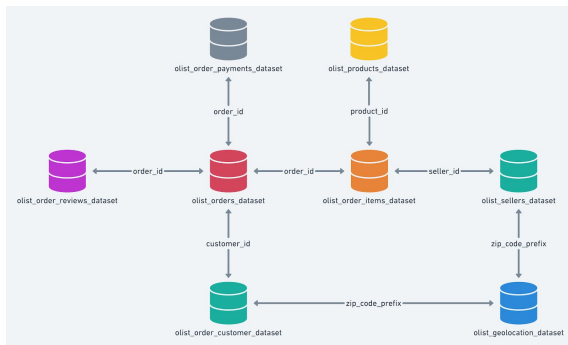


Data Science Project - Flow graph

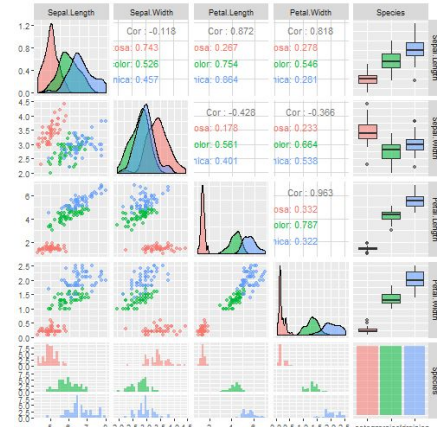
1. Data preparation



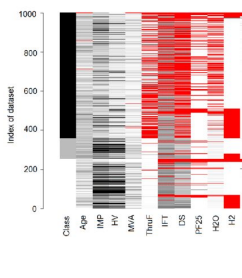
2. Flat-file generation Data enrichment and transformation

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	probeset	ITCC0000	ITCC0001	ITCC0002	ITCC0004	ITCC0007	ITCC0008	ITCC0009	ITCC0011	ITCC0012	ITCC0013	ITCC0014	ITCC0015	ITCC0016	ITCC0017	ITCC0018
2	1007_s_at	430.8	226.1	130.6	75	54.9	195.5	124.8	221.5	82.1	57.5	67.1	174	55.8	143.6	110.2
3	1093_at	79.5	95.7	178.8	185.8	144.7	111.1	150.2	157.5	118.4	140.7	198.1	123.7	254.8	130.4	161.5
4	117_at	65.9	14.9	1.8	20.7	4.7	35.2	1	25.5	19.3	17.1	23.7	39.7	4.1	29.8	64.1
5	121_at	203.3	79.4	74.1	60.3	79.7	54.5	60.1	113.3	95.6	85.8	80.1	81.9	72.7	91.9	112.1
6	1255_at	5	129.8	177.8	291.1	192.4	20	255.4	5.6	54.9	56.6	63.8	81.2	300	41.1	40.9
7	1294_at	114.6	57.1	40.1	37.1	44.2	50.9	40.4	28.3	56.3	51.6	38.6	127.8	44	79.4	83.6
8	1316_at	113	92.7	115.1	106.5	63.5	103.6	121.6	97.8	26.1	33.2	27	33.5	24	17.3	52.9
9	1320_at	12.1	17.9	1.7	1	2.2	11.2	17.2	16.1	33.2	58.1	27.4	14.3	46.6	37.9	40.4
10	1405_at	312.8	235.4	26.1	89.8	184.4	223.2	81	15.7	64.1	29	13	406.7	24.9	140.7	233.3
11	1431_at	18.8	13.7	16.2	8.6	8.2	66	10.3	15.1	33	14.2	79.8	12.5	15.1	7.4	20.9
12	1438_at	2.4	51.9	5.2	9.4	22	6.2	8.3	38.6	69.7	193.6	31.5	107.6	59.1	100.6	120.6
13	1487_at	94	88.9	96	52.1	96.3	67.8	32	91.4	63	99.1	160.7	106.4	54.8	84.4	101.6
14	1494_at	30	18.8	33	49	34	33.5	44.9	35.3	26.1	22.9	19.6	18.7	19.8	17	29.1
15	152256_at	33.4	55.5	57.9	53.1	52	60.2	63.7	117	182.2	79.1	185.9	187.1	137.1	143.8	251.8
16	152257_at	112.2	8.3	5.3	20.3	7	20.7	35	22.7	22.6	12.2	13.5	34.9	27.3	18.7	16.5
17	152258_at	11.2	8.3	5.3	20.3	7	20.7	35	22.7	22.6	12.2	13.5	34.9	27.3	18.7	16.5
18	152259_at	11.2	8.3	5.3	20.3	7	20.7	35	22.7	22.6	12.2	13.5	34.9	27.3	18.7	16.5

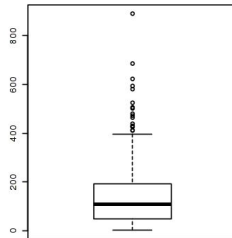
3. Exploratory Data Analysis (EDA)



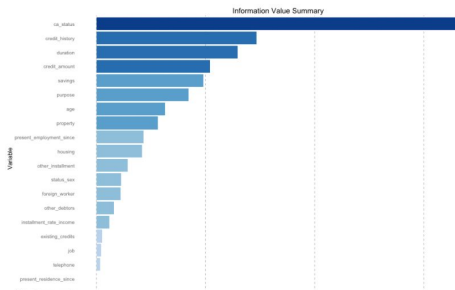
5. Missing values Imputation



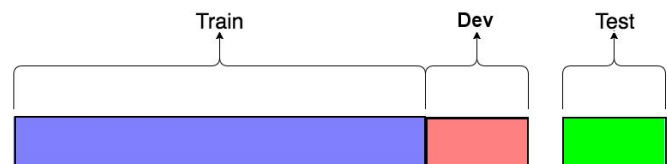
4. Outlier detection and treatment



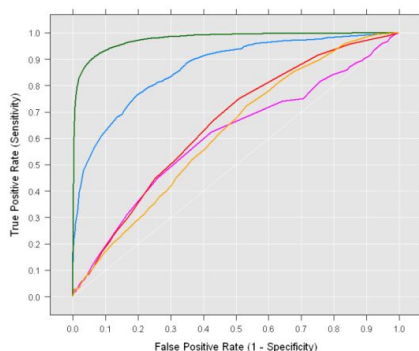
6. Variable Selection



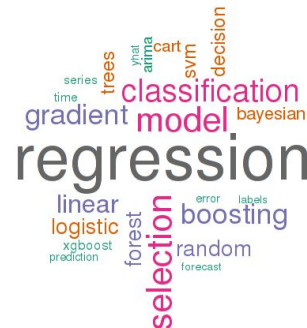
7. Train-Dev-Test preparation



9. Model Fine-tuning



8. Model Selection



Data Science Project - Flow graph

I. Data enrichment and transformation

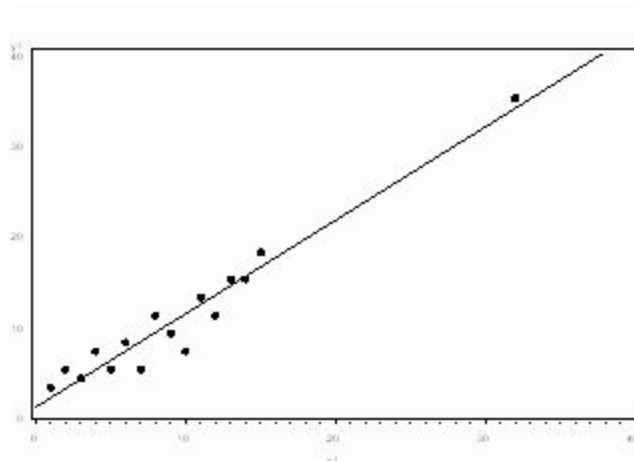
1. Combination of two or more variables: sum/difference/multiplication, division
2. Variable modification: Polynomial, Logarithm, Square root, exponential
3. Enrich with cluster analysis
4. Enrich with external data

II. Exploratory Data Analysis (EDA)

1. Statistical Analysis:
 - Mean and Standard Deviation
 - Median and IQR (25%-75%)
 - Minimum, Maximum, Count
2. Table One
3. Graphics:
 - Categorical variables: barplot
 - Numeric Variables: histogram (distplot)
 - Numeric vs Categorical: boxplot
 - Numeric vs Numeric: scatter plot
 - Pairs

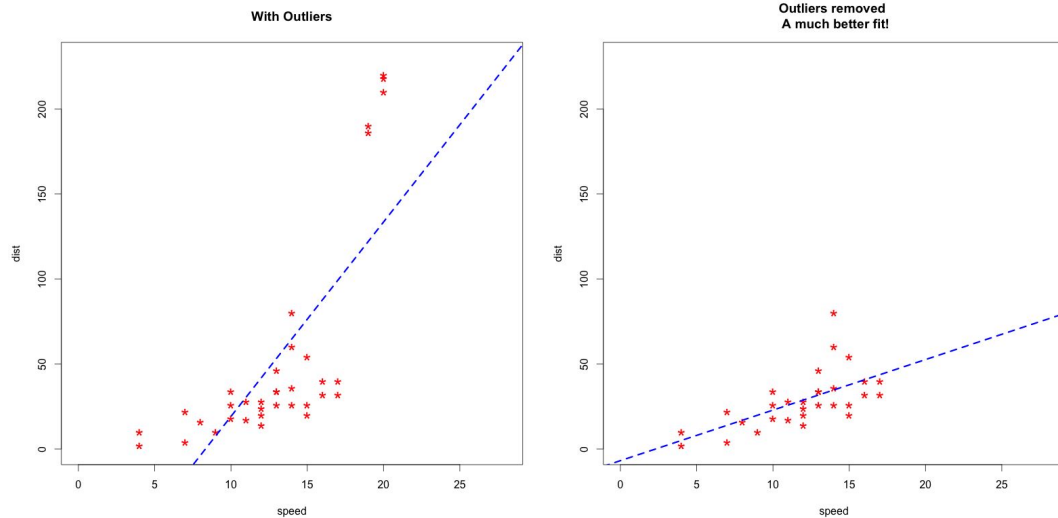
III. Data cleansing

1. Outlier detection - univariate:
 - Normal distribution: Standard Deviation (z-scores)
 - Non-normal distribution: IQR
2. Outlier detection - multivariate:
 - DBSCAN
3. Outlier treatment:
 - Obviously incorrect: convert value to NA
 - Example: age 450, body temperature 2°C, size -1.2
 - Does change the distribution but don't change relationships: convert value to NA and **report**.

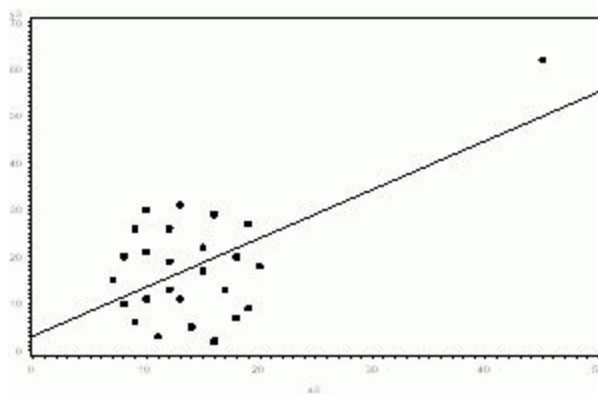


Data Science Project - Flow graph

- Does change both the distribution and the relationships: make analysis of the effect or influence in the presence and absence of the outliers. Report the findings and explain your decision of leaving or substituting the outliers with NAs.



- Makes a correlation where there is no correlation: convert value to NA



- Multivariate outliers: DBSCAN returns a list of rows containing mv-outliers. Add an variable indicator for outliers (0/1). Check if this indicator has a relation with the dependent variable (y or outcome). You can use correlation analysis or check the distribution of y between outlier and non outlier cases. You can drop it if there is no correlation or if the distribution is similar. If not, leave the outlier indicator in the dataset.

4. Missing values: Determinate the missingness mechanism.

- Missing Completely at Random (MCAR)*: missing values are generated randomly and there is no possible way to explain why those values are missing.
- Missing at Random (MAR)*: There is a variable that influence in the generation of missing values, but within subgroups of such a variable the missingness is completely at random. For example, women tend less to give write their age or weight than men. However, if we check the mechanism by separating women from men, there generative mechanism is MCAR.

Data Science Project - Flow graph

- *Missing not at Random (MNAR)*: There is a clear cause for the generation of missing values. An example is the visits to gynecologist (women's doctors). No men will have a visit in the dataset. Another example is a question made only to individuals with age 60 and more. No individual with 55 years old will have the question fulfilled.

5. Missingness treatment: If the missing mechanism is MCAR or MAR, we can use imputation techniques, otherwise we have to decide if dropping the rows or dropping the columns.

- Detecting Missingness Mechanism
 - For each variable (feature) having missing values, generate a dummy variable indicating for each value if it is full (0) or is missing (1).
 - Using only the variables which have no missing data, run a logistic regression (in R: glm; in Python: OLS) where the X is the dataset with only those variables and the outcome (y) is the missing dummy variable.
 - If all the p-values for the coefficients (betas) are non-significant ($p \geq 0.05$), we can assume that the mechanism is MCAR.
 - If there is one variable that was significant ($p\text{-value} < 0.05$) then we have to inspect the relationship between this variable and the missing indicator. We can check it using a boxplot and determinate a value which can divide the dataset into two subsets. Each subset will be analyzed separately as we did before. If in both groups we can demonstrate that the mechanism is MCAR, we say that for this variable original mechanism is MAR.
 - In the case we can't demonstrate MCAR or MAR, we say that the mechanism is MNAR.
- Imputation techniques:
 - Statistical imputation: using mean, median or mode (depending the data scale). **Not recommended!**
 - Model based imputation: we can use a predictive model to impute the missing values
 - KNN: this algorithm is very popular because it easily and fast impute data. We divide the dataset into complete data (including the variable that we want to impute) and train the model. Then we use the predict function to complete the data in the incomplete subset.
 - Random Forest (use as we explained with KNN)
 - Decision trees
 - Multiple imputation: The most common method is the Multiple Imputation by Chained Equations (MICE) which uses all the dataset for imputation. It begins with the variable with less missing values and impute it, Then uses the second with less missings and impute it, and so on. MICE generate a number of imputed datasets (default or most common is 5), that has different values for each imputed variables. The mean, median, standard deviation and IRQ for each variable on each imputed dataset are constant. Any further analysis must be made using all the imputed datasets (5 different analysis each time, the final result is the calculated as the mean of the outcome (y) of all the models.
 - When having lot of data (>20,000 rows), a simple model is recommended. With smaller datasets, multiple imputation is recommended.

Data Science Project - Flow graph

IV. Feature Selection

1. Univariate Analysis: Each variable on the dataset is analyzed by comparing its relationship with the dependent variable (outcome or y). The analysis depends on the independent (x) and dependent (y) data types:

- If x is nominal and y is nominal or ordinal: use chi-square
- If x is nominal and y is continuous: use spearman correlation
- If x is continuous and y is binomial (0/1, yes/no, true/false): use independent t-test
- If x is continuous and y is multinomial (more than two categories) or ordinal: use anova
- If both x and y are continuous, use pearson or spearman correlation

In case that y is nominal or ordinal with less than 6 categories, you can use tableone (R and Python) setting the y as the grouping category. Tableone makes the whole analysis automatically. You will want to include in your analysis those variables with a significant p-value (< 0.05).

2. Multivariate analysis: Using the whole dataset and running predictive models that are able to return a list recommended features by defining their influence in the model. For this step you don't need to have to partition the dataset into train, dev and test.

- **LASSO (L1 penalization)**: Lasso penalizes the growing of the values of the coefficients (betas) and can get them down to zero. Those variables with a zeroed coefficients are excluded from the analysis, giving to this algorithm the capacity of feature selection.
- **Random Forest**: this algorithm is able to generate a list of the most influential variables and their lift.
- **Gradient Descent**: same as Random forest.
- **Support Vector Machine (SVM)**: It can be used for feature selection by using L1 penalization.
- **Principal Component Analysis (PCA)**: can be used for the selection of variables comparing the variables with the highest correlation with the principal components that catch the much of the variability (80% cumulative variance).

3. Selection based on voting: using many of the techniques (univariate and multivariate), we make a table with all the variables on the dataset and indicate the recommended variables for each technique, then we select a threshold for the total votings and on this basis we select the variables that will be used to train our models.

Variable	Univarable	Lasso	RandomForest	GradientBoost	SVM	Sum
Fz_N100	1	1	1	1	1	5
FCz_N100	1	0	1	0	1	3
Cz_N100	1	1	1	0	1	4
FC4_N100	1	1	1	1	1	5
FC3_B0	1	0	0	0	0	1
FC4_B0	1	0	0	0	0	1
C3_B0	0	0	0	0	0	0
C4_B0	0	0	0	0	0	0
FC4_B1	0	1	0	0	1	2
FC3_B1	1	1	0	1	1	4