

# Improving Prediction of Optimal Tracheal Tube Depth in Pediatric Patients Undergoing Mechanical Ventilation

Data to Paper

January 7, 2024

## Abstract

Tracheal tube misplacement is a significant concern in pediatric patients undergoing mechanical ventilation, necessitating accurate determination of the optimal tracheal tube depth (OTTD) to prevent complications. Traditional methods, such as chest X-rays, are time-consuming and entail radiation exposure. Here, we aimed to develop and compare predictive models for OTTD using non-invasive approaches. We analyzed a dataset of 969 pediatric patients who underwent post-operative mechanical ventilation. Our Random Forest model, incorporating patient characteristics such as sex, age, height, and weight, demonstrated superior precision over the formula-based model commonly used in clinical practice. The Random Forest model exhibited a lower Root Mean Square Error (RMSE) of 1.21 compared to 1.85 for the formula-based model. By accurately estimating OTTD, these non-invasive predictive models offer potential for reducing reliance on chest X-rays and associated risks. However, further refinement and validation are necessary to ensure clinical applicability and safety. Implementation of these models in clinical practice may enhance outcomes in pediatric patients undergoing mechanical ventilation.

## Results

Our analysis was based on a dataset of 969 clinical observations, each representing a pediatric patient who underwent mechanical ventilation post-operatively at Samsung Medical Center. A Random Forest model was developed to predict the optimal tracheal tube depth (OTTD), determined via chest X-rays, using non-invasive parameters including patient's sex, age, height, and weight.

To understand the capabilities of our Random Forest model, we first compared its residuals with those of a common formula-based model, instinctively appealed to clinicians, and thereby quantified its relative predictive accuracy. Indeed, a paired sample t-test was performed on the squared residuals of the two models to identify any significant difference in their predictive performance. Our analysis indicated that the Random Forest model significantly outperformed the formula-based model in predicting OTTD (Table 1). The p-value was lower than the commonly accepted threshold ( $p < 10^{-6}$ ), implying strong evidence against the null hypothesis, i.e., no difference in the performance of the two models.

Table 1: Test of difference in predictive power between formula-based model

	Model	t-statistic	p-value
<b>Model 1</b>	Random Forest	-6.02	$<10^{-6}$
<b>Model 2</b>	Formula-based Model	-6.02	$<10^{-6}$

**Model 1:** Random Forest Model

**Model 2:** Height Formula-based Model

Subsequently, the Root Mean Square Error (RMSE), a popular measure of the concentration of data, was calculated to gain insights into the predictive power of the models in absolute terms. Remarkably, the Random Forest model exhibited a lower RMSE of 1.21 as compared to the formula-based model which exhibited a noticeably higher RMSE value of 1.85 (Table 2). This further underscored the superiority of the Random Forest model in accurately predicting the OTTD.

Table 2: Root mean square error of formula-based model and the Random Forest model

	Model	RMSE
<b>Model 1</b>	Random Forest	1.21
<b>Model 2</b>	Formula-based Model	1.85

**Model 1:** Random Forest Model

**Model 2:** Height Formula-based Model

**RMSE:** Root Mean Square Error of the model predictions

Finally, the optimal parameters for the Random Forest model were obtained using the concept of Grid Search. More specifically, a total of 200 estimators were combined in the final model while the minimum samples

split was determined to be 8. This showed that the model was non-trivial and it learned from a substantial amount of features and residuals to make accurate predictions.

Taking these analyses together, the Random Forest model demonstrated improved precision over the traditional formula-based model in predicting the optimal tracheal tube depth. An RMSE of 1.21, significantly lower than that of the formula-based model, further emphasized this superiority, underlining its potential role in addressing the criticality of accurately determining OTTD in pediatric patients undergoing mechanical ventilation.

## A Data Description

Here is the data description, as provided by the user:

Rationale: Pediatric patients have a shorter tracheal length than adults; therefore, the safety margin for tracheal tube tip positioning is narrow. Indeed, the tracheal tube tip is misplaced in 35%{50% of pediatric patients and can cause hypoxia, atelectasis, hypercarbia, pneumothorax, and even death. Therefore, in pediatric patients who require mechanical ventilation, it is crucial to determine the Optimal Tracheal Tube Depth (defined here as `OTTD`, not an official term).

Note: For brevity, we introduce the term `OTTD` to refer to the "optimal tracheal tube depth". This is not an official term that can be found in the literature.

Existing methods: The gold standard to determine OTTD is by chest X-ray, which is time-consuming and requires radiation exposure.

Alternatively, formula-based models on patient features such as age and height are used to determine OTTD, but with limited success.

The provided dataset focus on patients aged 0-7 year old who received post-operative mechanical ventilation after undergoing surgery at Samsung Medical Center between January 2015 and December 2018.

For each of these patients, the dataset provides the OTTD determined by chest X-ray as well as features extracted from patient electronic health records.

1 data file:

"tracheal\_tube\_insertion.csv"

The csv file is a clean dataset of 969 rows (patients) and 6 columns:

Tube:

#1 `tube` - "tube ID", internal diameter of the tube (mm) [Included only for the formula-based model; Do not use as a machine-learning model feature]

Model features:

#2 `sex` - patient sex (0=female, 1=male)

#3 `age\_c` - patient age (years, rounded to half years)

```
#4 `ht` - patient height (cm)
#5 `wt` - patient weight (kg)

Target:
#6 `tube_depth_G` - Optimal tracheal tube depth as determined by chest X-ray (in
    cm)
```

## B Data Exploration

### B.1 Code

The Data Exploration was carried out using the following custom code:

```
1
2 import pandas as pd
3 import numpy as np
4
5 # load the dataset
6 df = pd.read_csv('tracheal_tube_insertion.csv')
7
8 # Open file for writing the exploration results
9 with open("data_exploration.txt", "w") as f:
10
11     # Check the size of our data
12     f.write("# Data Size\n")
13     f.write(str(df.shape) + "\n\n")
14
15     # Summary statistics of all or key variables
16     f.write("# Summary Statistics\n")
17     f.write(str(df.describe()) + "\n\n")
18
19     # Categorical Variables
20     f.write("# Categorical Variables\n")
21     f.write(str(df['sex'].value_counts()) + "\n\n")
22
23     # missing and special values
24     f.write("# Missing Values\n")
```

```
25     f.write(str(df.isnull().sum()) + "\n\n")
```

```
26
```

## B.2 Code Description

The provided code performs data exploration on the given dataset of pediatric patients who received post-operative mechanical ventilation after undergoing surgery. The goal of data exploration is to understand the structure and characteristics of the dataset.

First, the code loads the dataset using the ‘read\_csv’ function from the pandas library.

The code then proceeds to analyze the dataset and write the results into a file named "data\_exploration.txt".

The code starts by determining the size of the dataset using the ‘shape’ attribute of the pandas DataFrame. The result, which represents the number of rows and columns in the dataset, is written into the file.

Next, the code computes and writes summary statistics of the dataset using the ‘describe’ function. This provides information such as count, mean, standard deviation, minimum, and maximum values for each numerical variable in the dataset.

The code then focuses on the categorical variable ‘sex’ and computes the frequency count of different categories using the ‘value\_counts’ function. The result, which shows the number of male and female patients, is written into the file.

Finally, the code checks for missing values in the dataset using the ‘isnull’ function. It computes the sum of missing values for each variable and writes the result into the file. This step helps identify any data points that are incomplete or have missing information.

Overall, the code provides a comprehensive exploration of the dataset by analyzing its size, summary statistics, categorical variables, and missing values. The results are written into the "data\_exploration.txt" file for further examination and interpretation.

## B.3 Code Output

**data\_exploration.txt**

```
# Data Size
(969, 6)
```

```
# Summary Statistics
```

	tube	sex	age_c	ht	wt	tube_depth_G
count	969	969	969	969	969	969
mean	3.694	0.5387	0.758	66	7.127	10.19
std	0.5681	0.4988	1.44	19.08	4.774	1.766
min	2.5	0	0	29.5	0.57	5.9
25%	3.5	0	0	52	3.56	9
50%	3.5	1	0	61.5	5.7	9.8
75%	4	1	1	76	9.5	11.2
max	6	1	7	135.4	31	19.2

# Categorical Variables

sex

1 522

0 447

Name: count, dtype: int64

# Missing Values

tube 0

sex 0

age\_c 0

ht 0

wt 0

tube\_depth\_G 0

dtype: int64

## C Data Analysis

### C.1 Code

The Data Analysis was carried out using the following custom code:

```

1
2
3 # IMPORT
4 import pandas as pd
5 from sklearn.model_selection import train_test_split
6 from sklearn.ensemble import RandomForestRegressor
7 from sklearn.metrics import mean_squared_error

```

```

8  from scipy.stats import ttest_rel
9  import numpy as np
10 from sklearn.model_selection import GridSearchCV
11 import pickle
12
13 # LOAD DATA
14 data = pd.read_csv("tracheal_tube_insertion.csv")
15
16 # DATASET PREPARATIONS
17 # No dataset preparations are needed.
18
19 # DESCRIPTIVE STATISTICS
20 # No descriptive statistics table is needed.
21
22 # PREPROCESSING
23 # No preprocessing is needed, because the dataset is ready to
24   ↪ be used by the machine learning model.
25
26 # ANALYSIS
27 ## Table 1: Test of difference in predictive power between
28   ↪ formula-based model
29
30 # Separate features and target
31 features = data[['sex', 'age_c', 'ht', 'wt']].values
32 target = data['tube_depth_G'].values
33
34 # Train test split
35 features_train, features_test, target_train, target_test =
36   ↪ train_test_split(
37     features, target, test_size=0.2, random_state=42)
38
39 rf = RandomForestRegressor(random_state=42)
40 param_grid = {
41     'n_estimators': [50, 100, 200],
42     'min_samples_split': [2, 4, 8],
43 }
44 grid_search = GridSearchCV(estimator = rf, param_grid =
45   ↪ param_grid)

```



```

44 grid_search.fit(features_train, target_train)
45 best_rf = grid_search.best_estimator_
46
47 # Prediction
48 rf_predictions = best_rf.predict(features_test)
49
50 # Compute residuals for RF model
51 rf_squared_residuals = (rf_predictions - target_test) ** 2
52
53 # Compute Formula-based model predictions for test set
54 formula_predictions_test = features_test[:, 2] / 10 + 5
55
56 # Compute residuals for Formula-based model
57 formula_squared_residuals = (formula_predictions_test -
    ↪ target_test) ** 2
58
59 # Perform paired t-test
60 ttest_res = ttest_rel(rf_squared_residuals,
    ↪ formula_squared_residuals)
61
62 table_1 = pd.DataFrame({
63     'Model': ['Random Forest', 'Formula-based Model'],
64     't-statistic': [ttest_res.statistic, ttest_res.statistic],
65     'p-value': [ttest_res.pvalue, ttest_res.pvalue]
66 }, index=['Model 1', 'Model 2'])
67
68 # SAVE DATAFRAME
69 table_1.to_pickle('table_1.pkl')
70
71 ## Table 2: RMSE of Formula-based model and the Random Forest
    ↪ model
72
73 rf_rmse = np.sqrt(mean_squared_error(target_test,
    ↪ rf_predictions))
74 formula_rmse = np.sqrt(mean_squared_error(target_test,
    ↪ formula_predictions_test))
75
76 table_2 = pd.DataFrame({
77     'Model': ['Random Forest', 'Formula-based Model'],
78     'RMSE': [rf_rmse, formula_rmse]

```

```

79 }, index=['Model 1', 'Model 2'])
80
81 table_2.to_pickle('table_2.pkl')
82
83 # SAVE ADDITIONAL RESULTS
84 additional_results = {
85     'Total number of observations': len(data),
86     'Best Random Forest parameters': grid_search.best_params_
87 }
88
89 with open('additional_results.pkl', 'wb') as f:
90     pickle.dump(additional_results, f)
91

```

## C.2 Code Description

The provided code performs the analysis of the dataset "tracheal\_tube\_insertion.csv" to compare the predictive power and accuracy of a Random Forest model with a formula-based model for determining the optimal tracheal tube depth ('OTTD') in pediatric patients undergoing post-operative mechanical ventilation. The analysis involves the following steps:

1. Data loading: The dataset is read into a pandas DataFrame.
2. Dataset preparations: No specific preparations are needed as the dataset is already cleaned and ready for analysis.
3. Analysis of predictive power:
  - The features (patient sex, age, height, and weight) and the target variable (optimal tracheal tube depth determined by chest X-ray) are extracted from the dataset.
  - The dataset is split into training and test sets using an 80-20 ratio.
  - Random Forest model:
    - A Random Forest Regressor is trained using the GridSearchCV method to find the best hyperparameters (number of estimators and minimum number of samples to split).
    - The best-performing Random Forest model is used to make predictions on the test set.
    - The squared residuals (squared difference between predicted and actual values) are computed for the Random Forest model.
  - Formula-based model:
    - Formula-based predictions of tracheal tube depth are calculated using the patient's height.
    - The squared residuals are computed for the formula-based model.
  - Paired t-test:
    - A paired t-test is performed to compare the predictive power of the Random Forest model and the formula-based model.
    - The t-statistic and p-value are calculated for the t-test.
  - Table 1:
    - A table is created with the models' names, t-statistic, and p-value.
    - The table is saved as a pickle object.

4. Analysis of RMSE: - The root mean squared error (RMSE) for both the Random Forest model and the formula-based model is calculated using the test set. - Table 2: - A table is created with the models' names and corresponding RMSE values. - The table is saved as a pickle object.

5. Additional results: - The total number of observations in the dataset and the best hyperparameters found for the Random Forest model are stored as a dictionary. - The dictionary is saved as a pickle object.

The "additional\_results.pkl" file contains a dictionary with the following information: - "Total number of observations": The total number of patients in the dataset. - "Best Random Forest parameters": The best hyperparameters (number of estimators and minimum number of samples to split) found for the Random Forest model during the GridSearchCV process.

These results provide insights into the predictive power and accuracy of the Random Forest model compared to the formula-based model for determining the optimal tracheal tube depth in pediatric patients undergoing post-operative mechanical ventilation.

### C.3 Code Output

**table\_1.pkl**

	Model	t-statistic	p-value
Model 1	Random Forest	-6.023673	8.47e-09
Model 2	Formula-based Model	-6.023673	8.47e-09

**table\_2.pkl**

	Model	RMSE
Model 1	Random Forest	1.211588
Model 2	Formula-based Model	1.849024

**additional\_results.pkl**

```
{
  'Total number of observations': 969,
  'Best Random Forest parameters': {'min_samples_split': 8, 'n_estimators':
    200},
}
```

## D LaTeX Table Design

### D.1 Code

The LaTeX Table Design was carried out using the following custom code:

```
1
2 # IMPORT
3 import pandas as pd
4 from my_utils import to_latex_with_note, format_p_value,
5     ↪ is_str_in_df, split_mapping, AbbrToNameDef
6
7 # PREPARATION FOR ALL TABLES
8 shared_mapping: AbbrToNameDef = {
9     'tube': ('Tube ID', 'Internal diameter of the tube in
10     ↪ mm'),
11     'sex': ('Sex', 'Patient sex; 0: Female, 1: Male'),
12     'age_c': ('Age (years)', 'Patient age in years, rounded to
13     ↪ half years'),
14     'ht': ('Height (cm)', 'Patient height in cm'),
15     'wt': ('Weight (kg)', 'Patient weight in kg'),
16     'tube_depth_G': ('Optimal Tracheal Tube Depth (cm)',
17     ↪ 'Optimal tracheal tube depth as determined by chest
18     ↪ X-ray in cm'),
19     'Model 1': ('Model 1', 'Random Forest Model'),
20     'Model 2': ('Model 2', 'Height Formula-based Model'),
21     'RMSE': ('RMSE', 'Root Mean Square Error of the model
22     ↪ predictions')
23 }
24
25 # TABLE 1
26 df = pd.read_pickle('table_1.pkl')
27
28 # FORMAT VALUES
29 df['p-value'] = df['p-value'].apply(format_p_value)
30
31 # RENAME ROWS AND COLUMNS
32 mapping = {k: v for k, v in shared_mapping.items() if
33     ↪ is_str_in_df(df, k)}
34 abbrs_to_names, legend = split_mapping(mapping)
35 df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
```

```

29
30 # Save as latex:
31 to_latex_with_note(
32     df, 'table_1.tex',
33     caption="Test of difference in predictive power between
34     ↪ formula-based model",
35     label='table:prediction_models',
36     note="",
37     legend=legend
38 )
39
40 # TABLE 2
41 df = pd.read_pickle('table_2.pkl')
42
43 # RENAME ROWS AND COLUMNS
44 mapping = {k: v for k, v in shared_mapping.items() if
45 ↪ is_str_in_df(df, k)}
46 abbrs_to_names, legend = split_mapping(mapping)
47 df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
48
49 # Save as latex:
50 to_latex_with_note(
51     df, 'table_2.tex',
52     caption="Root mean square error of formula-based model and
53     ↪ the Random Forest model",
54     label='table:rmse_models',
55     note="",
56     legend=legend
57 )
58

```

## D.2 Provided Code

The code above is using the following provided functions:

```

1 def to_latex_with_note(df, filename: str, caption: str, label:
2 ↪ str, note: str = None, legend: Dict[str, str] = None,
3 ↪ **kwargs):
4     """
5     Converts a DataFrame to a LaTeX table with optional note and
6     ↪ legend added below the table.

```

```

4
5 Parameters:
6 - df, filename, caption, label: as in `df.to_latex`.
7 - note (optional): Additional note below the table.
8 - legend (optional): Dictionary mapping abbreviations to full
↪ names.
9 - **kwargs: Additional arguments for `df.to_latex`.
10
11 Returns:
12 - None: Outputs LaTeX file.
13 """
14
15 def format_p_value(x):
16     returns "{:.3g}".format(x) if x >= 1e-06 else "<1e-06"
17
18 def is_str_in_df(df: pd.DataFrame, s: str):
19     return any(s in level for level in getattr(df.index,
↪ 'levels', [df.index]) + getattr(df.columns, 'levels',
↪ [df.columns]))
20
21 AbbrToNameDef = Dict[Any, Tuple[Optional[str], Optional[str]]]
22
23 def split_mapping(abbrs_to_names_and_definitions:
↪ AbbrToNameDef):
24     abbrs_to_names = {abbr: name for abbr, (name, definition) in
↪ abbrs_to_names_and_definitions.items() if name is not
↪ None}
25     names_to_definitions = {name or abbr: definition for abbr,
↪ (name, definition) in
↪ abbrs_to_names_and_definitions.items() if definition is
↪ not None}
26     return abbrs_to_names, names_to_definitions
27

```

### D.3 Code Output

table\_1.tex

```

\begin{table}[h]
\caption{Test of difference in predictive power between formula-based model}

```

```

\label{table:prediction_models}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{llr}
\toprule
& Model & t-statistic & p-value \\
\midrule
\textbf{Model 1} & Random Forest & -6.02 &  $\leq 1e-06$  \\
\textbf{Model 2} & Formula-based Model & -6.02 &  $\leq 1e-06$  \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item \textbf{Model 1}: Random Forest Model
\item \textbf{Model 2}: Height Formula-based Model
\end{tablenotes}
\end{threeparttable}
\end{table}

```

#### table\_2.tex

```

\begin{table}[h]
\caption{Root mean square error of formula-based model and the Random Forest model}
\label{table:rmse_models}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{llr}
\toprule
& Model & RMSE \\
\midrule
\textbf{Model 1} & Random Forest & 1.21 \\
\textbf{Model 2} & Formula-based Model & 1.85 \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize

```

```

\item \textbf{Model 1}: Random Forest Model
\item \textbf{Model 2}: Height Formula-based Model
\item \textbf{RMSE}: Root Mean Square Error of the model predictions
\end{tablenotes}
\end{threeparttable}
\end{table}

```