

# Predicting Optimal Tracheal Tube Depth in Pediatric Patients on Mechanical Ventilation

Data to Paper

February 20, 2024

## Abstract

Accurate placement of the tracheal tube is crucial for pediatric patients undergoing mechanical ventilation. Despite the potential risks of misplacement, determining the Optimal Tracheal Tube Depth (OTTD) remains challenging. Existing methods rely on chest X-ray or formula-based models but lack precision. To address this, we developed predictive models using patient features to estimate the OTTD. Our dataset comprised 969 pediatric patients aged 0-7 years who received post-operative mechanical ventilation. Leveraging Random Forest and Elastic Net algorithms, we evaluated the predictive performance based on sex, age, height, and weight. The Elastic Net model demonstrated superior accuracy in estimating the OTTD compared to the Random Forest model. Paired analysis further confirmed the superiority of the Elastic Net model. The development of accurate predictive models based on patient features offers an alternative to chest X-ray assessment for determining the OTTD. These findings have the potential to enhance patient safety and reduce complications associated with tracheal tube misplacement. However, further refinement and validation of the models are necessary for widespread implementation in clinical practice.

## Results

In this section, we present the results of our analysis aimed at determining the Optimal Tracheal Tube Depth (OTTD) in pediatric patients undergoing mechanical ventilation. We utilized a dataset of 969 pediatric patients aged 0-7 years who received post-operative mechanical ventilation after undergoing surgery. Our analysis focused on developing predictive models using patient features to estimate the OTTD and reducing the reliance on chest X-ray assessment.

To achieve this, we employed two different predictive models, Random Forest and Elastic Net, to estimate the OTTD based on patient features including sex, age, height, and weight. The Random Forest model utilized a maximum number of features of "sqrt" and a number of estimators of 500, while the Elastic Net model utilized an alpha value of 0.1 and an l1 ratio of 0.1.

First, we examined the descriptive statistics of patient features and OTTD stratified by sex (Table 1). The mean and standard deviation were calculated for patient features such as age, height, weight, and OTTD. We observed potential differences in these features based on patient sex. For instance, the mean OTTD was found to be 10.1 cm (SD=1.65) for female patients and 10.3 cm (SD=1.86) for male patients. Other features such as age, height, and weight also exhibited slight variations between the two groups.

Table 1: Descriptive statistics of patient features and OTTD stratified by sex

	sex	0	1
<b>Tube ID</b>	<b>mean</b>	3.68	3.7
	<b>std</b>	0.552	0.582
<b>Age</b>	<b>mean</b>	0.732	0.781
	<b>std</b>	1.4	1.47
<b>Height</b>	<b>mean</b>	65.4	66.5
	<b>std</b>	18.7	19.4
<b>Weight</b>	<b>mean</b>	6.84	7.37
	<b>std</b>	4.57	4.94
<b>Optimal Tracheal Tube Depth (OTTD)</b>	<b>mean</b>	10.1	10.3
	<b>std</b>	1.65	1.86

**Age:** Patient age (in years)

**Height:** Patient height (in cm)

**Weight:** Patient weight (in kg)

**Tube ID:** Internal diameter of the tube (in mm)

**Optimal Tracheal Tube Depth (OTTD):** As determined by chest X-ray (in cm)

Next, we evaluated the performance of the Random Forest and Elastic Net models in predicting the OTTD. The models were assessed based on mean squared error (MSE) values calculated on the test set (Table 2). The Elastic Net model outperformed the Random Forest model, achieving an MSE of 1.24 compared to the Random Forest model's MSE of 1.57. These findings indicate that the Elastic Net model demonstrated superior predic-

tive performance in estimating the OTTD compared to the Random Forest model.

Table 2: Performance of Random Forest and Elastic Net models in predicting OTTD

Model	Mean Squared Error
<b>Random Forest</b>	1.57
<b>Elastic Net</b>	1.24

Finally, we conducted a paired t-test to compare the squared residuals of the two models (Table 3). The analysis revealed a significant difference between the squared residuals of the Elastic Net and Random Forest models, with a t statistic of 2.36 (p-value=0.0194). This further supports the superior performance of the Elastic Net model in estimating the OTTD.

Table 3: Paired t-test comparing the squared residuals of the two models

Model 1 residuals	Model 2 residuals	t statistic	p value
<b>Random Forest</b>	Elastic Net	2.36	0.0194

In summary, our analysis suggests that the Elastic Net model shows promise in accurately estimating the Optimal Tracheal Tube Depth in pediatric patients on mechanical ventilation. Furthermore, patient features, such as sex, age, height, and weight, may play a crucial role in determining the OTTD. The results highlight the potential of utilizing predictive models based on patient features to minimize the need for chest X-ray assessment, providing an alternative approach in clinical practice.

## 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 header = ['Data Size', 'Summary Statistics', 'Categorical
    ↳ Variables', 'Missing Values']
9
10 # count of rows and columns
11 data_size = "Number of Rows: {}\nNumber of Columns:
    ↳ {}".format(df.shape[0], df.shape[1])
12
13 # summary statistics
14 summary_stats =
    ↳ df.describe(include='all').transpose().to_string()
15
16 # unique values in categorical variable
17 cat_vars = "sex: {}".format(df.sex.value_counts().to_string())
18
19 # check count of missing values in all columns
20 missing_values = df.isnull().sum().to_string()
21

```

```

22 exploratory_outputs = [data_size, summary_stats, cat_vars,
    ↪ missing_values]
23
24 with open("data_exploration.txt", 'w') as f:
25     for i,section in enumerate(exploratory_outputs):
26         f.write("# " + header[i] + "\n" + section + "\n\n")
27

```

## B.2 Code Description

The provided code performs an initial data exploration on the given dataset.

First, the code reads the "tracheal\_tube\_insertion.csv" file using the pandas library and stores it in a dataframe called "df".

The code then proceeds to perform the following analysis steps:

1. Data Size: The code determines the number of rows and columns in the dataset, and stores this information as "Number of Rows: <number>" and "Number of Columns: <number>".

2. Summary Statistics: The code computes the summary statistics for each column in the dataset using the describe() function. It generates statistics such as count, mean, standard deviation, minimum, and quartiles for numerical columns, and provides the count, number of unique values, top value, and frequency of top value for categorical columns. The summary statistics are stored as a table in the "summary\_stats" variable.

3. Categorical Variables: The code calculates the count of unique values in the 'sex' column, which represents the patient's sex. It provides the count of female patients and male patients in a formatted string.

4. Missing Values: The code checks for missing values in each column of the dataset using the isnull() function. It counts the number of missing values for each column and stores the result as a string in the "missing\_values" variable.

Finally, the code writes the results of the data exploration analysis to a text file named "data\_exploration.txt", using the "header" list as section titles. Each section includes the corresponding result or analysis step.

The "data\_exploration.txt" file will contain the following sections:

1. Data Size: The number of rows and columns in the dataset.
2. Summary Statistics: The summary statistics for each column in the dataset.
3. Categorical Variables: The count of unique values in the 'sex' column.
4. Missing Values: The count of missing values in each column.

This information can be used to gain a preliminary understanding of the dataset and identify any potential issues or patterns before proceeding with

further analysis.

### B.3 Code Output

#### data\_exploration.txt

##### # Data Size

Number of Rows: 969

Number of Columns: 6

##### # Summary Statistics

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

##### # Categorical Variables

sex: sex

1 522

0 447

##### # Missing Values

tube 0

sex 0

age\_c 0

ht 0

wt 0

tube\_depth\_G 0

## 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,
   ↪ GridSearchCV
6  from sklearn.ensemble import RandomForestRegressor
7  from sklearn.linear_model import ElasticNet
8  from sklearn.metrics import mean_squared_error
9  from scipy.stats import ttest_rel
10 import pickle
11 import warnings
12 warnings.filterwarnings('ignore')
13
14 # LOAD DATA
15 df = pd.read_csv('tracheal_tube_insertion.csv')
16
17 # DATASET PREPARATIONS
18 # No dataset preparations are needed.
19
20 # DESCRIPTIVE STATISTICS
21 # Table 0: "Mean and standard deviation of patient features
   ↪ and tube_depth_G stratified by sex"
22 df0 = df.groupby('sex').agg(['mean', 'std']).transpose()
23 df0.to_pickle('table_0.pkl')
24
25 # PREPROCESSING
26 # No preprocessing is needed, because all the features are
   ↪ numerical and there are no missing values.
27
28 # ANALYSIS
29 # Determining feature and target variables
30 X = df[["sex", "age_c", "ht", "wt"]]
31 y = df["tube_depth_G"]
32
33 # Splitting the data into training and testing sets
34 X_train, X_test, y_train, y_test = train_test_split(X, y,
   ↪ test_size=0.2, random_state=42)
35

```



```

36 # Table 1: "Performance (mean squared error) of Random Forest
    ↪ and Elastic Net models on the test set"
37 # Hyperparameter tuning for Random Forest
38 rf = RandomForestRegressor(random_state=42)
39 param_grid = {
40     'n_estimators': [200, 500],
41     'max_features': ['sqrt', 'log2']
42 }
43 CV_rf = GridSearchCV(estimator=rf, param_grid=param_grid, cv=
    ↪ 2)
44 CV_rf.fit(X_train, y_train)
45 best_rf = CV_rf.best_estimator_
46 rf_pred = best_rf.predict(X_test)
47 rf_mse = mean_squared_error(y_test, rf_pred)
48
49 # Hyperparameter tuning for Elastic Net
50 en = ElasticNet(random_state=42)
51 param_grid = {
52     'alpha' : [0.1, 0.5, 1, 2, 5],
53     'l1_ratio' : [0.1, 0.5, 1]
54 }
55 CV_en = GridSearchCV(en, param_grid=param_grid, cv=2)
56 CV_en.fit(X_train, y_train)
57 best_en = CV_en.best_estimator_
58 en_pred = best_en.predict(X_test)
59 en_mse = mean_squared_error(y_test, en_pred)
60
61 df1 = pd.DataFrame({'Model': ['Random Forest', 'Elastic Net'],
    ↪ 'Mean Squared Error': [rf_mse, en_mse]})
62 df1.set_index('Model', inplace=True)
63 df1.to_pickle('table_1.pkl')
64
65 # Table 2: "Paired t-test comparing the squared residuals of
    ↪ the two models"
66 rf_res = (rf_pred - y_test)**2
67 en_res = (en_pred - y_test)**2
68 ttest_results = ttest_rel(rf_res, en_res)

```

```

69 df2 = pd.DataFrame({'Model 1 residuals': ['Random Forest'],
    ↪ 'Model 2 residuals': ['Elastic Net'], 't statistic':
    ↪ [ttest_results.statistic], 'p value':
    ↪ [ttest_results.pvalue]})
70 df2.set_index('Model 1 residuals', inplace=True)
71 df2.to_pickle('table_2.pkl')
72
73 # SAVE ADDITIONAL RESULTS
74 additional_results = {
75     'Total number of observations': df.shape[0],
76     'Best parameters of Random Forest model': CV_rf.best_params_,
77     'Best parameters of Elastic Net model': CV_en.best_params_
78 }
79 with open('additional_results.pkl', 'wb') as f:
80     pickle.dump(additional_results, f)
81

```

## C.2 Code Description

The code performs data analysis on a dataset of pediatric patients who require mechanical ventilation after surgery. The goal is to determine the optimal tracheal tube depth (OTTD), which is crucial for patient safety.

The code first loads the dataset and prepares it for analysis. It then calculates descriptive statistics, such as the mean and standard deviation of patient features and OTTD stratified by sex, and stores the results in a table.

Next, the code performs preprocessing on the dataset, which is not needed in this case as all the features are numerical and there are no missing values.

The code then proceeds to the analysis step. It splits the dataset into training and testing sets and defines the feature variables (patient sex, age, height, and weight) and the target variable (OTTD).

Two machine learning models, Random Forest and Elastic Net, are trained on the training set and evaluated on the testing set. The code uses grid search to find the best hyperparameters for each model. The mean squared error (MSE) is calculated for both models on the testing set, and the results are stored in a table.

To compare the performance of the two models, a paired t-test is conducted on the squared residuals of the predictions from each model. The t-statistic and p-value of the t-test are then saved in another table.

Finally, the code saves additional results, including the total number of observations in the dataset, and the best parameters found for the Random Forest and Elastic Net models.

The "additional\_results.pkl" file contains a pickled dictionary with these additional results, including the total number of observations and the best parameters of each model.

Overall, the code performs data analysis on the pediatric patient dataset, determines the OTTD using machine learning models, compares the performance of the models, and saves the results for further reference.

### C.3 Code Output

#### table\_0.pkl

sex		0	1
tube	mean	3.681208	3.704598
	std	0.551846	0.582023
age_c	mean	0.731544	0.780651
	std	1.402500	1.472808
ht	mean	65.400447	66.514368
	std	18.701462	19.403722
wt	mean	6.841902	7.370556
	std	4.568146	4.935102
tube_depth_G	mean	10.062416	10.298276
	std	1.645478	1.857778

#### table\_1.pkl

	Mean Squared Error
Model	
Random Forest	1.572975
Elastic Net	1.239276

#### table\_2.pkl

	Model 2 residuals	t statistic	p value
Model 1 residuals			
Random Forest	Elastic Net	2.357123	0.01942

additional\_results.pkl

```
{
    'Total number of observations': 969,
    'Best parameters of Random Forest model': {'max_features': 'sqrt',
    'n_estimators': 500},
    'Best parameters of Elastic Net model': {'alpha': 0.1, 'l1_ratio': 0.1},
}
```

## 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
6
7 # PREPARATION FOR ALL TABLES
8 shared_mapping: Dict[str, Tuple[Optional[str], Optional[str]]]
    ↪ = {
9     'sex': ('Sex', 'Patient sex (0 for female, 1 for male)'),
10    'age_c': ('Age', 'Patient age (in years)'),
11    'ht': ('Height', 'Patient height (in cm)'),
12    'wt': ('Weight', 'Patient weight (in kg)'),
13    'tube': ('Tube ID', 'Internal diameter of the tube (in mm)'),
14    'tube_depth_G': ('Optimal Tracheal Tube Depth (OTTD)', 'As
    ↪ determined by chest X-ray (in cm)'),
15 }
16
17 # TABLE 0
18 df = pd.read_pickle('table_0.pkl')
19
20 # RENAME ROWS AND COLUMNS
21 mapping = {k: v for k, v in shared_mapping.items() if
    ↪ is_str_in_df(df, k)}
22 abbrs_to_names, legend = split_mapping(mapping)
```

```

23 df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
24
25 # Save as latex
26 to_latex_with_note(
27     df, 'table_0.tex',
28     caption="Descriptive statistics of patient features and OTTD
29     ↪ stratified by sex",
30     label='table:patient_features_ottd_by_sex',
31     legend=legend)
32
33 # TABLE 1
34 df = pd.read_pickle('table_1.pkl')
35
36 # RENAME ROWS AND COLUMNS
37 mapping = {k: v for k, v in shared_mapping.items() if
38     ↪ is_str_in_df(df, k)}
39 abbrs_to_names, legend = split_mapping(mapping)
40 df = df.rename(columns=abbrs_to_names, index=abbrs_to_names)
41
42 # Save as latex
43 to_latex_with_note(
44     df, 'table_1.tex',
45     caption="Performance of Random Forest and Elastic Net models
46     ↪ in predicting OTTD",
47     label='table:model_performance',
48     legend=legend)
49
50 # TABLE 2
51 df = pd.read_pickle('table_2.pkl')
52
53 # FORMAT VALUES
54 df['p value'] = df['p value'].apply(format_p_value)
55
56 # RENAME ROWS AND COLUMNS
57 mapping = {k: v for k, v in shared_mapping.items() if
58     ↪ is_str_in_df(df, k)}
59 abbrs_to_names, legend = split_mapping(mapping)
60 df.rename(columns=abbrs_to_names, index=abbrs_to_names,
61     ↪ inplace=True)

```

```

58 # Save as latex
59 to_latex_with_note(
60     df, 'table_2.tex',
61     caption="Paired t-test comparing the squared residuals of the
        ↪ two models",
62     label='table:paired_ttest',
63     legend=legend)
64

```

## 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:
  ↪ 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}
26     return abbrs_to_names, names_to_definitions
27

```

### D.3 Code Output

table\_0.tex

```

\begin{table}[h]
\caption{Descriptive statistics of patient features and OTTD stratified by sex}
\label{table:patient_features_ottd_by_sex}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{llrr}
\toprule
& sex & 0 & 1 \\\
\midrule
\multirow[t]{2}{*}{\textbf{Tube ID}} & \textbf{mean} & 3.68 & 3.7 \\\
\textbf{} & \textbf{std} & 0.552 & 0.582 \\\
\cline{1-4}
\multirow[t]{2}{*}{\textbf{Age}} & \textbf{mean} & 0.732 & 0.781 \\\
\textbf{} & \textbf{std} & 1.4 & 1.47 \\\
\cline{1-4}
\multirow[t]{2}{*}{\textbf{Height}} & \textbf{mean} & 65.4 & 66.5 \\\
\textbf{} & \textbf{std} & 18.7 & 19.4 \\\
\cline{1-4}
\multirow[t]{2}{*}{\textbf{Weight}} & \textbf{mean} & 6.84 & 7.37 \\\
\textbf{} & \textbf{std} & 4.57 & 4.94 \\\
\cline{1-4}
\multirow[t]{2}{*}{\textbf{Optimal Tracheal Tube Depth (OTTD)}} & \textbf{mean}

```

```

& 10.1 & 10.3 \\
\textbf{} & \textbf{std} & 1.65 & 1.86 \\
\cline{1-4}
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item \textbf{Age}: Patient age (in years)
\item \textbf{Height}: Patient height (in cm)
\item \textbf{Weight}: Patient weight (in kg)
\item \textbf{Tube ID}: Internal diameter of the tube (in mm)
\item \textbf{Optimal Tracheal Tube Depth (OTTD)}: As determined by chest X-ray
(in cm)
\end{tablenotes}
\end{threeparttable}
\end{table}

```

table\_1.tex

```

\begin{table}[h]
\caption{Performance of Random Forest and Elastic Net models in predicting OTTD}
\label{table:model_performance}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{lr}
\toprule
& Mean Squared Error \\
Model & \\
\midrule
\textbf{Random Forest} & 1.57 \\
\textbf{Elastic Net} & 1.24 \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item
\end{tablenotes}
\end{threeparttable}

```



```
\end{table}
```

**table\_2.tex**

```
\begin{table}[h]
\caption{Paired t-test comparing the squared residuals of the two models}
\label{table:paired_ttest}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{%
\begin{tabular}{llrl}
\toprule
& Model 2 residuals & t statistic & p value \\
Model 1 residuals & & & \\
\midrule
\textbf{Random Forest} & Elastic Net & 2.36 & 0.0194 \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item
\end{tablenotes}
\end{threeparttable}
\end{table}
```