# Data Science / Analytics network meeting - 13/3-2024

# "Good decisions require good data"

Pia Skare Rønnevik, Customer Success Manager, FANS
pia.roennevik@sas.com

§sas

# Copilot – What to think of in data preparation for analytics?

1. Data Collection and Sourcing:
- Gather relevant data from various sources. Ensure data quality by addressing missing values, duplicates, and inconsistencies.
- Verify the data's accuracy and reliability.

2. Data Cleaning and Preprocessing:
- Handle missing data: Impute missing values using techniques like mean, median, or regression.
- Remove duplicates and outliers.
- Standardize or normalize features to a common scale.
- Address data inconsistencies (e.g., typos, variations in formatting).

3. Feature Engineering:
- Create new features from existing ones (e.g., aggregations, transformations).
- Select relevant features based on domain knowledge and statistical significance.

4. Exploratory Data Analysis (EDA):
- Visualize data distributions, correlations, and patterns.
- Identify relationships between variables.
- Detect anomalies or unexpected behavior.

5. Data Transformation:
- Apply necessary transformations (e.g., log, square root) to achieve normality.
- Encode categorical variables (one-hot encoding, label encoding).
- Consider dimensionality reduction techniques (PCA, t-SNE) if needed.

6. Data Splitting:
- Divide the dataset into training, validation, and test sets.
- Ensure that the split maintains the distribution of target variables.

7. Handling Imbalanced Data:
- Address class imbalance using techniques like oversampling, undersampling, or synthetic data generation.

8. Feature Selection:
- Use statistical tests (e.g., chi-squared, ANOVA) or feature importance scores (from models) to select relevant features.

9. Data Validation and Quality Assurance:
- Validate data against business rules and domain knowledge.
- Monitor data quality over time.

10. Document the Process:
- Maintain clear documentation of data preparation steps.
- Record assumptions, transformations, and any decisions made during the process

§sas

# Copilot – is it preferable to use categorical variables instead of continuous

1. Continuous Variables:

- Continuous data represents measurements that can take any value within a range (e.g., height, temperature, weight).

- Advantages:
- Provides more information due to the fine-grained nature of measurements.
- Allows for precise modeling of relationships.
- Well-suited for linear regression and other models that assume linearity.

- Disadvantages:
- Assumes a linear relationship, which may not hold in all cases.
- Sensitive to outliers.
- Requires more data points for accurate estimation.

- When to use:
- Use when the relationship with the response variable is expected to be linear.
- Suitable for regression models, correlation analysis, and ANOVA.

2. Categorical Variables:

- Categorical data represents distinct categories or groups (e.g., gender, color, program type).

- Advantages:
- Simplifies complex relationships by grouping data into meaningful categories.
- Handles non-linear relationships more effectively.
- Robust to outliers.

- Disadvantages:
- May lose information due to grouping.
- Requires encoding (e.g., one-hot encoding) for modeling.
- Limited to the defined categories.

- When to use:
- Use when the relationship is non-linear or when there are clear groupings.
- Useful for decision trees, random forests, and logistic regression.
- Especially valuable when capturing interactions or nonlinear effects.

Considerations:

- Interaction Effects: Categorical variables allow modeling of interactions (e.g., gender and age).

- Data Context: Consider the context and domain knowledge. Some variables may naturally be categorical (e.g., program type).

- Data Exploration: Explore relationships visually (scatter plots, histograms) to decide how to treat variables.

§sas

# Copilot - is it best to have continuous or binary target?

1. Continuous Target (Regression):

- Definition: A continuous target variable can take any numeric value within a range. It has a natural order and can be measured with precision.

- Use Cases:
  - When you want to predict a numerical value, such as predicting house prices, temperature, or stock prices.
  - In scenarios where the relationship between predictors and the response is linear or can be approximated by linear models.
  - When you need to estimate the mean change in the dependent variable given changes in independent variables.

- Advantages:
  - Provides more information due to the greater number of possible values.
  - Allows for a wider variety of analyses.
  - Can handle non-linear relationships using techniques like polynomial regression1.

2. Binary Target (Classification):

- Definition: A binary target variable has only two possible outcomes (e.g., yes/no, true/false, 1/0).

- Use Cases:
  - When you want to perform classification tasks, such as spam detection, disease diagnosis, or customer churn prediction.
  - In scenarios where you need to assign instances to predefined classes.
  - When the outcome of interest is binary (e.g., whether a customer will buy a product or not).

- Advantages:
  - Simplicity: Binary targets are straightforward to interpret.
  - Commonly used in practical applications.
  - Suitable for algorithms like logistic regression, support vector machines, and decision trees

3. Hybrid Approaches:

Sometimes, you might encounter situations where a continuous target can be transformed into a binary one. For instance:
- Threshold-based classification: Convert a continuous target (e.g., probability) into binary outcomes by setting a threshold (e.g., 0.5).
- Risk prediction: Predicting the likelihood of an event (e.g., defaulting on a loan) using continuous scores and then categorizing into low/high risk.
- Ordinal regression: When the target has ordered categories (e.g., low/medium/high), you can use ordinal regression.

Consideration: If you can collect continuous data, it's often better to do so because it provides more information.

§sas

# VIMO!

- Valid
- Invalid
- Missing
- Outlier

§sas

# Missing values and outliers

Missing values:

- is due to information loss, dropouts and none-responses.
- give smaller sample size, could compromises the study and entail bias.
- could be ignored or imputed based on the present values.

Outliers:

- are extreme values that's outside the overall pattern of a distribution of variables.
- arises from respons and data errors.
- Could lead to under- or over-estimated values
- If you are sure the outliers are wrong, then ignore them or impute.

§sas

# How to check VIMO!

## Information catalog on Viya



## Visual Analytics

- Could check for missing values by just looking on histogram for the measure variables, and bar chart for the categorical variables. Can easily impute missing variables with calculated item.

- Could check for outliers by using box-plot

# The data - medical charges of a specific patient

- **age:** age of primary beneficiary

- **sex:** insurance contractor gender, female, male

- **bmi:** body mass index, providing an understanding of body, weights that are relatively high or low relative to height objective index of body weight (kg / m ^ 2) using the ratio of height to weight, ideally 18.5 to 24.9

- **children:** number of children covered by health insurance / Number of dependents

- **smoker:** smoking

- **region:** the beneficiary's residential area in the US, northeast, southeast, southwest, northwest

- **charges (insurance price):** Individual medical costs billed by health insurance

https://www.kaggle.com/datasets/nanditapore/medical-cost-dataset

§sas

# The Analytics Life Cycle



QUESTION

What is the medical expenses of specific patient?

**DataOps**

**Develop Models**

**Deploy Models**

Value

Given certain factors, we could say what the medical expenses is for a patient

- Access
- Prepare
- Visualize
- Govern

- Build
- Train
- Evaluate
- Select

- Validate
- Deploy
- Monitor
- Embed

§sas

# Data preparation in VA –
## How does the target look like?



Frequency of charges



Frequency of charges_log

# Data preparation in VA –
## How does the measure variables look like?

# Data preparation in VA –
## How does the categorical variables look like?
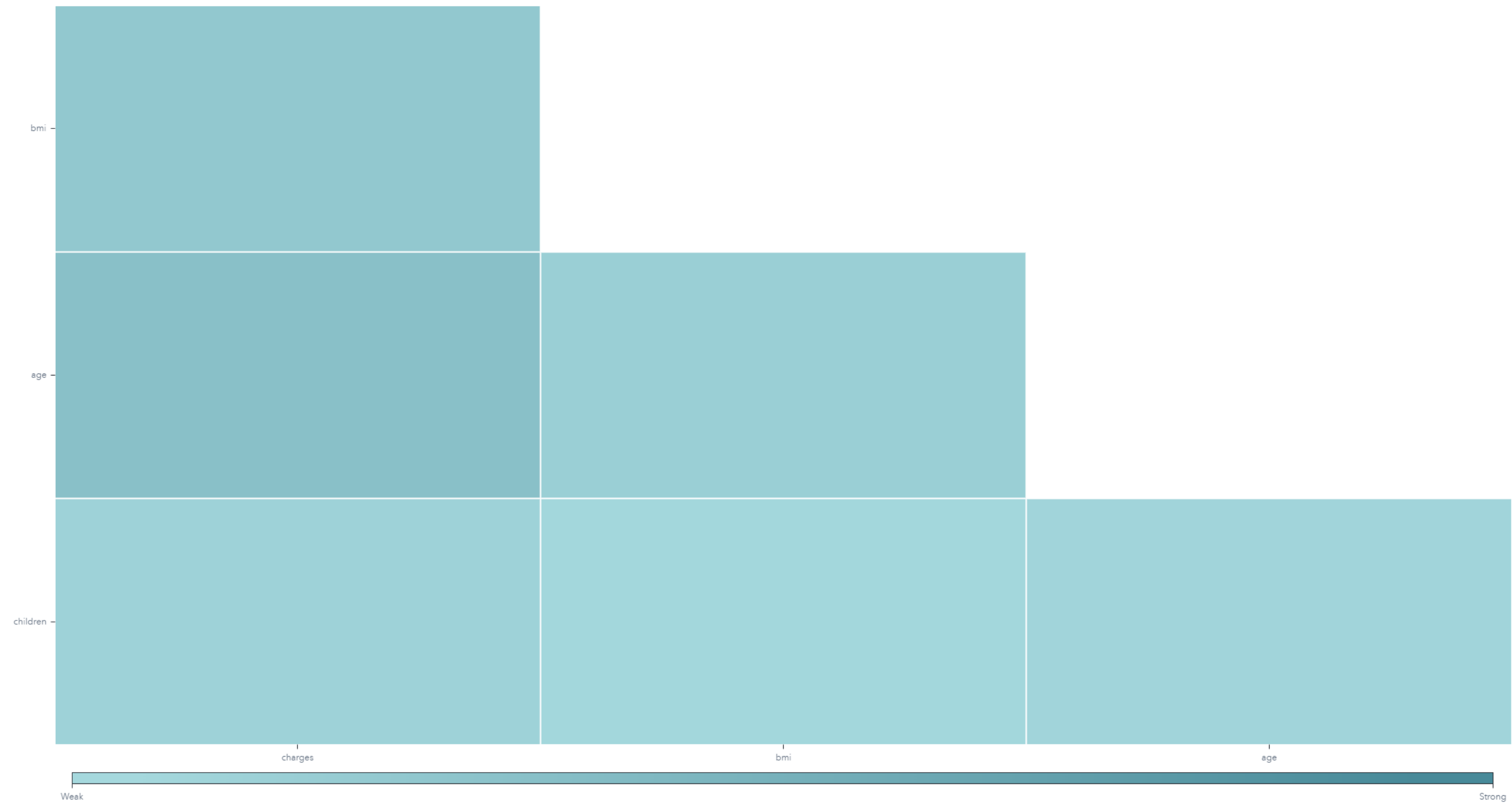
# Data preparation in VA –
## The correlation between the measure variables?



Correlation of Selected Measures

# Data preparation in VA –
## Re-define the measure variables as categorical

# References

- **Process for Data Quality Assurance at Manitoba Centre for Health Policy (MCHP):** [https://umanitoba.ca/manitoba-centre-for-health-policy/sites/manitoba-centre-for-health-policy/files/2021-11/data-quality-framework-document.pdf](https://umanitoba.ca/manitoba-centre-for-health-policy/sites/manitoba-centre-for-health-policy/files/2021-11/data-quality-framework-document.pdf)

- **Statistical data preparation: management of missing values and outliers:** [https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5548942/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5548942/)

§sas