

DigitalHouse >
Coding School

DATA SCIENCE

UNIDAD 1
MÓDULO 2

Introducción a Pandas

Agosto 2017

INTRODUCCIÓN A PANDAS

1 Limpieza y preparación de datos

2 Qué es Pandas?

3 Para qué sirve Pandas?

4 Práctica guiada

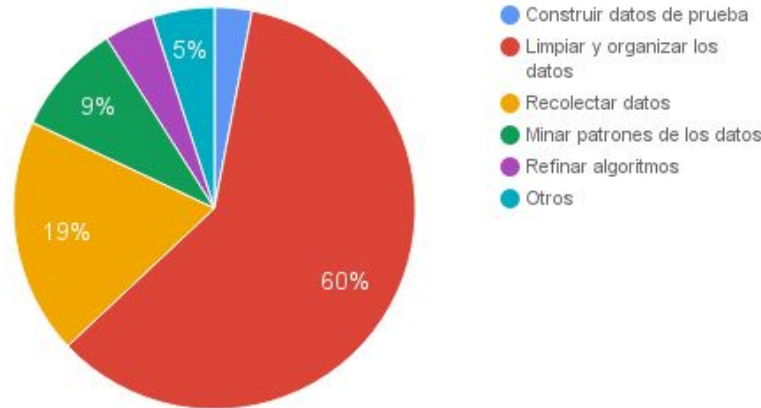
5 Práctica independiente

Limpieza y preparación de datos



- En 2009 Mike Driscoll (data scientist y CEO de Metamarkets) popularizó el término **“data munging”** para referirse al **arduo proceso de limpiar, preparar y validar los datos**

Como invierte su tiempo un data scientist?



Fuente:

<http://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

- **En 2013, Josh Wills** (ex director de Data Science de Cloudera y actual Director of Data Engineering en Slack) comenta: **"I'm a data janitor.** That's the sexiest job of the 21st century. It's very flattering, but it's also a little baffling."



Big Data Borat
@BigDataBorat

+ Follow

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

RETWEETS
506

LIKES
272



6:47 PM - 26 Feb 2013



12



506



272

Traducción:

En Data Science, se invierte un 80% del tiempo en preparar los datos y el 20% restante en quejarse de la necesidad de preparar los datos

- Proceso reproducible
- Control de versiones
- Creación de tests automáticos
- Facilita el mantenimiento
- Lenguaje dinámico que permite alta productividad
- Posibilidad de interfacear con C para performance
- Preferido en ambientes de data science

<http://lemire.me/blog/2014/05/23/you-shouldnt-use-a-spreadsheet-for-important-work-i-mean-it/>

Pandas



- **Colección de funciones y estructuras de datos** que facilitan el trabajo con datos estructurados
- Construido **en base a Numpy** inicialmente por Wes McKinney
- Nombre derivado de "**Panel Data System**" (término econométrico para datasets multidimensionales)
- Brinda capacidades flexibles de manipulación de datos similares a spreadsheets y bases de datos relacionales

- Combina la alta performance de las **operaciones sobre arrays de NumPy con la flexibilidad en la manipulación de datos** de un spreadsheet o una base de datos relacional
- Provee **funcionalidades de indexación avanzadas** para facilitar la manipulación, agregación y selección de partes de un dataset
- Provee **operaciones de agrupación por columnas, filtros y sumalizaciones**

- Veremos los siguientes objetos:
 - Series
 - DataFrames
 - Index

- En Pandas se utilizan dos estructuras de datos fundamentales ***Series*** y ***DataFrames***
- Ambas estructuras usan arrays de **Numpy** como base
- Una *Series* es un array unidimensional capaz de guardar cualquier tipo de datos (enteros, strings, floats, objetos Python, etc.)
- Un *DataFrame* es una matriz bidimensional. Puede verse como un conjunto de Series que comparten todos los mismos valores en el índice.

Index		Animales	Values	Name
	0	Perro		
	1	Oso		
	2	Jirafa		
	3	Tigre		
	4	Serpiente		
	5	Ratón		

- Una Series es un objeto similar a un **vector uni-dimensional**
- Contiene **un array de valores y un array asociado de etiquetas** de estos valores denominado como índice
- Una colección Series también puede ser considerado como un **diccionario de tamaño fijo con sus claves ordenadas**
- Comparados con los arrays de NumPy, permiten pasar una lista de valores del índice para seleccionar un subconjunto de valores

Eje 1 (columnas) →

Eje 0 (filas) ↓

df.iloc(2)

	Animales	Dueños
0	Perro	Juan
1	Oso	Pedro
2	Jirafa	Cristian
3	Tigre	Esteban
4	Serpiente	Pablo
5	Ratón	Claudio

df.iloc(5)['Animales'] **df['Dueños']**

- Representa una **estructura de datos tabular** que contiene una **colección de columnas**, cada una de las cuales tiene un tipo determinado (number, string, boolean, etc.)
- Inspirados en el paquete data.frame de R
- Permiten operaciones “ricas” sobre índices como los JOIN y GROUP BY en SQL
- Ideales para organizar el resultado de un análisis en un formato útil para graficar el resultado o mostrarlo

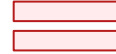
Práctica Guiada 1

Objetos en Pandas

	Nombre	Apellido
0	Juan	Perez
1	Pedro	García
2	Matías	Zabala



True
False
True



	Nombre	Apellido
0	Juan	Perez
1	Matías	Zabala

- Para acceder y subsetting a los objetos en Pandas es necesario introducir los siguientes métodos:
 - o .loc
 - o .iloc

Práctica Guiada 2

Subsetting e Indexing de Objetos en Pandas

LAB: Pandas con datos de la Encuesta Permanente de Hogares

Conclusiones

- En el día a día de un data scientist la limpieza, preparado y normalización de los datos con los que trabaja es la tarea que más tiempo insume
- Python es un lenguaje que numerosos features que nos pueden facilitar estas tareas
- Pandas es la librería principal al momento de tener que realizar estas tareas
- Los DataFrame de Pandas son la herramienta fundamental que vamos a utilizar a lo largo del curso y de nuestro día de trabajo a partir de ahora