# Scraping Job Posts

Domingos Lopes

Bootcamp #8
New York City Data Science Academy

# Introduction

- Task:

    Scrape job posts from multiple sources, and use the data to build a job recommendation system.

# First Approach

- Get a list of the most important companies.
- Scrape each one individually.

Main advantage:

- Cleaner data

Main disadvantages:

- A lot of spiders to create
- Potentially requires high maintenance

# Second Approach

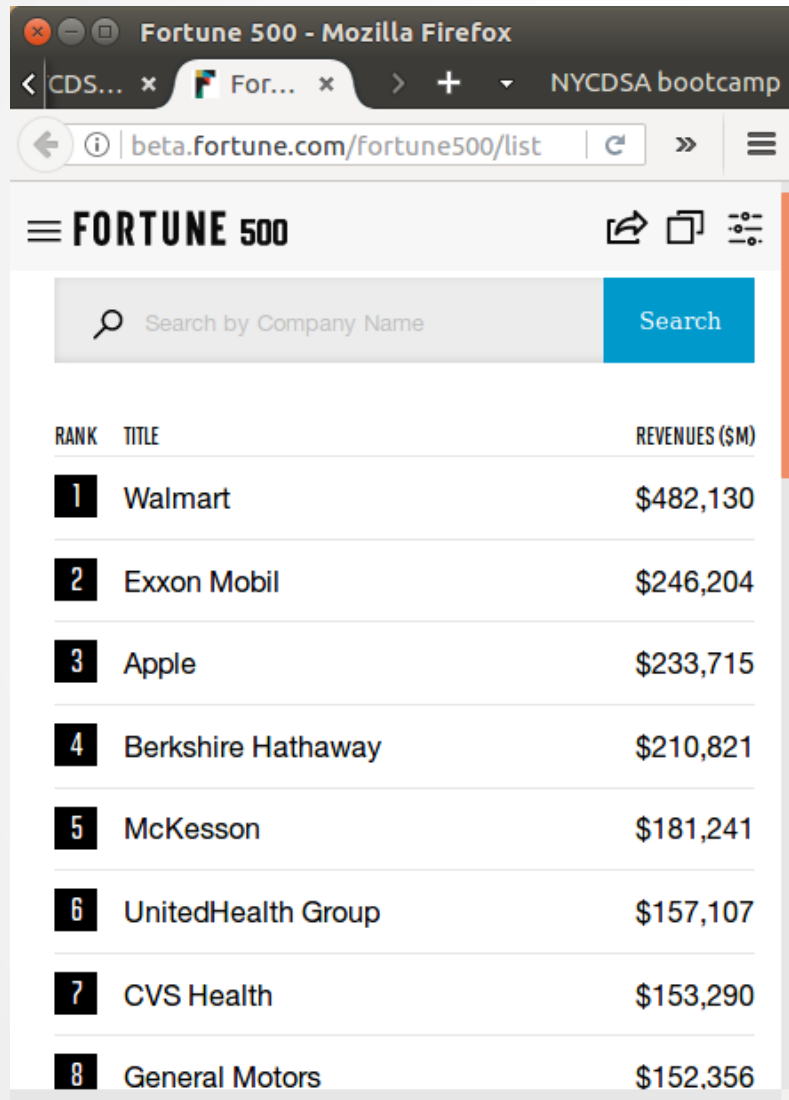- Scrape job posts aggregators

Main advantages:

- – Can obtain great amounts of data using a single spider

- – Many more companies available

- – Often, an API will be available
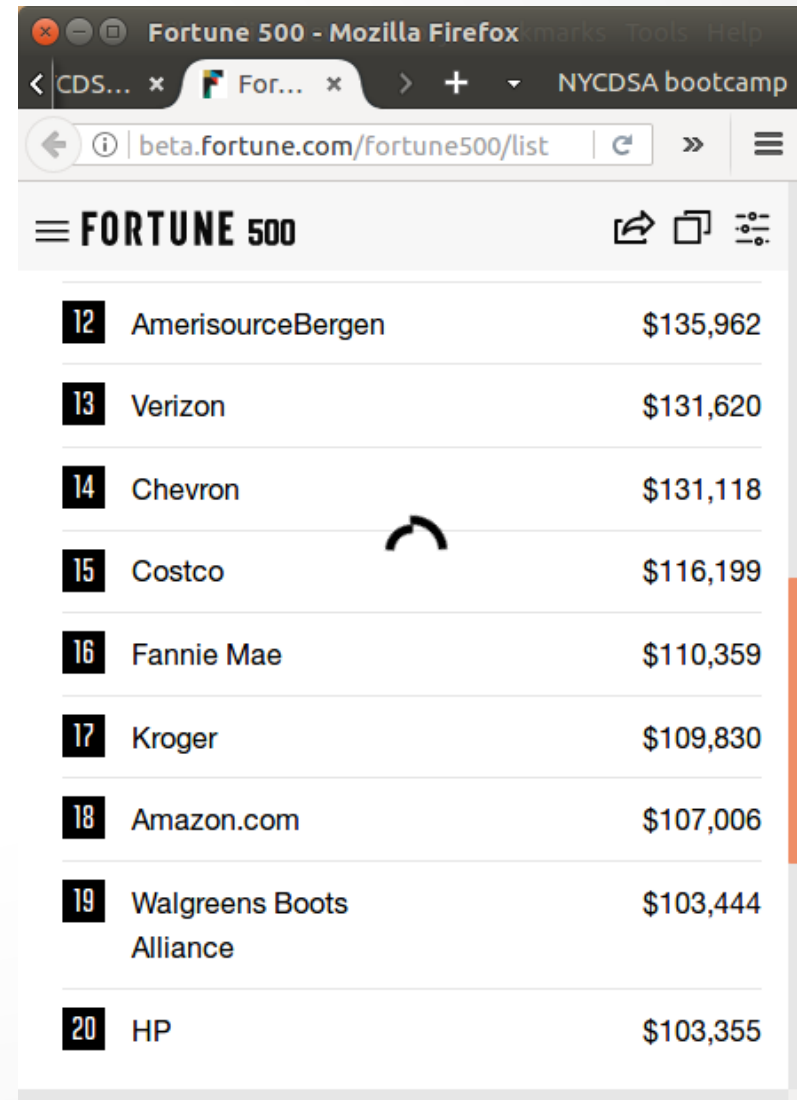
Main disadvantage:

- – Dirty data

# First Approach Challenges

- Get list of top companies: http://beta.fortune.com/fortune500/list

# First Approach Challenges

- List starts with 20 rows.

- Dynamically expands as the user scrolls down.

- Not very scrapy-friendly.

    Solution:

    Listen for AJAX (Asynchronous Javascript And XML) requests and identify the ones in which the important data is transmitted. Then replicate those within scrapy.
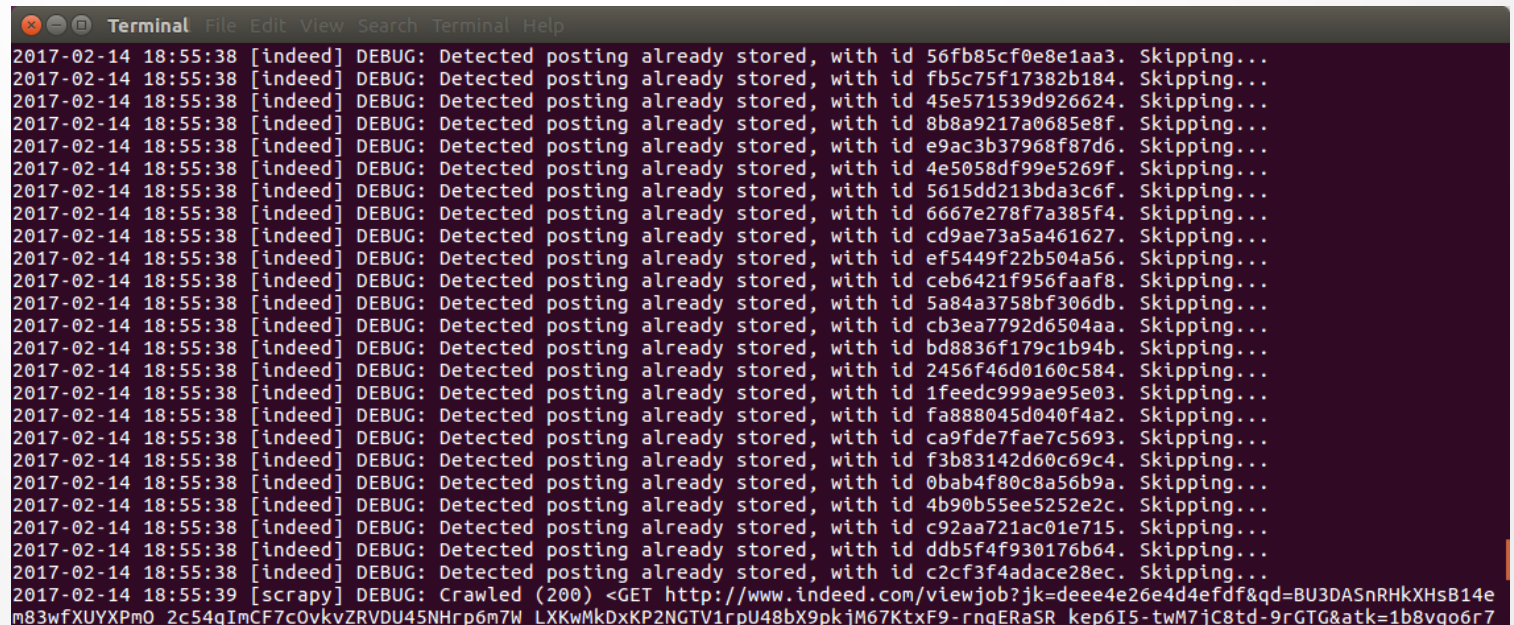
- Apple and Amazon had a similar obstacle.

# First Approach Challenges

- Other challenges:
  - Add new job posts incrementally (not scraping old jobs again),
  - Take advantage of a common structure among posts
  - Location filter
  - Amazon uses the same url for html and json responses
- Results:
  - List of companies: 500 entries with all sorts of key indicators
  - Posts by company: Amazon (727/5862), Apple (84/877), Facebook (44/79)
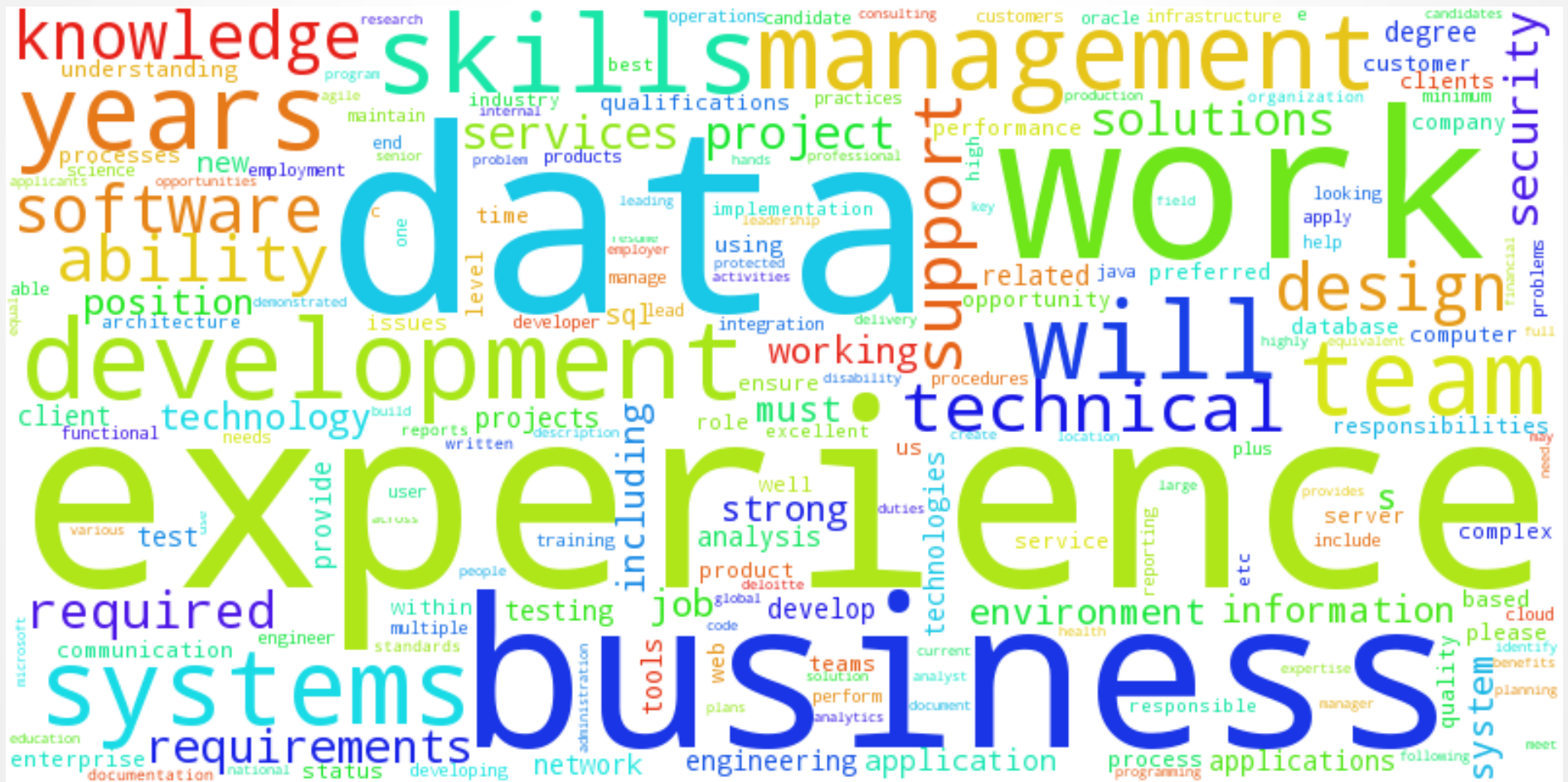
# Second Approach Challenges

- Scraped Indeed and Dice

- Both provide an API to search positions, but one must scrape their page for each individual position.

- A lot of duplicate results within the same search.
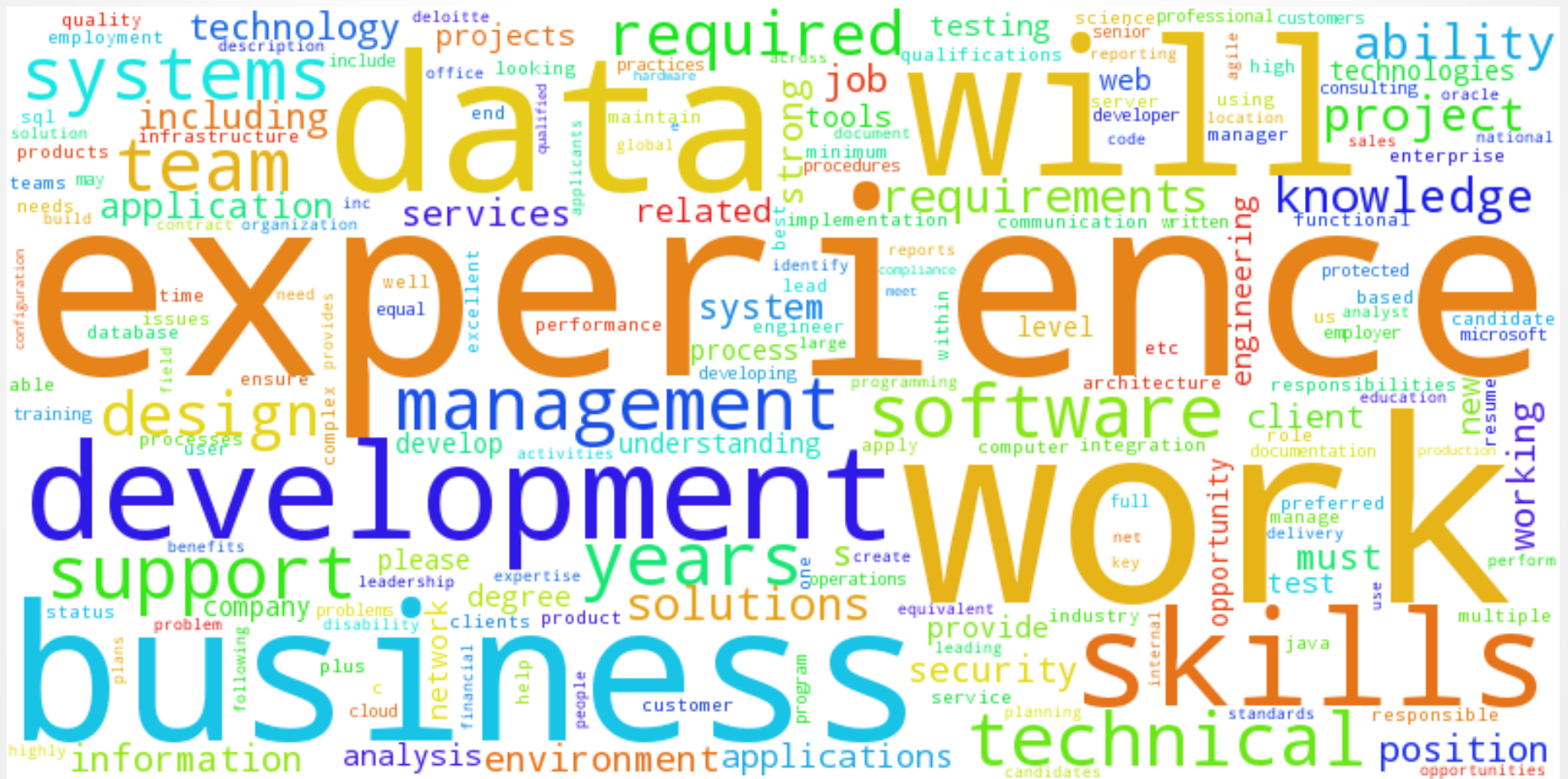
# Indeed and Dice data

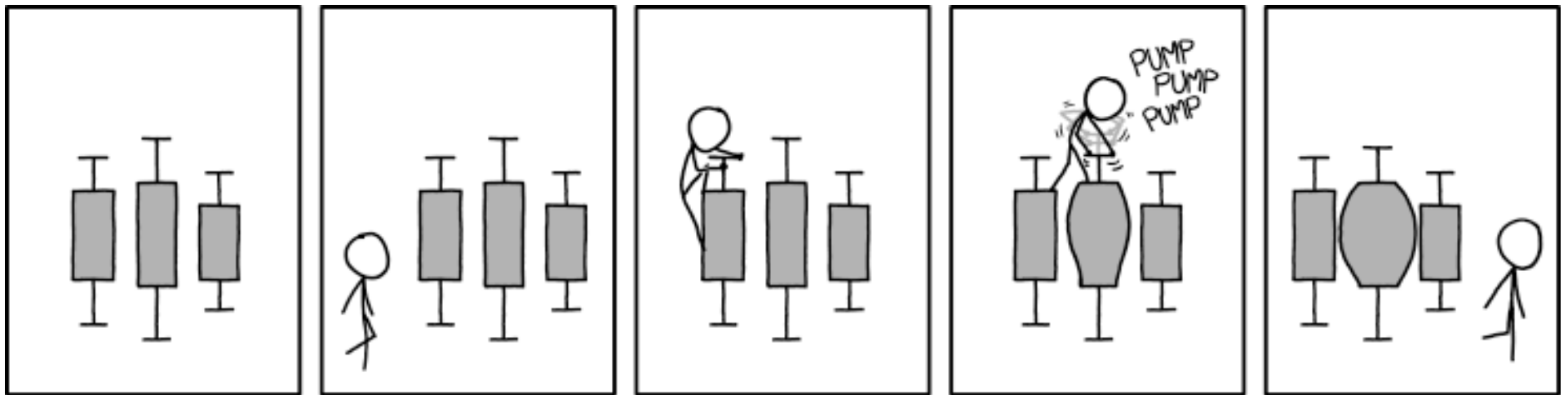From "data scientist" search query

# Indeed and Dice data

Also using "data engineer" and "business analyst"

# Moving forward

- Use even more search queries
- Possibly scrape completely unrelated results – for classification purposes
- Be able to classify sentences as either part of a description or as part of requirements
- Proper classification of job postings (into several dimensions)

# Thank you!



Credit: xkcd