# DATA PREPARATION FOR OPTIMIZATION IN DATA MINING

Ashish Bhagchandani,
Bachelor of Engineering Student:
Information Technology Department
Gandhinagar Institute of Technology, Gandhinagar, India.

*Abstract:* Data can be structured or unstructured, it totally depends on how data has been generated. A typical database which is well organized into formatted and brief repository, so that the elements of database can be made addressable and used for more effective processing and analysis. Data preparation is the fundamental stage for data analysis. As there is a lot of poor quality or dirty data available at many resources, so as to remove all the unwanted data sets, data preparation is must. This paper shows various techniques which should be kept in consideration for data preparation.

*IndexTerms* - **Data mining, data preparation, data analysis, datasets.**

## 1. INTRODUCTION

As a clear amber fluid, gasoline—the power behind the transportation industry barely resembles the adhesive black ooze pumped-up up through oil wells. The distinction between the main two liquids is that the results of multiple steps of refinement that distill helpful product from the staple. Knowledge preparation could be a terribly similar method. The staple comes from operational systems that have usually accumulated crud, within the kind of eccentric business rules and layers of system enhancements and fixes, over the course of your time. Fields within the knowledge square measure used for multiple functions. Values become obsolete. Errors square measure fastened on associate degree current basis, therefore interpretations amendment over time. Of getting ready knowledge is just like the process of refinement oil. Valuable stuff lurks within the slime of operational knowledge. 0.5 The battle is refinement. The opposite 0.5 is changing its energy to a helpful form—the equivalent of running associate degree engine on gas.

The proliferation of knowledge could be a feature of recent business. Our challenge is to create sense of the information, to refine the information in order that the engines of knowledge mining will extract worth. One in all the challenges is that the sheer volume of knowledge. A client could decision the decision center many times a year, pay a bill once a month, flip the phone on once every day, build and receive phone calls many times every day. Over the course of your time, many thousands or countless customers square measure generating many countless records of their behavior. Even on today's computers, this is often tons of knowledge process. as luck would have it, laptop systems became powerful enough that the matter is admittedly one in all having associate degree adequate take into account shopping for hardware and software; technically, process such Brobdingnag Ian quantities of knowledge is feasible [9].

Data comes in several forms, from several systems, and in many alternative varieties. Knowledge is usually dirty, incomplete, typically incomprehensible and incompatible. This is, alas, the world. And yet, knowledge is that the staple for data processing. Oil starts out as a thick tarry substance, mixed with impurities. It's solely by prying varied stages of refinement that the staple becomes usable—whether as clear gas, plastic, or fertilizer. Even as the foremost powerful engines cannot use fossil fuel as a fuel, the foremost powerful algorithms, the engines of knowledge of information mining square measure unlikely to search out attention-grabbing patterns in unprepared data.

After over a century of experimentation, the steps of refinement oil square measure quite well understood—better understood than the processes of getting ready knowledge. This chapter illustrates some tips and principles that, supported expertise, ought to build the method more practical. It starts with a discussion of what knowledge ought to appear as if once it's been ready, describing the client signature. It then dives into what knowledge truly feels like, in terms of knowledge varieties and column roles. Since a serious a part of victorious data processing is within the derived variables, ideas for these square measure given in some detail. The chapter ends with a glance at a number of the difficulties given by dirty knowledge and missing values, and also the process challenge of operating with giant volumes of economic knowledge [9].

## 2. LITERATURE SURVEY

### 2.1. What Data Should Look Like

The place to start out the discussion on information is at the end: what the info ought to seem like. All data processing algorithms wish their inputs in tabular form—the rows and columns thus common in spreadsheets and databases. In contrast to spreadsheets, though, every column should mean an equivalent issue for all the rows. Some algorithms want their information in a very specific format. For example, market basket analysis sometimes appearance at solely the merchandise purchased at any given time. Also, link analysis wants references between records so as to attach them. However, most algorithms, and particularly call trees, neural networks, clustering, and statistical procedure, square measure searching for information in a very specific format known as the client signature.

### 2.2. The Customer Signature

The customer signature is consider as brief introduction of the customer behavior which gives importance to all the current attributes and also customer behavior overtime. Customer signature is unique in characteristic like a signature on check. Often customer signature have no more identifying information than a string of seemingly random digits representing a household, individual, or account number, therefore unlike a signature on a check, though, the customer signatures used for analysis and not identification. Figure 1 shows that a customer signature is simply a row of data that represents the customer and whatever might be useful for data mining [6].
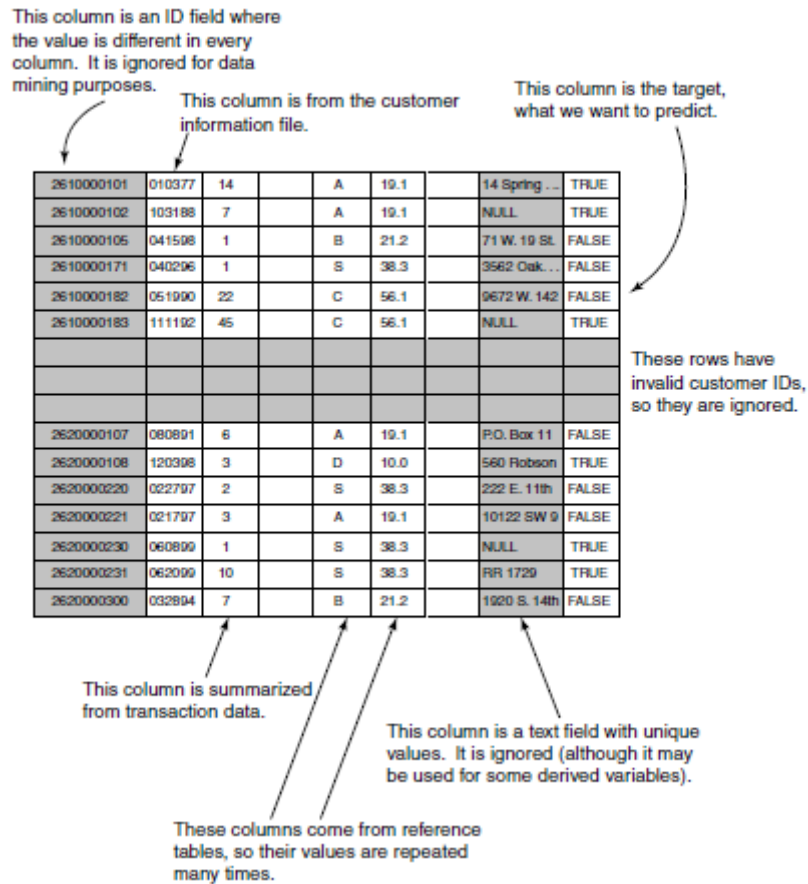


Figure 1: Each row in the customer signature represents one customer (the unit of data mining) with fields describing that customer.

It is maybe unfortunate that there's no giant info with up-to-date client signatures prepared for all modeling applications. Such a system could seem terribly helpful initially sight. However, the dearth of such a system is a chance as a result of modeling efforts need understanding the information [11]. No single consumer signature works for all modeling efforts, though some client signatures work well for several applications "Customer" may be a unit of knowledge mining in client signatures. This focuses totally on customers, that the unit of knowledge mining is usually associate degree account, a private, or a unit. There are a unit different prospects. Acquisition modeling usually happens at geographic levels, census block teams, or nada codes. And applications outside client relationship management area unit even additional heterogeneous.

## 3. VARIOUS CONSIDERATION FOR DATA PREPARATION

### 3.1. The Columns

The columns within the information contain values that describe aspects of the client. In some cases, the columns come back directly from existing business systems; additional typically, the columns are the results of some calculation—so referred to as derived variables.

Each column contains values. The vary refers to the set of allowable values for that column. Table 1 shows range characteristics for typical types of data used for data mining [5].

Table 1: Range Characteristics for Typical Types of Data Used for Data Mining

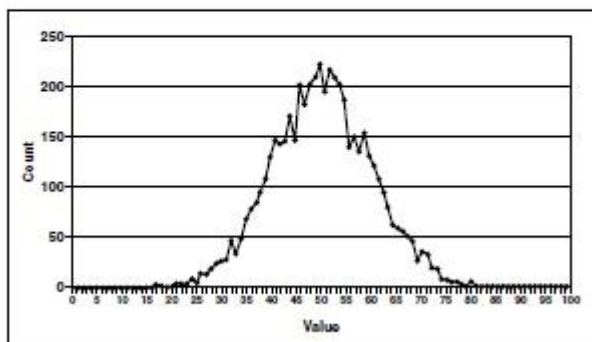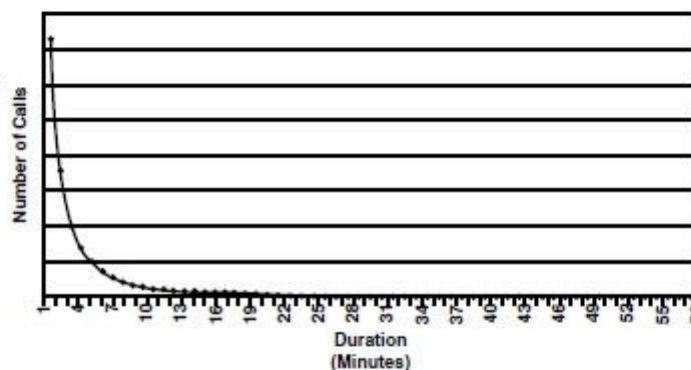| Variable Type | Typical Range Characteristics |
|---|---|
| Categorical variables | List of acceptable values |
| Numeric | Minimum and maximum values |
| Dates | Earliest and latest dates, often latest date is less than or equal to current date |
| Monetary amounts | Greater than or equal to 0 |
| Durations | Greater than or equal to 0 (or perhaps strictly greater than 0) |
| Binned or quantile values | The number of quantiles |
| Counts | Greater than or equal to 0 (or perhaps greater than or equal to 1) |



This histogram is for the month of claim for a set of insurance claims.

This is an example of a typically uniform distribution. That is, the number of claims is roughly the same for each month.

This histogram shows the number of telephone calls made for different durations.

This is an example of an exponentially decreasing distribution.

This histogram shows a normal distribution with a mean of 50 and a standard deviation of 10. Notice that high and low values are very rare.

Figure 2: Histograms show the distribution of data values

Histograms, such as those in Figure 2, shows however usually every worth or vary of values happens in some set of information. The vertical axis could be a count of records, and also the horizontal axis is that the values within the column. The form of this bar chart shows the distribution of the values (strictly speaking, in a very distribution, the counts area unit divided by the overall variety of records therefore the space below the curve is one). If we tend to area unit operating with a sample, and also the sample is arbitrarily chosen, then the distribution of values within the set ought to be concerning constant because the distribution within the original information.

The distribution of the values provides necessary insights into the information. It shows that values area unit common and that area unit less common. Simply staring at the distribution of values brings up questions—such as why AN quantity is negative or why some categorical values don't seem to be gift. Though statisticians tend to be a lot of involved with distributions than information miners, it's still necessary to appear at variable values [3]. Here, we tend to illustrate some special cases of distributions that area unit necessary for data processing functions, likewise because the special case of variables synonymous with the target.

### 3.1.1. Columns with One Value

The most degenerate distribution could be a column that has just one price. Unary valued columns, as they're additional formally illustrious, don't contain any data that helps totally differentiate to tell apart between different rows. As a result of they lack any data content, they ought to be unheeded for data processing functions.

Having just one price is typically a property of the info. It's not uncommon, for example, for an information to possess fields outlined within the information that don't seem to be nevertheless inhabited. The fields are only placeholders for future values, so all the values are uniformly something such as "null" or "no" or "0." Before throwing out unary variables, check that NULLs are being counted as values. Appended demographic variables sometimes have only a single value or NULL when the value is not known [5].

Unary-valued columns additionally arise once the information mining effort is concentrated on a set of shoppers, and also the field accustomed filter the records is maintained within the ensuing table. The fields that outline this set could all contain an equivalent worth. If we tend to are building a model to predict the loss-ratio (an insurance measure) for automobile customers in New Jersey, then the state field can continuously shave "NJ" crammed in. This field has no info content for the sample getting used, therefore it ought to be neglected for modeling functions.

### 3.1.2. Columns with Almost Only One Value

In "almost-unary" columns, almost all the records have the same value for that column. There may be a few outliers, but there are very few. For example, retail this chart shows an almost-unary column. The column was created by binning telephone call durations into 10 equal-width bins. Almost all values, 9988 out of 9995, are in the first bin. If variable width bins had been chosen, then the resulting column would have been more useful.

Data could summarize all the purchases created by every client in every department. Very few customers could build a buying deal from the automotive department of a food market or the tobacco department of a retail store. So, most customers can have a $0 for total purchases from these departments. Purchased information usually comes in Associate in Nursing "almost-unary" format, as well. Fields like "people collect ceramic ware dolls" or "amount spent on greens fees" can have a null or $0 price for nigh only a few individuals. Or, some data, like survey information, is just offered for a really tiny set of the purchasers [5].
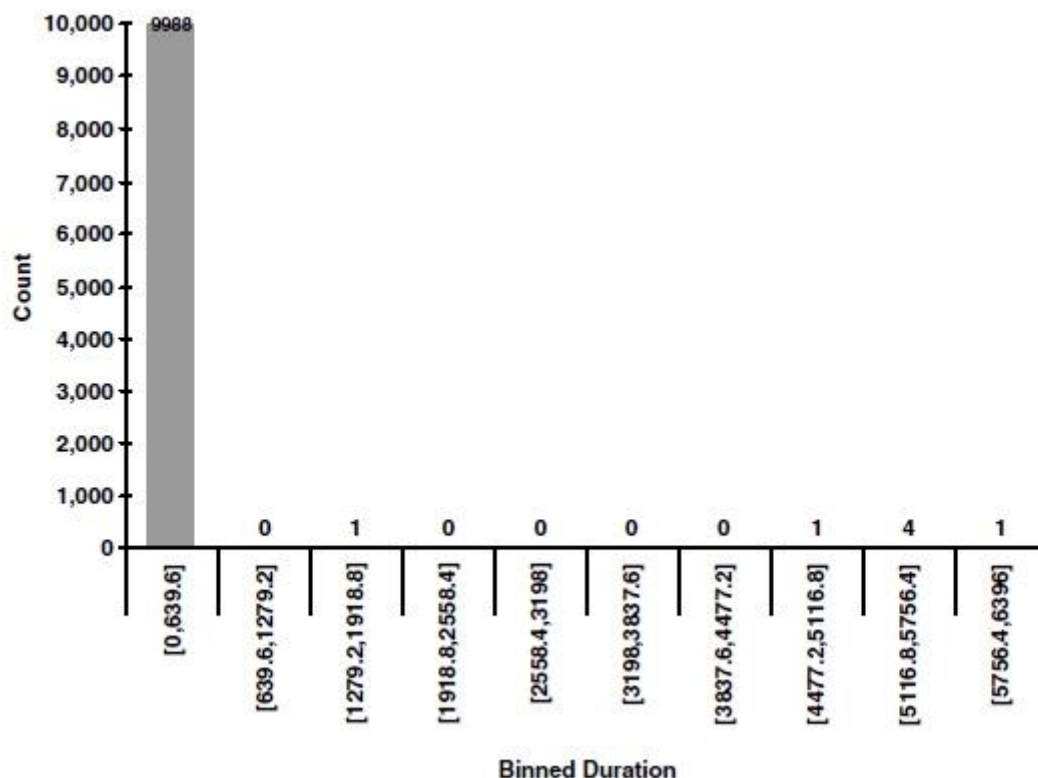


Figure 3: An almost-unary field, such as the bins produced by equal-width bins in this case, is useless for data mining purposes.

These are all extreme examples of data skew, shown in Figure 3. This chart shows an almost-unary column. The column was created by binning telephone call durations into 10 equal-width bins. Almost all values, 9988 out of 9995, are in the first bin. If variable width bins had been chosen, then the resulting column would have been more useful. The big question with "almost-unary" columns is, "When will they be ignored?" To justify ignoring them, the values should have two main characteristics.

First, the majority the records should have an equivalent worth. Second, there should be thus few records with a special worth, that they represent a negligible portion of the info. What's a negligible portion of information the info the information? It's a gaggle thus little that notwithstanding the data mining algorithms known it utterly, the cluster would be too little to be important.

Before ignoring a column, though, it's vital to grasp why the values are thus heavily inclined. What will this column tell America concerning the business? Maybe few folks ever get automotive merchandise as a result of solely a couple of the stores in question even sell them. Characteristic customers as "automotive product- consumers," during this case, might not be helpful. In alternative cases, an occasion may well be rare for alternative reasons. The quantity of individuals who cancel their utility on any given day is negligible, however over time the numbers accumulate. That the cancellations ought to be accumulated over an extended period of time, like a month, quarter, or year [2]. Or, the quantity of individuals collect ceramic ware dolls could also be terribly rare in itself, however once combined with alternative fields, this may counsel a vital phase of collectors. The rule of thumb is that, notwithstanding a column proves to be terribly informative, it's unlikely to be helpful for data processing if it's almost-unary. That is, absolutely understanding the rows with completely different values doesn't yield unjust results. As a general rule of thumb, if ninety five to ninety nine p.c of the values within the column are identical, the column—in isolation—is seemingly to be useless while not some work. As an example, if the column in question represents the target variable for a model, then proportional sampling will produce a sample wherever the rare values are additional extremely inhabited. Another approach is to mix many such columns for making derived variables that may encourage be valuable. As associate example, some census fields are sparsely inhabited, like those for explicit occupations.

However, combining a number of these fields into one field such as "high standing occupation"—can prove helpful for modeling functions.

### 3.1.3. Columns with Unique Values

At the opposite extreme square measure categorical columns that combat a special worth each single row—or nearly every row [1]. These columns establish every client unambiguously (or shut enough), for example:
1) Customer name
2) Address
3) Telephone number
4) Customer ID
5) Vehicle identification number

These columns are not terribly useful. Why? They are doing not have prophetic price, as a result of the unambiguously establish every row. Such variables cause overfitting. Sometimes these columns contain a wealth of knowledge. Lurking within phone numbers and addresses is vital geographical info. Customers' 1st names offer a sign of gender. Client numbers could also be consecutive assigned, telling America that customers are additional recent—and thence show up as vital [12].

### 3.1.4. Columns Correlated with Target

When a column is just too extremely correlate with the target column, it will mean that the column is simply an equivalent word. Here are two examples:
1) "Account range is NULL" is also synonymous with failure to reply to a selling campaign. Solely responders opened accounts and were appointed account numbers.
2) "Date of churn isn't NULL" is synonymous with having churned.

Another danger is that the column reflects previous business practices. As an example, the information could show that each one customers with telephony even have telephony. This is often a results of product bundling; telephony is sold-out in a very product bundle that perpetually includes telephony. Or the information could show that nearly all customers reside within the wealthiest areas, as a result of this wherever client acquisition campaigns within the past were targeted. This illustrates that knowledge miners have to be compelled to recognize historical business practices. Columns synonymous with the targets ought to be unnoticed.

### 3.2. Model Roles in Modeling

Columns contain data with data types. In addition, columns have roles with respect to the data mining algorithms [5]. Three important roles are:
1) **Input columns:** These are columns that are used as input into the model.
2) **Target column:** This column or set of columns is merely used once building prognostic models. These are what's fascinating, like propensity to shop for a selected product, chance to retort to a proposal, or chance of remaining a client. Once building afloat models, there doesn't have to be a target.
3) **Ignored columns:** These are columns that aren't used. Different tools have different names for these roles. Figure 17.4 shows how a column is removed from consideration in Agnos Knowledge Studio.

There are a unit some additional advanced roles further, that area unit used beneath specific circumstances. There are numerous model roles out there in SAS Enterprise jack. These model roles include:
1) **Identification column**: These area unit columns that unambiguously determine every row.
In general, these columns area unit unnoticed for data processing functions, however area unit vital for rating.
2) **Weight column**: this is often a column that specifies a "weight" to be applied to every row. This is often some way of making a weighted sample by as well as the burden within the knowledge.
3) **Cost column**: the value column specifies a value related to a row. For example, if we tend to area unit building a client retention model, then the "cost" would possibly embody associate estimate of every customer's price. Some tools will use this info to optimize the models that they're building. The extra model roles out there within the tool area unit specific to SAS Enterprise Miners.

### 3.3. Variable Measures

Variables seem in knowledge and have some necessary properties. Though databases square measure involved with the kind of variables (and we'll come to the present topic during a moment), data processing worries with the live of variables [10]. It's the live that determines however the algorithms treat the values. The subsequent measures square measure necessary for knowledge mining:

1) **Categorical variables**: This may be compared for equality however there's no important ordering. For instance, state abbreviations square measure categorical. The very fact that Alabama is next to AK alphabetically doesn't mean that they're nearer to every aside from Alabama and Tennessee, that share a geographic border however seem abundant any apart alphabetically.

2) **Ordered variables**: This may be compared with equality and with larger than and fewer than. Schoolroom grades, that vary from A to F, square measure associate example of ordered values.

3) **Interval variables**: The square measure ordered and support the operation of subtraction
(Although not essentially the other computing like addition and multiplication). Dates and temperatures square measure samples of intervals.

4) **True numeric variables**: The square measure interval variables that support addition and different mathematical operations. Financial amounts and client tenure (measured in days) square measure samples of numeric variables. The distinction between true numeric and intervals is delicate. However, data processing algorithms treat each of those an equivalent means. Also, note that these measures type a hierarchy. Any ordered variable is additionally categorical, any interval is additionally categorical, and any numeric is additionally interval.

There is a distinction between live and knowledge kind. A numeric variable, as an example, would possibly represent a secret writing scheme, say for account standing or maybe for state abbreviations. Though the values seem like numbers, they're extremely categorical. Nothing codes square measure a typical example of this development. Some algorithms expect variables to be of a particular live. Simple regression and neural networks, as an example, expect their inputs to be numeric [7]. So, if a zipper code field is enclosed and keep as variety, then the algorithms treat its values as numeric, typically not an honest approach. Call trees, on the opposite hand, treat all their inputs as categorical or ordered, even after they square measure numbers.

Measure is one necessary property. In observe, variables have associated varieties in databases and file layouts. The subsequent sections bring up knowledge varieties and measures in additional detail.

### 3.3.1. IDs and Keys

The purpose of some variables is to produce links to alternative records with a lot of data. IDs and keys are typically hold on as numbers, though they will even be hold on as character strings. As a general rule, such IDs and keys shouldn't be used directly for modeling functions.

A good example of a field that ought to typically be neglected for data processing functions are account numbers. The irony is that such fields might improve models, as a result of account numbers aren't assigned every which way. Often, they're assigned consecutive, therefore older accounts have lower account numbers; probably they're supported acquisition channel, therefore all internet accounts have higher numbers than alternative accounts. It's higher to incorporate the relevant data expressly within the client signature, instead of counting on hidden business rules. In some cases, IDs do write in code purposeful data. In these cases, the data ought to be extracted to form it a lot of accessible to the information mining algorithms.

Table 2: Data mining expects information to be in a very explicit format

| | |
|---|---|
| 1. | All information ought to be in a very single table. |
| 2. | Each row ought to correspond to Associate in nursing entity, like a client, that's relevant to the business. |
| 3. | Columns with one worth ought to be neglected. |
| 4. | Columns with a unique worth for each column ought to be ignored— though their info is also enclosed in derived columns. |
| 5. | For prognostic modeling, the target column ought to be known and every one substitutable columns removed. |

Alas, this can be not however information is found within the planet. Within the planet, information comes from supply systems, which can store every field in a very explicit means [13]. Often, we would like to exchange fields with values keep in reference tables, or to extract options from additional difficult information varieties. Consecutive section talks concerning putt this information along into a client signature.

## 4. CONSTRUCTING THE CUSTOMER SIGNATURE

Building the client signature, particularly the primary time, may be a terribly progressive method. At a minimum, client signatures got to be engineered a minimum of two times—once for building the model and once for grading it. In follow, exploring information and building models suggests new variables and transformations, that the method is recurrent repeatedly. Having a repeatable method simplifies the info mining work [8].

| | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SCORE | | | | | | 4 | 3 | 2 | 1 | | P |
| MODEL SET | | | | 4 | 3 | 2 | 1 | | P | | |
| MODEL SET | | | 4 | 3 | 2 | 1 | | P | | | |

Figure 4: Building customer signatures is an iterative process; start small and work through the process step-by-step, as in this example for building a customer signature for churn prediction.

The first step within the method, shown in Figure 4, is to spot the on the market sources of knowledge. After all, the client signature may be an outline, at the client level, of what's better-known regarding every client. The outline relies on the market information. This information could reside in a very information warehouse. It would equally well reside in operational systems and a few may be provided by outside vendors. When doing prognosticative modeling, it's significantly vital to spot wherever the target variable is returning from.

The second step is distinguishing the client. In some cases, the client is at the account level. In others, the client is at the individual or manage level. In some cases, the signature could don't have anything to try and do with someone in any respect. We've used signatures for understanding merchandise, zip codes, and counties, as an example, though the foremost common use is for accounts and households.

Once the client has been known, information sources got to be mapped to the client level. This could need extra operation tables—for instance, to convert accounts into households. It's going to not be attainable to search out the shoppers within the on the market information. Such a state of affairs needs revisiting the client definition. The key to putting together client signatures is to begin easy and build up. Rank the info sources by the convenience with that they map to the client. Begin with the best one, and build the signature victimization it [13]. You'll be able to use a signature before all the info is place into it. Whereas awaiting a lot of difficult information transformations, get your feet wet and perceive what's on the market. Once building client signatures out of transactions, take care to urge all the transactions related to a specific client.

### 4.1. Cataloging the Data

The data mining cluster at a mobile telecommunications company needs to develop a churn model in-house. This churn model can predict churn for one month, given a one-month lag time. So, if the information is accessible for Gregorian calendar month, then the churn prediction is for Apr. Such a model provides time for gathering the information and rating new customers, since the Gregorian calendar month information is accessible someday in March. At this company, there square measure many potential sources of knowledge for the client signatures. All of those square measure unbroken during an information repository with eighteen months of history. Each file is AN end-of-the-month snapshot—basically a dump of AN operational system into a knowledge repository.

The UNIT_MASTER file contains an outline of each variety phone number signaling in commission and an exposure of what's far-famed regarding the phone number at the top of the month. samples of fields during this file square measure the phone range, request account, request arrange, telephone set model, last beaked date, and last payment.

The TRANS_MASTER file contains each dealing that happens on a selected signal throughout the course of the month. These square measure account level transactions that embrace connections, disconnections, telephone set upgrades, and so on.

The BILL_MASTER file describes request data at the account level [9].

Multiple handsets may well be connected to constant request account—particularly for business customers and customers on family request plans. Though alternative sources of knowledge were offered within the company, these weren't instantly highlighted to be used for the client signature. One source, as an example, was the decision detail records—a record of each phone call—that is helpful for predicting churn. Though this information was eventually employed by the information mining cluster, it absolutely was not a part of this primary effort.

### 4.2. First Attempt

The first plan to build the client signature has to target the only knowledge supply. During this case, the only knowledge supply is that the UNIT_MASTER file, that handily stores knowledge at the phone range level, the extent being employed for the client signature.

It is value mentioning main two issues with this file and therefore the client definition:
1) Customers could modification their sign.
2) Telephone numbers is also reassigned to new customers.

These issues are going to be addressed later; the primary client signature is at the phone range level to induce started. The method accustomed build the signature has four steps: distinguishing the time frames, making a recent pic, pivoting columns, and calculative the target.
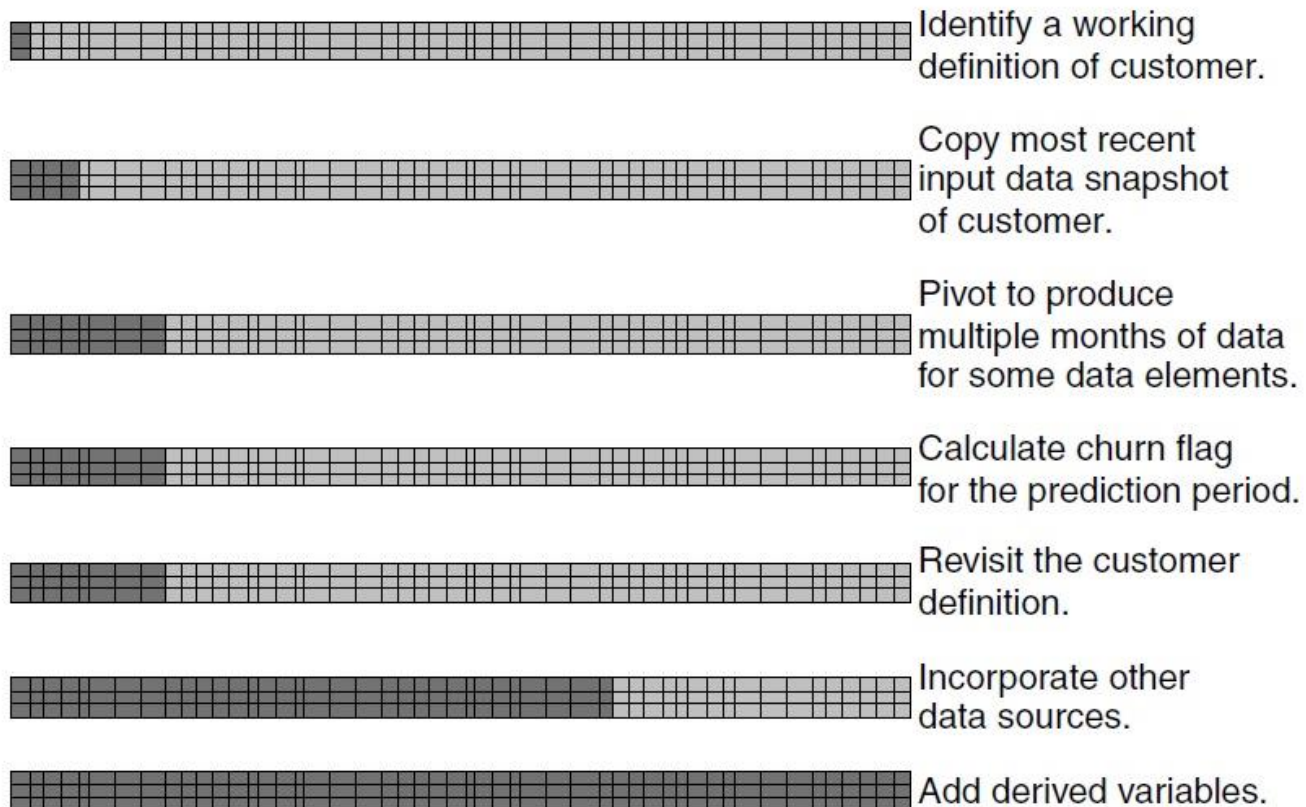
Figure 5: A model time chart shows the time frame for the input columns and targets when building a customer signature

Identifying the Time Frames, the first try at building the client signature has to take into consideration the timeframe for the info. Figure 5 shows a model time chart for this knowledge. The last word model set ought to have over just once frame it. However, the primary try focuses on only 1 timeframe. The timeframe outlined churn throughout one month—August. All of the input file return from a minimum of one month before. The cutoff date is June thirty, so as to produce one month of latency. Taking a Recent pic, the most recent pic of knowledge is outlined by the cutoff date. These fields within the signature describe the foremost recent data better-known a couple of client before he or she churned (or didn't churn).

This is a collection of fields from the UNIT_MASTER file for June—fields like the telephone set kind, asking arrange, and so on. It's necessary to stay the timeframe in mind once filling the client signature. It's an honest plan to use a naming convention to avoid confusion. During this case, all the fields may need a suffix of "_01," indicating that they're from the foremost recent month of input file. TI P Use a naming convention once building the client signature to point the timeframe for every variable. For example, the foremost recent month of input file would have a "_01" suffix; the month before, "_02"; so on. At now, presumptively not a lot of is understood concerning the fields, therefore descriptive data is beneficial. For example, the asking arrange may need an outline, monthly base, per-minute price, and so on. All of those options square measure fascinating and of potential worth for modeling—so it's cheap to bring them into the model set. Though descriptions aren't aiming to be used for modeling (codes square measure a lot of better), they assist the info miners perceive the info.

## 4.3. Pivoting Columns
Some of the fields in UNIT_MASTER represent knowledge that's rumored in a very regular statistic. For example, bill quantity contains a worth for each month, and every of those values has to be place into a separate column. These columns return from totally different UNIT_MASTER records, one for June, one for could, one for April, and so on. Employing a naming convention, the fields would be, for example:
Last_billed_amount_01 for June (which could already be within the snapshot)
Last_billed_amount_02 for could
Last_billed_amount_03 for Gregorian calendar month

At now, the client signature is beginning to form. Though the input fields solely return from one supply, the suitable fields are chosen as input and aligned in time [4].

Calculating the Target, a client signature for prophetic modeling wouldn't be helpful while not a target variable. Since the client signature goes to be used for churn modeling, the target has to be whether or not or not the client churned in August. This is often within the account standing field for the August UNIT_MASTER record. Note that solely customers UN agency were active on or before June thirty square measure enclosed within the model set. A client that starts in July and cancels in August isn't enclosed.

## 5. CONCLUSION

Data is that the gas that powers data processing. The goal of knowledge preparation is to supply a clean fuel, therefore the analytic engines work as with efficiency as doable. For many algorithms, the most effective input takes the shape of client signatures, one row of knowledge with fields describing numerous aspects of the client. Several of those fields square measure input fields, a number of square measure targets used for prophetical modeling. Sadly, client signatures aren't the means information is found in obtainable systems—and permanently reason, since the signatures amendment over time. In fact, they're perpetually being engineered and restored, with newer information and newer concepts on what constitutes helpful info.

Source fields are available in many totally different varieties, like numbers, strings, and dates. However, the foremost helpful values square measure sometimes those who square measure supplemental in. making derived values is also as straightforward as taking the ad of 2 fields. Or, they will need rather more subtle calculations on terribly massive amounts of knowledge. This can be significantly true once making an attempt to capture client behavior over time, as a result of statistic, whether or not regular or irregular, should be summarized for the signature.

Data conjointly suffers (and causes USA to suffer beside it) from problems— missing values, incorrect values, and values from totally different sources that disagree. Once such issues square measure known, it's doable to figure around them. The largest issues square measure the unknown ones—data that appears correct however is wrong for a few reason. Many data processing efforts need to use information that's but good. Therefore, data preparation plays a major role in detecting the wrong data by performing the troubleshot among the given dataset, which will help in optimizing output for the given requirement.

## REFERENCES

[1] Cooley, R., Mobasher, B. & Srivastava, J. Knowledge and Information Systems (1999) 1: 5. https://doi.org/10.1007/BF03325089

[2] David J. Hand, Statistical and Numerical Computing, 2013, doi: https://doi.org/10.1002/9780470057339.vad002.pub2

[3] E. Begoli and J. Horey, "Design Principles for Effective Knowledge Discovery from Big Data," *2012 Joint Working IEEE/IFIP Conference on Software Architecture and European Conference on Software Architecture*, Helsinki, 2012, pp. 215-218.
doi: 10.1109/WICSA-ECSA.212.32

[4] H. Michael Chung and Paul Gray, Special Section: Data Mining, Journal of Management Information Systems, Routledge, pp. 11-16

[5] Kantardzic, M., Data Mining: Concepts, Models, Methods, and Algorithms, Second Edition: Concepts, Models, Methods, and Algorithms, Wiley, 2011

[6]Kohavi, R., Mason, L., Parekh, R. et al. Machine Learning (2004) 57: 83. https://doi.org/10.1023/B:MACH.0000035473.11134.83

[7] Larose, D.T. and Larose, C.D., Discovering Knowledge in Data: An Introduction to Data Mining, Wiley, 2014

[8] McLellan, E., MacQueen, K. M., & Neidig, J. L. (2003). Beyond the Qualitative Interview: Data Preparation and Transcription. *Field Methods*, *15*(1), 63–84. https://doi.org/10.1177/1525822X02239573

[9] Michael J., A. Berry, Gordon, S. Linoff, Data Mining Techniques, pp, 539-596

[10] Pyle, D., Data Preparation for Data Mining, Elsevier Science 1999

[11] Rud, O.P., Data Mining Cookbook: Modeling Data for Marketing, Risk, and Customer Relationship Management, Wiley, 2001

[12] Stefan Stieglitza, Milad Mirbabaiea, Björn Rossa, Christoph Neubergerb, Social media analytics – Challenges in topic discovery, data collection, and data preparation

[13] Tan, P.N., Introduction to Data Mining, Pearson, 2018