# Formulation and Applications of
# a Probabilistic Pareto Chart

Kenneth A. Hart,[*] Bradley A. Steinfeldt,[†] and Robert D. Braun[‡]

*Georgia Institute of Technology, Atlanta, GA, 30332-1510, USA*

**A probabilistic treatment of the Pareto chart can provide benefits to the fields of quality control, sensitivity analysis, and conceptual design. The probabilistic Pareto chart can inform a decision-maker about the relative significance of distributed factors and highlight anomalies in a dataset. This investigation provides a framework for creating a probabilistic Pareto chart, as well as examples to enable a discussion of the information provided by both the deterministic and probabilistic Pareto charts. The applications presented in this investigation demonstrate the probabilistic Pareto chart's ability to highlight anomalous trends and to determine the significance of variables in non-linear functions.**

## Nomenclature

| | |
|---|---|
| $\mathcal{R}$ | Universal gas constant |
| $C$ | Coefficient |
| $M$ | Molar mass |
| $T$ | Temperature |
| $V$ | Velocity |

*Subscripts*

| | |
|---|---|
| $\infty$ | Freestream |
| $D$ | Drag |
| $N$ | Normal |
| $T$ | Tangential |
| $W$ | Wall |

*Symbols*

| | |
|---|---|
| $\alpha$ | Angle of attack |
| $\beta$ | Sideslip angle |
| $\sigma$ | Accommodation coefficient |

## I.   Introduction

THE Pareto chart enables users to determine the relative contributions of factors to a whole and provides quantitative insights in a wide variety of disciplines.[1] The chart first appeared in 1951 in Juran's Quality Control Handbook as a technique for visualizing the Pareto principle.[2] The Pareto principle, also known as the 80/20 rule, states that 80% of the total contributions are generally due to 20% of the individual factors. Its origin is in early-20th century, when Vilfredo Pareto analyzed the distribution of wealth in Italy and it has since proliferated through the field of quality control and other industries as well.[2,3] A generic example of the Pareto chart is reproduced in Figure 1. As seen in Figure 1, the contribution of three of the defect categories (G, C, and A) contribute to 90% of the response. Similarly, the Pareto plot can also be used to show contributions to the variability in the response and the choice between examining the response or the variability in the response is situational.

---

[*]Graduate Research Assistant, Daniel Guggenheim School of Aerospace Engineering, AIAA Student Member.
[†]Research Engineer II, Daniel Guggenheim School of Aerospace Engineering, AIAA Member.
[‡]David & Andrew Lewis Professor of Space Technology, Daniel Guggenheim School of Aerospace Engineering, AIAA Fellow.

American Institute of Aeronautics and Astronautics

**Figure 1:** A Pareto chart for nine defect categories.

While there have been a few modifications to the Pareto chart in the past, it has usually remained in the same form since its inception.[1] The chart is comprised of a bar chart, sorted by bar height, and a line plot showing the cumulative total. Often there is a second, horizontal line indicating an acceptable level for the total contribution of the "vital few," a term coined by Juran.[2] One implicit assumption in the creation of a Pareto chart is that the contribution of each variable is deterministic, however in many cases it is probabilistic. For example, a chart illustrating defects from one factory would not be the same as a chart for another factory. Each of the factories could instead be considered a sample, so the bar heights would be distributions instead of deterministic levels. The distributions associated with each defect could be considered the error or the uncertainty in that defect's contribution to the response or to the variability in the response, depending on the application.

## II.   Motivation

Consider assessing relative significance of probabilistic or uncertain terms in an analytic equation. If these terms are insignificant, fixing them to constant values can reduce the complexity of the equation. This significance test is a sensitivity analysis and it is commonly performed using Monte Carlo simulation,[4–6] though other methods, such as Bayesian analysis and complex step, exist to assess this data.[7–10] Scatter plots are used to show the trend in the output of the equation with a single variable and if that trend appears to be non-random, then that variable is considered significant. Appropriate analysis of each individual scatter plot can be a time-consuming process and it provides no information as to which variables are the "vital few" and which are not. The probabilistic Pareto chart can be used to gain similar insights as the scatter plot simultaneously for many distributed variables.

Whether the data is sampled from the same source at different times, from different sources at the same time, or from different sources at different times is irrelevant to the process of generating a probabilistic Pareto chart. In the example of the factories, consider that the data on defects is collected across several factories and totaled on a monthly basis. If it is uncertain whether the relative contributions of each defect are time-dependent, a probabilistic Pareto chart could provide that insight. The data across all the factories could be summed for each month, and then each month would be considered a sample of the probabilistic Pareto chart. If the variances of each distribution are small, then deterministic Pareto charts for each month would vary only slightly and it could be concluded that the relative contributions of the defects are time-invariant. The same process could be applied to the data to determine whether there is significant variability between factories. Finally, if no summation across factories or across months is applied, then the probabilistic Pareto chart could show that the relative contributions of each defect do not depend on either factory or time.

In the factory example previously described, the probabilistic Pareto chart can show whether the distribution of defects is common across all factories, or whether each factory is significantly different where non-systematic causes should be investigated. If you consider Table 2 as a report of the defects from 30 factories, then generating a deterministic Pareto chart showing which defects are most frequent will show what the overall trend is, but not whether this general trend is followed by each of the factories. This overall trend can also lead to a false conclusion by way of Simpson's Paradox, which occurs when the trends in partitioned data differ from those in aggregated data.[11] Historically, failure to understand or to identify

American Institute of Aeronautics and Astronautics

instances of Simpson's Paradox has lead to poor decision-making, which can be prevented by making both the aggregated and the partitioned trends available to the decision-maker.[12] One possible way of analyzing the partitioned data would be to create a Pareto chart for each factory and compare them against each other. This process has been investigated previously,[13, 14] however 30 factories can be grouped into 435 pairs, so a pairwise comparison is not the ideal solution when there are many samples.

A probabilistic Pareto chart can eliminates the need for pairwise data comparison as it simultaneously illustrates the relative contributions of each variable and the portion of the population for which a given trend applies. Additionally, the chart does so in a manner that scales well with both the number of samples and the number of variables. Furthermore, outliers in the data can be readily identified when it is organized in a probabilistic Pareto chart.

## III.   Formulation

The data required to create a deterministic Pareto chart can take any form, so long as it can be tabulated. The information in the table could be displayed in the form of a pie chart. However, the goal of a Pareto chart is to distinguish the "vital few" from the "trivial many," so the data is represented in a bar chart sorted by bar height. Similarly the probabilistic Pareto chart requires data on the contributions of individual factors, however the key distinction is that the contributions are sampled more than once.

With any table of data similar to Table 2, a probabilistic Pareto chart can be created. The chart has a column for each possible factor and a row for each sample. There are ten steps required to visualize this table of data in a probabilistic Pareto chart which are outlined below. The sum in Step 4 provides the metric for sorting the columns in Step 5. To check that Steps 1-7 have been followed correctly, all the values in the last column of the second table should be ones.

1. Create a table of data in the same form as Table 2

2. Sum across the columns

3. Normalize the rows with respect to the total

4. Sum down the rows

5. Sort the columns by their totals in descending order

6. Create a second table that is equal in size to the original

7. Populate the table by calculating the running total across each row

8. Calculate the desired quantiles of the columns of the new table

9. Generate the box and whisker plots from the data in the original table

10. Generate the percentile line plots from the second table

The user of the probabilistic Pareto chart has several options and ways to vary the chart shown in Figure 2(b). For example, the percentiles to display on the probabilistic Pareto chart are the user's choice. If the goal is to determine whether the data depends on a specific variable then the user would be more likely to choose the mean and a 95% confidence interval. Another attribute that can be changed is showing the distributed data in box and whisker charts. Box and whisker charts are demonstrated in this investigation because they do not assume a form of the distributions, make outliers apparent, and do not obfuscate trends in the cumulative distribution lines. In other cases it may be more advantageous to display the distributions in other ways such as histograms, point density plots, violin plots, or as parameter-estimated distributions such as Beta distributions. Another aspect that can be varied is the color of the box and whisker charts. If the factors are sources of uncertainty, for example, then the box and whisker charts associated with epistemic uncertainties could be coded with one color and aleatory uncertainties can be coded with another color. In the case that motivated this investigation, the charts were colored based on whether a term in the rarefied aerodynamic coefficients of a cone had a closed-form or non-closed form solution.

American Institute of Aeronautics and Astronautics

(a) Deterministic Pareto Chart



(b) Probabilistic Pareto Chart

**Figure 2:** Comparison of deterministic and probabilistic Pareto charts for the same set of data.

American Institute of Aeronautics and Astronautics

# IV. Applications

## A. Identify Simpson's Paradox

The phenomenon of Simpson's Paradox arises when trends in partitioned data disappear when the data are combined.[11] Cases of Simpson's Paradox are often found in the social sciences, such as the 1973 investigation of gender bias in graduate admissions to the University of California, Berkeley. When admission was compared at the university-level, a greater portion of male applicants were admitted compared to female applicants; however, at the department-level there was no significant bias against female applications. In studying the data and comparing trends at the department- and university-level, it was concluded that men generally applied to less-competitive departments and women to more-competitive departments.[15] Though the original intention of the investigation was to determine gender discrimination, the conclusions from the data offered deeper insight into gender and graduate admissions because both the trends at the partition-level and aggregate-level were studied and reconciled.

Appendix B contains a table of data that demonstrates Simpson's Paradox as it applies to Pareto analysis. If the data from each trial is aggregated, then the deterministic Pareto chart in Figure 3(a) is generated. Partitioning the data for each trial results in Figure 3(b), which shows the distributions of each factor's significance. The vital few factors in the deterministic chart are $\{4, 5, 14, 15, 16\}$, while in the probabilistic chart the vital few are $\{13, 14, 15, 16, 17\}$. Re-partitioning the data and examining trials 1-27 separately from trials 28-30 shows that there are two separate trends. The frequency of defects is greater for trials 28-30 than for the other trials and in separate categories, indicating that these three trials should be analyzed separately from trials 1-27. The differences in the vital few between the deterministic and probabilistic Pareto charts indicates further investigation of datasets demonstrating Simpson's Paradox.

## B. Significance of Variables in Nonlinear Functions

Another application for the probabilistic Pareto chart is in determining the significance of input variables in a non-linear function. The significance of variables in a linear equation can be determined from the gradient and the ranking of the variables will remain constant no matter which combination of variables is input to the function. For a non-linear equation the terms in the gradient are functions of the inputs, so the ranking of the input variables can change from one set of inputs to the next. For most problems, the individual input variables are bounded in magnitude, and within that hypercube domain of inputs there may be variables that show consistent membership in the "trivial many".

$$
\begin{aligned}
C_D = {} & \frac{2\mathcal{R}\cos(\alpha)\cos(\beta)T_\infty}{MV_\infty^2}\left(\left(\operatorname{erf}\left(\frac{\cos(\alpha)\cos(\beta)V_\infty}{\sqrt{2}\sqrt{\frac{\mathcal{R}T_\infty}{M}}}\right)+1\right)\left(\frac{\sqrt{\frac{\pi}{2}}\cos(\alpha)\cos(\beta)\sigma_N V_\infty\sqrt{\frac{T_W}{T_\infty}}}{2\sqrt{\frac{\mathcal{R}T_\infty}{M}}}\right.\right. \\
& +(2-\sigma_N)\left(\frac{2\mathcal{R}\cos^2(\alpha)\cos^2(\beta)T_\infty}{MV_\infty^2}+\frac{1}{2}\right)\Bigg) \\
& +\frac{MV_\infty^2 e^{-\frac{M\cos^2(\alpha)\cos^2(\beta)V_\infty^2}{2\mathcal{R}T_\infty}}\left(\frac{1}{2}\sigma_N\sqrt{\frac{T_W}{T_\infty}}-\frac{\cos(\alpha)\cos(\beta)(\sigma_N-2)V_\infty}{\sqrt{2\pi}\sqrt{\frac{\mathcal{R}T_\infty}{M}}}\right)}{2\mathcal{R}T_\infty}\Bigg) \\
& +\sigma_T\left(\sin^2(\alpha)\cos(\beta)-\cos(\alpha)\sin^2(\beta)\right) \\
& \left(\cos(\alpha)\cos(\beta)\left(\operatorname{erf}\left(\frac{\cos(\alpha)\cos(\beta)V_\infty}{\sqrt{2}\sqrt{\frac{\mathcal{R}T_\infty}{M}}}\right)+1\right)+\frac{\sqrt{\frac{2}{\pi}}\sqrt{\mathcal{R}T_\infty}e^{-\frac{M\cos^2(\alpha)\cos^2(\beta)V_\infty^2}{2\mathcal{R}T_\infty}}}{\sqrt{M}V_\infty}\right)
\end{aligned}
\tag{1}
$$

The drag coefficient of a flat plate in rarefied flow is one example of a multivariable nonlinear function. The equation for the drag coefficient, given by Eq. 1 is derived from gas kinetic theory and is highly nonlinear in eight independent variables.[16] Not all of these variables may be known in application. For instance, an object in a low Earth orbit (LEO) may be large enough to determine its attitude, but it might not be possible to know the freestream temperature about that object. To determine which variables are significant over

American Institute of Aeronautics and Astronautics

## Deterministic Pareto Chart



(a) Deterministic Pareto Chart

## Probabilistic Pareto Chart



(b) Probabilistic Pareto Chart

**Figure 3:** An example of Simpson's Paradox as applied to Pareto analysis.

American Institute of Aeronautics and Astronautics

the hypercube domain of input variables, the inputs were normalized and a Monte Carlo simulation was performed on the partial derivatives of $C_D$. The ranges on these variables are provided in Table 1.

**Table 1:** Upper and lower bounds for variables in Eq. 1.

| Bound | $\alpha$ (deg) | $\beta$ (deg) | $V_\infty$ (m/s) | $T_\infty$ (K) | $M$ (g/mol) | $\sigma_N$ | $\sigma_T$ | $T_W$ (K) |
|---|---|---|---|---|---|---|---|---|
| | | | Input Variable | | | | | |
| Lower | -90 | -90 | 5500 | 200 | 2 | 0 | 0 | 100 |
| Upper | 90 | 90 | 9500 | 2000 | 46 | 1 | 1 | 500 |

In both of the Pareto charts shown in Figure 4, it is clear that the angle of attack and sideslip angle dominate the variability in the response. The next two terms are the accomodation coefficients, after that none of the variables significantly contribute to the variability in $C_D$. This indicates that fixing $M$, $T_\infty$, $T_W$, and $V_\infty$ would not significantly change the resulting $C_D$ value. This is evident in both charts by examination of the bar/box heights and the cumulative lines. Since $\alpha$ and $\beta$ have the same trigonometric effect on drag, the wide spread in their relative significance in Figure 4(b) can be attributed to the two angles trading off between each other. The accommodation coefficients are also related to each other, trading off the remainder of the variability.



(a) Deterministic Pareto Chart



(b) Probabilistic Pareto Chart

**Figure 4:** Relative significance of inputs to flat plate drag in rarefied flow.

American Institute of Aeronautics and Astronautics

## V.    Summary

This investigation shows how a probabilistic Pareto chart can be created from data that are aggregated into deterministic Pareto charts. The primary advantages demonstrated are the identification of anomalous data and illustration of the relative significance of variables in nonlinear functions. In order to provide insight into whether a global trend is consistent with the individual trends and to improve the outcomes of decisions made when Simpson's Paradox may apply, both the deterministic and the probabilistic Pareto charts are necessary to the analyst.

One aspect of the probabilistic Pareto chart that can be explored in the future is the representation of the distributed data. The charts presented above use box-and-whisker plots to show this distribution, however other representations may be more revealing. Options for displaying this data include mixed Gaussians or violin charts, histograms, and beta distributions. The mixed Gaussians and histograms would reveal modality in the data, though beta distributions and mixed Gaussians assume a distribution *a priori*. Another related investigation would be to determine which cumulative percentile lines to include on the chart. The probabilistic charts shown above have three cumulative percentile lines drawn, though which three are drawn varies. It may be more useful to draw lines at 0%, 80%, and 100% to show the range of the data and a relevant cutoff. The usefulness of the probabilistic Pareto chart has been shown above, though how it should be integrated with existing quality control tools should be investigated.

## References

[1]Wilkinson, L., "Revising the Pareto Chart," *The American Statistician*, Vol. 60, No. 4, Nov. 2006, pp. 332–334.

[2]Juran, J. M. and Gryna, F. M., *Quality Control Handbook*, McGraw-Hill, Inc, New York, 1951.

[3]Pareto, V., *Manuale di economia politica*, Societa Editrice, Milan, 1906.

[4]Claxton, K., Sculpher, M., McCabe, C., Briggs, A., Akehurst, R., Buxton, M., Brazier, J., and O'Hagan, T., "Probabilistic sensitivity analysis for NICE technology assessment: not an optional extra," *Health economics*, Vol. 14, No. 4, April 2005, pp. 339–347.

[5]Boyle, P., Broadie, M., and Glasserman, P., "Monte Carlo methods for security pricing," *Journal of Economic Dynamics & Control*, Vol. 21, 1997, pp. 1267–1321.

[6]Critchfield, G. C., Willard, K. E., and Connelly, D. P., "Probabilistic sensitivity analysis methods for general decision models," *Computer and Biomedical Research*, Vol. 19, 1986, pp. 254–265.

[7]Oakley, J. E. and O'Hagan, A., "Probabilistic sensitivity analysis of complex models: a Bayesian approach," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, Vol. 66, No. 3, Aug. 2004, pp. 751–769.

[8]Wu, Y. T., "Computational methods for efficient structural reliability and reliability sensitivity analysis," *AIAA Journal*, Vol. 32, No. 8, 1994, pp. 1717–1723.

[9]Homma, T. and Saltelli, A., "Importance measures in global sensitivity analysis of nonlinear models," *Reliability Engineering & System Safety*, Vol. 52, 1996, pp. 1–17.

[10]Martins, J. R. R. A., Kroo, I. M., and Alonso, J. J., "An Automated Method for Sensitivity Analysis Using Complex Variables," *Proceedings of the 38th AIAA Aerospace Sciences Meeting*, Reno, NV, Jan. 2000.

[11]Blyth, C. R., "On Simpson's Paradox and the Sure-Thing Principle," *Journal of the American Statistical Association*, Vol. 67, No. 338, 1972, pp. 364–366.

[12]Pearl, J., Causality: Models, Reasoning, and Inference, Vol. 29, Cambridge University Press, Sept. 2009.

[13]Grzegorzewski, P., "Statistical Tests for Comparing Pareto Charts," *Communications in Numerical Analysis*, Vol. 2012, 2012, pp. 12.

[14]Kenett, R. S., "Two methods for comparing Pareto charts," *Journal of quality technology*, Vol. 23, No. 1, 1991, pp. 27–31.

[15]Bickel, P. J., Hammel, E. A., and O'Connell, J. W., "Sex Bias in Graduate Admissions: Data from Berkeley," *Science*, Vol. 187, No. 4175, 1975, pp. 398–404.

[16]Sentman, L. H., "Free Molecule Flow Theory and Its Application to the Determination of Aerodynamic Forces," Tech. Rep. LMSC-448514, Lockheed Missiles & Space Company, Sunnyvale, CA, Oct. 1961.

# Appendix A: Data from Figure 2

**Table 2:** Number of defects associated with nine unique categories

| Trial | A | B | C | D | E | F | G | H | I |
|-------|------|-----|-----|---|----|---|-----|----|----|
| 1 | 3 | 1 | 10 | 0 | 1 | 0 | 72 | 0 | 1 |
| 2 | 99 | 17 | 114 | 0 | 30 | 1 | 201 | 1 | 3 |
| 3 | 71 | 12 | 183 | 0 | 56 | 1 | 216 | 7 | 8 |
| 4 | 1 | 1 | 59 | 0 | 1 | 0 | 60 | 0 | 0 |
| 5 | 72 | 10 | 85 | 0 | 51 | 0 | 203 | 5 | 8 |
| 6 | 36 | 8 | 96 | 0 | 20 | 3 | 332 | 4 | 5 |
| 7 | 4 | 0 | 5 | 0 | 0 | 0 | 63 | 0 | 0 |
| 8 | 12 | 2 | 24 | 0 | 2 | 0 | 192 | 0 | 0 |
| 9 | 7 | 1 | 16 | 0 | 2 | 0 | 61 | 0 | 0 |
| 10 | 29 | 26 | 39 | 0 | 27 | 0 | 112 | 1 | 1 |
| 11 | 19 | 2 | 20 | 0 | 7 | 0 | 58 | 0 | 1 |
| 12 | 23 | 11 | 47 | 0 | 17 | 0 | 55 | 0 | 9 |
| 13 | 67 | 5 | 83 | 0 | 11 | 1 | 107 | 3 | 4 |
| 14 | 23 | 11 | 29 | 0 | 19 | 0 | 46 | 5 | 10 |
| 15 | 19 | 2 | 19 | 0 | 2 | 0 | 68 | 0 | 0 |
| 16 | 60 | 10 | 100 | 0 | 10 | 0 | 261 | 0 | 0 |
| 17 | 24 | 2 | 26 | 0 | 2 | 0 | 120 | 0 | 0 |
| 18 | 15 | 2 | 29 | 0 | 3 | 0 | 567 | 0 | 0 |
| 19 | 40 | 1 | 90 | 0 | 16 | 0 | 315 | 0 | 0 |
| 20 | 19 | 10 | 26 | 0 | 16 | 0 | 75 | 2 | 6 |
| 21 | 4 | 3 | 18 | 0 | 3 | 0 | 23 | 0 | 0 |
| 22 | 18 | 2 | 29 | 0 | 13 | 0 | 54 | 0 | 0 |
| 23 | 19 | 3 | 43 | 0 | 4 | 0 | 99 | 0 | 1 |
| 24 | 8 | 1 | 22 | 0 | 3 | 0 | 55 | 0 | 1 |
| 25 | 85 | 7 | 138 | 0 | 15 | 0 | 199 | 0 | 1 |
| 26 | 6 | 4 | 6 | 0 | 5 | 0 | 46 | 0 | 4 |
| 27 | 10 | 2 | 13 | 0 | 9 | 0 | 131 | 0 | 0 |
| 28 | 126 | 67 | 296 | 0 | 87 | 0 | 309 | 1 | 33 |
| 29 | 96 | 25 | 221 | 0 | 41 | 1 | 320 | 10 | 11 |
| 30 | 34 | 11 | 106 | 0 | 17 | 0 | 245 | 0 | 7 |

American Institute of Aeronautics and Astronautics

# Appendix B: Data from Figure 3

**Table 3:** Number of defects associated with twenty unique categories

| | Factor | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Trial | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 26 | 283 | 520 | 161 | 8 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 20 | 255 | 528 | 184 | 10 | 0 | 0 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 167 | 522 | 276 | 24 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 134 | 504 | 319 | 34 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 13 | 202 | 531 | 234 | 17 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 38 | 335 | 496 | 124 | 5 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 46 | 362 | 479 | 107 | 4 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 15 | 222 | 531 | 214 | 14 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 90 | 457 | 391 | 56 | 1 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 80 | 442 | 408 | 63 | 1 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 15 | 218 | 531 | 218 | 15 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 79 | 439 | 411 | 65 | 1 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 88 | 454 | 394 | 57 | 1 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 15 | 222 | 531 | 215 | 14 | 0 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 47 | 365 | 477 | 105 | 3 | 0 | 0 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 39 | 338 | 494 | 121 | 5 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 13 | 203 | 531 | 233 | 17 | 0 | 0 | 0 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 132 | 503 | 322 | 34 | 0 | 0 | 0 | 0 |
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 164 | 521 | 279 | 25 | 0 | 0 | 0 |
| 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 20 | 253 | 528 | 186 | 11 | 0 | 0 |
| 21 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 26 | 285 | 519 | 159 | 8 | 0 | 0 |
| 22 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 15 | 222 | 531 | 215 | 14 | 0 | 0 | 0 |
| 23 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 162 | 520 | 281 | 25 | 0 | 0 | 0 | 0 |
| 24 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 10 | 183 | 527 | 256 | 21 | 0 | 0 | 0 | 0 |
| 25 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 23 | 272 | 523 | 170 | 9 | 0 | 0 | 0 |
| 26 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 33 | 317 | 505 | 136 | 6 | 0 | 0 |
| 27 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 17 | 235 | 531 | 202 | 12 | 0 | 0 |
| 28 | 0 | 0 | 27 | 632 | 2491 | 1658 | 186 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 29 | 0 | 20 | 543 | 2405 | 1798 | 227 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 30 | 0 | 72 | 1070 | 2659 | 1116 | 79 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |