# Data Exploration and Preprocessing Report

Team ID: SWTID1720067156

July 10, 2024

## Data Collection and Preprocessing Phase

**Team ID:** SWTID1720067156
**Project Title:** Lymphography Classification Tool
**Maximum Marks:** 6 Marks

## Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

## Section Description

| Section | Description |
|---|---|
| Data Overview | Dimension: 614 rows × 13 columns<br>Descriptive statistics:<br><br>```<br>Attribute information:<br>--- NOTE: All attribute values in the database have been entered as<br>            numeric values corresponding to their index in the list<br>            of attribute values for that attribute domain as given below.<br>1. class: normal find, metastases, malign lymph, fibrosis<br>2. lymphatics: normal, arched, deformed, displaced<br>3. block of affere: no, yes<br>4. bl. of lymph. c: no, yes<br>5. bl. of lymph. s: no, yes<br>6. by pass: no, yes<br>7. extravasates: no, yes<br>8. regeneration of: no, yes<br>9. early uptake in: no, yes<br>10. lym.nodes dimin: 0-3<br>11. lym.nodes enlar: 1-4<br>12. changes in lym.: bean, oval, round<br>13. defect in node: no, lacunar, lac. marginal, lac. central<br>14. changes in node: no, lacunar, lac. margin, lac. central<br>15. changes in stru: no, grainy, drop-like, coarse, diluted, reticular,<br>                     stripped, faint,<br>16. special forms: no, chalices, vesicles<br>17. dislocation of: no, yes<br>18. exclusion of no: no, yes<br>19. no. of nodes in: 0-9, 10-19, 20-29, 30-39, 40-49, 50-59, 60-69, >=70<br><br>Missing Attribute Values: None<br><br>Class Distribution:<br>Class:          Number of Instances:<br>normal find:  2<br>metastases:   81<br>malign lymph: 61<br>fibrosis:      4<br>``` |

| | |
|---|---|
| Univariate Analysis |  |
| Bivariate Analysis |  |
| Multivariate Analysis |  |
| Data Preprocessing Code Screenshots | |
| Loading Data | ```python
data_file_path = './data/lymphography.data'
data = pd.read_csv(data_file_path, header=None)

print("Given data")
print(data.head())
```<br> |
| Handling Missing Data | No missing attributes |
| Data Transformation | ```python
column_names = [
    "class", "lymphatics", "block_of_affere", "bl_of_lymph_c", "bl_of_lymph_s", "by_pass", "extravasates",
    "regeneration_of", "early_uptake_in", "lym_nodes_dimin", "lym_nodes_enlar", "changes_in_lym",
    "defect_in_node", "changes_in_node", "changes_in_stru", "special_forms", "dislocation_of",
    "exclusion_of_no", "no_of_nodes_in"
]
data.columns = column_names
``` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | Done |