# Machine Learning: Geotechnical Engineering

## Predicting UCS and RQD using Random Forest

**A Presentation for Operations Technical Support - Ok Tedi Mining Limited**

Presented by Jesse Gabriel

October 27, 2025

# Outline

- What & Why Machine Learning (ML)?
- Machine Learning Models and Workflow
- Raw Data and EDA
- Preprocessing
- ML Model Training & Evaluation Approach
- Results: UCS (ESTUCS)  RQD Predictions
- Discussion: Further Work
- Discussion: Related Work on Application

# What & Why Machine Learning (ML)?

- Machine Learning (ML), a subfield of AI, involves development of algorithms that enable computers to learn and improve task performance from data and experience without explicit programming.
- Types: Supervised, Unsupervised, Reinforcement, Semi-supervised
- Data inputs: Spreadsheet, Image, Text, Software files (e.g. GIS shapefiles), etc.
- ML can handle large datasets and generate useful prediction outputs

$$\boxed{\text{INPUT DATA}} \rightarrow \boxed{\text{ML MODEL}} \rightarrow \boxed{\text{OUTPUT}}$$



Figure: Sample spreadsheet data (geotechnical log)

# Machine Learning Models and Workflow

EDA → PREPROCESSING → MODEL TRAINING → RESULTS → POSTPROCESSING

$$\hat{y} = f_\theta(X) + \varepsilon \qquad \text{with} \qquad X = \{x_1, \ldots, x_p\}$$

$\hat{y}$ : prediction $\quad f_\theta$ : model $\quad \theta$ : parameters $\quad X$ : features $\quad x_i$ : feature $\quad p$ : dimensionality $\quad \varepsilon$ : noise

```python
from sklearn.ensemble import RandomForestRegressor
# Example: X is feature matrix, y target (UCS or RQD)
# X = data[feature_cols].values
# y = data['UCS'].values
model = RandomForestRegressor(n_estimators=200, random_state=42,....)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```
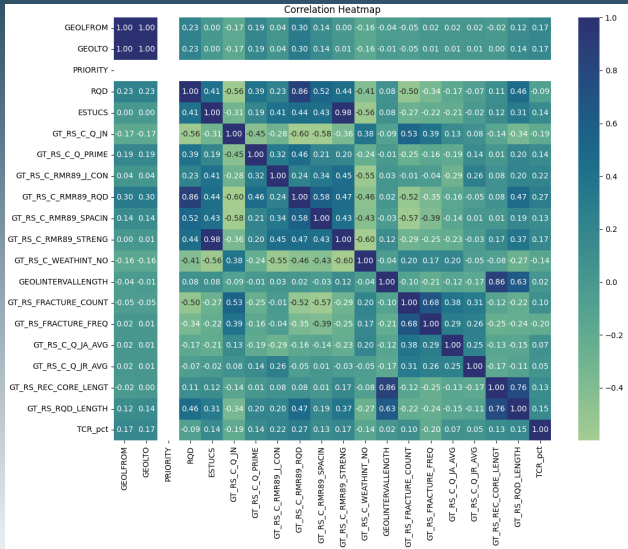
# Raw Data and EDA

- Geotechnical Log data (13236 rows 31 columns)

Table: Data types and missing values summary

| # | Variable | Data Type | Missing | Missing (%) |
|---|---|---|---|---|
| 1 | HOLEID | object | 0 | 0.00 |
| 2 | PROJECTCODE | object | 0 | 0.00 |
| 3 | GEOLFROM | float64 | 0 | 0.00 |
| 4 | GEOLTO | float64 | 0 | 0.00 |
| 5 | PRIORITY | int64 | 0 | 0.00 |
| 6 | GT_RockStrength | object | 1998 | 15.10 |
| 7 | RQD | float64 | 194 | 1.47 |
| 8 | GT_Weathering | object | 14 | 0.11 |
| 9 | ESTUCS | float64 | 0 | 0.00 |
| 10 | GT_RS_FABRIC | object | 5010 | 37.85 |
| 11 | GT_RS_C_Q_JN | float64 | 65 | 0.49 |
| 12 | GT_RS_C_Q_PRIME | float64 | 11 | 0.08 |
| 13 | GT_RS_C_RMR89_DESC | object | 167 | 1.26 |
| 14 | GT_RS_C_RMR89_J_CON | float64 | 89 | 0.67 |
| 15 | GT_RS_C_RMR89_RQD | float64 | 174 | 1.31 |
| 16 | GT_RS_C_RMR89_SPACIN | float64 | 183 | 1.38 |
| 17 | GT_RS_C_RMR89_STRENG | float64 | 71 | 0.54 |
| 18 | GT_RS_C_WEATHINT_NO | float64 | 15 | 0.11 |
| 19 | GEOLINTERVALLENGTH | float64 | 0 | 0.00 |
| 20 | GT_Alteration | object | 5011 | 37.86 |
| 21 | GT_RS_FRACTURE_COUNT | float64 | 2 | 0.02 |
| 22 | GT_RS_FRACTURE_FREQ | float64 | 130 | 0.98 |
| 23 | GT_RS_AVG_DEFECT_RGH | object | 5522 | 41.72 |
| 24 | GT_RS_C_Q_JA_AVG | float64 | 609 | 4.60 |
| 25 | GT_RS_C_Q_JR_AVG | float64 | 634 | 4.79 |
| 26 | GT_RS_REC_CORE_LENGT | float64 | 52 | 0.39 |
| 27 | GT_RS_RQD_LENGTH | float64 | 186 | 1.41 |
| 28 | GT_RS_LITH_ROCKTYPE | object | 6 | 0.05 |
| 29 | GT_Logger | object | 111 | 0.84 |
| 30 | GT_LoggedDate | object | 79 | 0.60 |
| 31 | TCR_pct | float64 | 12 379 | 93.53 |

# Raw Data and EDA

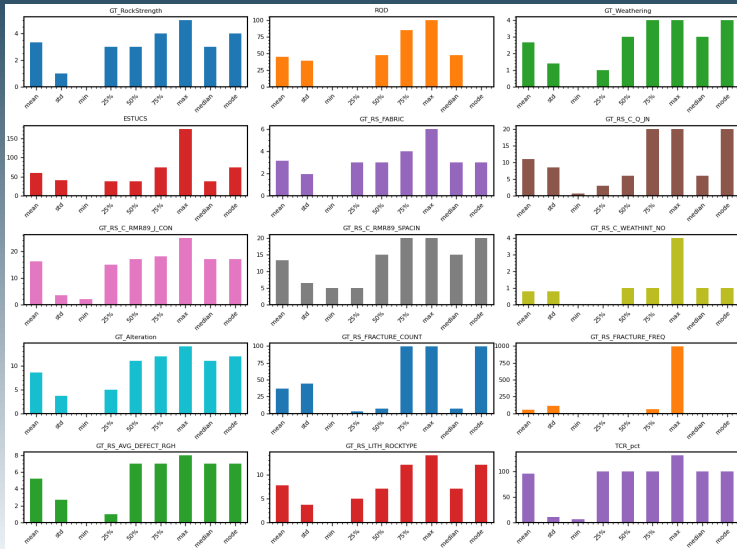- Correlation matrix heatmap

# Preprocessing

- Feature selection
- Handling outliers (e.g. RQD above 100%)
- Handling missing values
- Computing summary statistics
- Encoding of categorical data

Table: Sample encoded data: `GT_RS_AVG_DEFECT_RGH`

| Original Value | Encoded Value | Count |
|---|---|---|
| UR | 7 | 4876 |
| PR | 1 | 1879 |
| PS | 2 | 400 |
| US | 8 | 264 |
| SI | 3 | 58 |
| UP | 6 | 51 |
| PP | 0 | 40 |
| SP | 4 | 19 |
| SS | 5 | 10 |

# Preprocessing

- Summary statistics of selected & encoded data (features & targets)

# ML Model Training & Evaluation Approach

- **Random Forest Regressor**
- **Hyperparameter Optimization with Optuna**: We define an objective function that tests different hyperparameter combinations for a Random Forest using cross-validated $R^2$ as the evaluation metric.
- **Parameter Search Space**: Key parameters tuned include n_estimators, max_depth, min_samples_split, min_samples_leaf, and max_features, while bootstrap and random_state are fixed.
- **Automated Study Execution**: Optuna runs 50 trials to maximize the model's $R^2$ score, automatically identifying the best-performing hyperparameters.
- **Final Model Training & Evaluation**: The Random Forest is retrained on the full training data using the best parameters; predictions on the test set are evaluated with $R^2$ and RMSE metrics.

# ML Model Training & Evaluation Approach

- $R^2$ measures the proportion of variance in the target explained by the model. $RMSE$ quantifies the average magnitude of prediction errors.

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2}$$

$y_i$ : actual value   $\hat{y}_i$ : predicted value   $\bar{y}$ : mean of actual values   $n$ : number of samples

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}$$

$y_i$ : actual value   $\hat{y}_i$ : predicted value   $n$ : number of samples

# ML Model Training & Evaluation Approach

```python
def objective(trial):
    params = {
        'n_estimators': trial.suggest_int('n_estimators', 200, 800),
        'max_depth': trial.suggest_int('max_depth', 5, 30),
        'min_samples_split': trial.suggest_int('min_samples_split', 5, 20),
        'min_samples_leaf': trial.suggest_int('min_samples_leaf', 5, 15),
        'max_features': trial.suggest_categorical('max_features', ['sqrt', 'log2']),
        'bootstrap': True,
        'random_state': 42
    }
    rf = RandomForestRegressor(**params)
    score = cross_val_score(rf, X_train, y_train, cv=5, scoring='r2').mean()
    return score
study = optuna.create_study(direction="maximize")
study.optimize(objective, n_trials=50, show_progress_bar=True)
best_params = study.best_params
best_params_dict[target] = best_params
```

# Results: UCS (ESTUCS) & RQD Predictions

- Both models show solid performance on the test set.
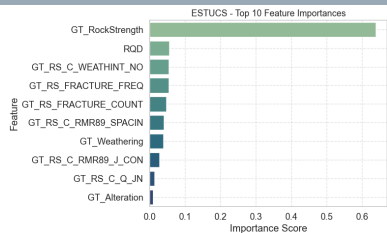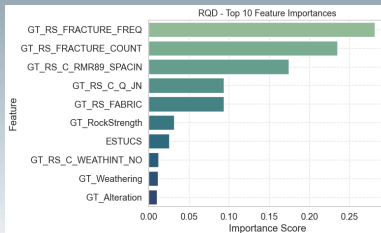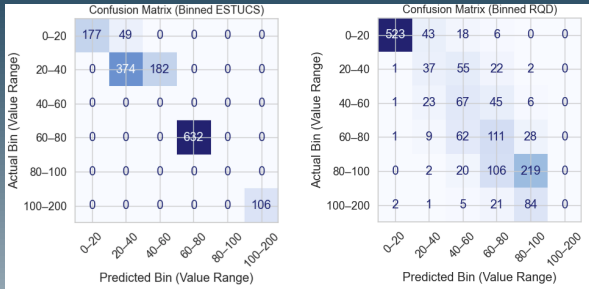- Predicted values of UCS and RQD exported as csv/excel files

Table: Random Forest Model Summary

| Target | Best Trial | Best Hyperparameters | Train Size | Test Size | $R^2$ (Test) | RMSE (Test) |
|--------|-----------|----------------------|------------|-----------|--------------|-------------|
| RQD | 29 | n_estimators=458, max_depth=26, min_samples_split=7, min_samples_leaf=5, max_features='log2' | 6077 | 1520 | 0.8594 | 14.7149 |
| ESTUCS | 41 | n_estimators=535, max_depth=16, min_samples_split=8, min_samples_leaf=5, max_features='sqrt' | 6077 | 1520 | 0.9818 | 5.3558 |

# Results: UCS (ESTUCS) & RQD Predictions

- Confusion Matrix and Feature Importances

# Discussion: Further Work

- Limitation: Proper selection of features required for this work
- Limitation: Need to handle imbalanced data to improve prediction
- Other derived/measured parameters apart from RQD & UCS can be predicted using appropriate models (FoS, displacement, etc.)
- Trained ML model can be deployed within existing/industry software or isolated web/desktop apps to complement standard workflows

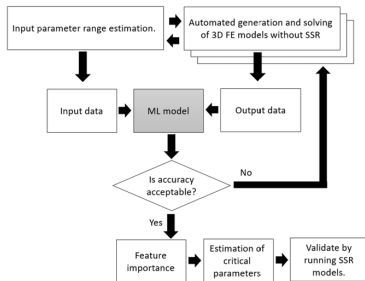| INPUT DATA | PREDICT | PREDICTED OUTPUT |
|---|---|---|

- Related work from Rocscience



Figure 5. Flowchart of proposed back-analysis methodology.

# Acknowledgments: THANK YOU

- Mr. Dauba Dauba (Snr. Geotechnical Engineer), for organizing and facilitating this presentation.
- OTML Operations Technical Support Team for your support and facilitation.
- All participants and attendees, for your engagement and valuable contributions.
- Further questions / discussions: Email: jessegabriel11@gmail.com