# Data Science -Challenges and Impact of Statistics

Author1- Wasim Khan, Lecturer , Government  Women's Polytechnic College, Indore
EmailID-wasukhan1982@gmail.com
Author2 -Mr. Dharmendra Gupta, Asst prof, IPS Academy , IES Indore
dharmendragupta86@gmail.com

## Abstract-

In this paper, focus on different essential motives why analyzing the area of statistics is integral in contemporary society. First, statisticians are guides for studying from statistics and navigating frequent troubles that can lead you to incorrect conclusions. Second, given the growing importance of decisions and opinions based on data, it's quintessential that you can significantly determine the great of analyses that others present to you.We substantiate our premise that information is one of the maximum crucial disciplines to provide equipment and strategies to locate structure in and to offer deeper insight into information, and the maximum critical subject to analyze and quantify uncertainty. We supply a top level view over special proposed systems of Data Science and cope with the effect of information on such steps as information acquisition and enrichment, information exploration, records analysis and modelling, validation and representation and reporting. Also, we suggest fallacies when neglecting statistical reasoning. The subject of information is the technological know-how of learning from information. Statistical information enables you operate the proper techniques to gather the statistics, appoint the perfect analyses, and efficaciously gift the consequences. Statistics is a important technique in the back of how we make discoveries in technological know-how, make selections based on information, and make predictions. Statistics permits you to recognize a topic much greater deeply.
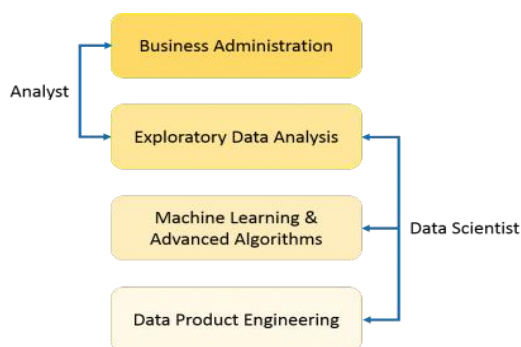
**Keywords** Structures of data science · Impact of statistics on data science · Fallacies in data science.

## Introduction-

In the present, everything has become data-driven. the quantity of knowledge made each second within the world multiples into terabytes, and this suggests that the sector of knowledge science has additionally mature at the same pace at the same time. Analyzing such giant amounts of information need capable data scientists. Therefore it's no surprise that there's an enormous demand for excellent information analysts which it's become such a profitable field nowadays. Statistics may be a vary of procedures for gathering, organizing, analyzing and presenting quantitative information. 'Data' is that the term for facts that are obtained and later recorded, and, for statisticians, 'data' sometimes refers to quantitative information that square measure numbers. Basically so, statistics may be a scientific approach to analyzing numerical information. Using statistics, going to gain deeper and extra fine grained insights into but exactly our information is structured and supported that structure however we will

optimally apply various knowledge science techniques to urge even additional information. The main aim of information science is to analyze the unstructured data being created nowadays, however typically often not possible to try to qualitatively – it's to be done quantitatively. Once analyzing this information, organizations ought to get real insights regarding their customers and their wants, in order that these insights may be translated into correct business price quickly. Therefore, the encumbrance is on information scientists to hold out their analyses properly, thus on improve and optimize the method business is conducted. Organizations in an exceedingly form of fields, starting from health care to recreation presently follow this model.

 "Unstructured data" will incorporate messages, features, images, on-line networking, and different shopper made substance. Informatics of times obliges managing associate awe-inspiring live info and composing calculations to concentrate bits of data from this information.



## Importance of data Statistics-

Data is a technological relation to show how that helps us make decisions and draw conclusions within the presence of variability. Number of Data statistics

technique is one of the most vital disciplines to provide equipment and methods to locate structure in and to present deeper perception into information, and the maximum essential field to investigate and quantify uncertainty. For instance, civil engineers running in the transportation subject are involved approximately the capability of regional highway structures. A typical trouble would involve records at the range of network, home-primarily based trips, the wide variety of folks consistent with household, and the number of cars per household, and the goal might be to produce a journey-generation version relating journeys to the quantity of humans in keeping with family and the variety of automobiles per family. A statistical technique called regression evaluation can be used to construct this model.
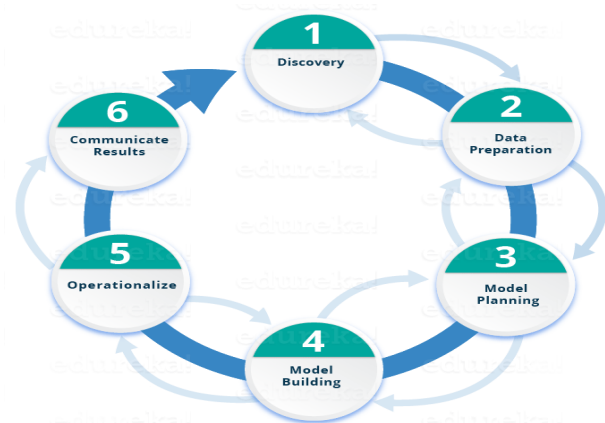
## Data science Predictive Technique-

**Predictive causal analytics –**The purpose of predictive analytics is to come across issues before they even arise. With each era of planes being extra linked, it's miles easier than it's ever been to gather huge amounts of actual-time performance data..If you want a model which could are expecting the opportunities of a selected occasion within the destiny, you want to use predictive causal analytics. Say, if you are presenting money on credit, then the possibility of customers making future credit score bills on time is a matter of concern for you. Here, you may construct a model which can perform predictive analytics on the price records of the patron to are expecting if the future bills may be on time or now not.

**Prescriptive analytics:-**Prescriptive analytics -takes predictive analytics one step in addition by using supplying unique and actionable subsequent steps for the way to clear up the troubles introduced up inside the predictive facts evaluation. While predictive analytics can inform you what is going to take place, while it will take place and why, prescriptive analytics applies many layers of device getting to know to suggest options for taking advantage of future opportunities or mitigating future risks and the capability results of each choice. In exclusive terms, it no longer completely predicts however indicates pretty a couple of prescribed actions and associated consequences. The most effective example for that is often Google's self-the use of vehicle that I had cited prior to too. The data focused with the help of cars is also accustomed educate self-using automobiles. You will be able to run algorithms in this information to carry intelligence thereto. This might exchange your automobile to require options like as soon as to expose, that course to require, whereas to gradual down or accelerate.

**Machine learning for making predictions** - if you have transactional data of a finance enterprise and need to build a model to decide the future trend, then device gaining knowledge of algorithms are the pleasant bet. This falls below the paradigm of supervised studying. It's far called supervised because you already have the data based on which you can teach your machines. As an example, a fraud detection model can be skilled the usage of a ancient file of fraudulent purchases.

**Machine learning for pattern discovery** - if you don't have the parameters based on which you may make predictions, then you definitely need to find out the hidden patterns within the dataset that allows you to make significant predictions. This is not anything but the unsupervised model as you don't have any predefined labels for grouping. The most common algorithm used for sample discovery is clustering.

**Lifecycle of Data Science-**Here is a brief overview of the main phases of the Data Science Lifecycle:-
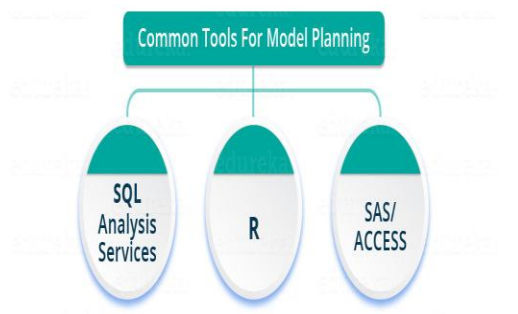


**Phase 1—discovery:** earlier than you begin the project, it is crucial to apprehend the diverse specifications, requirements, priorities and required price range. You should possess the ability to invite the right questions. Here, you assess when you have the specified assets found in phrases of human beings, generation, time and information to aid the mission. In this section, you furthermore might want to frame the commercial enterprise problem and formulate preliminary hypotheses (ih) to test.

**Phase 2**—records practice: in this segment, you require analytical sandbox in which you

could carry out analytics for the entire duration of the assignment. You need to discover, preprocess and condition records prior to modeling. In addition, you will perform etlt (extract, rework, load and remodel) to get records into the sandbox. Let's have a observe the statistical evaluation waft below.

You can use R programming for statistics cleansing, transformation, and visualization. This could assist you to identify the outliers and establish a relationship between the variables. Once you have wiped clean and prepared the facts, it's time to do exploratory analytics on it. Permit's see how you can attain that.

**Phase 3**—Model planning-Here , you'll decide the methods and techniques to attract the relationships among variables. Those relationships will set the base for the algorithms which you'll implement within the subsequent segment. You'll observe exploratory statistics analytics (eda) using various statistical formulas and visualization gear.



R has a complete set of modeling capabilities and provides a good environment for building interpretive models.

**SQL** Analysis services can perform in-database analytics using common data mining functions and basic predictive models.

**SAS/ACCESS** can be used to access data from Hadoop and is used for creating repeatable and reusable model flow diagrams.

**Phase 4**—Model building:on this phase, you may increase datasets for training and testing functions. You will don't forget whether your present gear will suffice for strolling the fashions or it's going to want a better environment (like speedy and parallel processing). You will examine various gaining knowledge of strategies like category, association and clustering to build the version.

**Phase 5**—Operationalize-in this phase, you supply very last reports, briefings, code and technical documents. Similarly, now and again a pilot task is also implemented in a actual-time production surroundings. This could provide you a clean photo of the overall performance and other related constraints on a small scale before full deployment.

**Phase 6**—Communicate results: now it is essential to evaluate if you have been able to obtain your intention which you had deliberate inside the first segment. So, in the closing section, you discover all of the key findings speak to the stakeholders and determine if the results of the venture are a success or a failure primarily based at the criteria evolved in phase 1.

Now, I will take a case study to explain you the various phases described above.

**Case Study: Diabetes Prevention**

What if we could predict the occurrence of diabetes and take appropriate measures beforehand to prevent it?

In this use case, we will predict the occurrence of diabetes making use of the entire lifecycle that we discussed earlier. Let's go through the various steps.

**Step 1:** First, we will collect the data based on the medical history of the patient as discussed in

Phase 1. You can refer to the sample data below

```
.
;npreg;glu;bp;skin;bmi;ped;age,income
1;6;148;72;35;33.6;0.627;50
2;1;85;66;29;26.6;0.351;31
3;1;89;80;23;28.1;0.167;21
4;3;78;50;32;31;0.248;26
5;2;197;70;45;30.5;0.158;53
6;5;166;72;19;25.8;0.587;51
7;0;118;84;47;45.8;0.551;31
8;1;103;30;38;43.3;0.183;33
9;3;126;88;41;39.3;0.704;27
10;9;119;80;35;29;0.263;29
11;1;97;66;15;23.2;0.487;22
12;5;109;75;26;36;0.546;60
13;3;88;58;11;24.8;0.267;22
14;10;122;78;31;27.6;0.512;45
15;4;97;60;33;24;0.966;33
16;9;102;76;37;32.9;0.665;46
17;2;90;68;42;38.2;0.503;27
18;4;111;72;47;37.1;1.39;56
19;3;180;64;25;34;0.271;26
20;7;106;92;18;39;0.235;48
21;9;171;110;24;45.4;0.721;54
```

**Attributes:**

npreg – Number of times pregnant
glucose – Plasma glucose concentration
bp – Blood pressure
skin – Triceps skinfold thickness
bmi – Body mass index

ped – Diabetes pedigree function
age – Age
income – Income

**Step 2:** Now, once we have the data, we need to clean and prepare the data for data analysis.

This data has a lot of inconsistencies like missing values, blank columns, abrupt values and incorrect data format which need to be cleaned.

Here, we have organized the data into a single table under different attributes – making it look more structured.
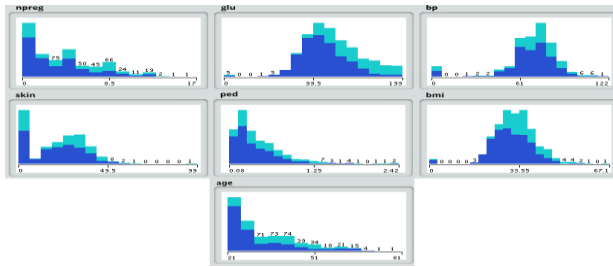
Let's have a look at the sample data below.

| | npreg | glu | bp | skin | bmi | ped | age | income |
|---|---|---|---|---|---|---|---|---|
| 1 | 6 | 148 | 72 | 35 | 33.6 | 0.627 | 50 | |
| 2 | 1 | 85 | 66 | 29 | 26.6 | 0.351 | 31 | |
| 3 | 1 | 89 | 6600 | 23 | 28.1 | 0.167 | 21 | |
| 4 | 3 | 78 | 50 | 32 | 31 | 0.248 | 26 | |
| 5 | 2 | 197 | 70 | 45 | 30.5 | 0.158 | 53 | |
| 6 | 5 | 166 | 72 | 19 | 25.8 | 0.587 | 51 | |
| 7 | 0 | 118 | 84 | 47 | 45.8 | 0.551 | 31 | |
| 8 | one | 103 | 30 | 38 | 43.3 | 0.183 | 33 | |
| 9 | 3 | 126 | 88 | 41 | 39.3 | 0.704 | 27 | |
| 10 | 9 | 119 | 80 | 35 | 29 | 0.263 | 29 | |
| 11 | 1 | 97 | 66 | 15 | 23.2 | 0.487 | 22 | |
| 12 | 5 | 109 | 75 | 26 | 36 | 0.546 | 60 | |
| 13 | 3 | 88 | 58 | 11 | 24.8 | 0.267 | 22 | |
| 14 | 10 | 122 | 78 | 31 | 27.6 | 0.512 | 45 | |
| 15 | 4 | | 60 | 33 | 24 | 0.966 | 33 | |
| 16 | 9 | 102 | 76 | 37 | 32.9 | 0.665 | 46 | |
| 17 | 2 | 90 | 68 | 42 | 38.2 | 0.503 | 27 | |
| 18 | 4 | 111 | 72 | 47 | 37.1 | 1.39 | 56 | |
| 19 | 3 | 180 | 64 | 25 | 34 | 0.271 | 26 | |
| 20 | 7 | 106 | 92 | 18 | | 0.235 | 48 | |
| 21 | 9 | 171 | 110 | 24 | 45.4 | 0.721 | 54 | |

**Step 3:**

Now let's do some analysis as discussed earlier in Phase 3.

First, we will load the data into the analytical sandbox and apply various statistical functions on it. For example, R has functions like describe which gives us the number of missing values and unique values. We can also use the summary function which will give us statistical information like mean, median, range, min and max values. Then, we use visualization techniques like histograms, line graphs, box plots to get a fair idea of the distribution of data.

**Step 4:**
Now, based on insights derived from the previous step, the best fit for this kind of problem is the decision tree. Let's see how? Since, we already have the major attributes for analysis like npreg, bmi, etc., so we will use supervised learning technique to build a model here.

Further, we have particularly used decision tree because it takes all attributes into consideration in one go, like the ones which have a linear relationship as well as those which have a non-linear relationship. In our case, we have a linear relationship between npreg and age, whereas the nonlinear relationship between npreg and ped.

**Step 5:**
In this phase, we will run a small pilot project to check if our results are appropriate. We will also look for performance constraints if any. If the results are not accurate, then we need to replan and rebuild the model.

**Conclusion**
The role data in information technology is under-envisioned as, e. G., compared to engineering science. This yields, specifically, for the areas of statistics acquisition and enrichment yet as for superior modeling required for prediction. Stirred with the help of this finish, statisticians area unit nicely-recommended to further obnoxiously play their role during this trendy and properly frequently occurring space of information technology. solely complementing and/or combining mathematical strategies and machine algorithms with applied mathematics reasoning, notably for big data, can cause scientific consequences based mostly all on acceptable processes. Ultimately, solely a balanced interaction of all sciences involved can cause winning solutions in facts technology.

**References**

1.Jeff Leek (2013-12-12). "The key word in 'Data Science' is not Data, it is Science". Simply Statistics.
2.Data Munging with Perl. DAVID CROSS. MANNING.
3.Allen, E., Allen, L., Arciniega, A., Greenwood, P.: Construction of equivalent stochastic differential equation models. Stoch. Anal. Appl. 26, 274–297 (2008)MathSciNetCrossRefGoogle Scholar
4. Anderson, C.: The End of Theory: The Data Deluge Makes the Scientific Method Obsolete. Wired Magazine
5. Aue, A., Horváth, L.: Structural breaks in time series. J. Time Ser. Anal. 34(1), 1–16 (2013)MathSciNetCrossRefGoogle Scholar
6. Statistical Modelling: the two cultures Leo Breiman. Statistical Science. Vol. 16 No.3 (August 2001) 199-215
7. Bischl, B., Mersmann, O., Trautmann, H., Weihs, C.: Resampling methods for meta-model validation with recommendations for evolutionary computation. Evol. Comput. 20(2), 249–275 (2012)CrossRefGoogle Scholar
8.Data Munging with Perl. DAVID CROSS. MANNING. Chapter 1 Page 4.
9.What is Data Science? http://www.datascientists.net/what-is-data-science
10.The Data Science Venn Diagram. http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram