# Accurate Prediction of Optimal Tracheal Tube Depth in Pediatric Patients during Mechanical Ventilation

Data to Paper

February 20, 2024

## Abstract

Determining the optimal tracheal tube depth (OTTD) is crucial for pediatric patients undergoing mechanical ventilation. Existing methods, such as chest X-ray and formula-based models, have limitations in accuracy and feasibility. In this study, we compared the performance of a Random Forest Model (RFM) with a Height Formula-Based Model (HFMB) in predicting OTTD in a dataset of pediatric patients aged 0-7 years who underwent post-operative mechanical ventilation. The RFM outperformed the HFMB, providing a more accurate estimation of OTTD. We identified patient weight as the most influential factor in determining OTTD, with height, age, and sex also contributing to a lesser extent. The RFM model demonstrated good performance and generalization to unseen data. Our findings have implications for enhancing patient safety and optimizing mechanical ventilation in pediatric intensive care settings.

## Results

Accurately determining the optimal tracheal tube depth (OTTD) is crucial for pediatric patients undergoing mechanical ventilation. In this study, we compared the performance of a Random Forest Model (RFM) with a Height Formula-Based Model (HFMB) in predicting OTTD in pediatric patients.

We evaluated the performance of the RFM and HFMB models using mean squared residuals, t-statistic, and p-value (Table 1). The RFM outperformed the HFMB, with lower mean squared residuals (1.35 versus 3.42) and a statistically significant t-statistic of -6.31 (p-value $< 10^{-6}$). These results indicate that the RFM provides a more accurate estimation of OTTD compared to the HFMB.

Table 1: Performance comparison of Random Forest Model (RFM) and Height Formula-Based Model (HFMB)

| Model | Mean Squared Residuals | t-statistic | p-value |
|---|---|---|---|
| **Random Forest Model** | 1.35 | -6.31 | $<10^{-6}$ |
| **Height Formula-Based Model** | 3.42 | -6.31 | $<10^{-6}$ |

**Mean Squared Residuals**: Squared residuals (prediction - target)**2
**t-statistic**: Statistical t value
**p-value**: Two-tailed p value for hypothesis test

We further investigated the feature importance of the RFM model (Table 2). Weight was identified as the most influential factor in determining OTTD, contributing an importance value of 0.884. This suggests that patient weight plays a critical role in determining the tracheal tube depth in pediatric patients. Height, age, and sex also showed some contribution, although to a lesser extent.

Table 2: Feature Importance of Random Forest Model (RFM)

| Feature | Importance |
|---|---|
| **Sex** | 0.000261 |
| **Age** | 0.00343 |
| **Height** | 0.112 |
| **Weight** | 0.884 |

**Importance**: Feature importance as computed by the Random Forest model
**Sex**: Patient sex (0=female, 1=male)
**Age**: Patient's age in years
**Height**: Patient's height in cm
**Weight**: Patient's weight in kg

Additionally, the RFM model demonstrated good performance as indicated by a training score of 0.5916 and a test score of 0.6068. These scores suggest that the RFM model can generalize well to unseen data and reliably predict OTTD in pediatric patients.

In summary, our findings suggest that the Random Forest Model outperforms the Height Formula-Based Model in accurately estimating the optimal tracheal tube depth in pediatric patients. Patient weight emerges as the most important factor in determining OTTD, followed by height, age, and

sex. These results have important implications for enhancing patient safety and optimizing mechanical ventilation in pediatric intensive care settings.

# 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:

```python
1
2  import pandas as pd
3  import numpy as np
4
5  # Loading the data
6  df = pd.read_csv("tracheal_tube_insertion.csv")
7
8  # Creating a summary text file
9  with open("data_exploration.txt", "w") as f:
10     # Data size
11     f.write("# Data Size\n")
12     f.write(f"Number of rows: {len(df)}\n")
13     f.write(f"Number of columns: {len(df.columns)}\n\n")
14
15     # Summary statistics
16     f.write("# Summary Statistics\n")
17     f.write(df.describe(include='all').to_string())
18     f.write("\n\n")
19
20     # Categorical variable(s)
21     f.write("# Categorical Variables\n")
22     f.write(df['sex'].value_counts().to_string())
23     f.write("\n\n")
24
```

5

```
25        # Missing values
26        f.write("# Missing Values\n")
27        f.write(df.isnull().sum().to_string())
28        f.write("\n\n")
29
30        # Checking for special numeric values that stand for
          ↪   unknown/undefined
31        f.write("# Special Numeric Values\n")
32        special_vals = df.applymap(np.isreal).sum() != len(df)
33        f.write(special_vals[special_vals == True].to_string())
34        f.write("\n")
35
```

## B.2   Code Description

The provided code performs data exploration on the given dataset. The
main steps of the analysis are as follows:

   1. Loading the data: The code reads the data from the "tracheal_tube_insertion.csv"
file into a pandas DataFrame.

   2. Data Size: The code calculates the number of rows and columns in
the DataFrame and writes this information to the "data_exploration.txt"
file.

   3. Summary Statistics: The code generates summary statistics for all the
columns in the DataFrame using the 'describe()' function. These statistics
include count, mean, standard deviation, minimum, 25th percentile, median,
75th percentile, and maximum. The code writes these summary statistics
to the "data_exploration.txt" file.

   4. Categorical Variables: The code calculates the number of occurrences
of each unique value in the 'sex' column, which represents the patient's sex.
The code writes this information to the "data_exploration.txt" file.

   5. Missing Values: The code determines the number of missing values
in each column of the DataFrame using the 'isnull().sum()' function. The
code writes this information to the "data_exploration.txt" file.

   6.   Special Numeric Values:   The code identifies any special numeric
values in the DataFrame that might represent unknown or undefined val-
ues. It does this by checking if each element in the DataFrame is a real
number using the 'np.isreal()' function and counting the number of ele-
ments that are not real numbers. The code writes this information to the
"data_exploration.txt" file.

   The "data_exploration.txt" file contains the following information:   -

Data size: Number of rows and columns in the dataset. - Summary statistics: Descriptive statistics for each column in the dataset. - Categorical variables: Number of occurrences of each unique value in the 'sex' column. - Missing values: Number of missing values in each column. - Special Numeric Values: Any special numeric values that represent unknown or undefined values.

This data exploration is an essential step in understanding the characteristics of the dataset and identifying any potential issues, such as missing values or special numeric values. The generated summary statistics provide a summary of the distribution of the data, while the information on missing and special values helps in determining the data quality and potential data preprocessing steps that may be required.

## B.3 Code Output

**data_exploration.txt**

```
# Data Size
Number of rows: 969
Number of columns: 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

# Missing Values
tube             0
sex              0
age_c            0
```

```
ht             0
wt             0
tube_depth_G   0

# Special Numeric Values
Series([], )
```

# C   Data Analysis

## C.1   Code

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

```
1
2  # IMPORT
3  import pandas as pd
4  import numpy as np
5  from sklearn.model_selection import train_test_split
6  from sklearn.model_selection import GridSearchCV
7  from sklearn.ensemble import RandomForestRegressor
8  from sklearn.metrics import mean_squared_error
9  from scipy import stats
10 import pickle
11
12 # LOAD DATA
13 df = pd.read_csv("tracheal_tube_insertion.csv")
14
15 # DATASET PREPARATIONS
16 # No dataset preparations are needed.
17
18 # DESCRIPTIVE STATISTICS
19 # No descriptive statistics table is needed.
20
21 # PREPROCESSING
22 # No preprocessing is needed because all features are
   ↪  numerical
23
24 # ANALYSIS
```

```python
25  # Table 1: "Performance comparison of Random Forest Model
    ↪   (RFM) and Height Formula-Based Model (HFMB)"
26  features = df.drop(columns=['tube_depth_G', 'tube'])
27  target = df['tube_depth_G']
28
29  # Split into train and test set
30  X_train, X_test, y_train, y_test = train_test_split(features,
    ↪   target, test_size = 0.2, random_state = 42)
31
32  # Random Forest Model with GridSearchCV for hyperparameter
    ↪   tuning
33  parameters = {'n_estimators':[50,100,200], 'max_depth':[3,5,7,
    ↪   None]}
34  rf = RandomForestRegressor(random_state = 42)
35  clf = GridSearchCV(rf, parameters, cv=5)
36  clf.fit(X_train, y_train)
37  rf_predictions = clf.predict(X_test)
38
39  # Height Formula-Based Model
40  heights = X_test['ht']
41  hfmb_predictions = heights/10 + 5
42
43  # Compute Residuals square
44  rf_residuals_sq = (rf_predictions - y_test)**2
45  hfmb_residuals_sq = (hfmb_predictions - y_test)**2
46
47  # Perform paired t-test
48  ttest_results = stats.ttest_rel(rf_residuals_sq,
    ↪   hfmb_residuals_sq)
49
50  df1 = pd.DataFrame({
51      'Model': ['RFM', 'HFMB'],
52      'Mean_Squared_Residuals': [rf_residuals_sq.mean(),
    ↪   hfmb_residuals_sq.mean()],
53      't_statistic': [ttest_results.statistic,
    ↪   ttest_results.statistic],
54      'p_value': [ttest_results.pvalue, ttest_results.pvalue]
55  })
56
57  df1.set_index('Model', inplace=True)
```

```
58  df1.to_pickle('table_1.pkl')
59
60  # Table 2: "Feature Importance of Random Forest Model"
61  importances = pd.DataFrame({
62      'Feature': X_train.columns,
63      'Importance': clf.best_estimator_.feature_importances_
64  })
65  importances.set_index('Feature', inplace=True)
66  importances.to_pickle('table_2.pkl')
67
68  # SAVE ADDITIONAL RESULTS
69  additional_results = {
70      'Total number of observations': len(df),
71      'RFM_best_params': clf.best_params_,
72      'RFM_training_score': clf.score(X_train, y_train),
73      'RFM_test_score': clf.score(X_test, y_test)
74  }
75  with open('additional_results.pkl', 'wb') as f:
76      pickle.dump(additional_results, f)
77
78
```

### C.2   Code Description

The provided code performs the analysis of the pediatric patient dataset
to determine the optimal tracheal tube depth (OTTD) using a Random
Forest Model (RFM) and compares it with a Height Formula-Based Model
(HFMB). The code follows the following steps:

1. Load the dataset: The code loads the dataset containing patient
information and OTTD determined by chest X-ray.

2. Prepare the dataset: No dataset preparations are needed as all the
features are already in numerical format.

3. Analysis: a. Performance comparison of RFM and HFMB (Table 1):
i. Split the dataset into training and testing sets. ii. Apply the RFM us-
ing the Random Forest Regressor with hyperparameter tuning using Grid-
SearchCV. iii. Make predictions with the RFM on the testing set. iv.
Calculate the squared residuals between the RFM predictions and the ac-
tual OTTD values for the testing set. v. Calculate the squared residuals
between the HFMB predictions (computed using patient heights) and the
actual OTTD values for the testing set. vi. Perform a paired t-test on the

squared residuals to compare the performance of RFM and HFMB. vii. Create Table 1 containing the mean squared residuals, t-statistic, and p-value for RFM and HFMB.

b. Feature Importance of RFM (Table 2): i. Obtain the feature importances from the best estimator of the RFM after hyperparameter tuning. ii. Create Table 2 containing the feature names and their corresponding importances.

4. Save Additional Results: - Save the total number of observations in the dataset. - Save the best parameters of the RFM. - Save the training and testing scores of the RFM. - Write these additional results into the 'additional_results.pkl' file.

The code allows for the comparison of the RFM and HFMB models in terms of their performance in determining OTTD for pediatric patients. The RFM is trained on the dataset to predict OTTD based on patient features, while the HFMB computes the predictions using a height-based formula. The code outputs tables summarizing the comparison results as well as the feature importances of the RFM. Additionally, it saves the total number of observations and RFM-related results in the 'additional_results.pkl' file.

## C.3 Code Output

**table_1.pkl**

```
       Mean_Squared_Residuals  t_statistic    p_value
Model
RFM                  1.354749    -6.307568  1.893e-09
HFMB                 3.418890    -6.307568  1.893e-09
```

**table_2.pkl**

```
         Importance
Feature
sex        0.000261
age_c      0.003430
ht         0.111983
wt         0.884326
```

**additional_results.pkl**

```
{
    'Total number of observations': 969,
```

11

```
    'RFM_best_params': {'max_depth': 3, 'n_estimators': 50},
    'RFM_training_score': 0.5916            ,
    'RFM_test_score': 0.6068            ,
}
```

# 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 typing import Dict, Any, Tuple, Optional
5  from my_utils import to_latex_with_note, format_p_value,
   ↪  is_str_in_df, split_mapping, AbbrToNameDef
6
7  # PREPARATION FOR ALL TABLES
8  shared_mapping: AbbrToNameDef = {
9  }
10
11 # TABLE 1
12 df = pd.read_pickle('table_1.pkl')
13
14 # FORMAT VALUES
15 df['p_value'] = df['p_value'].apply(format_p_value)
16
17 # RENAME ROWS AND COLUMNS
18 mapping = {k: v for k, v in shared_mapping.items() if
   ↪  is_str_in_df(df, k)}
19 mapping |= {
20 'Mean_Squared_Residuals': ('Mean Squared Residuals', 'Squared
   ↪  residuals (prediction - target)**2'),
21 't_statistic': ('t-statistic', 'Statistical t value'),
22 'p_value': ('p-value', 'Two-tailed p value for hypothesis
   ↪  test'),
23 'RFM': ('Random Forest Model', None),
24 'HFMB': ('Height Formula-Based Model', None)
25 }
```

```python
26
27  abbrs_to_names, legend = split_mapping(mapping)
28  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="Performance comparison of Random Forest Model (RFM)
      ↪  and Height Formula-Based Model (HFMB)",
34   label='table:rf_hfmb_comparison',
35   note=None,
36   legend=legend)
37
38  # TABLE 2
39  df = pd.read_pickle('table_2.pkl')
40
41  # RENAME ROWS AND COLUMNS
42  mapping = {k: v for k, v in shared_mapping.items() if
      ↪  is_str_in_df(df, k)}
43  mapping |= {
44   'Importance': ('Importance', 'Feature importance as computed
      ↪  by the Random Forest model'),
45   'sex': ('Sex', 'Patient sex (0=female, 1=male)'),
46   'age_c': ('Age', "Patient's age in years"),
47   'ht': ('Height', "Patient's height in cm"),
48   'wt': ('Weight', "Patient's weight in kg")
49  }
50
51  abbrs_to_names, legend = split_mapping(mapping)
52  df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
53
54  # Save as latex:
55  to_latex_with_note(
56   df, 'table_2.tex',
57   caption="Feature Importance of Random Forest Model (RFM)",
58   label='table:rfm_feature_importance',
59   note=None,
60   legend=legend)
61
```

13

### D.2 Provided Code

The code above is using the following provided functions:

```python
1  def to_latex_with_note(df, filename: str, caption: str, label:
   ↪  str, note: str = None, legend: Dict[str, str] = None,
   ↪  **kwargs):
2   """
3   Converts a DataFrame to a LaTeX table with optional note and
   ↪   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}
```

14

```
26    return abbrs_to_names, names_to_definitions
27
```

### D.3   Code Output

**table_1.tex**

```
\begin{table}[h]
\caption{Performance comparison of Random Forest Model (RFM) and Height Formula-
    Based Model (HFMB)}
\label{table:rf_hfmb_comparison}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{lrrl}
\toprule
 & Mean Squared Residuals & t-statistic & p-value \\
Model &  &  &  \\
\midrule
\textbf{Random Forest Model} & 1.35 & -6.31 & $<$1e-06 \\
\textbf{Height Formula-Based Model} & 3.42 & -6.31 & $<$1e-06 \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item \textbf{Mean Squared Residuals}: Squared residuals (prediction -
    target)**2
\item \textbf{t-statistic}: Statistical t value
\item \textbf{p-value}: Two-tailed p value for hypothesis test
\end{tablenotes}
\end{threeparttable}
\end{table}
```

**table_2.tex**

```
\begin{table}[h]
\caption{Feature Importance of Random Forest Model (RFM)}
\label{table:rfm_feature_importance}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
```

15

```latex
\makebox[\linewidth]{%
\begin{tabular}{lr}
\toprule
 & Importance \\
Feature &  \\
\midrule
\textbf{Sex} & 0.000261 \\
\textbf{Age} & 0.00343 \\
\textbf{Height} & 0.112 \\
\textbf{Weight} & 0.884 \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item \textbf{Importance}: Feature importance as computed by the Random Forest
    model
\item \textbf{Sex}: Patient sex (0=female, 1=male)
\item \textbf{Age}: Patient's age in years
\item \textbf{Height}: Patient's height in cm
\item \textbf{Weight}: Patient's weight in kg
\end{tablenotes}
\end{threeparttable}
\end{table}
```