

Data preparation phase - Overview

referred to as "data wrangling," or "data munging" or "data janitor work".

Includes everything from list verification to removing commas in the data, and debugging databases.

Cleaning data is very time consuming

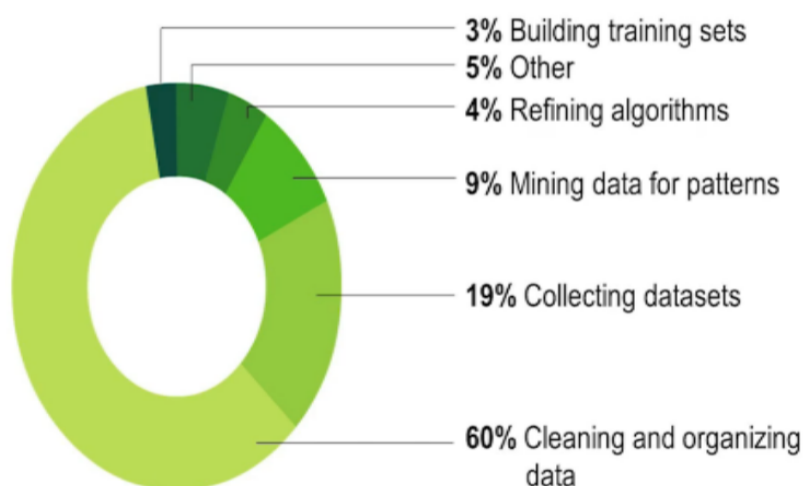
Introduction to data preparation

The chart below shows that 3 out of every 5 data scientists spend most time during their working day cleaning and organizing data.

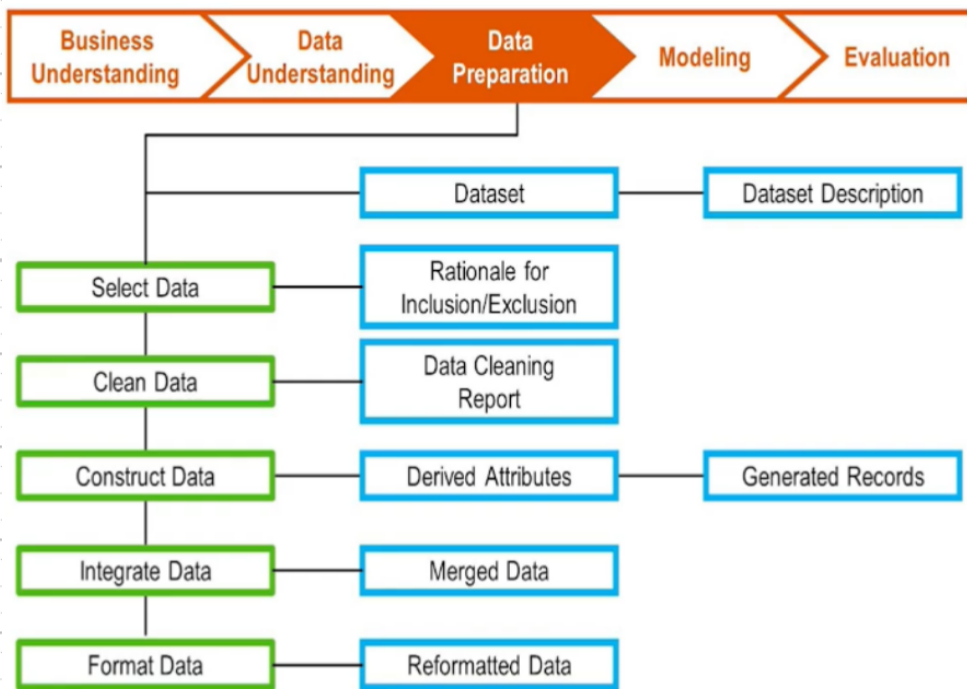
A New York Times article reported that data scientists spend from **50% to 80%** of their time mired in the more mundane task of collecting and preparing unruly digital data before it can be explored for useful nuggets.

For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insights.
New York Times. STEVE LOHR.
AUG. 17, 2014

What data scientists spend most time doing



CrowdFlower Data Science Report 2016



Data prep Phase -
This phase covers all the activities to construct the final analytical dataset.

from the initial raw data.

Data preparation tasks are likely to be performed multiple times and not in any prescribed order.

Tasks Include:

- Table
- Record
- Attribute selection
- Transformations
- Data cleaning for the modeling tool

Select Data task - Decides on the data that can be used for analysis.

The criteria to choose data includes the:

- Relevance to the data mining goals
- Data quality
- Technical constraints.

And this covers Attribute selection as well as the selection of records in a table.

The clean data task - raises the data quality to the level that's required for the selected analysis

The construct Data task - Includes the data preparation operations, such as the production of derived attributes, entire new records, or transformed values for existing attributes.

The integrate Data task combines information from multiple tables or records to create new records or values

Format Data task - produces the transformations that are primarily syntactic modifications made to the data, that don't actually change its meaning, but might be required by the modeling tool.

Data Preparation-Phase 3: Outputs

- Dataset
 - This is the dataset (or datasets) produced by the Data Preparation phase, which will be used for modeling or the major analysis work of the project.
- Dataset description
 - Describes the dataset (or datasets) that will be used for the modeling or the major analysis work of the project.

Phase 3.1: Select Data

- Tasks
 - Decide on the data to be used for analysis
 - Criteria include relevance to the data science goals and quality and technical constraints such as limits on →

data volume or data types.

- Note that data selection covers selection of attributes (columns) as well as selection of records (rows) in a table

• Output - Rationale for inclusion/exclusion

- List the data to be included/excluded and the reasons for these decisions.

Phase 3.2 - Clean Data

• Task

- Raise the data quality to the level required by the selected analysis techniques.

- This may involve selection of clean subsets of the data, the insertion of suitable defaults or more ambitious techniques such as the estimation missing data by modeling.

• Output - Data cleaning report

- Describe what decisions and actions were taken to address the data quality problems reported during the Verify Data Quality task of the Data Understanding phase.

Phase 3.3 - Construct Data

• Task

- This task includes constructive data preparation operations such as the production of derived attributes, entire new records, or transformed values for existing attributes

• Output - Derived Attributes

- Derived attributes are new attributes that are constructed from one or more existing attributes in the same record. Examples: $\text{area} = \text{length} * \text{width}$.

• Output - Generated records

- Describe the creation of completely new records

Phase 3.4 - Integrate Data

• Task

- These are methods whereby information is combined from multiple tables or records to create new records or values.

• Output - Merged Data

- Merging data tables refers to joining together two or

more tables that have different information about the same objects

- Merged data also covers aggregations.

↳ Aggregation refers to operation where new values are computed by summarizing together info from multiple records.

Example:

An aggregation converts a table of customer purchases, where there are multiple purchases per customer with one record for each purchase.

And you aggregate this into a new table where there is one record for each customer, and that will include fields such as:

- The number of purchases made
- The average purchase amount
- % of orders that have been charged to a credit card

Phase 3.5 - Format Data

- Task

- Formatting transformations refer to primarily syntactic modifications made to the data that do not change its meaning, but might be required by the modeling tool.

- Output - Reformatted data

- Some tools have requirements on the order of the attributes, such as the first field being a unique identifier for each record or the last field being the outcome field the model is to predict.

Example -

- Some tools have requirements on the order of the attributes, such as the first field being the unique identifier for each record, or the last field, the target field, the thing that we're trying to predict. It might be important to change the order of the records in the dataset, perhaps the modeling tool requires that the records be sorted

ted according to the value of the outcome attribute.

A common situation is that the records of the dataset are initially ordered in some way but the modeling tool itself needs them to be in a fairly randomized order.

Example

When using neural networks it's generally best for the records to be presented in a random order. And some data science tools, machine learning tools do this randomizing for you automatically. Additionally, there are purely syntactic changes made to satisfy the requirements of the specific modeling tool.

Example

Removing commas from within text fields in comma-delimited data-files is really important when you're doing text analysis

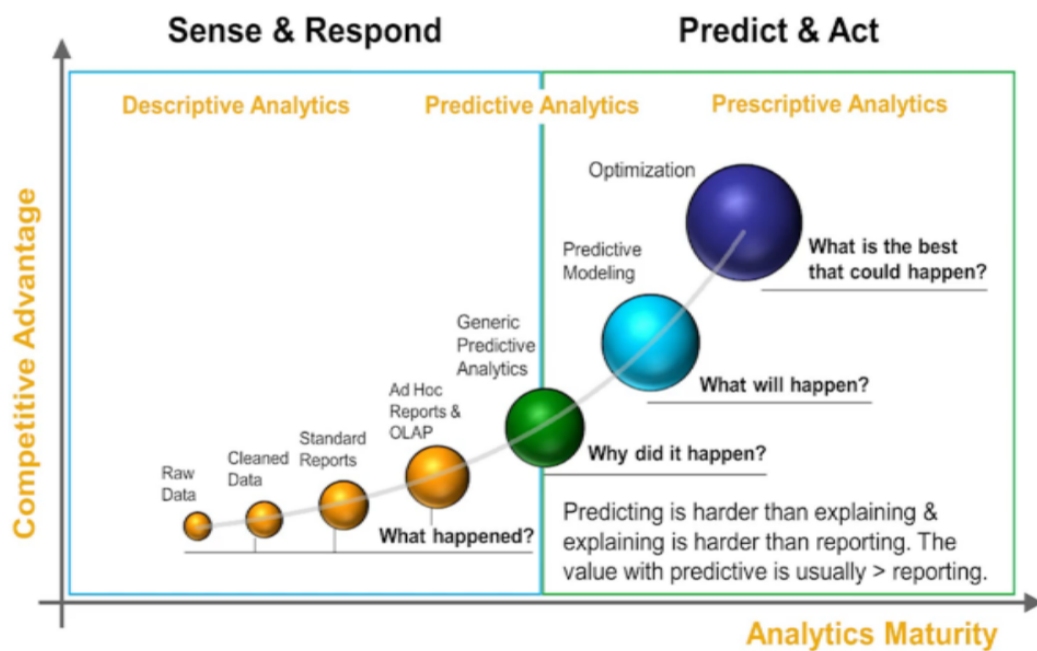


Predictive Modeling Methodology -

Predictive modeling encompasses a variety of techniques that analyze current and historical facts to make predictions about future or otherwise unknown events.

The output of a predictive model is a score or a probability of the targeted event occurring in a specified time frame in the future. Although, most often the unknown event of interest is in the future, predictive analytics can also be applied to any type of unknown, whether in the past, present or future

Introduction



The key is unlocking data to move decision making from sense & respond to predict & act

Example

- identifying suspects after a crime has been committed or credit card fraud as it occurs.

Descriptive analytics uses data visualization to provide insight into the past and answer the question "What has happened?".

Descriptive analytics are useful for company reports giving total stock inventory, average dollars spent per customer, and Year-over-Year changes in sales.

Predictive analytics uses statistical models and forecasting techniques to understand the future and answer: "What could happen?".

Predictive analytics analyses patterns in historical data and develops models that predict the probability of events happening in the future.

Prescriptive analytics, which use optimization and simulation algorithms enable us to answer: "What should we do" →

Prescriptive analytics predicts not only what will happen, but also why it will happen, providing recommendations regarding actions that will take advantage of the predictions.

Prescriptive analytics uses a combination of techniques and tools such as:

- Business rules
- Algorithms
- Optimization
- Machine learning
- Mathematical modeling and processing



Use predictive analytics to solve a variety of business challenges



- Churn Reduction
- Customer Acquisition
- Lead Scoring
- Product Recommendation
- Campaign Optimization
- Customer Segmentation
- Next Best Offer/Action



- Predictive Maintenance
- Load Forecasting
- Inventory/Demand Optimization
- Product Recommendation
- Price Optimization
- Manufacturing Process Optimization
- Quality Management
- Yield Management



- Fraud and Abuse Detection
- Claims Analysis
- Collection and Delinquency
- Credit Scoring
- Operational Risk Modeling
- Crime Threat
- Revenue and Loss Analysis



- Cash Flow and Forecasting
- Budgeting Simulation
- Profitability and Margin Analysis
- Financial Risk Modeling
- Employee Retention Modeling
- Succession Planning



- Life Sciences
- Healthcare
- Media
- Higher Education
- Public Sector/Social Sciences
- Construction and Mining
- Travel and Hospitality
- Big Data and IoT



Common use cases for machine learning and data analysis across industries

There are two phases to the predictive modeling process:

Model build

(The training phase)

Predictive models are built or "trained" on historic data with a known outcome

The input variables are called explanatory or independent variables

For model building the outcome, which is called target or dependent variable, is known.

Model apply

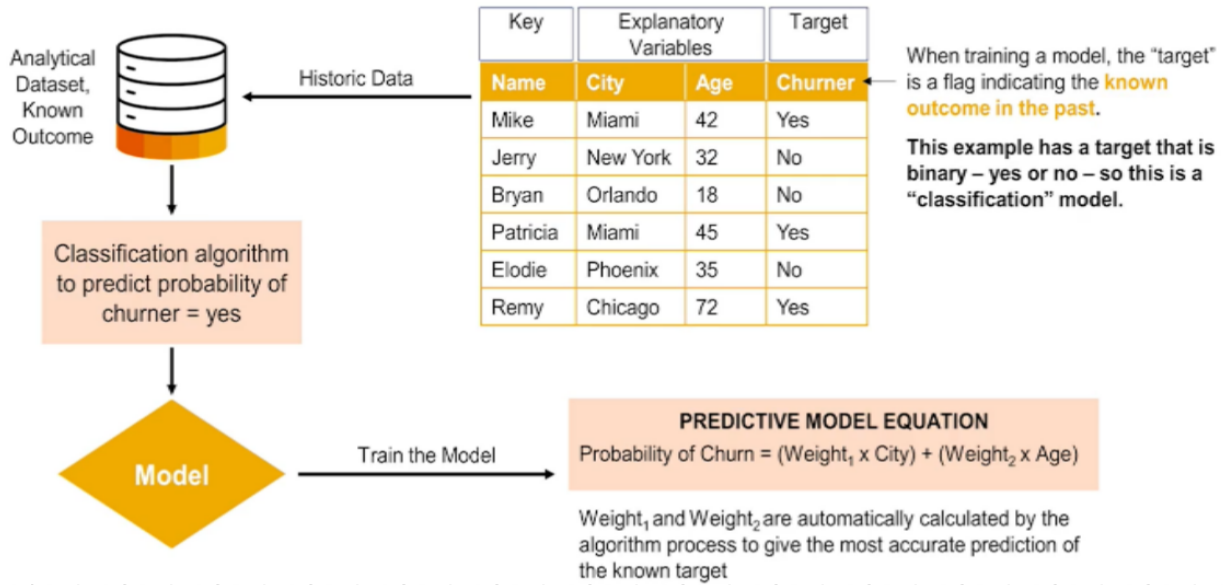
(The applying phase)

Once the model has been built, it is applied on more recent data, which has an unknown outcome (because the outcome is in the future)

The model calculates the score or probability of the target category occurring

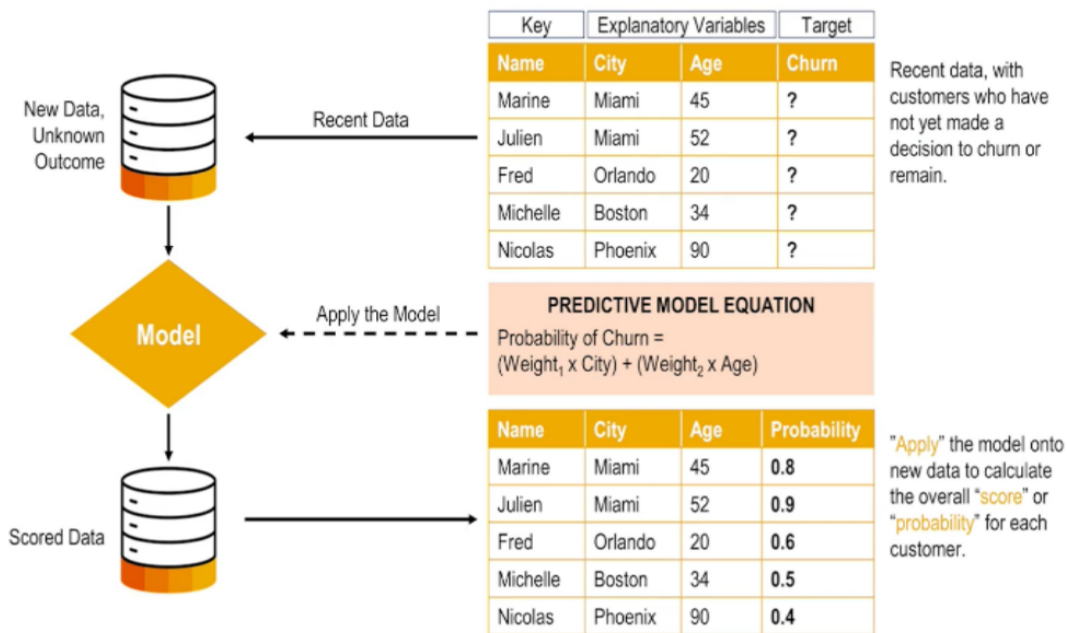
Eg. Calculates the probability of a customer responding to a market campaign

Building the model – Training phase



Identifies the key characteristics between a churmer and a non churmer

Using the model – Applying phase

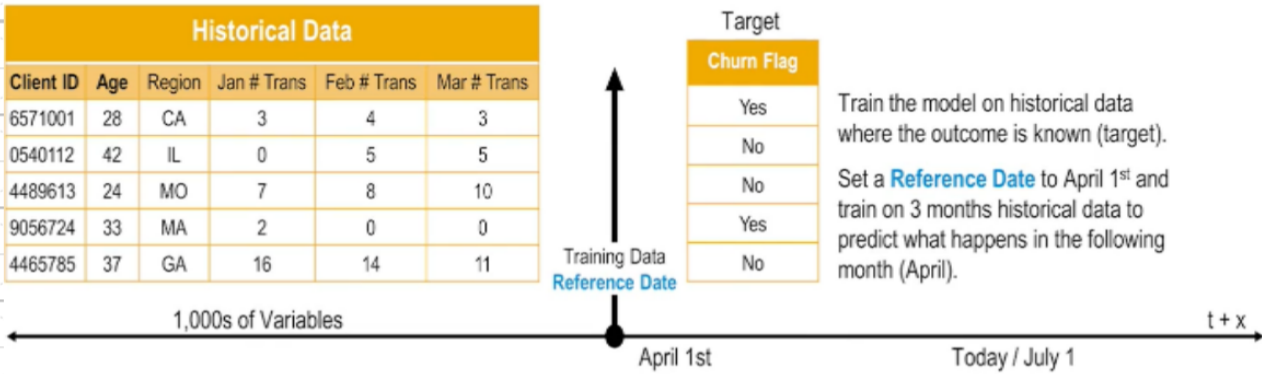


When we apply the model onto new data we don't know the target value →

Example moving data through time – Training phase

Imagine a Telco company on July 1st that wants to predict customer churn using a classification model

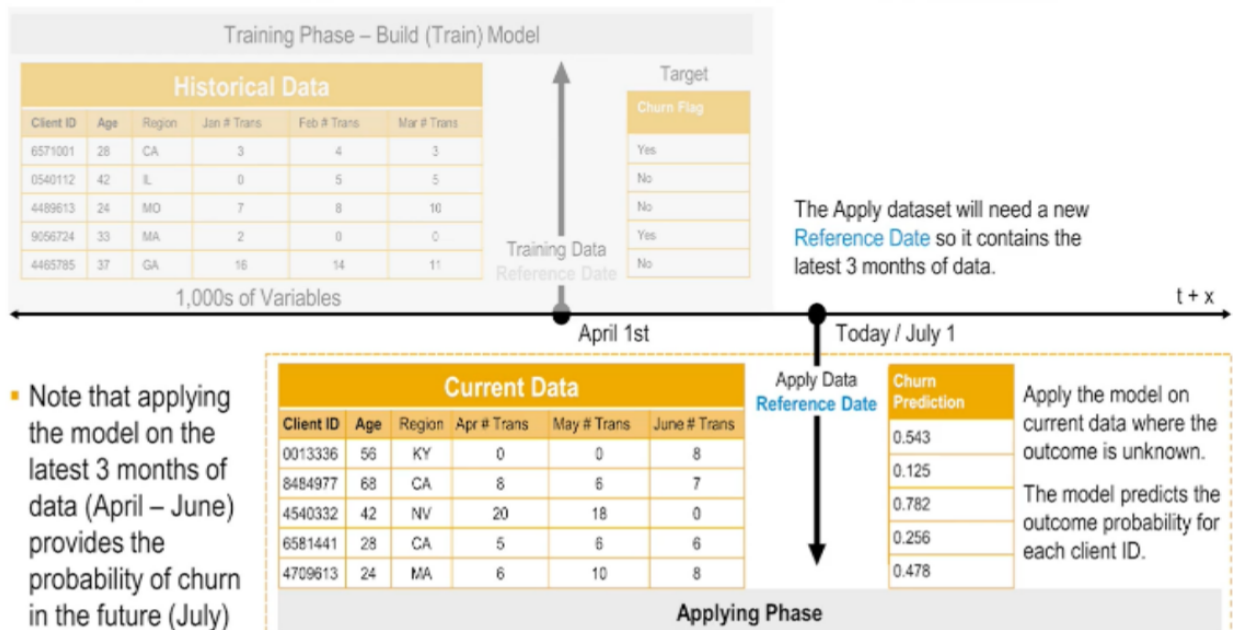
Training Phase – Build (Train) Model



- Note that the target data time frame (April) occurs AFTER the historical data time frame (January to March)
- The model is trained to identify patterns in the data in the past to predict the target in the following or later months

Example moving data through time – Applying phase

Imagine a Telco company on July 1st that wants to predict customer churn using a classification model



- Note that applying the model on the latest 3 months of data (April – June) provides the probability of churn in the future (July)

Depending on the business requirements models may need to be applied every month, week, day or even minutes or seconds

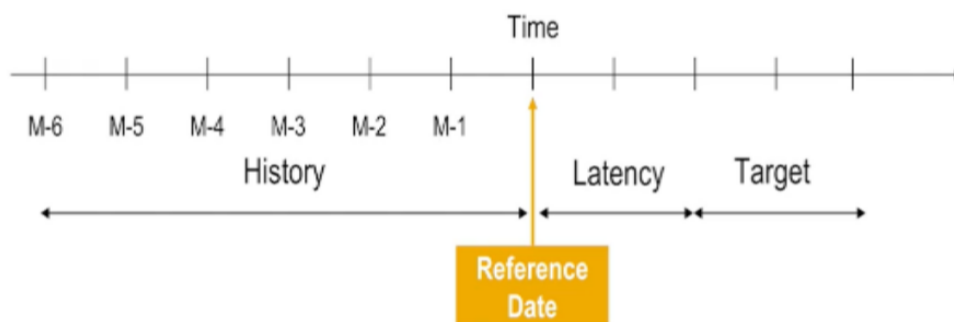
When the data-set changes the reference date also changes so any derived variables like a customer's age must be updated relative to the new reference date

Latency

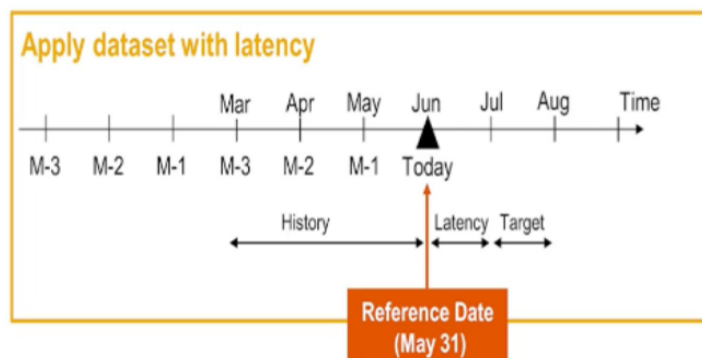
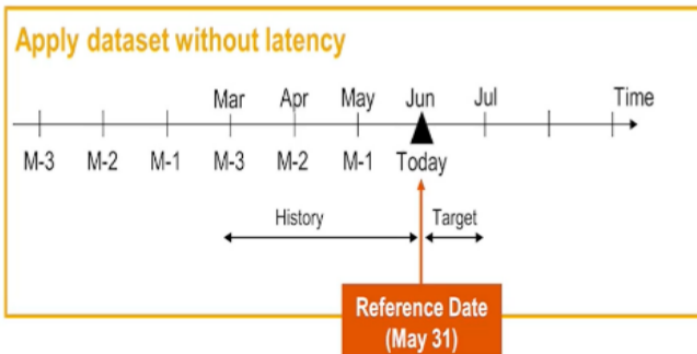
Sometimes when you design a model, you might want to build in a "latency" period.

Many datasets used for predictive modeling have the following structure:

- **Historic Data:** (in the past, compared to the reference date) with dynamic data computed in relation to the reference date. Usually short-term, mid-term, and long-term indicators.
- **Latency Period:** (starting after the reference state) a period where no data is collected. This is used to represent the time required by the business to collect new data, apply the model, produce the scores, and define the campaign. Not all predictive models require a latency period, although many churn models will.
- **Target:** (starting after the reference state + latency period) a period where the targeted behavior is observed.

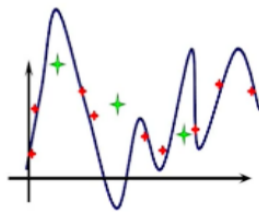


Why is latency needed?

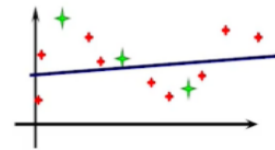


latency
 Periods can
 give the
 business
 time to
 convince the
 customer
 to stay

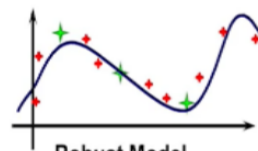
Model fitting



Over-Fit Model/Low Robustness
 (No Training Error, High Test Error)



Under-Fit Model/High Robustness
 (High Training Error = High Test Error)



Robust Model
 (Low Training Error \approx Low Test Error)



Dont overfit - accurate on training data but not on "real data".

Overfitting occurs when a model is too complex

A model that's been overfitted over-reacts to minor fluctuations in the data

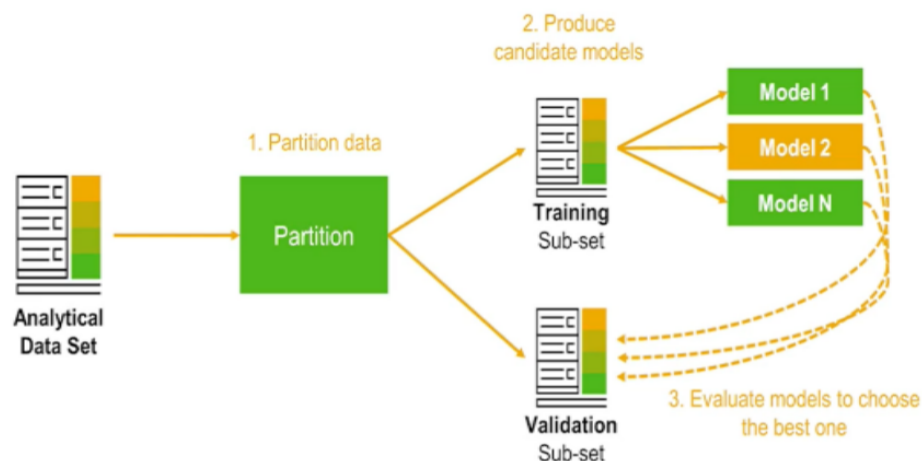
Additional techniques for avoiding overfitting -

- Cross-validation
- Regularization
- Early stopping
- Pruning decision trees

Dont Underfit

Hold-out sample

- During the Model Build phase, a "hold-out" sample is created.
- This is a sample of observations withheld from the model learning, so that the model's ability to predict future probabilities can be estimated by its ability to predict the data in the hold-out sample.



↖ a hold out sample is data that's separated before training to test and validate the model and its robustness after training

Data manipulation -

- Most data mining activities will require the data to be "prepared" before analysis
- Data manipulation is often driven by domain knowledge.
- This is a process where database tables are merged and aggregated, new variables and transformations are created to try and improve model quality, IF/THEN conditions are created, filters are applied etc...
- The first step is to identify the entity for analysis
 - An entity is the object targeted by the planned analytical task
 - It may be a customer, product, store, etc..., and it is usually identified by a unique identifier

- The entity defines the "granularity" of the analysis

Items of significance to an enterprise
are data entities

customer material product

Entity makes business sense

It can be characterized in terms of attributes

Can be associated with predictive metrics related to the
tasks you want to perform

Defining the entity is quite difficult

Feature engineering -

A feature is an attribute or property that is shared by all the entities on which an analysis or predictions are to be done

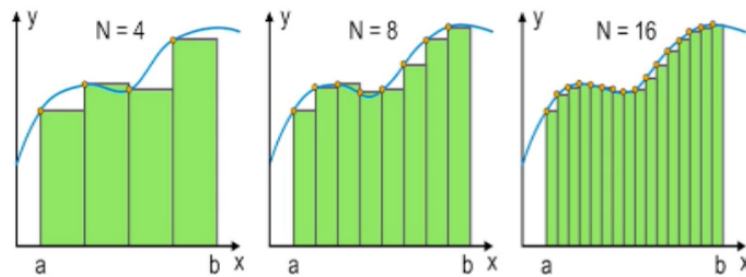
- In feature engineering, new features are created to extract more information from existing features.
- These new features may have an improved ability to explain the variance in the training data and improve model accuracy
- Feature engineering is highly influenced by business understanding

Main goals of feature engineering

- preparing the proper input dataset that's compatible with the machine learning algorithm requirements.
- try and improve the model performance

Binning

- Binning is one of the fundamental feature engineering techniques.
- The original data values which fall into a given small interval, a bin, are replaced by a value representative of that interval, often the central value.



Binning illustration of numerical data

To learn more about data binning see:

Binning is a technique used to reduce the effects of minor observation errors.

A bin is the original data value which falls into a given small interval, the bin is usually replaced by a value representative of that interval, and that's often the central value.

Binning "smooths" the input data, it might lower the chance of overfitting in small data-sets →

there are two general methods of dividing data into bins.

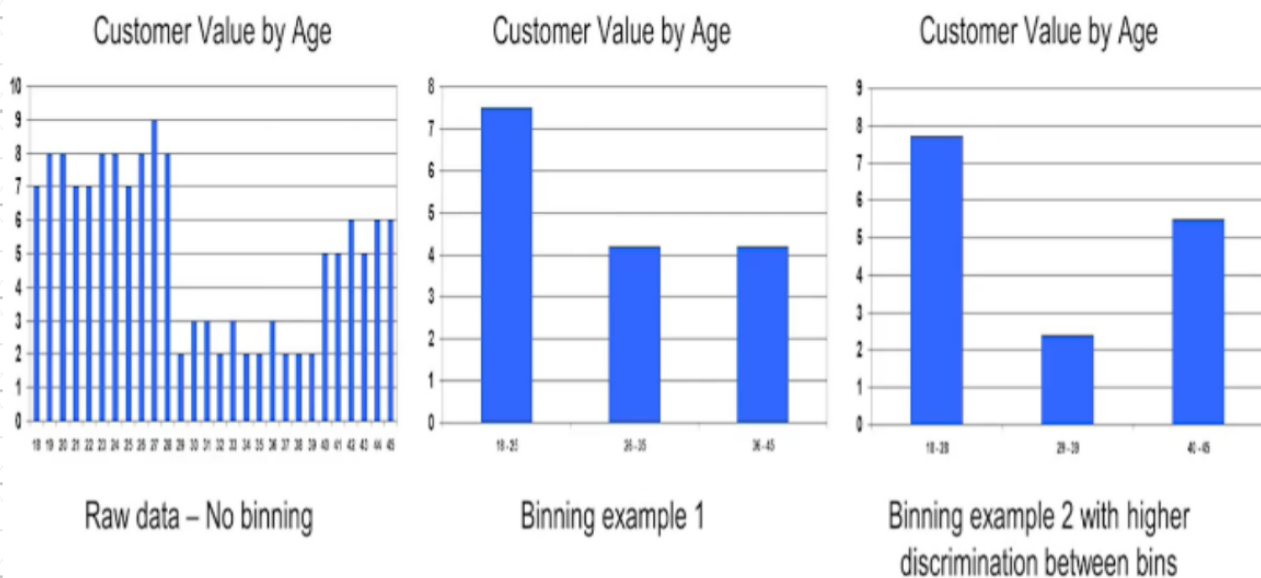
equal frequency binning

- basically, the bins have equal frequency

equal width binning

- bins have equal widths with the cut-off points for each calculated based on the maximum and minimum values of the attribute and the number of bins being created

Continuous variable binning – variable "AGE"



Binning helps to improve model performance, and it does this by capturing non-linear behaviour of continuous variables, it minimizes the impact of outliers, and removes noise from large numbers of distinct values.

Grouped values are easier to display and understand

Improves model build speed - predictive algorithms build faster as the number of distinct values decreases

Merge tables -

Merge tables

A	B
1	A_NUMBER
2	7809702612
3	6139214653
4	7809538328
5	7783183499
6	7788829500
7	6132919446

Table 1: A_NUMBER_FACT

A	B
1	CUSTOMER_ID A_NUMBER
2	1000172 2042930441
3	1000198 2502048322
4	1000210 2502264353
5	1000213 2502280241
6	1000258 2503672523
7	1000260 2503889993

Table 2: CUSTOMER_ID_LOOKUP

A	B	C	D	E	F	G	H	
1	CUSTOMER_ID	GENDER	AGE	ZIP_CODE	DISTRIBUTION_CHANNEL_ID	DEVICE_BRAND_NAME	DEVICE_MODEL_NAME	TENURE_MTHS
2	1000272	M	49	85394	RC11001	Samsung	Galaxy S7 Edge	12
3	1000286	F	36	30832	PH00001	Apple	iPhone 7	12
4	1000210	F	24	11208	PD50001	Samsung	Galaxy S7	12
5	1000213	F	33	10625	PD50001	Samsung	Galaxy S7 Edge	12
6	1000258	F	35	30009	PD50001	Huawei	Honor 8	12
7	1000260	F	36	48643	PH00001	OnePlus	3T	12
8	1000281	M	48	2135	WNV0001	Huawei	Honor 8	12

Table 3: CUSTOMER

A	B	C	D	E	F	G	H	I	
1	A_NUMBER	CUSTOMER_ID	GENDER	AGE	ZIP_CODE	DISTRIBUTION_CHANNEL_ID	DEVICE_BRAND_NAME	DEVICE_MODEL_NAME	TENURE_MTHS
2	77809702612	1000172	M	49	85394	RC11001	OnePlus	3T	2
3	6139214653	1000198	M	62	49509	ALCO001	Google	Pixel	5
4	7809538328	1000210	F	53	11203	PD50001	Apple	iPhone 7	5
5	7783183499	1000213	M	35	91105	WNV0001	Google	Pixel XL	10
6	7809538328	1000215	F	48	30546	SPR0001	Google	Pixel	1
7	6132919446	1000260	F	55	90656	WLV0001	Apple	iPhone 7	9
8	7809538328	1000177	M	21	90805	PD50001	Apple	iPhone 7	1
9	6132919446	1000128	F	29	49623	DIAC001	Apple	iPhone 7	7
10	6132919446	1000129	F	57	23862	PH00001	Apple	iPhone 7	7
11	6132919446	1000120	M	32	52847	PD50001	Apple	iPhone 7	8
12	6132919446	1000121	M	22	44887	WAP0001	Apple	iPhone 7	6

Merged Table

To learn more about merging tables, see:

<https://365datascience.com/what-are-joins/>

Left-outer Join.

A - NUMBER - unique line number associated to each account entity
 merged
 CUSTOMER_ID_LOOKUP - associates with the unique line number.
 merge key

Customer_ID comes from the Customer_ID in table 3

Data aggregation

- When you know that there can be multiple entries for one individual in a transaction table, for example, you have to compute an aggregate to avoid creating duplicates in your dataset
- For example, when you are counting the number of purchases in a store - a customer can make several purchases within a month - so there will be multiple entries for that customer

Create variables such as the number of purchases per week or per month, total amount spent on purchases per week, per month or per quarter, date of last and first purchase.

is the process where raw data is gathered and expressed

There are a range of aggregation functions you can consider, for example:

Functions	Description	Returned Values
<i>Count</i>	computes the number of occurrences	number of occurrences
<i>Sum</i>	compute the sum	sum
<i>Average</i>	compute the mean	mean
<i>Min</i>	identifies the minimum value	minimum value
<i>Max</i>	identifies the maximum value	maximum value
<i>Exists</i>	checks if at least one event exists for the current reference	0 if no event has been found 1 if at least one event has been found
<i>NotExists</i>	checks if no event exists for the current reference	0 if at least one event has been found 1 if no event has been found
<i>First</i>	identifies the first occurrence note that this function needs a date column	value of the first chronological occurrence for the current reference
<i>Last</i>	identifies the last occurrence note that this function needs a date column	value of the last chronological occurrence for the current reference

Data Encoding

- Data encoding is an essential part of the data preparation process.
- The encoding process prepares missing values in the data, deals with outliers, and creates data bins or bands to transform raw data into a "machineable" source of information.

Nominal variable -

• A nominal variable is a discrete (categorical), qualitative variable that characterizes, describes, or names an element of a population.

Examples:

- Hair Color
- Make of car
- Gender
- Postal code
- Residence city

the order of the categories
does not matter

Ordinal value

• An ordinal variable is a discrete (categorical), qualitative variable that has order

• Examples:

- Gold, Silver, bronze

- Satisfaction level

- Pain level

The order of the categories DOES matter.

Continuous Variable

• A continuous variable is a quantitative variable.

• It is a real number that can take any value (with fractions/decimal places) between two specific numbers.

• It accommodates all basic math operations

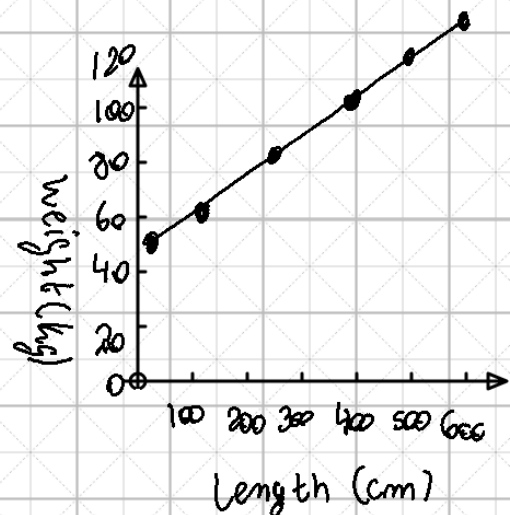
Examples:

• Income (\$)

• Age (years)

• Running time (minutes)

• Bank account balance (\$)



Missing values

- A missing value is an empty cell in your dataset
- Missing values in a dataset can be due to error or simply not available
- They can be removed from the dataset, estimated or kept.
- The analysis could also be stopped so that further investigation of the missing values can be undertaken

Outliers

- For a continuous variable - An outlier is a single or low frequency occurrence of the value of a variable.
- For a categorical variable (nominal or ordinal) - an outlier is a single or very low frequency occurrence of a category of a variable

Ex: ○ ○ ○
○ △ ○
○ ○ ○

Feature Selection

- IS the process of selecting a subset of relevant explanatory variables or predictors for use in data science model construction.
- it is also known as variable selection, attribute selection, or variable subset selection.
- Often, data contains many features that are either redundant or irrelevant, and can be removed without incurring much loss of information.
- Remember that domain knowledge can be the best selection criterion

Traditional approaches to data selection

- traditional approaches can be very time consuming, especially when there are 1000's of variables to analyse
- the most popular form of feature selection is stepwise regression
this algorithm adds the best feature (or deletes the worst feature) in a series of iterative steps. The main control issue is deciding when to stop the algorithm
- Other automated selection processes are back ward elimination and forward selection.

Backward Elimination

1. Backward elimination starts with all candidate features
2. Test the deletion of each feature using the chosen model comparison criterion, deleting the feature (if any) that improves the model the most by being deleted
3. Repeat this process until no further improvement is possible

Forward selection

1. Forward selection starts with no features in the model
2. Test the addition of each feature using the chosen model comparison criterion
3. Add the feature (if any) that improves the model the most
4. Repeat this process until no other feature additions improve the model



Stepwise Regression

- This is a combination of backward elimination and forward selection
- At each stage in the process, after a new variable is added, a test is made to check if some variables can be deleted without appreciably increasing the error
- The procedure terminates when the measure is maximized, or when the available improvement falls below some critical value.
- One of the main issues with stepwise regression is that it is prone to overfitting the data.

Modern approaches to variable selection

