

UNIT - II

1) Exploratory Data Analysis:

Exploratory data analysis (EDA) is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those we believe to be there.

EDA is the first step towards building a model. It's traditionally presented as a bunch of histograms & stem & leaf plots. But EDA is a critical part of the data science process, & also represents a philosophy & way of doing stat's practiced by a strain of statisticians coming from the Bell labs tradition.

(1941) John Tukey, a mathematician at Bell labs, developed exploratory data analysis in contrast to confirmatory data analysis, which concerns itself with modeling & hypotheses as described in the previous section. In EDA, there is no hypothesis & there is no model. The "exploratory" aspect means that our understanding of the problem we're solving, & might solve, is changing as we go.

The basic tools of EDA are plots, graphs & summary stat's. Generally speaking, it's a method of systematically going thru the data, plotting distributions of all var's (using box plots), plotting time series of data, transforming var's, looking at all pairwise relationships b/w var's using scatterplot matrices & generating summary stat's for all of them. At the very least that would mean computing their mean, min, max, the upper & lower quartiles, & identifying outliers.

But as much as EDA is a set of tools, it's also a mindset. And that mindset is abt our relationship with the data. We want to understand the data - gain intuition, understand the shape of it, & try to connect our understanding of the process that generated the data to data itself. EDA happens b/w us & the data & isn't abt providing anything to anyone else.

In stat's, EDA is an approach to analyzing data sets to summarize their main characteristics, often with visual methods. A statistical model can be used or not, but primarily EDA is for seeing what the data can tell us beyond the formal modeling / hypothesis testing task.

EDA is different from Initial data analysis, which focuses more narrowly on checking assumptions required for model fitting & hypothesis testing.

2) Philosophy of Exploratory Data Analysis:

"Long before worrying abt how to convince others, u have to understand what's happening urself"

- Andrew Gelman.

While at Google, Rachel was fortunate to work along side 2 former Bell labs / AT&T statisticians - Darly Pregibon & Diane Lambert, who also work in this vein of applied stats - and learned from them to make EDA a part of her best practices.

Yes, even with very large scale data, they did EDA. In the context of data in an Internet-lengthy Company, EDA is done for some the same reasons it's done with smaller datasets, but there are additional reasons to do it with data that has been generated from logs.

There are 8 imp reasons anyone working with data should do EDA. Namely, to gain intuition ^{immediate understanding} abt the data; to compare ^{assumptions} distributions; for sanity checking ^{knowledge} (making sure the data is on the scale u expect, in the format u thought it should be); to find out where data is missing / if there are outliers; & to summarize the data.

In the context of data generated from logs, EDA also helps with debugging the logging process. For eg, "Patterns" u find in the data could actually be something wrong in the logging process that needs to be fixed. If u never go to the trouble of debugging, u'll continue to think ur patterns are real. The ~~eggs~~ engineers we've worked with are always grateful for help in this area.

In the end, EDA helps u make sure the product is performing as intended.

Although there's lots of visualization involved in EDA, we distinguish b/w EDA & data visualization in that EDA is done toward the beginning of analysis, & data visualization, as it's used in our vernacular, is done towards the end to communicate one's findings. With EDA, the graphics are solely done for u to understand what's going on.

With EDA, u can also use the understanding u get to inform & improve the development of alg's. For eg, Suppose u are trying to develop a ranking alg that ranks content that u are showing to users. To do this u might want to develop a

Before u decide how to ^{measure} Quantify Popularity (which could be, for eg, highest frequency of clicks, or the post with most no. of comments, or comments above some threshold, or some weighted avg of many metrics), u need to understand how the data is behaving, & the best way to do that is looking at it & getting ur hands dirty.

Plotting data & making comparisons can get u extremely far, & it's far better to do than getting dataset & immediately running a regression just bcoz u know how. It's been a service to analysts & data scientists that EDA has not been enforced as a critical part of the process of working with data. Take this opportunity to make it part of your process.

2. The Data Science Process:-

Let's put it all together into what we define as the data science process. The more eg's u see of people doing data science, the more you'll find that they fit into the general framework shown below:

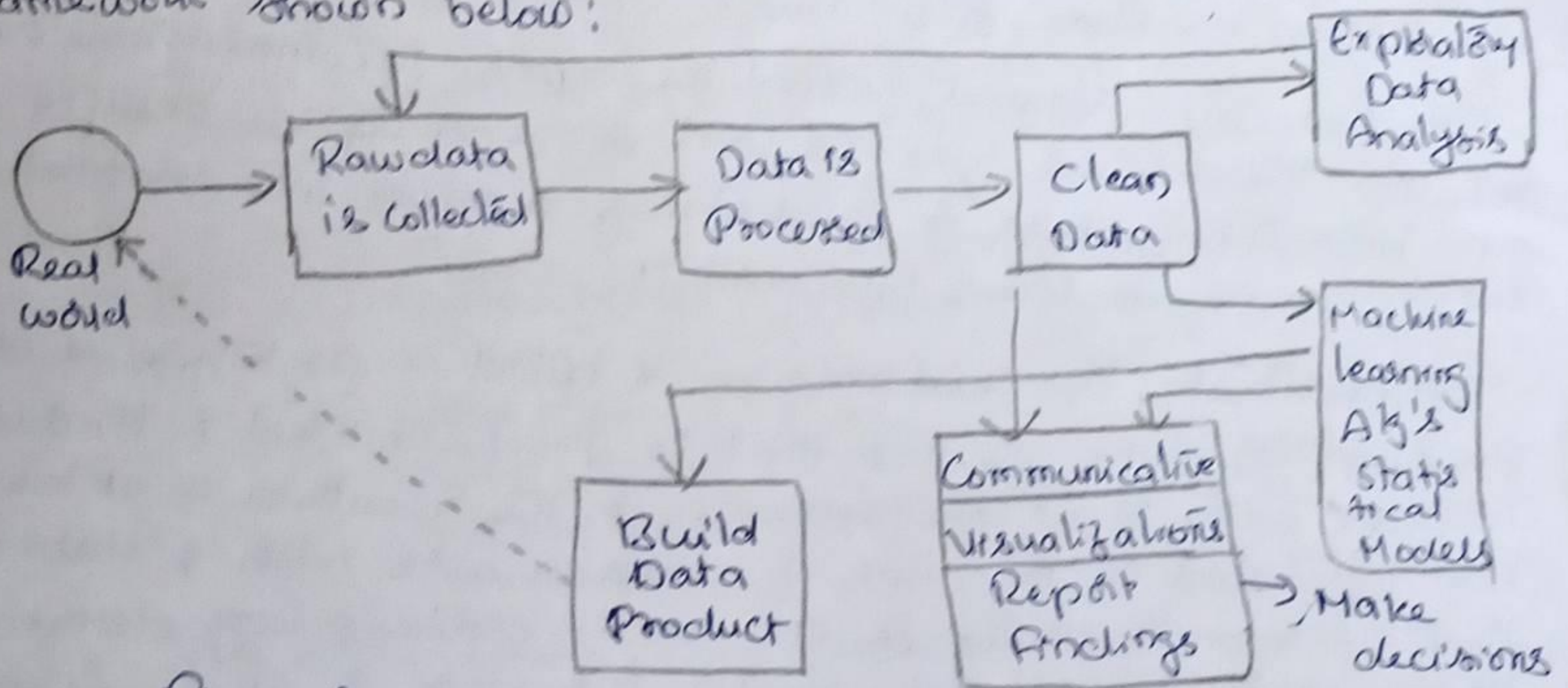


Fig: The data Science process.

First we have the real world. Inside the real world is lots of people busy at various activities. Some people using Google+, others competing in the olympics, there is spammers sending spam, & there is people getting their blood drawn. Say we have data on one of these things.

Specifically, we'll start with raw data - logs, olympics results, Enron emp emails, or sequenced genetic material. We want to process this to make it clean for analysis. So

we build & use pipelines of data munging; joining, scrubbing, wrangling, or whatever u want to call it. To do this we use tools such as python, shell scripts, R & SQL & all of the above.

Eventually we get the data down to a nice format, like something with columns:

name | event | year | gender | event time

Once we have this clean dataset, we should be doing some kind of EDA. In the course of doing EDA, we may realize that it isn't actually clean bcoz of duplicates, missing values, absurd outliers & data that wasn't actually logged or incorrectly logged. If that's the case, we may have to go back to collect more data, or spend more time cleaning the dataset.

Next, we design our model to use some alg's like KNN, linear regression, Naive Bayes, or Bayes something else. The model we choose depends on the type of problem we're trying to solve, of course, which could be a classification problem, prediction problem, or a basic description problem.

We then can interpret, visualize, report, & communicate our results. This could take the form of reporting the results up to our boss or coworkers, or publishing a paper in a journal & going out & giving academic talks abt it.

Alternatively, our goal may be to build & prototype a "data product". A data product that is productionized & that users interact with is at one extreme & the weather is at the other, but regardless of the type of data u work with & "data product" that gets built on top of it - be it public policy determined by a statistical model, health insurance, or election polls that gets built on top of it - be it public policy determined widely reported & perhaps influence viewer opinions - u should consider the extent to which ur model is influencing the very phenomenon that u r trying to observe & understand.

3) A data scientist role in this process: (3)

This model so far seems to suggest this will all magically happen without human intervention. By "human" here, we mean "data scientist". Someone has to make the decisions about what data to collect, & why. That person needs to be formulating questions & hypotheses & making a plan for how the problem will be attacked. And that someone is the data scientist in our beloved data science team.

Let's revise & at least add an overlay to make clear that the data scientist to be involved in this process throughout, meaning they're involved in the actual coding as well as in the higher level process as shown below:

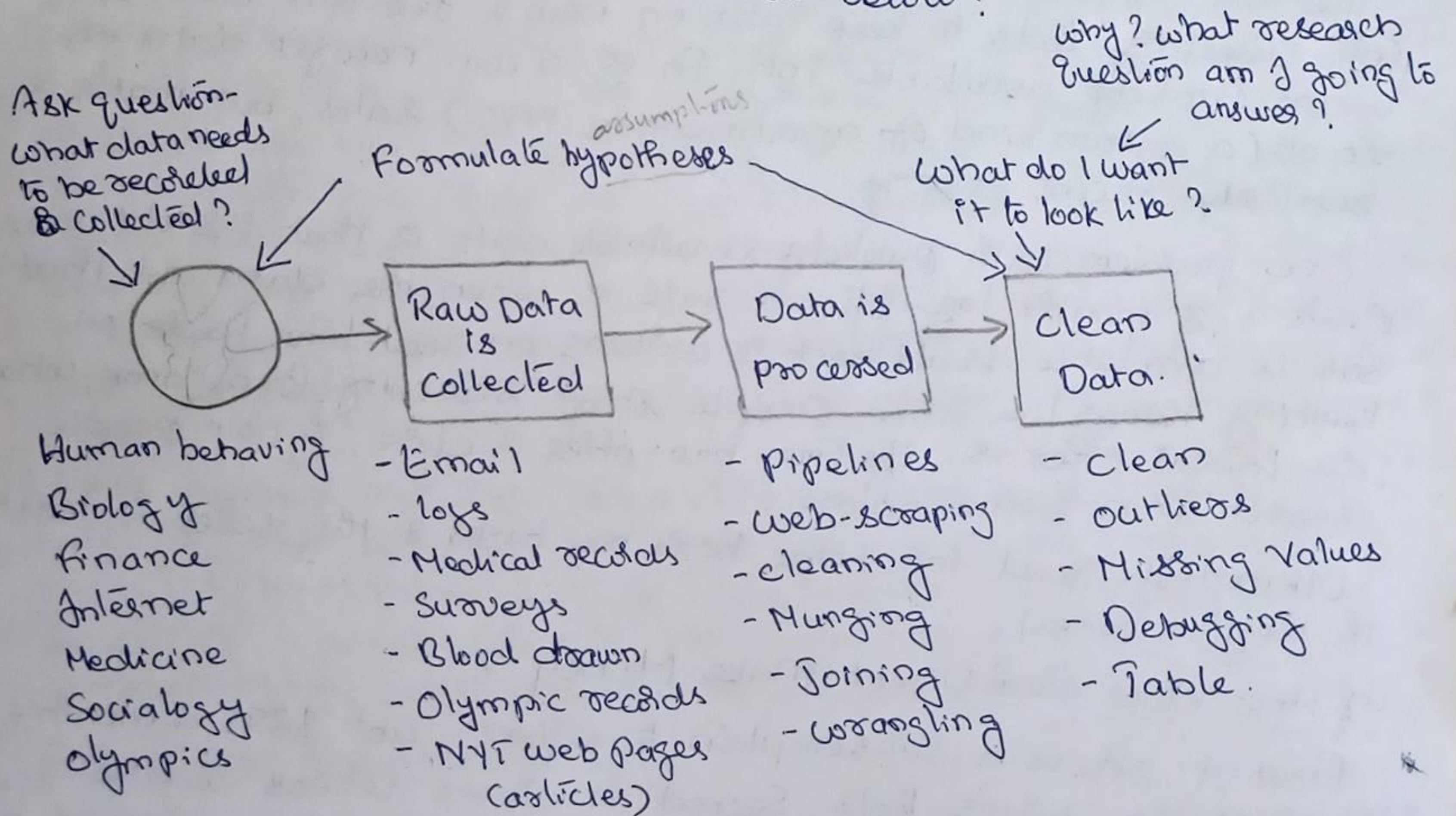


Fig: The data-scientist is involved in every part of this process.

4) Case Study: RealDirect:

Doug Pearson, the CEO of RealDirect, has a background in real estate law, startups & online advertising. His goal with RealDirect is to use all the data he can access about real estate to improve the way people sell & buy houses.

Normally, people sell their homes about once every 7 yrs, & they do so with the help of professional brokers & current data. But there's a problem both with the broker system & the data quality. RealDirect addresses both of them.

First, the brokers. They're typically "free agents" operating on their own - think of them as home sales consultants. This means that they guard their data aggressively, & the really good ones have lots of experience. But in the grand schemes of things, that really means they have only slightly more data than the inexperienced brokers.

RealDirect is addressing this problem by hiring a team of licensed real estate agents who work together & pool their knowledge. To accomplish this, it built an interface for sellers, giving them useful data-driven tips on how to sell their house. It also uses interaction data to give real-time recommendations on what to do next.

The team of brokers also become data experts, learning to use info-collecting tools to keep tabs on new & relevant data & to access publicly available info. For eg, you can now get data on Co-op (a certain kind of apartment in NYC) sales, but that's a relatively recent change.

One problem with publicly available data is that it's old news - there's a 3 month lag b/w a sale & when the data abt that sale is available. RealDirect is working on real-time feeds on ~~things~~ things like when people start searching for a home, what the initial offer is, the time b/w offer & close, & how people search for a home online.

Ultimately, good info helps both the buyer & the seller. At least if they're honest.

(i) How does RealDirect make Money?

First it offers a subscription to sellers - abt \$395 a month - to access the selling tools. Second, it allows sellers to use RealDirect's agents at a reduced commission, typically 2% of the sale instead of the usual 2.5% & 3%. This is where the magic of data pooling comes in: it allows RealDirect to take a smaller commission bcoz it's more optimized, & therefore gets more volume.

The site itself is best thought of as a platform for buyers & sellers to manage their sale & purchase process. There're statuses for each person on site: active, offer rejected, showing, in contract etc. Based on user status, different actions are suggested by the software.

There are some challenges they have to deal with as well, of course. First off, there's a law in New York that says you can't show all the current housing listings unless those listings are behind a registration wall, so RealDirect requires registration. On the one hand, this is an obstacle for buyers, but seriously buyers are likely willing to do it. Moreover, places that don't require registration, like Zillow, aren't true competitors to RealDirect because they're merely showing listings without providing additional service. Doug pointed out that you also need to register to use Pinterest, & it has lots of users in spite of this.

RealDirect comprises licensed brokers in various established realtor associations, but even so it has had its share of hate mail from realtors who don't appreciate its approach to cutting commission costs. In this sense, RealDirect is breaking directly into a guild. On the other hand, if a realtor refuses to show houses because they're being sold on RealDirect, the potential buyers would see those listings elsewhere & complain. So the traditional brokers have little choice but to deal with RealDirect even if they don't like it. In other words, the listings themselves are sufficiently transparent so that the traditional brokers can't get away with keeping their buyers away from these houses.

Doug talked about key issues that a buyer might care about: nearby parks, subway, & schools, as well as the comparison of prices per square foot of apartments sold in the same building or block. This is the kind of data they want to increasingly cover as part of the service of RealDirect.

Exercise: Real Direct Data Strategy:

You have been hired as chief data scientist at RealDirect.com, & report directly to CEO. The company doesn't yet have its data plan in place. It's looking to you to come up with a data strategy.

→ You can use any or all of the datasets.

→ First challenge: load in & clean up data. Next, conduct EDA to find out missing values.

→ Once the data is in good shape, conduct EDA to visualize & make comparisons. If you have time, start looking for meaningful patterns in this dataset.

→ Summarize your findings in a brief report aimed at the CEO.

Data scientist role (continue):

We can think of the data science process as an extension of variation of the scientific method:

- Ask a question → Do background research
- Construct a hypothesis
- Test the hypothesis by doing an experiment
- Analyze the data & draw a conclusion
- Communicate the results.

In both the data science process & the scientific method, not every problem requires one to go through all steps, but almost all problems can be solved with some combination of stages.

Exploratory Data Analysis Tools:

Now a days, ample of tools are available in the market which are free & quite interesting to work with. These tools don't require you to code explicitly but simple drag-drop clicks does the job.

1) Weka : Data mining sw in java:

Weka is a collection of machine learning algorithms for data mining tasks. The algorithms can either be applied directly to a dataset or called from your own Java code. Weka contains tools for data pre-processing, classification, regression, clustering and visualization. It is an open source sw issued under GNU General Public license.

2) R : R is a lang & environment for statistical computing & graphics. R provides a wide variety of statistics & graphical techniques & is highly extensible.

3) Gephi : The Open Graph Viz Platform:

Gephi is the leading visualization & exploration sw for all kinds of graphs & networks. Gephi is open source & free. It runs on windows, macOS & Linux. Gephi is a tool for data analysts & scientists keen to explore & understand graphs.

The goal is to help data analysts to make hypotheses, & discover patterns.

4) OpenRefine:

OpenRefine (formerly Google Refine) is a powerful tool for working with messy data | cleaning it, transforming it from one format into another, & extending it with web services & external data.

5) Orange:

It is an open source machine learning & data visualization tool for novice & expert. It is an interactive data analysis tool with a large toolbox. It performs simple data analysis with clever data visualization.

6) Trifacta:

Trifacta's wrangler tool is challenging the traditional methods of data cleaning & manipulation. Since excel possess limitations on data size, this tool has no such boundaries & u can securely work on big data sets.

7) Rapid Miner:

It is more than just a data cleaning tool. It extends its expertise in building ML models. Not just a GUI, it also extends support to people using Python & R for model building.

8) Qlikview:

It is one of the most popular tool in business intelligence industry around the world. With its state of art visualization capabilities, you'd be amazed by the amt of ctrl u get while working on data. It has an inbuilt recommendation engine to update u from time to time about best visualization methods while working on datasets.

9) DataCracker:

It's a data analysis tool which specializes on survey data. Many companies do survey but they struggle to analyze it statistically. Survey data is never clean. It comprises of

lot of missing & inappropriate value. This tool reduce our
agony & enhances our experience of working on messy data.