_week_viewed_products**: porcentaje de la cantidad de productos vistos de esa semana sobre el total de eventos

4.0.10. Porcentaje de la actividad del ultimo mes

Una lógica análoga a la sección precedente se sigue en esta parte. Los motivos de este feature son simplemente una ampliación de la idea anterior. Si el usuario ingresó muchas veces a la página en mayo es muy probable que compre en la primera semana de junio. De la misma manera, si el usuario compró en mayo es algo probable que no compre por las siguientes dos.

De más está decir que se pensaron los siguientes features:

- **percentage_last_month_activity**: porcentaje de la cantidad de eventos de ese mes sobre el total de eventos
- **percentage_last_month_conversions**: porcentaje de la cantidad de compras de ese mes sobre el total de eventos

- `percentage_last_month_checkouts`: porcentaje de la cantidad de checkouts de ese mes sobre el total de eventos
- `percentage_last_month_viewed_products`: porcentaje de la cantidad de productos vistos de ese mes sobre el total de eventos

4.0.11. Días entre el último checkout y última actividad

La intención de este feature es medir la diferencia de días que tiene cada usuario entre la compra de un celular y la ultima vez que visualizo el producto comprado. De esta forma poder predecir en base a los productos vistos si es posible que se haga una compra.

4.0.12. Estados de celulares

Utilizando la lógica de que hay empresas que compran celulares en mal estado con el único fin de usar sus partes como respuestos se plantea agregar una columna que indique porcentaje de celulares en estado Bom - Sem Touch ID vs Bom sobre todos los celulares vistos.

4.0.13. Varianza logarítmica de productos vistos

Se propone analizar la varianza en los precios de los productos visitados. Es decir, si los usuarios ven telefonos de un rango pequeño de precio o, por el contrario, articulos de precios muy variados. Se utilizó una escala logarítmica para seguir manteniendo las proporciones sin tener una gran diferencia entre la varianza de un usuario y la de otro.

4.0.14. ¿El usuario compró más de la media?

Se propone como feature evaluar si el usuario compró un celular por encima de la media de precios.

4.0.15. ¿Cuántas veces vio el último modelo que compró?

La idea de este atributo es evaluar una cierta correlación entre los usuarios de la cantidad de veces que se ve un modelo antes de comprarlo. Si un usuario ve una cantidad significativamente grande de veces un modelo es muy probable que lo compre. Esto no quiere decir que sea imposible que un usuario pueda ver una vez un modelo y no comprarlo o verlo muchas veces y no comprarlo pero se busca analizar el caso más general.

4.0.16. ¿Cuántas veces vio la última marca que compró?

Una lógica similar a la sección precedente es la que se sigue con este feature. La idea sería también analizar si un usuario se restringe a un modelo en particular o si puede ver distintos celulares de la misma marca y elegir uno de ellos.

4.0.17. Comportamiento en sesiones de las últimas semanas

Se presenta una idea similar a la de las secciones 3.9 y 3.10:

4.1. Feature selection

En este notebook se utilizan distintas formas de seleccionar los features de manera de eliminar aquellos atributos que resulten ruidosos con el modelo de **Random Forest**. Se eligió dicho modelo por su popularidad para la selección de features ya que los árboles que crea el algoritmo toman distintos subconjuntos de atributos tomados al azar y arrojan distintos resultados. De esta manera, con cientos o miles de árboles el algoritmo adopta una amplia capacidad predictora de cada atributo.

Se utiliza como métrica el AUC debido a que es la utilizada por la plataforma de Kaggle para evaluar la eficiencia de los distintos modelos utilizados.

Las distintas formas de selección utilizadas se describen como sigue:

- **Cumulative importance**

El nombre del método fue inventado por el grupo de Trabajo y denota el método más intuitivo para la selección de features. Con el uso del algoritmo de **Random Forest** se parte de una lista de todos los features ordenados según importancia y se genera una lista de listas que agrega un feature a la vez. Por ejemplo siendo a,b,c features se parte de una lista como [a,b,c] y luego se obtiene [[a], [a,b], [a,b,c]]. El objetivo de este método es encontrar el *codo*, es decir, los features que incrementan el AUC local.

- **Forward Selection**

Con este método se comienza con ningún atributo y en cada paso se agrega el atributo que genere mejor resultado. Se agregan atributos siempre y cuando los resultados mejoren. El algoritmo termina cuando el resultado no se puede mejorar o cuando ya se han agregado todos los atributos.

- **Backward Selection**

Este método funciona a la inversa de **Forward Selection**. Se comienza con todos los atributos y se quita en cada iteración el atributo que aumente el resultado de la métrica. De esta manera el algoritmo termina cuando al quitar un atributo el resultado empeora o cuando ya no hay más atributos por quitar.

- **Stepwise Selection**

Este método es una variante que combina los dos métodos anteriores. En cada paso se considera agregar o quitar una variable de manera de aumentar el AUC local.

5. Algoritmos utilizados

5.1. Parameter tuning

6. Desarrollo

6.1. Submission framework

6.2. Encontrando el mejor submit

TO-DO metricassss

En el notebook principal llamado TP2 se realiza el desarrollo principal del Trabajo Práctico. A grandes rasgos en un primer lugar se utilizan los dataframes con todos los atributos seleccionados. Luego se definen y aplican los algoritmos de clasificación para realizar los entrenamientos y posteriores predicciones de conversiones. Finalmente se arman las postulaciones de labels.

La elección del algoritmo para realizar el *submit* se hace en base a todos los algoritmos y a combinaciones duales de ellos. Las combinaciones se realizan con el algoritmo **Voting Classifier** que es de la librería de **sklearn**. Para cada una de estas combinaciones se utiliza un set de features diferente a fin de elegir el que arroje mejor resultado. Dichos set de features se obtienen a partir de la selección que fue detallada previamente.

Finalmente, una vez elegido el algoritmo o el ensamble de algoritmos predilecto se entrenó con todo el dataframe (o mejor dicho **X train**) para enviar el submit.

TO-DO plotsssss

7. Resultados obtenidos

8. Conclusiones

its a kind of magic