

Contents

1	Introduction	1
1.1	Data Analysis	1
1.1.1	Procedure for analyzing data	2
1.1.2	What is Data?	4
1.1.3	Connection to the Scientific Method	5
1.2	Operational environment: Maintenance	6
1.2.1	Preventive Maintenance	7
1.2.2	Condition-based maintenance & monitoring	9
1.2.3	A possible solution	11
2	The hosting company: Zensor	13
2.1	About the company	13
2.1.1	Philosophy	14
2.2	Zensor Approach	14
2.2.1	What are the stages	14
2.2.2	Advanced features and metrics	15
2.3	Workflow	16
2.3.1	Structure of a Deployable Script	16
3	Tools for data analysis	19
3.1	Pandas	19
3.1.1	Series and DataFrame	19
3.1.2	Core Features	20
3.2	InfluxDB	23
3.2.1	TSDB: time series database	23
3.2.2	Influx solution	24
3.3	Grafana	26
3.3.1	Dashboard what is it?	27
3.3.2	Key strengths	28
	Bibliography	33

Chapter 1

Introduction

In this opening chapter, we take a close look at two major research areas data-analysis and maintenance, and try to highlight where they intersect. This chapter is divided into two parts:

1. First, to set the scene, we will give a definition of data analysis and briefly review the main steps involved in the process.
2. Second, we will talk about maintenance, how it has changed over time, and what benefits data analysis techniques can bring to the industry.

Shall we start?

1.1 Data Analysis

Data analysis is the act of analyzing, cleansing, manipulating, and modeling data in order to identify usable information, generate conclusions, and help decision-makings [1]. In today's corporate world, data analysis plays an important part in making decisions more scientific and assisting firms in operating more efficiently.

Data analysis has several dimensions and approaches, including a wide range of techniques known by various names and applied in a variety of business, science, and social science sectors [2]. Let's give some examples. *Data mining* (DM) is a type of data analysis technique that focuses on statistical modeling and knowledge discovery for predictive rather than purely descriptive purposes, whereas *Business Intelligence* (BI) is a type of data analysis that focuses on aggregation and is primarily concerned with business information. Furthermore, Data analysis, in statistical applications, can be separated into descriptive statistics, *Exploratory data analysis* (EDA), and *Confirmatory data analysis* (CDA) [3]. EDA is concerned with finding new features in data, whereas CDA is concerned with validating or refuting current assumptions [4]. Finally, predictive analytics focuses on the application of statistical models for predictive forecasting or classification, while text

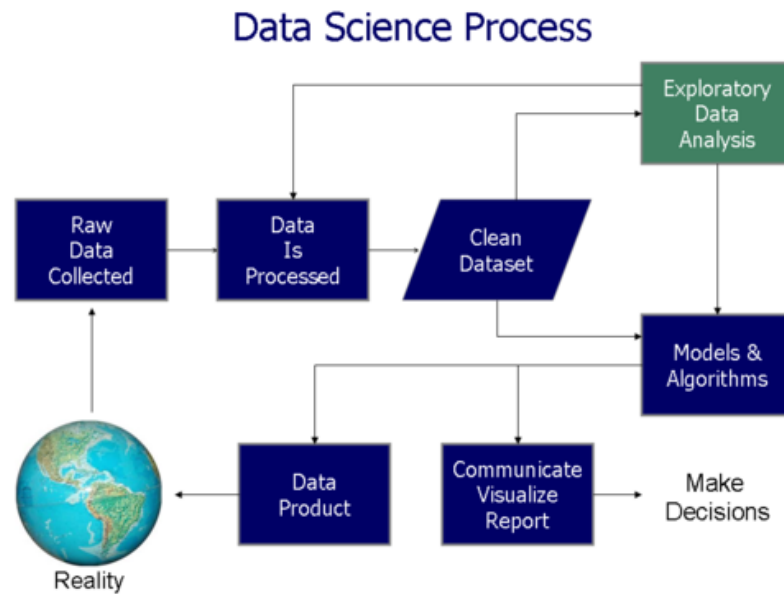


Figure 1.1: Data science process flowchart (Source: [3])

analytics, applies statistical, linguistic, and structural techniques to extract and classify information from textual sources, a species of unstructured data. All of the aforementioned are examples of data analysis [5].

1.1.1 Procedure for analyzing data

The term “*analysis*” refers to the process of breaking down a whole into its constituent parts for closer evaluation. Data analysis is the act of getting raw data and then transforming it into information that users can utilize to make decisions [1]. Data is gathered and processed in order to answer questions, test hypotheses, or refute theories. Statistician Tukey, defined data analysis in 1961 [6], as following:

“Procedures for analyzing data, techniques for interpreting the results of such procedures, ways of planning the gathering of data to make its analysis easier, more precise or more accurate, and all the machinery and results of (mathematical) statistics which apply to analyzing data.”

There are various distinct phases that can be identified as show in 1.1, these are iterative in the sense that input from later phases may lead to further effort in earlier ones [3]. And now we are going to present them one by one, in a slightly more detailed manner, stressing that similar stages can be found

in the *Cross-industry standard process* (CRISP) framework, which is used in Data mining.

Data Collection

Data is collected from a broad variety of sources, like sensors in the environment, including traffic cameras, satellites, recording devices, etc. and it may also be obtained through interviews, downloads from online sources, or reading documentation [3].

Data Processing

Data, when initially obtained, must be processed or organized for analysis. For instance, these may involve placing data into rows and columns in a table format (known as structured data) for further analysis, often through the use of spreadsheet or statistical software [3].

Data Cleaning & Cleansing

After it has been processed and structured, data may be missing, duplicated, or contain errors. Data cleaning will be essential as a result of challenges with the way data is entered and stored. The process of detecting and correcting (or removing) corrupt or inaccurate records from a record set, table, or database is, indeed, known as data cleansing (or *cleaning*), and it entails identifying incomplete, incorrect, inaccurate, or irrelevant parts of the data and then replacing, modifying, or deleting the dirty or coarse data [7]. The actual data cleaning workflow may include removing typographical errors and/or validating and correcting values against a known list of entities.

Exploratory data analysis

Once the datasets are cleaned, they can then be analyzed. EDA is the first step toward building a model, since it is a critical part of the data science process, and also represents a philosophy developed by Tukey in contrast to CDA, which concerns itself with modeling and hypotheses [6]. Indeed in EDA, there is no hypothesis and there is no model. The “exploratory” aspect means that your understanding of the problem you are solving, or might solve, is changing as you go.

The process of data exploration may result in additional data cleaning or additional requests for data; thus, the initialization of the iterative phases mentioned in the lead paragraph of this section. Data visualization is also a technique used, in which the analyst is able to examine the data in a graphical format in order to obtain additional insights, regarding the messages within the data [3].

Modelling and algorithms

Mathematical formulas or models (known as algorithms) can be applied to data to identify relationships between variables, such as correlation or causation [2]. In general, models can be created to evaluate a specific variable based on other variables in the dataset, with some residual error depending on the accuracy of the implemented model (e.g., $Data = Model + Error$) [8].

Data product

A data product is a computer application that takes data inputs and generates outputs, feeding them back into the environment. As such it may be based on a model or algorithm. For instance, an application that analyzes data about customer purchase history, and uses the results to recommend other purchases the customer might enjoy [3].

Data visualization & Communication

Once data is analyzed, it may be reported in many formats to the users of the analysis to support their requirements [4]. The users may have feedback, which results in additional analysis. As such, much of the analytical cycle is iterative, as stated before. A company dashboard, for instance, that visualizes some company's *Key performance indicators* (KPIs), is both a data product and a decision support system, with a nice user interface that allows to access multiple visualization 3.3.

1.1.2 What is Data?

We discussed a lot about data in the previous subsection; certainly it is important that one have the ability to analyse and investigate data, but understanding data is an important prerequisite for being able to use it properly, and perhaps the single most important element. Fortunately the NOIR system, commonly used, can help us. It defines the type of data as nominal, ordinal, interval or ratio as show in table 1.1.

So we can confirm that valuable data is no longer only a collection of numbers and classified variables and a strong data scientist needs to be versatile and comfortable dealing with a variety of types of data, including:

- Traditional: numerical, categorical, or binary
- Text: emails, tweets, New York Times articles
- Records: user-level data, timestamped event data, log files
- Complex: Geo-based location data (GIS), Network & Images
- Sensor data, my use-case (see Chapter [ref])

Variable	Type(s)	Description	Examples
Categorical	Nominal	Named categories with no implied order	Blood groups, breed, gender, neuter status
	Ordinal	Ordered categories where the differences between categories are not necessarily equal	Scoring systems, cancer staging, onset of disease (peracute, acute, chronic)
Continuous	Interval	Equal distances between values but the zero point is arbitrary	IQ, ordinal data with equal-appearing categories
	Ratio	Above as for interval and a meaningful zero; data usually obtained by measurement	

Table 1.1: NOIR system of classification of types of data (Source: [4])

1.1.3 Connection to the Scientific Method

In both the data science process and the scientific method, not every problem requires one to go through all the steps, but almost all problems can be solved with *some* combination of previously mentioned stages [3]. In fact, We can think of the data science process as an extension of or variation of the scientific method:

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

As an example, if your end goal is a **data visualization** (which itself could be thought of as a data product), it's possible you might not do any machine learning or statistical modeling, but you would want to get all the way to a clean dataset, do some Exploratory data analysis, and then create the visualization. This has happened to me many times during my internship, as we will see in [Chapter 4].

Data-driven context There are numerous areas in which data analysis shines, but in this text we will focus on one specific area: maintenance. For complex systems such as airplanes, railways, power plants, is a big issue (and challenge) as it ensures the system reliability and safety during their life cycles.

But what does the term maintenance mean and which and how many types exist?

1.2 Operational environment: Maintenance

In industrial, business, and domestic settings, maintenance entails functioning checks, maintaining, repairing, or replacing necessary devices, equipment, machinery, building structures, and supporting utilities. This has evolved over time to encompass a variety of terms that indicate various cost-effective techniques for keeping equipment operating; these actions might occur before or after a failure [9].

Types The marine and air transportation, offshore structures, industrial plant and facility management industries depend on maintenance, repair and overhaul MRO including scheduled or preventive paint maintenance programmes to maintain and restore coatings applied to steel in environments subject to attack from erosion, corrosion and environmental pollution [10].

Architectural conservation employs MRO to preserve, rehabilitate, restore, or reconstruct historical structures with stone, brick, glass, metal, and wood which match the original constituent materials where possible, or with suitable polymer technologies when not. The basic types of maintenance falling under MRO include:

- **Corrective maintenance**, where equipment is repaired or replaced after wear, malfunction or break down.
- **Preventive maintenance**, where equipment is checked and serviced in a planned manner (in a scheduled points in time or continuously).
- **Reinforcement**, where equipment is reinforced and hardened to prevent failure.

These are huge topics and we are going to focus mostly on corrective and preventive maintenance. **1.DG:Expand this connection?**

Corrective maintenance is a type of reactive maintenance that is performed on equipment after it has broken down or malfunctioned, sometimes referred to as “fighting fires“. Not only can worn equipment damage other parts and cause multiple damages, but it can also result in significant repair and replacement costs as well as lost revenue due to downtime during



Figure 1.2: *USS Ronald Reagan (CVN 76) Undergoes Preventive Maintenance* [11]

overhaul. Traditional procedures like welding and metal flame spraying, as well as designed solutions with thermoset polymeric materials, are used to rebuild and resurface equipment and infrastructure damaged by erosion and corrosion as part of corrective or preventive maintenance programs.

1.2.1 Preventive Maintenance

Preventive maintenance (PM) is, according to Decourcy Hinds [12]:

... a routine for periodically inspecting with the goal of noticing small problems and fixing them before major ones develop.

Ideally, nothing breaks down!

The main objectives of preventive maintenance are as follows:

1. **Enhance** capital equipment productive life.
2. **Reduce** critical equipment breakdown.
3. **Minimize** production loss due to equipment failures.

Many people, me included, confuse the phrases *preventive*, *predictive*, and *prescriptive* maintenance, and while they are distinct, the latter two might be considered kinds of preventive maintenance. Preventative maintenance in all forms aids manufacturers in transitioning from a repair-and-replace to a preventive maintenance approach [13]. Let's look at the different sorts of preventative maintenance.

- **Planned preventive maintenance (PPM)**, more commonly referred to as planned maintenance or scheduled maintenance, is any variety of scheduled maintenance to an object or item of equipment. Specifically, planned maintenance is a scheduled service visit carried out by a competent and suitable agent, to ensure that an item of equipment is operating correctly and to therefore avoid any unscheduled breakdown and downtime. It can be further split in two subcategories:

- **Calendar-based maintenance** is performed on equipment according to a calendar timetable. In other words, a maintenance activity is triggered by the passage of time. Calendar-based maintenance includes things like: cleaning your air conditioner and replacing the air filter in your heating, ventilation, and air conditioning equipment every three months.

When a scheduled task is due, *Computerized Maintenance Management Systems* (CMMS) are frequently used to keep schedules straight and issue recurring work orders.

- **Usage-based maintenance** uses triggers based on how much each piece of equipment is used. Maintenance managers can create a preventative maintenance schedule based on predefined parameters by tracking usage using equipment monitors and the operating hours of each piece of machinery.

For example, when machine X reaches a certain number of hours of operation, this creates a trigger to schedule a booking for a service technician to perform on-demand maintenance.

- **Predictive maintenance (PdM)** are intended to assist in determining the state of in-service equipment so that maintenance can be scheduled. Because actions are performed only when warranted, this strategy promises cost reductions over routine or time-based preventive maintenance. As a result, it is viewed as condition-based maintenance which is carried out in accordance with estimations of an item's degradation status. Predictive maintenance's key benefit is that it allows for easy scheduling of corrective maintenance and prevents unexpected equipment breakdowns [14].

It is particularly useful when coupled with CMMS software; Logging work requirements data allows managers to review data and notice failure patterns over time. This information allows to predict when outages will occur based on historical data and plan maintenance tasks to avoid them [13].

- **Condition-based maintenance (CBM)** to put it succinctly, is *maintenance when it is needed*. Although being chronologically much older, it is considered one sector or practice within the broader and younger

predictive maintenance field, where new *Artificial Intelligence* (AI) technology and *Internet Of Things* (IoT) connectivity abilities are put to use, to help schedule preventive maintenance tasks. CBM is performed after one or more indicators show that equipment is going to fail or that equipment performance is deteriorating; this concept is applicable to both mission-critical systems that incorporate active redundancy and fault reporting, and non-mission-critical systems that have a more limited budget; let's explore this in a little more detail.

1.2.2 Condition-based maintenance & monitoring

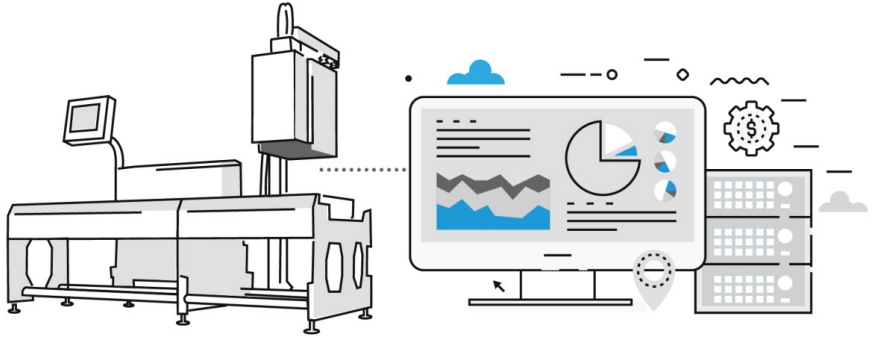


Figure 1.3: Condition-based maintenance schema (Source: [15])

CBM was developed to ensure that the proper equipment was maintained at the right time. It is focused on prioritizing and optimizing maintenance resources using real-time data, through *Condition monitoring* (CM), which is the process of observing the system's state; how? Through monitoring parameter(s) of condition in machinery (vibration, voltage, temperature etc.), in order to identify a significant change which is indicative of a developing fault [16]. The equipment's health will be determined by this system, and maintenance will be performed only when it is genuinely required.

Recent technological advancements have enabled widespread equipment instrumentation, and, when combined with improved tools for analyzing condition data, maintenance personnel is more than ever able to determine when it is appropriate to undertake maintenance on a piece of equipment. CBM, in essence, should allow maintenance professionals to focus on only the necessary tasks, reducing spare parts costs, system downtime, and maintenance time. Furthermore using *Machine Learning* (ML) based scheduled

maintenance could not only predicts a possible failure, but also attempts to make result-oriented maintenance recommendations based on that machine’s analysis.

Challenges Despite its usefulness, there are several challenges to the use of Condition-based maintenance.

1. First and most important of all, the initial cost can be high since it requires improved instrumentation of the equipment. Often the cost of sufficient instruments can be quite large, especially on equipment that is already installed, even though wireless systems have reduced the initial cost. Therefore, it is important for the installer to decide the importance of the investment before adding CBM to the equipment.

For instance, a result of this cost is that the first generation of CBM in the oil and gas industry has only focused on vibration in heavy rotating equipment [17].

2. Secondly, introducing this process will invoke a major change in how maintenance is performed, and potentially to the whole maintenance organization in a company. As we well know organizational changes are in general difficult.
3. Lastly, the technical side of it is **not** simple! Even if some types of equipment can easily be observed by measuring simple values such as vibration (displacement, velocity or acceleration), temperature or pressure, it is not trivial to turn this measured data into actionable knowledge about the health of the equipment.

Advantages and disadvantages Condition-based maintenance has some advantages and disadvantages over Planned preventive maintenance:

Pros	Cons
Improved system reliability	High installation costs, for minor equipment items often more than the value of the equipment
Decreased maintenance costs	Unpredictable maintenance periods cause costs to be divided unequally
Decreased number of maintenance operations causes a reduction of human error influences	Increased number of parts (the CBM installation itself) that need maintenance and checking

1.2.3 A possible solution

Today, most often, technical data sheets coupled with the knowledge of a number of unique experienced individuals are used to determine when the asset requires maintenance. Product quality only becomes an issue when customers start complaining and repairs are done when it's already far too late. All of these puts tremendous strain on the people responsible for the asset, while it could be avoided. Unexpected shutdowns are costly and very demanding for the workforce involved; a possible solution would be making assets smart in order to increase the availability. The best way to validate the actual health is by having a continuous look at a broad data set and adding a specific multi-aspect monitoring setup consisting of different sensor types that follow the behavior of the asset's general state-of-health [18].

Infrastructure Analytics (Infralytics®) Zensor platform is designed to make PdM as easy as possible. It graphically shows you the current state of your asset, it alarms you when a certain component is about to break down and the collection and aggregation of all kinds of sensor and operational data makes it possible to go beyond the symptom level of a mechanical failure and detect the underlying reason for the breakdown. Relying on Infralytics® will not only lower your maintenance costs significantly, it will also have a serious impact on your plant's efficiency by increasing uptime and lowering the occurrence of unexpected mechanical breakdowns [19].

2. not all data are the same, treated differently depending on the context. Increase the awareness of the problem

Chapter 2

The hosting company: Zensor

2.1 About the company

Zensor [18] provides full, integrated and intelligent monitoring solutions for the industrial production, renewable energy and infrastructure sectors. This allows to Zensor's customers to convert the potential contained in the world of IoT and Industry 4.0 into their reality. The company enable digitalization, but with an interface oriented towards the real human: an expert solution without the need for internal experts. It provides not only monitoring devices or data analysis, but offers a full, standardized and standalone end-to-end product that leverages the value contained in a subset of the following aspects:

- operational efficiency
- predictive maintenance
- energy efficiency
- ageing and degradation
- safety

A standard offering consists in several aspects (if required) ordered logically below; as such Zensor takes end-to-end responsibility in monitoring the health and efficiency of structures and processes.

1. Hardware {sensors and acquisition units}
2. Installation and Commissioning {engineering and CAD}
3. Data Management {data transfer, storage, coupling to existing data sources (SCADA, weather, operational...), data cleaning and treatment}
4. Analysis and Reporting {predictions, trend and event detection, real-time reporting through online dashboards.}

Table 2.1: Assets, Industries, and Infrastructure for which Zensor has specific products

Asset	Industries	Infrastructure
Rolling cranes	Metal Production	Offshore wind
Grinders and crushers	Mining and Materials	Rail
Flattener rollers	Food Production	Civil Infrastructure
Rolling mills	Glass Production	Energy
Conveyor belts	Discrete Manufacturing	
Tunnels	Textiles	
Chain transporters		
Sieves		
Bridges		

2.1.1 Philosophy

Vision Technological advance can only bring real value to society when the user-facing component is driven by Simplicity and Clarity for the end-user. Zensor sees this as the fastest and most certain route to a world where man-made structures affect the sustainability of our planet in the least possible way:

1. they are intrinsically safe;
2. their useful life is optimized to the maximum;
3. their impact on environment and society is quantified and communicated.

Mission Translate technological innovations in monitoring and analysis into easy-to-understand, tangible and relevant information that we share in the way tailored for either production managers, management, or maintenance professionals ... As such, we are the knowledgeable and easy-to-reach companion for owners and operators in making their assets increasingly safe, efficient and sustainable. If up to us, till eternity.

2.2 Zensor Approach

2.2.1 What are the stages

Hardware \rightsquigarrow *Installation* \rightsquigarrow *DataManagement* \rightsquigarrow *Analysis*

Hardware and/or other Existing Data sources Sometimes not enough data is available from the beginning. Deriving valuable insights from a monitoring system states from the data: identifying and locating the relevant data in existing databases/data warehouse or putting the right sensors, with appropriate settings, on the right positions and measurement conditions; afterwards reading them out in the optimal way. All of this are defined clearly in every asset-specific package.

Installation and Commissioning Push the button of Industry 4.0. Initially links to the existing data sources are established and data gets ingested. Where required a set of acquisition units and sensors is installed on the machine or structure. After a final verification on the spot (SAT) the monitoring system is launched: the assets enter the IoT.

Data Management Data is continuously streaming in from individual setups as well as historian sources. Structuring, verifying and cleaning the data sets is an essential prerequisite to allow for a profound analysis afterwards: on your way to an automated, continuous and smart follow-up.

Analysis and Clear Reporting Advanced insights are unlocked using algorithms based on physics as well as big data approaches. Clear dashboards, warnings and periodic reports inform the owner or line manager about the present state and upcoming issues. Surprises are avoided, standstills reduced.

2.2.2 Advanced features and metrics

Here is a non-exhaustive list of the main monitoring metrics available

Availability Have a continuous idea of availability, automatically as the platform combines different input streams and contextual information.

Performance Based on the data collected and machine-learning based methods for determining the operational condition the performance is calculated.

Warnings Whenever values start to deviate, or data streams stop, warnings are sent. This avoids 'black holes' in the insights of the production line or assets.

Quality Coupling to existing databases or using human input fields the product quality is linked to operational process parameters.

MTTF The *MeanTimeTillFailure* is tracked continuously, for each asset covered the overall 'disturbance free' operation is displayed.

MTBF As events and operating conditions are automatically detected the *MeanTimeBetweenFailures* is determined continuously, giving a good insight on where optimization is possible.

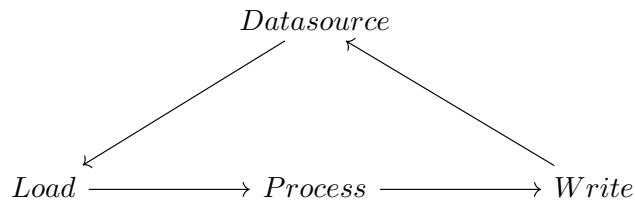
OEE Using all parameters cited above the *OverallEquipmentEffectiveness* is determined, a major parameter for optimizing asset management strategies and future investments, with a huge cost savings potential.

Factory information systems Such systems are crucial for obtaining operational excellence. When well managed they maximize efficiency and effectiveness. Automated data collection and advanced analysis makes this possible.

2.3 Workflow

Preamble: I expand on the stages I have seen, not on those I have not seen. I am part of this project; I mention what I have done ...

2.3.1 Structure of a Deployable Script



Most scripts that run on the Zensor platform have a very common structure. For a given time window, they:

1. Load some data (either raw or from InfluxDB).
2. Process it in some (clever!) way.
3. Write the results out to InfluxDB, to be shown in a dashboard.

What time window they operate on will depend on what the task is, but also on whether the script is being invoked automatically by cron, or manually. If a script is being invoked manually, this is usually to run it over historical data e.g. rerunning a script for the month of February 2020. We typically call this **backfilling**. Typically, if the script is running in cron, it's loading "recent" data, e.g. from the past hour or past day, ending at the time the

script started. Scripts on the Zensor platform need to support running in both modes, so there are a few guidelines to keep in mind when writing a script.

Chapter 3

Tools for data analysis

The idea of this chapter is to show the main tools used during the work as seen at [2.3.1], and to highlight the important pieces.

3.1 Pandas

Pandas is a data manipulation and analysis software library created for the Python programming language. It provides data structures and functions for manipulating numerical tables and time series, and it is free software distributed under the BSD three-clause licence [20]. The name derives from the word “panel data”, which is an econometrics term for data sets that comprise observations for the same individuals over several time periods and, at the same time, is a parody of the term “Python data analysis” [21].

3.1.1 Series and DataFrame

Pandas is primarily used to analyse data. It supports data import from a variety of file formats, including comma-separated values (CSV), JSON, SQL database tables or queries, and Microsoft Excel [22]. Further more Pandas supports a variety of data manipulation operations such as merging, reshaping, and selecting, as well as data cleaning and handling. To accomplish this, Pandas define and makes use of two important software *Classes*, Series and DataFrame, which we will now briefly introduce.

Series is a one-dimensional labeled array capable of holding any data type (integers, strings, floating point numbers, Python objects, etc.). The axis labels are collectively referred to as the *index*.

DataFrame is a 2-dimensional labeled data structure with columns of potentially different types; in some ways it is like a spreadsheet or SQL table (see figure: 3.1) and each individual column is a *Series*. It is generally the

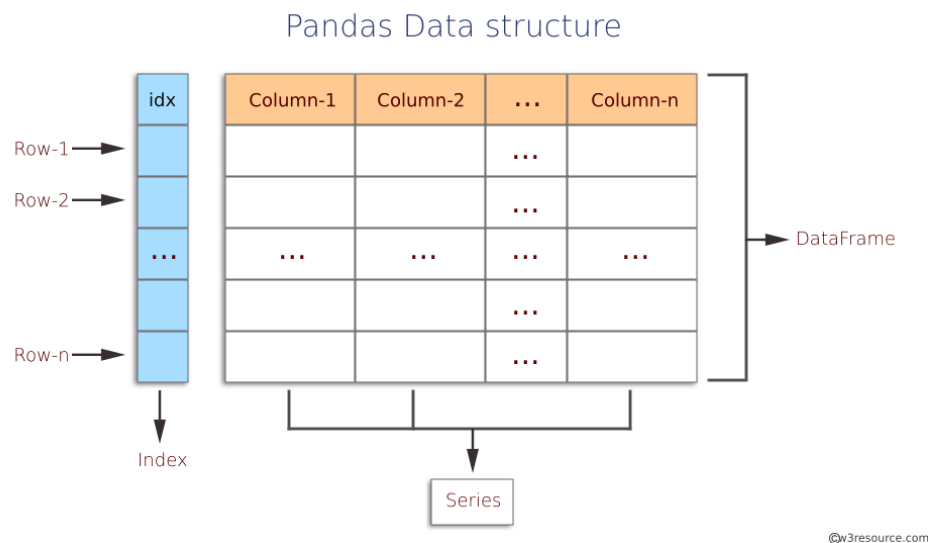


Figure 3.1: Pandas Data structure

most commonly used pandas object, since it accepts many different kinds of input at makes him very flexible [23]. Along with the data, one can optionally pass index (row labels) and columns (column labels) arguments. By doing so you are guaranteeing the index and/or columns of the resulting DataFrame. If axis labels are not passed, they will be constructed from the input data based on common sense rules. And this is just one way of building a dataframe, a foretaste of the flexibility of this library.

3.1.2 Core Features

The idea of this subsection is to give an outline of how many possible use-case this library can cover, and, at the same time, explore a couple of them that proved to be crucial during my internship experience; let's start with the idea of *grouping*.

Groupby The name **GroupBy** should be quite familiar to those who have used a SQL-based tool or worked with relation database. This "engine" allows split-apply-combine operations on heterogeneous data sets. by "group by" we are referring to a process involving one or more of the following steps:

1. Splitting the data into groups based on some criteria.
2. Applying a function to each group independently.
3. Combining the results into a data structure

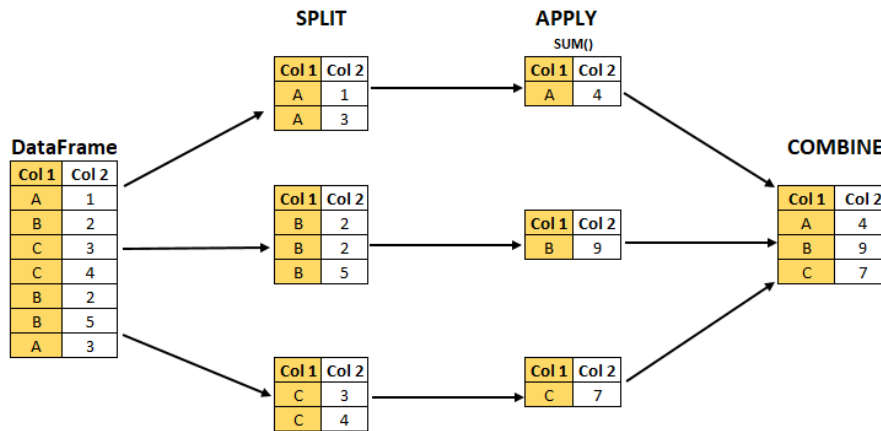


Figure 3.2: Shows the Split-Apply-Combine using an aggregation function (source: Analytics Vidhya [24])

Out of these three, the *split* step is the most straightforward. In fact, in many situations, we would like to split the data set into groups and do something with those groups [23]. In the *apply and combine* step, we might wish to do one of the following:

- Aggregation: compute a summary statistic (or statistics) for each group, some examples:
 - Compute group sums or means.
 - Compute group sizes / counts.
- Transformation: perform some group-specific computations and return a like-indexed object, for instance:
 - Standardize data (zscore) within a group.
 - Filling NAs (value that are not valid) within groups with a value derived from each group.
- Filtration: discard some groups, according to a group-wise computation that evaluates True or False, like:
 - Discard data that belongs to groups with only a few members.
 - Filter out data based on the group sum or mean.
- Some combination of the above: **GroupBy** will examine the results of the *apply* step and try to return a sensibly *combined* result if it doesn't fit into either of the above two categories.

Since the set of object instance methods on pandas data structures are generally rich and expressive, we often simply want to invoke, say, a `DataFrame` function on each group. With this engine you can try multiple different approaches, testing what suits more your necessities, even though often is hard to define this three separates step for badly shaped data.

Time series resampling Pandas contains extensive capabilities and features for working with time series data for all domains. Using the **NumPy** datetimes dtypes, it has consolidated a large number of features from other Python libraries (like **scikits.timeseries**) as well as created a tremendous amount of new functionality for manipulating time series data. As an example, Pandas supports:

- Parsing time series information from various sources and formats
- Manipulating and converting date times with timezone information
- Moving window statistics and linear regressions

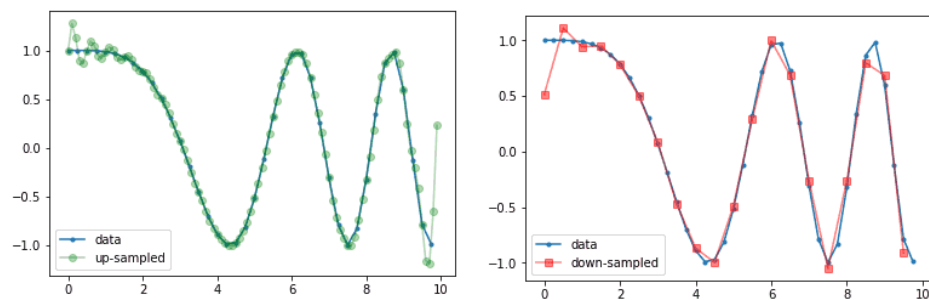


Figure 3.3: Examples of signal waveform data processed by resample

Let's now focus a bit about one key aspect of handling/manipulating sensor data, crucial in our context, the necessity of resampling. Fortunately, Pandas has, once again, a simple-to-use, powerful, and efficient functionality for performing resampling operations during frequency conversion (e.g., converting secondly data into 5-minutely data). This is also extremely common in, but not limited to, financial applications.

To keep things simple we could say that resample is a time-based **Groupby** followed by a reduction method on each of its groups; as a positive side, this method can be used directly from *DataFrameGroupBy* objects that we discussed in paragraph [3.1.2]. The resample function is very flexible and allows you to specify many different parameters to control the frequency conversion and resampling operation, both upsampling and downsampling.

Others functionality Furthermore, other time series features are available, and not only that notably:

- Date range generation and frequency conversions
- Data alignment, shifting and lagging
- Integrated missing data handling
- Data set reshaping, pivoting, merging and combining
- Label-based slicing, sophisticated indexing, and big data set sub-setting
- Insertion and deletion of columns in a data structure.

Conclusion Pandas provides a solid foundation upon which a very powerful data analysis ecosystem can be established, especially since the library is performance-optimized, with important code paths implemented in Cython or C.

3.2 InfluxDB

Time Series In mathematics, a time series is a series of data points which are indexed (or listed or graphed) in time order (see 3.4). More generally, a time series is a sequence taken at evenly spaced intervals over a period of time. As a result, it's a succession of discrete-time data. Ocean tidal heights, sunspot counts, and the Dow Jones Industrial Average's daily closing value are all examples of time series.

3.2.1 TSDB: time series database

It follows that a time series database TSDB is a software system that is designed to store and serve this peculiar type of data, time series, using time(s) and value(s) pairs(s). Timescale, popular TSDB, CEO *Ajay Kulkarni* [26] put it:

Time-series datasets track changes to the overall system as INSERTs, not UPDATEs.

This practice of recording each and every change to the system as a new, different row is what makes time-series data so powerful. It allows us to measure change: analyse how something changed in the past, monitor how something is changing in the present, predict how it may change in the future.

So here's how I like to define time-series data: data that collectively represents how a system/process/behaviour changes over time.

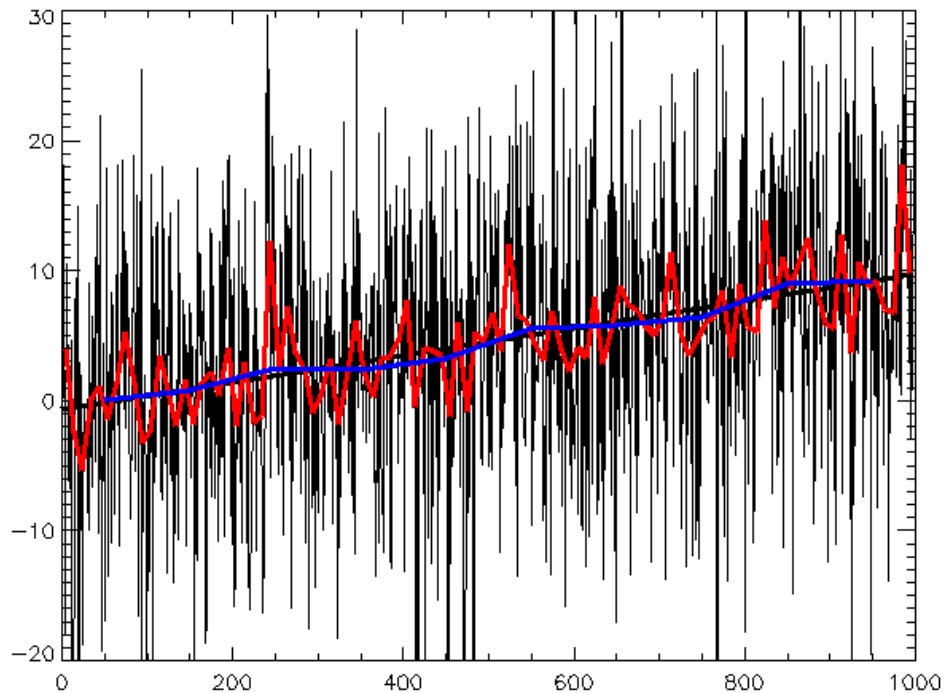


Figure 3.4: Random data plus trend, with best-fit line and different smoothings applied. (source: Wikimedia Commons [25])

Although it is possible to store time-series data in many diverse database types, the design of these systems with **time** as a **key** index is distinctly different from relational databases, which reduce discrete relationships through referential models. In many cases, the repositories of time-series data will utilize compression algorithms to manage the data efficiently [27]. Furthermore, time series databases can also be configured to regularly delete old data, unlike traditional databases which are designed to store data indefinitely.

3.2.2 Influx solution

InfluxDB is an open-source time series database TSDB created by the InfluxData organization [28]. It is written in the Go programming language and is used for time series data storage and retrieval in various sectors such as operations monitoring, application metrics, Internet of Things, and, especially important in our context, sensor data and real-time analytics. It also supports the processing of data from Graphite, a data logging and graphing tool for time series data. **3. Quest'ultima frase non è così rilevante**

Core Features InfluxDB has no external dependencies and provides a SQL-like vocabulary with built-in time-centric functions for querying a data

structure made up of measurements, series, and points, which listens on port 8086 [29].

time	location	scientist	butterflies	honeybees
2015-08-18T00:00:00Z	1	langstroth	12	23
2015-08-18T00:00:00Z	1	perpetua	1	30
2015-08-18T00:06:00Z	1	langstroth	11	28
2015-08-18T05:54:00Z	2	langstroth	2	11
2015-08-18T06:00:00Z	2	langstroth	1	10
2015-08-18T06:06:00Z	2	perpetua	8	23
2015-08-18T06:12:00Z	2	perpetua	7	22

Table 3.1: Sample time series dataset: number of butterflies and honeybees counted by two scientists

Each point is made up of a fieldset and a timestamp, which are key-value pairs. These form a series when they are grouped together by a set of key-value pairs known as a tagset. Finally, a measurement is created by grouping series together using a string identification. 64-bit integers, 64-bit floating points, strings, and booleans are among the possible values as shown in the table above [3.1]. The time and tagset are used to sort the points. As a side note it is important to know that data is downsampled and removed according to retention policies, which are set by measurement and that *Continuous Queries* are executed on a regular basis and the results are stored in a goal measurement.

Design Tradesoff InfluxDB is a time series database and optimizing for this use case involves a number of trade-offs, primarily to increase performance at the cost of functionality [29]. Here is a list of three design ideas that lead to compromises that I personally experienced during my experience:

1. Deletes are a rare occurrence. When they do occur, it is almost always against large ranges of old data that are cold for writes.
2. Updates to existing data are a rare occurrence, and contentious updates never happen. Time series data is predominantly new data that is never updated.
3. Many time series are ephemeral. There are often time series that appear only for a few hours and then go away, e.g., a new host that gets started and reports for a while and then gets shut down.

4. Potrebbero essere più di 3, mi sembrava un buon compromesso

Pros	Cons
Restricting access to deletes and updates allows for increased query and write performance	Delete and Update functionality is significantly restricted, since influxDB is not CRUD
InfluxDB is good at managing discontinuous data	Schema-less design means that some database functions are not supported e.g. there are no cross table joins

Table 3.2: Pros and cons of InfluxDB

Conclusion It is therefore no surprise to conclude that, when compared to a general purpose relational database like SQL Server, InfluxDB, using default single node configuration, outperformed both write speed, disk storage usage (by a factor of 27x) and query execution time, where InfluxDB is up to 20x faster with an average of 8x faster [30]. It can be seen that InfluxDB, tailored-made for Time Series data, is released after the other competitive technologies (like Graphite, TimescaleDB, Prometheus) and yet still among the top list. https://db-engines.com/en/ranking_trend/time+series+dbms
5.Non so se valga la pena citarlo The InfluxQL, SQL-like query language, helps make it easier to use and adapt for people who are used to working with relational databases such as MySQL, even though influx implementation is a smaller subset, not fully supporting *Create Read Update Delete* (CRUD) operations.

3.3 Grafana

Grafana is a web-based analytics and interactive visualization application that runs on a variety of platforms. When connected to supported data sources, it produces web-based charts, graphs, and alerts. Grafana Enterprise, a paid version with more features, is available as a self-hosted installation or as a Grafana Labs cloud service account [32]. Grafana is split into two parts: a front end and a back end, both of which are built in TypeScript and Go and, through a plug-in, system, it can be expanded, adding specific customizations to the existing platform [33]. Using interactive query builders, end users can develop complicated monitoring dashboards. But what is a dashboard anyway?

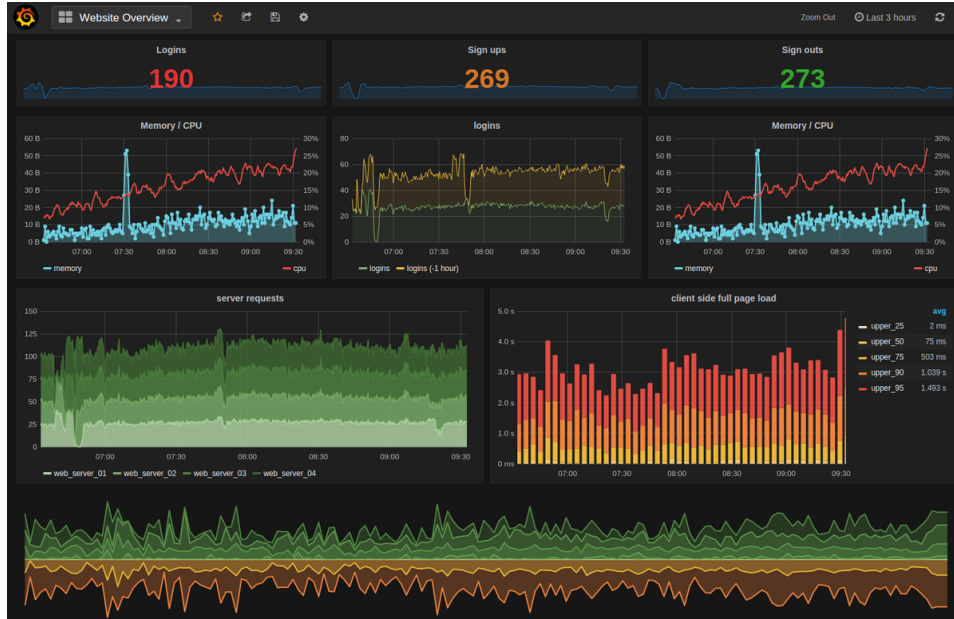


Figure 3.5: Grafana Graph Visualization (Source: Flickr [31])

3.3.1 Dashboard what is it?

A dashboard, in business, is a type of graphical user interface that often enables quick access to key performance indicators (KPIs) related to a certain goal or business activity. In our context, "dashboard" refers to a "progress report" or "report" and, as a type of data visualization, it is mostly accessible by a web browser and is usually linked to regularly updating data sources. The term dashboard is derived from the automotive dashboard, where drivers may monitor the primary functions at a glance using the instrument panel.

The success of dashboard projects depends on the relevancy/importance of information provided within the dashboard. This includes the metrics chosen to monitor and the timeliness of the data forming those metrics; data must be up-to-date and accurate. Well known dashboards include Google Analytics dashboards, used on 55% of all websites [34], which show user activity on a website, or the UK government, and similar for each country, coronavirus tracker, for the COVID-19 pandemic [35].

An interesting project is the GLAM Wiki dashboard, from Israel [36]. Its purpose is to assist GLAM institutions (galleries, libraries, archives, and museums) in tracking the use of their free-content files that they have submitted to Wikimedia projects. Based on multiple specified indices and several time frames, the dashboard visualizes statistical data that shows the extent of exposure and usage of these public-domain assets. The collecting data, which is presented in a variety of diagrams and graphs, allows the institutions to obtain insights, discover patterns and preferences, and understand the overall

impact of these free materials on the global audience of Wikimedia-project users.

3.3.2 Key strengths

Now that we are familiar with the "dashboard" concept, we can highlight why we should use Grafana; here are some of the key features [37]:

- **Visualize** fast and flexible visualization with a variety of options allows data to be displayed the way the user wants it;
- **Dynamic Dashboard** dynamic and reusable dashboards may be created using template variables.
- **Explore Metric** ad-hoc queries and dynamic drill-down permits for data discovery. View splits and side-by-side comparisons of different time ranges and data sources is easily achievable;
- **Explore Logs** fast switch from metrics to logs preserving label filters. Furthermore, searching the logs is rather quick and can be performed on live streams;
- **Alerting** most of the vital/operational metrics may be alerted visually and different types of notifications (SMS, mail, Slack) may be despatched with the aid of Grafana;
- **Mixed Data Source** the same chart can have different data source: these can be selected based on queries, with built-in support for most of the prominent data sources available in the market as well as custom ones;
- **Annotations** graphs may have events that can be annotated, two solutions are possible: use native annotation store, with the ability to add annotation events directly from the graph panel or via the HTTP API, or querying other data sources. Event metadata and tags can be seen when hovering over events.

Loading speed of Grafana dashboards 6.Potenzialmente da rimuovere/es-pandere Loading speed of a Grafana dashboard depends on 5 major things:

1. pre-selected and saved time window: the larger the time period you query, the longer it takes to open and display the contents;
2. data frequency in the panels: in case of the very high frequency, non aggregated data, even if selected time period is minutes, it will take time to load;
3. the number of panels with the data inside;

4. your database structure;
5. whether calculations have to happen inside the panel before the data is displayed.

Conclusion Grafana is the right choice when visualizing infrastructure, applications, network devices, sensors, and more. This is a great 24/7 monitoring solution for NOC and DevOps teams. It can also help to manage all data from other application monitoring tools like AppDynamics, New Relic, Splunk, Dynatrace and all-in-one web interface for data viewing, alerting, and reporting. A further comparison with other data visualization tools, such as Power BI and Tableau, might make interesting reading.

Acronyms

Infralytics [©]	Infrastructure Analytics, 11
AI	Artificial Intelligence, 8
BI	Business Intelligence, 1
CBM	Condition-based maintenance, 8
CDA	Confirmatory data analysis, 1
CM	Condition monitoring, 9
CMMS	Computerized Maintenance Management Systems, 8
CRISP	Cross-industry standard process, 2
CRUD	Create Read Update Delete, 25
DM	Data mining, 1
EDA	Exploratory data analysis, 1
IoT	Internet Of Things, 8
KPIs	Key performance indicators, 4
ML	Machine Learning, 9
MRO	Maintenance, repair and overhaul, 6
PdM	Predictive maintenance, 8
PM	Preventive maintenance, 7
PPM	Planned preventive maintenance, 7
TSDB	Time Series Database, 23

Bibliography

- [1] Meta S. Brown. “Transforming Unstructured Data into Useful Information”. In: *Big Data, Mining, and Analytics*. Auerbach Publications, Mar. 2014, pp. 227–246. DOI: 10.1201/b16666. (Visited on 08/26/2021).
- [2] Claude A. Pruneau. “The Multiple Facets of Correlation Functions”. In: *Data Analysis Techniques for Physical Scientists*. Cambridge University Press, 2017, pp. 526–576. DOI: 10.1017/9781108241922.013.
- [3] Schutt Rachel and O’Neil Cathy. *Doing data science*. O’Reilly Media, 2013, pp. 1–30, 120–140. ISBN: 9781449358655.
- [4] Vicki Adams. “Introduction to data analysis”. In: *The Journal of small animal practice* 49 (Sept. 2008), pp. 375–6. DOI: 10.1111/j.1748-5827.2008.00647.x.
- [5] James Goodnight. “The forecast for predictive analytics: hot and getting hotter”. In: *Statistical Analysis and Data Mining* 4 (Jan. 2011), pp. 9–10. DOI: 10.1002/sam.10106. (Visited on 06/16/2019).
- [6] John W. Tukey. “The Future of Data Analysis”. In: *The Annals of Mathematical Statistics* 33.1 (1962), pp. 1–67. DOI: 10.1214/aoms/1177704711. URL: <https://doi.org/10.1214/aoms/1177704711>.
- [7] Wikipedia Contributors. *Data Cleansing*. Wikipedia, Apr. 2019. URL: https://en.wikipedia.org/wiki/Data_cleansing (visited on 03/10/2022).
- [8] Charles M. Judd. *Data analysis: a model-comparison approach*. Harcourt Brace Jovanovich, 1989, pp. 10–15. URL: https://openlibrary.org/books/OL18760366M/Data_analysis (visited on 03/14/2022).
- [9] European Federation of Maintenance Societies. *EFNMS: What Does EFNMS Stand For?* <http://www.efnms.eu/>, 2016. URL: <http://www.efnms.eu/about-us/what-does-efnms-stand-for/> (visited on 03/07/2022).
- [10] *Paints and varnishes – Corrosion protection of steel structures by protective paint systems – Part 9: Protective paint systems and laboratory performance test methods for offshore and related structures*. Tech. rep. Geneva, CH: International Organization for Standardization, Jan.

2018. URL: <https://www.iso.org/standard/64832.html> (visited on 03/07/2022).
- [11] U.S. Indo-Pacific Command. *USS Ronald Reagan (CVN 76) Undergoes Preventive Maintenance*. flickr, Jan. 2016. URL: <https://www.flickr.com/photos/us-pacific-command/23569479804/> (visited on 03/07/2022).
 - [12] Michael Decourcy Hinds. “PREVENTIVE MAINTENANCE: A CHECKLIST”. In: *The New York Times* (Feb. 1985). URL: <https://www.nytimes.com/1985/02/17/realestate/preventive-maintenance-a-checklist.html> (visited on 02/28/2022).
 - [13] Jonathan Trout. *Preventive Maintenance: An Overview*. Reliableplant, Sept. 2019. URL: <https://www.reliableplant.com/Read/12494/preventive-maintenance> (visited on 02/28/2022).
 - [14] Daniel Penn. *What is Preventive & Predictive Maintenance?* Daniel Penn Associates, Jan. 2020. URL: <https://www.danielpenn.com/preventive-predictive-maintenance-2020/> (visited on 03/08/2022).
 - [15] Bizerba S.p.A. *Condition-based maintenance*. Bizerba.com. URL: https://www.bizerba.com/it_it/argomenti/digital-services/condition-based-maintenance/condition_based_maintenance.html (visited on 03/09/2022).
 - [16] Leith Hitchcock. “ISO Standards for Condition Monitoring”. In: *Engineering Asset Management*. Springer, 2012, pp. 606–613. DOI: 10.1007/978-1-84628-814-2_65. (Visited on 03/08/2022).
 - [17] Valerio Dida et al. *Manufacturing: Analytics unleashes productivity and profitability | McKinsey*. McKinsey, 2020. URL: <https://tinyurl.com/mckinseybusiness>.
 - [18] Yves Van Ingelgem and Marteen Durie. *Zensor: our approach*. Zensor, 2022. URL: <https://www.zensor.be/our-approach> (visited on 01/16/2022).
 - [19] Yves Van Ingelgem and Marteen Durie. *What Does The Future Of Maintenance Look Like?* Zensor, Dec. 2020. URL: <https://www.zensor.be/blog/futureofmaintenance> (visited on 03/07/2022).
 - [20] *License:OLDAP-2.7 – Free Software Directory*. 1999. URL: <https://directory.fsf.org/wiki/License:OLDAP-2.7> (visited on 01/31/2022).
 - [21] Wes McKinney. “Data Structures for Statistical Computing in Python”. In: *Proceedings of the 9th Python in Science Conference*. Ed. by Stéfan van der Walt and Jarrod Millman. 2010, pp. 56–61. DOI: 10.25080/Majora-92bf1922-00a.

- [22] *IO tools (text, CSV, HDF5, ...)* – *pandas 1.4.0 documentation*. URL: https://pandas.pydata.org/docs/user_guide/io.html (visited on 01/31/2022).
- [23] *pandas-dev/pandas: Pandas 1.4.0*. Jan. 2022. DOI: 10.5281/zenodo.5893288. URL: <https://doi.org/10.5281/zenodo.5893288>.
- [24] Anurag Pandey. *Split-Apply-Combine Strategy for Data Mining*. en. Oct. 2020. URL: <https://medium.com/analytics-vidhya/split-apply-combine-strategy-for-data-mining-4fd6e2a0cc99> (visited on 02/28/2022).
- [25] *File:Random-data-plus-trend-r2.png* - *Wikimedia Commons*. en. URL: <https://commons.wikimedia.org/wiki/File:Random-data-plus-trend-r2.png> (visited on 02/08/2022).
- [26] Matt Asay. *Why time series databases are exploding in popularity*. en-US. June 2019. URL: <https://www.techrepublic.com/article/why-time-series-databases-are-exploding-in-popularity/> (visited on 02/08/2022).
- [27] Michael Duffy. *DevOps Automation Cookbook*. en. Google-Books-ID: k_SoCwAAQBAJ. Packt Publishing, Nov. 2015. ISBN: 9781784398392.
- [28] *InfluxDB: Open Source Time Series Database*. URL: <https://www.influxdata.com/> (visited on 02/06/2022).
- [29] *InfluxDB OSS 1.8 Documentation*. URL: <https://docs.influxdata.com/influxdb/v1.8/> (visited on 01/31/2022).
- [30] Syeda Noor et al. *Université libre de Bruxelles Advanced Databases Time Series Databases and InfluxDB*. 2017. URL: https://cs.ulb.ac.be/public/_media/teaching/influxdb_2017.pdf (visited on 02/18/2022).
- [31] Linux Screenshots. *grafana dashboard*. Jan. 2016. URL: <https://www.flickr.com/photos/xmodulo/24311604930/> (visited on 02/14/2022).
- [32] *Grafana: The open observability platform*. en. URL: <https://grafana.com/> (visited on 02/06/2022).
- [33] *Grafana documentation*. en. URL: <https://grafana.com/docs/grafana/latest/> (visited on 01/31/2022).
- [34] W3Techs. *Usage Statistics and Market Share of Google Analytics for Websites*. en. W3techs. URL: <https://w3techs.com/technologies/details/ta-googleanalytics> (visited on 02/28/2022).
- [35] *UK Summary | Coronavirus (COVID-19) in the UK*. en. URL: <https://coronavirus.data.gov.uk> (visited on 02/14/2022).
- [36] *GLAM Wiki Dashboard*. URL: <https://glamwikidashboard.org/about> (visited on 02/14/2022).

- [37] Sunil Kumar and Prof. Saravanan. “A Comprehensive study on Data Visualization tool - Grafana”. English. In: 8.5 (May 2021), pp. 1–7. ISSN: ISSN-2349-5162. URL: <https://www.jetir.org/papers/JETIR2105788.pdf>.