

Área de Data Science

Projeto Final - Trainees 2020 - Grupo Turing

Introdução

Olá!! Seja bem-vindx ao projeto final da área de [Data Science](#) do Grupo Turing!

Data Science tem uma função muito importante no mundo em que vivemos hoje. Na era da informação, a humanidade produz milhares e milhares de dados a cada momento, sejam eles sobre produtos, textos, informações de clientes, dados climáticos, dentre outros. Sendo assim, Data Science surge como a ciência responsável por coletar dados, entendê-los, processá-los, e, principalmente, conseguir extrair valor para contar uma história com eles.



Até agora, você já deve ter usado várias ferramentas de ciência de dados nas aulas do Turing Academy. Nesse projeto, você irá utilizar algumas novas técnicas, além de várias que você já deve conhecer e que poderão ser aprimoradas durante esse processo. Vamos lá!

Enunciado

Você já deve ter trabalhado com vários datasets, sejam eles de cursos do DataCamp ou de aulas que você participou no Turing. Mesmo sendo diferentes, esses datasets tinham uma característica em comum: todos foram previamente coletados e organizados por outra pessoa.

Nesse projeto, queremos que você trabalhe com os dados desde o início, ou seja, você deve **coletar os dados** da internet e realizar uma **análise de dados**. O site que

you should use for this project is <https://scrapethissite.com/pages/forms/>. It contains data about the performance of various ice hockey teams in recent years and, as the name suggests, it was made to be scraped, that is, for didactic purposes.



Besides that, remember that part of the process of working with data is understanding its history, so don't forget to also study a little about the problem you will be analyzing. For this, try to research a little about ice hockey rules, how the competition structure has changed over the years, if new teams have emerged in that period, etc. This will help a lot in the process of extracting interesting insights.

Web Scraping

The process of collecting data from the internet is called **Web Scraping**. For this, we have different tools that we can use, but we recommend that you use two libraries in this project: **Requests** and **BeautifulSoup**.

The two have different functions: we use the *requests* library to navigate through pages and the *BeautifulSoup* library is used to extract data from pages. Even so, generally you will always use these two together.

To do this task, look at the notebook of BeautifulSoup from the Grupo Turing available at [árvore de habilidades](#).

At the end of this step, you should have your collected and organized data in a .csv file.

Análise e Visualização de Dados

Now that you have collected the data, it's time to analyze it. To make this more interesting, we have some more data for you! It's in the *DataSet ESPN.csv* file, located in the Data Science folder in the Grupo Turing Drive. It covers the period from 2001 to 2011, and contains the following columns:

- *Team*: O nome da equipe
- *Home Games*: Quantidade de jogos disputados em casa, na temporada regular daquele ano
- *Home Total*: Número total de torcedores para todos os jogos em casa da temporada
- *Home Average*: Média de torcedores por jogo em casa da temporada
- *Road Games*: Quantidade de jogos disputados fora de casa, na temporada regular daquele ano
- *Road Average*: Média de torcedores por jogo fora de casa da temporada
- *Overall Games*: Total de jogos disputados na temporada regular daquele ano
- *Overall Average*: Média de torcedores por jogo da temporada (considerando tanto jogos em casa quanto fora de casa)
- *Year*: O ano em que a temporada terminou (por exemplo, se a temporada em questão for a de 2000-2001, o valor de “Year” é 2001)
- *Save Percentage*: Porcentagem de finalizações defendidas
- *Penalty Minutes*: Minutos de penalidade que o time sofreu
- *Penalty Minutes Against*: Minutos de penalidade que o time adversário sofreu

Caso você tenha ficado curioso, os dados foram extraídos desse site: <http://www.espn.com/nhl/statistics>

Talvez você ainda esteja se perguntando o que exatamente fazer com esses dados, mas a ideia é que isso não seja nada rigoroso mesmo: siga sua curiosidade.

Levante perguntas sobre os dados, e procure respondê-las você mesmo. Algumas ideias possíveis são: Será que times que marcam mais, ganham mais? Ou é o oposto que vale, e uma defesa sólida é mais importante? Como que o desempenho do goleiro afeta o resultado do time como um todo? Algum time tem evoluído muito nas últimas temporadas? A torcida é importante para o desempenho da equipe? Como o número de torcedores evoluiu ao longo do tempo?

Para quaisquer dúvidas sobre o projeto ou sobre onde procurar conteúdos, converse com seu mentor ou com a equipe de Data Science (Victor, Azank e Julia). Esperamos que você se divirta e que aprenda com essa experiência. Quando terminar, não esqueça de subir tudo no seu GitHub.