→ Health Analytics Lab Notes :

- Split 60:40 between Notebook & Report
  - ↳ Non-technical
    - Characteristics
    - Meaning of the report
    - Explain what the output/results means
    - What is valuable within the notebook

↓ Technical
- Dataset :
  → Date & time features
  → Time series data
  → Cross sectional data
- Table of content on jupyter Notebook
- Problem statement (1-2 sentences)
- EDA
  ↳ overview
  Understand → ↳ shape of data
  the distribution → ↳ Data type
  ↳ Descriptive stats
  (df.info & df.des

↳ Visualisations
- Univariate Analysis :
  Histogram or bar graph
  Box plots

- Bivariate Analysis
  Scatter plot for continuous vs. categorical variables

- Correlation Analysis
  Heat Map or pair plot to visualize correlation

. Data Pre-processing :

※ Encoding Categorical Variables

※ Feature Engineering

※ Try Cat-boost in feature selection

Gradient feature also code.

. Modelling :

① Selection of model → select
 appropriate algorithm. For health
analytics = Logistic Regression, Random
forest, Gradient boost. Neural Network

②

③ ROC_AUC plot, PR_AUC plot
  ↳ Good submissions
  ↳ Used for im-balanced data

④ Hyper parameter Tuning

Grid Search CV or Randomized Search
CV to find best hyperparameters
for your model.

② Metrics :
Don't only use one metric to
evaluate performance
· Precision
· Recall
∘ F1 - score
· ROC - AUC
∘ Confusion matrix
· Regression . MSE ; RMSE . MAE,
                R - Squared

· Conclusion

↳ Summarize findings

↳ Discuss implications of
  model results in context of
  health analytics