# MAC0459/MAC5865 - Tópicos em Ciência e Engenharia de Dados

**Aula 02**

**Sejam bem-vindas, sejam bem-vindos!**

**Entre no link https://app.sli.do/event/kgwf8yjm ou e faça suas perguntas da aula.**

**R. Hirata Jr.**

# Objetivos de hoje

- Ao final da aula de hoje você deve:
  - Conhecer alguns datasets abertos
  - Conhecer algumas perguntas simples que podem ser feitas para os datasets
  - Conhecer algumas perguntas interessantes que podem ser feitas para os datasets
  - Conhecer os tipos de dados

# Relembrando a semana passada

# Simple pipeline – Scientific Method

1. Pose a question
2. Formulate a hypothesis
3. Formulate an experiment
4. Observe (data collecting)
5. Analyse the results
6. Go back to step 2 if the hypothesis is not correct/supported
7. Report results

# Simple pipeline – Engineering Method

1. Define a problem

2. Specify requirements

3. Brainstorm, evaluate, choose solution

4. Develop a prototype solution

5. Tests solution

6. Go back to step 3 if the results, or data, etc in case the solution does not meet requirements

# Simple pipeline – Data Science Method

1. Pose question

2. Get the data

3. Explore the data

4. Model the data

5. Report results

# Repeatability vs Reproducibility

- Repeatability
  - variability caused by the measurement device
- Reproducibility
  - variability caused by different labs/operators

# CS, DS, Real Science

- "Computer scientists, by nature, don't respect data"
- Real scientists will only use as many significant digits as the worst precision of any measurement in the process.

# Hypothesis vs data driven science

- HD science - ask specific questions of the world and then generating the specific data to confirm, or deny it
- Data driven science - a new paradigm to model the world

# Learning to ask questions

- Computer scientists students are not used to ask questions, why?
- Good data scientists develop an inherent curiosity about the world around them and have wide-ranging interests.

# What is Science?

- "We absolutely must leave room for doubt or there is no progress and there is no learning. There is no learning without having to pose a question. And a question requires doubt. People search for certainty. But there is no certainty. People are terrified — how can you live and not know? It is not odd at all. You only think you know, as a matter of fact. And most of your actions are based on incomplete knowledge and you really don't know what it is all about, or what the purpose of the world is, or know a great deal of other things. It is possible to live and not know."  Feynman

# What is Science?

- In our example, did we left room for doubt?
- The lemma of our disciple should be:

# De omnibus dubitandum

Pronto, pode acordar!

# Você está presente?

# Learning to ask questions

- The baseball encyclopedia
- The Internet Movie Database (IMDb)
- Google Ngrams
- New York Taxi Records

# Learning to ask questions

- ## The baseball encyclopedia

# The baseball encyclopedia

- Obvious questions:

  - How can we best measure an individual player's skill or value?

  - How fairly do trades between teams generally work out?

  - What is the general trajectory of player's performance level as they mature and age?

# The baseball encyclopedia

- Demographic and Social issues

  – Do left-handed people have shorter lifespans than right-handers?

  – How often do people return to live in the same place where were born?

  – Do player salaries generally reflect past, present or future performance?

  – To what extent have heights and weights been increasing in the population at large?

# The baseball encyclopedia

- Identifiers and reference tags (metadata) are as interesting, or more, than the statistical records.

- Statistical proxy: use data you have to substitute for the one you really want.

- The data set of your dreams likely does not exist.

- A good data scientist is a pragmatist!

# Learning to ask questions

- The Internet Movie Database (IMDb)

# The Internet Movie Database (IMDb)

- Some obvious questions:

  - Which actors appeared in most films? Earned the most money? Appeared in the lowest rated films? Had the longest career or the shortest lifespan?

  - What was the highest rated film each year, or the best in each genre? Which movies lost the most money, had the highest-powered casts, or got the least favorable reviews.
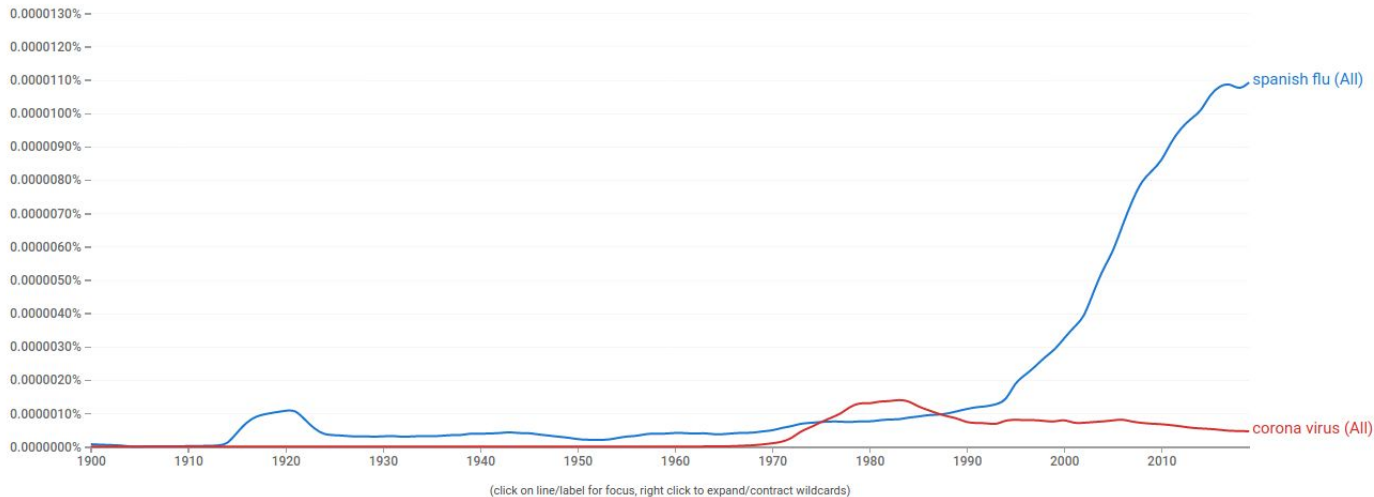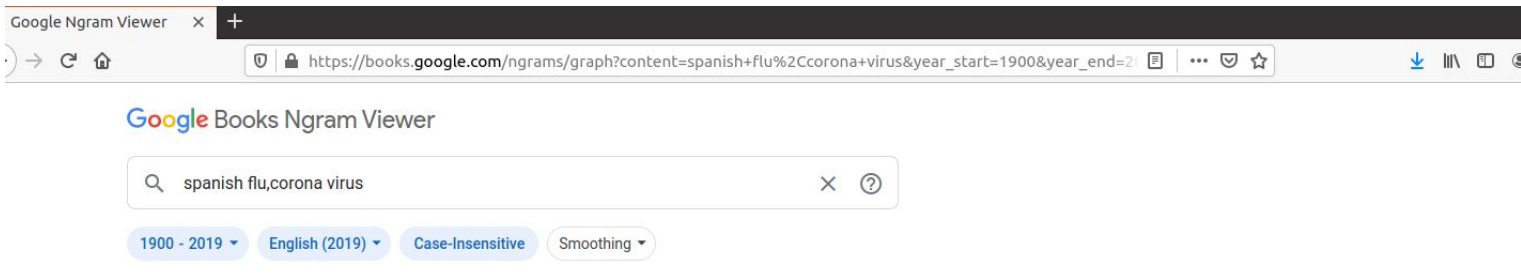
# The Internet Movie Database (IMDb)

- Larger scale questions

  - How well does movie gross correlate with viewer ratings and awards? Do customers instinctively flock to trash, or is virtue on the part of the creative team properly rewarded?

  - How do Hollywood movies compare to Bollywood movie, in terms of ratings, budget, and gross? Are American movies better received than foreign films, and how does this differ between between US and non-US reviewers?

# The Internet Movie Database (IMDb)

- Larger scale questions

  - What is the age distribution of actors and actresses in films? How much younger is the actress playing the wife, on average, than the actor playing the husband? Has this disparity been increasing or decreasing with time?

  - Live fast, die young, and leave a good-looking corpse? Do movie stars live longer or shorter lives than bit players, or compared to the general public?

# Learning to ask questions

Google Ngram Viewer

# Google Ngram Viewer

- Some interesting questions:

  – How often new words emerge and get popular? Do these words tend to stay in common usage, or rapidly fade away? Can we detect when words change meaning over time? Like the transition of gay from happy to

# Learning to ask questions

## New York Taxi Records

# New York Taxi Records

- ## Some interesting questions:

  – How much money do drivers make each night, on average? What is the distribution? Do drivers make more money on sunny days, or rainy days?

  – Where are the best spots in the city for drivers to cruise, in order to pick up profitable fares? How does this vary at different times of the day?

  – How much are drivers tipped, and why? Do faster drivers get tipped

# New York Taxi Records

# Obrigado!