

Machine Learning para el análisis de texto

Daniel Jiménez M.

Universidad Nacional de Colombia

16 -11 -2020

- Brett Lantz, Machine Learning with R, 2016, PACKT
- Bradley B & Brandon G, Hands-On Machine Learning with R, 2020, The R Series.
- Silge G, Kuhn M, Tidy Modeling with R, 2020, The R Series.
- Silge G, Hvitfeldt E, Supervised Machine Learning for Text Analysis in R, 2020, The R Series.

Las librerías para esta clase son :

```
library(tidyverse)
library(tidymodels)
library(skimr)
library(themis)
```

Algunas definiciones importantes

La ciencia de datos, el Machine Learning y la AI son areas diferentes de estudio y trabajo, a pesar que se ha popularizado el termino **Data Scientist** como una sumatoria de los anteriores y por fines técnicos esto se mantendrá así. Más a continuación se presentarán algunas definiciones importantes y el contexto en el cual se desenvuelven.

Algunas definiciones importantes

La siguiente imagen es una tesis de Robinson (2017) sobre las diferencias en acciones de los data scientist:

- **Data science produces insights**
- **Machine learning produces predictions**
- **Artificial intelligence produces actions**

El foco principal de la ciencia de datos es la producción de insights, esto se logra de la siguiente manera :

- Los insights son entendibles ;
- Se basan en el entendimiento descriptivo y prescriptivo de un problema;
- Se intenta entender las causalidades entre las variables exploratorias;

Una definición clásica dice que el data science combina : Estadística, ingeniería de sistemas y conocimiento de negocios.

Pero a diferencia de la definición clásica que es la que provee al nombre de Data Scientist a la persona que hace un poco de todo, es la siguiente:

Data Science

El data science es la disciplina que envuelve las siguientes corrientes:

- Inferencia estadística;
- Visualización de datos;
- Diseño de experimentos;
- Dominio del negocio y lo más importante;
- Comunicación.

El machine learning produce predicciones (llamese forecast o clasificación, detección de anomalías entre otros).

Se divide en :

- Análisis supervisado: los datos tienen la forma de $\{(x_1, y_1), \dots, (x_n, y_n)\}$
 - El objetivo es estudiar el comportamiento de la variable **y**, condicional a las variables respuestas.
 - Puntualmente se debe entender el comportamiento (distribución) de la variable **y** dada las explicativas.

- Análisis no supervisado: Los datos tiene la siguiente forma $\{x_1, \dots, x_n\}$
 - El objetivo es estudiar la variable x y los posibles conglomerados que se encuentren en ello.
 - Matemáticamente hay que estudiar la variable x

Son acciones, desarrolladas por algoritmos que entienden la lógica del funcionamiento de un entorno.

Una inteligencia artificial que todos usamos (o eso espero) es google maps, cuando selecciona la mejor ruta!

Machine Learning para la clasificación de textos

En esta sección veremos dos ejemplos, uno será a nivel de clasificación, donde hablaremos de las principales métricas para evaluar estos modelos.

Algunos de los algoritmos de clasificación son:

- Decision tree;
- Naive Bayes;
- Knn;
- Logistic Regression;
- SVM;
- Linear Discrimination Analysis;
- Random Forest

En esta oportunidad nos concentraremos en la regresión logística y el random forest

- Es un modelo basado en variables categóricas;
- Sus variables explicativas se dan en valores continuos
- No solo clasifica sino que genera probabilidad sobre los eventos

Condiciones para un problema logístico

- Un problema binario;
- Cuando se necesite la probabilidad de una predicción;
- Cuando exista un perímetro de decisión;
- Cuando se necesite saber el impacto de las features.

La función logística funciona a través de la sigmoide

Esta es la sigmoide

$$\sigma(\theta^t X) = \frac{1}{1 + e^{-\theta^t X}}$$

Notese que si $\theta^t X$ es muy grande entonces la función tiende a uno, en caso contrario a cero

###Proceso de la logística

- Se inicia en σ
- Se calcula $\hat{y} = \sigma(\theta^t X)$
- Se calculan los errores $y - \hat{y}$
- Se minimiza la función de costo

Función de costo : Mejora los parámetros de la regresión a través de interacciones del siguiente tipo

$$\sigma(w^T x + b) = \frac{1}{1 + e^{-(w^T x + b)}}$$

Donde

$$J = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(a^{(i)}) + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

cambiando así los pesos de los parámetros

Para mejorar la función de costo es necesario el gradiente descendiente

Gradiente descendiente

Función que interactivamente busca encontrar la aproximación de σ que minimice la función de costo.

- Precisión : Calidad del modelo

$$\frac{TP}{TP + FP}$$

Esto se puede leer como la precisión del modelo

- Recall: Porcentaje de positivos que tiene el modelo

$$\frac{TP}{TP + FN}$$

- F1-Score: Mide el rendimiento del modelo, basado en la siguiente relación

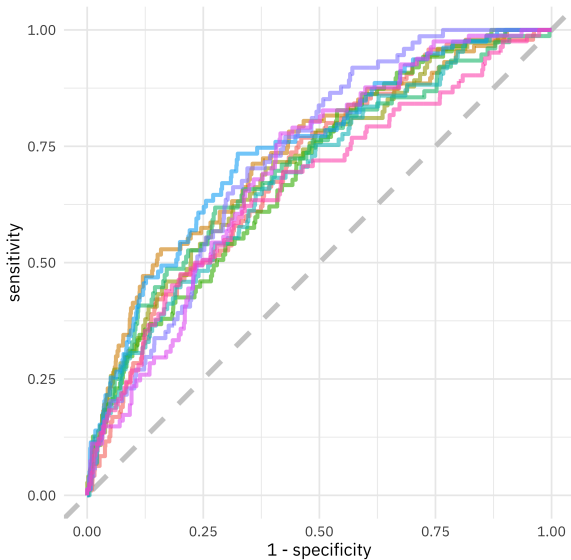
$$\frac{2 * (precision * recall)}{precision + recall}$$

Entre esta medida este más cercana a uno es mejor

- Accuracy : Mide el porcentaje de aciertos que tiene el modelo

$$\frac{(Tp + Tn)}{(Tp + Tn) + (Fp + Fn)}$$

Ejemplo



Random Forest

Es una adaptación de los árboles de decisión, en donde los valores de cada uno de sus árboles es independiente a nivel de interacción y dependiente de output del anterior.

Random Forest

- Es difícil de interpretar a nivel de salidas
- Es útil cuando los modelos primarios no tienen suficiente accuracy
- Su ventaja se basa en la aleatoriedad.

Random Forest

- Su ventaja primordial es que si se tiene un training set grande, el valor de sus predicciones será buenísimas
- Computacionalmente no es pesado
- Es super útil si se trabaja para estimar datos perdidos
- Posee buenos métodos experimentales

Observación : Tiende a sobre-ajustarse en tareas repetitivas de clasificación!!!!.