

Flexible Session #1

Exploratory Data Analysis

Ivan Corneillet

Data Scientist

Learning Objectives

After this lesson, you should be able to:

- Review Step ③ Parse the Data and more specifically
 - Descriptive Statistics and Exploratory Data Analysis
 - Apply *pandas* on a Kaggle dataset
- Have fun doing Data Science!



DS

Announcements and Exit Tickets

DS

Review

Python and *pandas*

<i>Measure of Centrality</i>	<code>.mean()</code>	<code>.median()</code>	<code>.mode()</code>
<i>Measure of Dispersion</i>	<code>.var()</code> , <code>.std()</code>	<code>.min()</code> , <code>.max()</code> <code>.quantile()</code>	
<i>Summary</i>	<code>.describe()</code>		
<i>Graphical Methods</i>		<code>.plot(kind = 'box')</code>	<code>.plot(kind = 'hist')</code>
<i>Correlation Matrix</i>	<code>.corr()</code>		
<i>Scatter plot</i>	<code>DataFrame.plot(kind = 'scatter', x = 'SeriesName', y = 'SeriesName')</code>		
<i>Scatter matrix</i>	<code>pd.tools.plotting.scatter_matrix(DataFrame)</code>		
<code>.shape</code> , <code>.columns</code> , <code>.set_index()</code> , <code>.drop()</code>	<code>.count()</code> , <code>.sum()</code> , <code>.unique()</code> <code>.value_counts()</code> , <code>.isnull()</code> , <code>.notnul()</code> , <code>.dropna()</code> , <code>pd.crosstab(Series, Series)</code>		<code>np.sort()</code> , <code>.apply()</code>



DS

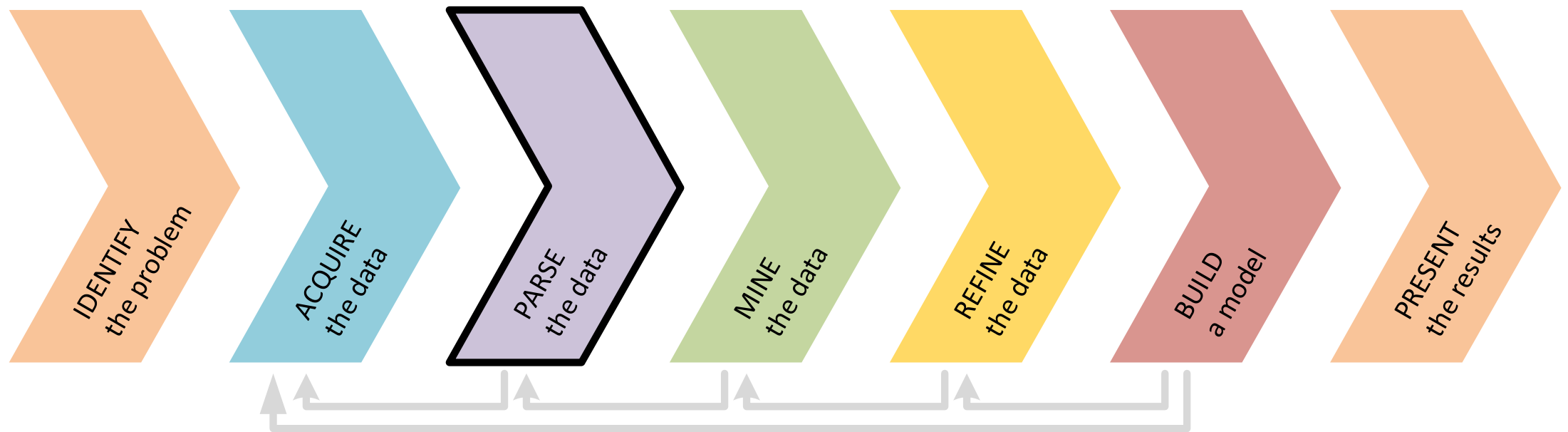
Q & A



DS

Today

Today we'll keep our focus on **PARSE** the data



And more precisely on the Exploratory Data Analysis using the *pandas* library

Research Design and Data Analysis	Research Design	Data Visualization in <i>pandas</i>	Statistics	Exploratory Data Analysis in <i>pandas</i>
Foundations of Modeling	Linear Regression	Classification Models	Evaluating Model Fit	Presenting Insights from Data Models
Data Science in the Real World	Decision Trees and Random Forests	Time Series Data	Natural Language Processing	Databases

Here's what's happening today:

- Announcements and Exit Tickets
- Review
- **③** Parse the Data
 - Kaggle – Exploratory Data Analysis
- Review
- Exit Tickets

DS

Kaggle

Exploratory Data Analysis

A black circle containing the white text "DS".

DS

Q & A

Next Class

Inferential Statistics for Model Fit

Learning Objectives

After this next lesson, you should be able to:

- Explain the difference between causation and correlation
- Identify a normal distribution within a dataset using summary statistics and visualization
- Test a hypothesis within a sample case study
- Validate your findings using statistical analysis (t-tests, p-values, t-values, confidence intervals)



DS

Exit Ticket

Don't forget to fill out your exit ticket [here](#)

Slides © 2016 Ivan Corneillet Where Applicable
Do Not Reproduce Without Permission