

1 Pre-analysis & Questions

The first objective is to observe what features may influence the cost of tubes. We start our observation with the features given in the training set where we define a variable for each feature.

Variable	Feature
x_1	tube_assembly_id
x_2	supplier
x_3	quote_date
x_4	annual_usage
x_5	min_order_quantity
x_6	bracket_pricing
x_7	quantity
x_8	cost

1.1 Tube Physical Properties

As a supplier, we have to think on which tube features the cost will be based. We know that a tube assembly is made with one or more components. Some numerical tube properties may be helpful to check.

- The weight
- The quantity
- The volume
- The number of bends used. Logically, it is more difficult to bend a tube than to keep it linear, so it should be more expensive.
- The component types used to assemble the tube.

1.2 Supplier Features

- The date when the supplier has quoted the price which is certainly less 20 years ago than today when not adjusted.
- The suppliers may use different mathematical models to quote their price.
- The supplier uses or not a bracket pricing which can influence what features to use in both cases.

1.3 Other Observations

- The tube assembly ID may be used at some points in the prediction. We need to investigate why and how these tubes are chosen in the train set.
- The costs have too many decimals to be a real cost. Maybe a conversion is needed in some way to get the real cost.
- The supplier may have decided to quote the tubes with very expensive or cheap price. These prices can be considered as anomalies.

1.4 Questions to Answer

To achieve our goal of predicting the cost with a good accuracy, we need to answer the following questions.

1. What features are used to determine the cost and what features to exclude from the analysis?

2. Do the costs presented in the training set are the real costs or do they need to be adjusted?
3. Is there a unique mathematical model describing the cost in function of the quantity for each supplier?
4. What particular features need to be classified? Why? How?
5. Are there anomalies to consider in the dataset? Why?

2 Preparing & Cleaning the Dataset

In this section, we will answer the first question: *What features are used to determine the cost and what features to exclude from the analysis?* We will also explain why we chose to keep and exclude features and how we will clean the dataset.

From the dataset, we note that there are a total of 2048 components. These components are spread among the `comp_[type].csv` files uniquely. This means that we can create a single table **Component** by merging those files together. To avoid to many columns, we will remove some features that we do not want in our analysis.

The file `bill_of_materials.csv` gives us the list of components with their respective quantity used to assemble a tube. Thus, to calculate the total weight for each tube, we use the formula

$$W_T = \sum_{i=0}^n W_i * Q_i$$

where W_T is the total weight of the tube T , $W = (W_1, \dots, W_n)$ is the vector of component weights, $Q = (Q_1, \dots, Q_n)$ the vector of component quantities and $n \leq 8$ the number of possible components used to assemble a tube T .

Let the total volume estimation of a tube assembly be denoted by V_T . The volume is function of the length, the wall thickness and the diameter of the tube and its formula is

$$V_T = \pi L t (d - t)$$

, where t is the wall thickness, d the outside diameter and L the developed length of the tube.

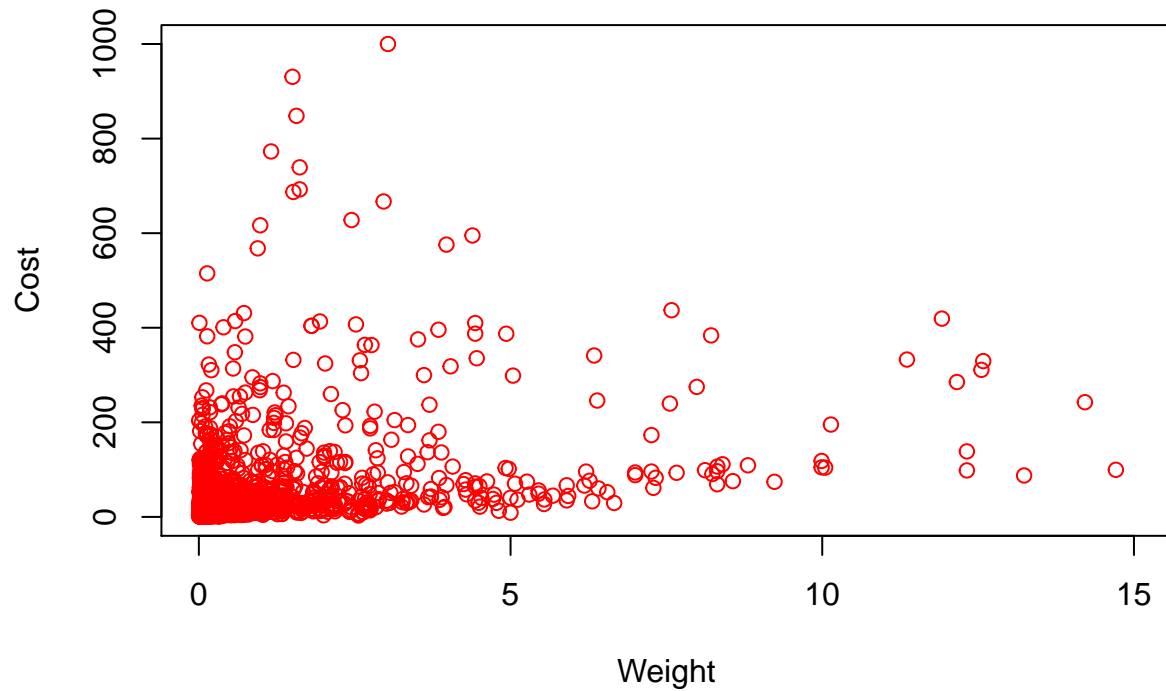
We use the Chi-Squared test to check if a feature is independent or dependent of the cost. We reject the hypothesis of independence if the p-value is greater than 0.05.

	tubeAssemblyID	totalWeight	volume	numberOfBends	cntQty	cost
1	2	0.02	1723.49	8	8.00	21.91
2	4	0.02	1723.49	9	8.00	21.97
3	5	0.21	7562.44	4	8.00	28.37
4	12	0.02	2604.10	7	8.00	22.42
5	13	0.21	20027.98	3	1.00	10.00
6	14	0.10	3170.01	4	8.00	21.99
7	21	0.05	2404.38	6	1.00	3.43
8	22	0.02	3335.12	6	1.00	8.56
9	24	0.03	1945.45	2	8.00	20.93
10	25	0.03	1277.31	3	8.00	20.78

Table 2:

Pearson's Chi-squared test

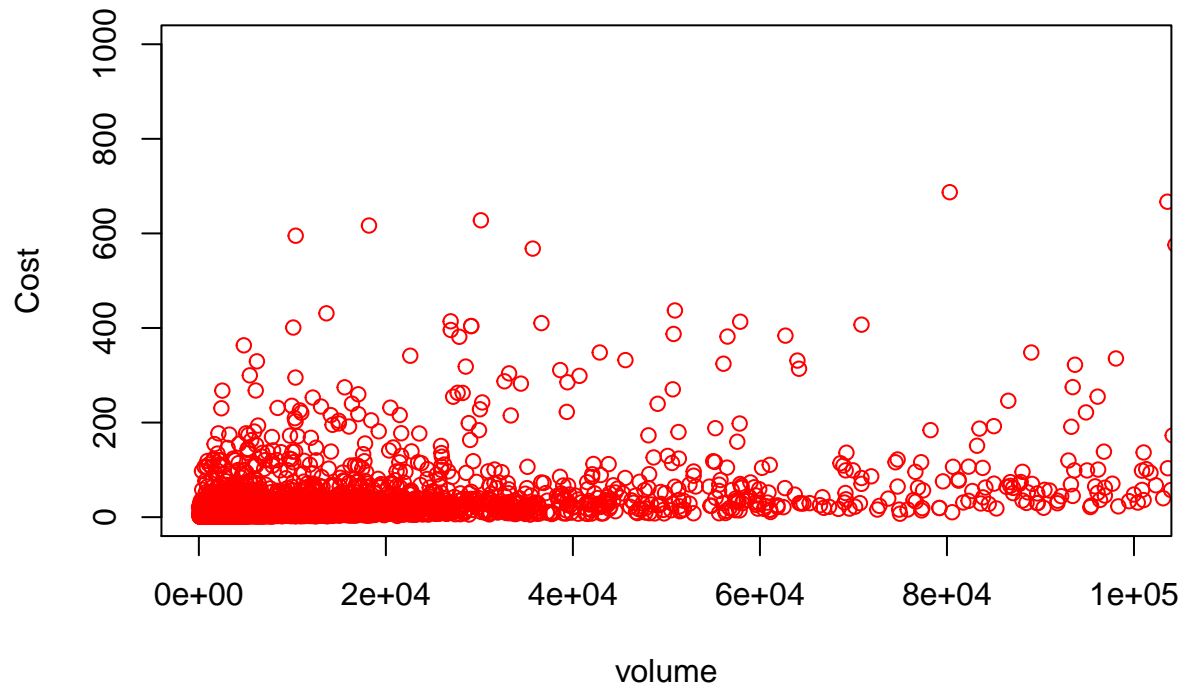
data: tbl X-squared = 5244100, df = 4262600, p-value < 2.2e-16



Following the Chi-Square test done and the plot, the cost is independent of the weight of a tube since the p-value is less than 0.05.

Pearson's Chi-squared test

data: tbl X-squared = 22286000, df = 20737000, p-value < 2.2e-16



Following the Chi-Square test done and the plot, the cost is independent of the volume of a tube since the p-value is less than 0.05.

Pearson's Chi-squared test

data: tbl X-squared = 87527, df = 93262, p-value = 1

Following the Chi-Square tests done, the cost is dependent of the number of bends since the p-value is greater than 0.05.

3 Real Cost vs Adjusted Cost

In this section, we will answer the question: *Do the costs presented in the training set are the real costs or do they need to be adjusted?*

4 Supplier's Model

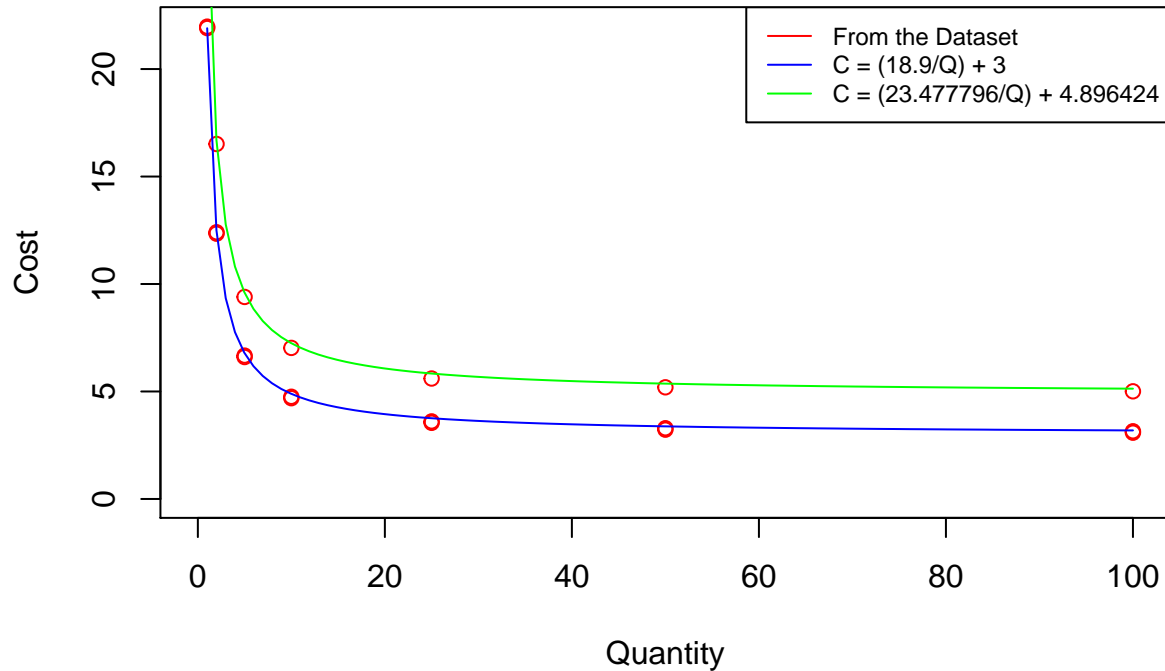
In this section, we will answer the question: *Is there a unique mathematical model describing the cost in function of the quantity for each supplier?* The first objective is to check the existence of a mathematical model representing the cost in function of the quantity. The second objective is to show if the model is applied by a unique supplier. The last objective is to show if each supplier has its own model. If the unicity does not hold, then we have to check if a model is applied by more than one supplier or if a supplier can apply more than one model depending of other features.

We denote $C_\beta(Q)$ our cost heuristic function of a tube assembly given by a supplier where β is our learning parameters and Q the vector of quantities.

4.1 Existence of a Mathematical Model

We start with the 3 first tube assemblies (TA-00002, TA-00004 and TA-00005) which are quoted by the supplier S-0066.

Cost of tubes 2,3,5 in function of the quantity (Supplier S-0066)



fkTubeAssembly	supplierID	quantity	cost
1	2 S-0066	1	21.91
2	2 S-0066	2	12.34
3	2 S-0066	5	6.60
4	2 S-0066	10	4.69
5	2 S-0066	25	3.54
6	2 S-0066	50	3.22
7	2 S-0066	100	3.08
8	2 S-0066	250	3.00
9	4 S-0066	1	21.97
10	4 S-0066	2	12.41
11	4 S-0066	5	6.67
12	4 S-0066	10	4.75
13	4 S-0066	25	3.61
14	4 S-0066	50	3.29
15	4 S-0066	100	3.15
16	4 S-0066	250	3.07
17	5 S-0066	1	28.37
18	5 S-0066	2	16.51
19	5 S-0066	5	9.40
20	5 S-0066	10	7.03
21	5 S-0066	25	5.60
22	5 S-0066	50	5.19
23	5 S-0066	100	5.01
24	5 S-0066	250	4.90

Table 3: Table built from Tube 2, 4 and 5

From the plot, we see that the curve representing the circles is clearly an hyperbola of equation

$$C_{\beta}(Q) = \frac{\beta_0}{Q} + \beta_1$$

where $Q \geq 1$, β_1 is the cost at the last level of purchase based on quantity and supplier, and $\beta_0 = C_{\beta}(1) - \beta_1$. This equation indicates that if Caterpillar buy more tubes, cheaper will be the cost per tube. This proves the existence of a mathematical model representing the cost in function of the quantity.

We need to find on which features depend β_0 and β_1 .

4.2 Unicity of the Model per Supplier

5 Classified Features

6 Anomalies Detection

fkTubeAssembly	supplierID	minQty	purchaseLevel	cost
5013	S-0054	1	8	7.135981
1243	S-0066	15	8	14.942439
18244	S-0054	1	8	16.781380
20621	S-0066	1	8	16.920483
20557	S-0066	1	8	16.953868
20558	S-0066	1	8	17.051240

fkTubeAssembly	supplierID	minQty	purchaseLevel	cost
8661	S-0066	1	8	17.056804
19148	S-0066	1	8	17.129138

	fkTubeAssembly	supplierID	minQty	purchaseLevel	cost
2194	19143	S-0066	1	8	55.65231
2195	20477	S-0066	1	8	59.63344
2196	20619	S-0066	1	8	60.37903
2197	20639	S-0066	1	8	63.07206
2198	18838	S-0066	1	8	66.51625
2199	20272	S-0066	1	8	68.06029
2200	20273	S-0066	1	8	68.11037
2201	5000	S-0066	1	8	141.65967

used (Mb) gc trigger (Mb) max used (Mb)

Ncells 543833 29.1 940480 50.3 843326 45.1 Vcells 797459 6.1 150343808 1147.1 187920775 1433.8