

LESSONS FROM THE GERMAN TANK PROBLEM

GEORGE CLARK, ALEX GONYE, AND STEVEN J. MILLER

ABSTRACT. During World War II the German army used tanks to devastating advantage. The Allies needed accurate estimates of their tank production and deployment. They used two approaches to find these values: spies, and statistics. This note describes the statistical approach. Assuming the tanks are labeled consecutively starting at 1, if we observe k serial numbers from an unknown number N of tanks, with the maximum observed value m , then the best estimate for N is $m(1+1/k) - 1$. This is now known as the German Tank Problem, and is a terrific example of the applicability of mathematics and statistics in the real world. The first part of the paper reproduces known results, specifically deriving this estimate and comparing its effectiveness to that of the spies. The second part presents a result we have not found in print elsewhere, the generalization to the case where the smallest value is not necessarily 1. We emphasize in detail why we are able to obtain such clean, closed-form expressions for the estimates, and conclude with an appendix highlighting how to use this problem to teach regression and how statistics can help us find functional relationships.

CONTENTS

| | |
|---|----|
| 1. Introduction | 2 |
| 2. Derivation with a known minimum | 4 |
| 2.1. The probability that the sample maximum is m | 5 |
| 2.2. The best guess for \hat{N} | 6 |
| 3. Derivation with an unknown minimum | 7 |
| 3.1. The probability that the spread is s | 7 |
| 3.2. The best guess for \hat{N} | 8 |
| 4. Comparison of Approaches | 10 |
| Appendix A. The German Tank Problem and Linear Regression | 11 |
| A.1. Theory of Regression | 11 |
| A.2. Issues in Applying to the German Tank Problem | 12 |
| A.3. Resolving Implementation Issues | 14 |
| A.4. Determining the Functional Form | 15 |
| References | 16 |

Date: January 22, 2021.

2010 Mathematics Subject Classification. 62J05, 60C05 (primary), 05A10 (secondary).

Key words and phrases. German Tank Problem, Binomial Identities, Regression.

The third named author was partially supported by NSF Grant DMS1561945. Parts of this paper were given by the first two named authors as their senior colloquium; we thank our colleagues in the department for comments. The third author has presented this topic at several conferences and programs, and thanks the participants for many suggestions; for video see <https://www.youtube.com/watch?v=quV-MCB8Ozs> or <https://youtu.be/I3ngtIYjw3w>.

We thank our colleagues at Williams and audiences at several talks over the years, and the referee, for valuable feedback that improved the exposition, and Jason Zhou for pointing out some typos in an earlier draft.

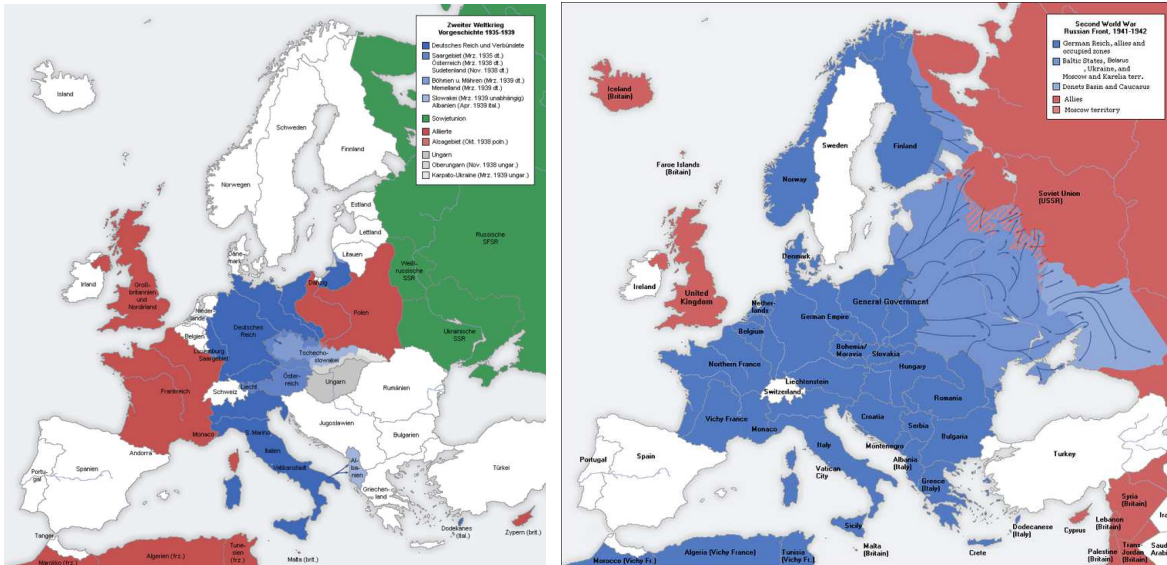


FIGURE 1. Left: Europe before the start of major hostilities. Right: Europe in 1942. Images from Wikimedia Commons from author San Jose.

1. INTRODUCTION

In this paper we revisit a famous and historically important problem which has since become a staple in many probability and statistics classes: the German Tank Problem. This case study illustrates that one does not need to use the most advanced mathematics to have a tremendous impact on real world problems; the challenge is frequently in creatively using what one knows.¹

By the end of 1941 most of continental Europe had fallen to Nazi Germany and the other Axis powers, and by 1942 their forces had begun their significant advance in the Eastern Front deep into the Soviet Union; see Figure 1 for an illustration of their rapid progress, or <https://www.youtube.com/watch?v=WOVEy1tC7nk> for an animation of territorial gains day by day. A key component to their rapid conquests lay with their revolutionary use of tanks in modern warfare. While other militaries, most notably France, used tanks as a modern, armored form of cavalry, the Germans were the first to fully utilize tanks' speed and strength to their advantage. Tanks would move rapidly and punch through enemy lines creating gaps which German infantry would stream through. Once through the holes in the line, the Germans would wreak havoc on lines of communication, creating logistical nightmares for the combatants left on the front lines. This lightning fast warfare has become dubbed Blitzkrieg (or lighting war) by many historians.

With the Nazis utilizing tanks with such devastating results, it was essential for the Allies to stop them. A key component to the solution was figuring out how many tanks the Germans were building, or had deployed in various theaters, in order to properly allocate resources. As expected, they tried spying (both with agents and through decrypting intercepted messages) to accurately estimate these numbers. It was essential that the Allies obtain accurate values, as there is a tremendous danger in both under- and over-estimating the enemy's strength. The consequence of underestimating is clear, as one could suddenly be outnumbered in battle. Overestimating is also bad, as it can lead

¹Another great example is the famous Battle of Midway and the role the cryptographers played in figuring out the Japanese target; see for example [1].

to undo caution and failure to exploit advantages, or to committing too many resources in one theater and thus not having enough elsewhere. The U.S. Civil War provides a terrific example of these consequences, where General George McClellan would not have his Union army take the field against the Confederates because, in his estimation, his forces were greatly outnumbered though in fact they were not. This situation led to one of President Abraham Lincoln's many great quips: *If General McClellan isn't going to use his army, I'd like to borrow it for a time.*² Considering how close Pickett's charge came to succeeding at Gettysburg, or what would have happened if Sherman hadn't taken Atlanta before the 1864 elections (where McClellan, now as the Democrat nominee for president, was running against Lincoln on a platform of a negotiated peace), the paralysis from incorrect analysis could have changed the outcome of the war.

Returning to World War II and the problem of determining the number of tanks produced and facing them in the field, the Allies understandably wanted to have a way to evaluate the effectiveness of their estimates. During a battle, they realized that the destroyed and captured tanks had serial numbers on their gearboxes which could help with this problem. Assuming that the serial numbers are sequential and start with 1, given a string of observed numbers one can try to estimate the largest value. This discovery contributed to the birth of the statistical method, the use of a population maximum formula, and changed both the war and science.³

The original variant of the problem assumes the first tank is numbered 1, there are an unknown number N produced (or in theater), that the numbers are consecutive, and that k values are observed with the largest being m . Given this information, the goal is to find \hat{N} , the best estimate of N . The formula they derived is

$$\hat{N} = m \left(1 + \frac{1}{k} \right) - 1.$$

For the reader unfamiliar with this subject, we deliberately do not state here in the introduction how well this formula does versus what was done by spies and espionage.⁴ While this story has been well told before (see for example [3, 4, 5, 6]), our contribution is to extend the analysis to consider the more general case, namely what happens when we do not know the smallest value. To our knowledge this result has not been isolated in the literature; we derive in §3 that if s is the spread between the smallest and the largest of the observed k serial numbers, then

$$\hat{N} = s \left(1 + \frac{2}{k-1} \right) - 1.$$

In Appendix A we show how to use regression to show that these are reasonable formulas, and thus the German Tank Problem can also be used to introduce some problems and subtleties in regression analysis, as well as serve as an introduction to mathematical modeling.

As it is rare to have a clean, closed-form expression such as the ones above, we briefly remark on our fortune. The key observation is that we have a combinatorial problem where certain binomial identities are available, and these lead to tremendous simplifications.

²There are many different phrasings of this remark; this one is taken from <https://thehistoriansmanifesto.wordpress.com/2013/05/13/best-abraham-lincoln-quotes/>.

³There are advantages in consecutive labeling; for example, it simplifies some maintenance issues as it is clear which of two tanks is older.

⁴That said, as this paper is appearing in a mathematics journal and not a bulletin of a spy agency, we invite the reader to conjecture which method did better.

2. DERIVATION WITH A KNOWN MINIMUM

In this section we prove

$$\widehat{N} = m \left(1 + \frac{1}{k} \right) - 1$$

when we observe k tanks, the largest labeled m , and knowing that the smallest number is 1 and the tanks are consecutively numbered. Before proving it, as a smell test we look at some extreme cases. First, we never obtain an estimate that is less than the largest observed number. Second, if there are many tanks and we observe just one (so $k = 1$), then \widehat{N} is approximately $2m$. This is very reasonable, and essentially just means that if we only have one data point, it's a good guess that it was in the middle. Further, as k increases the amount we must inflate our observed maximum value decreases. For example, when $k = 2$ we inflate m by approximately a factor of $3/2$, or in other words this is saying our observed maximum value is probably about two-thirds of the true value. Finally, if k equals the number of tanks N , then m must also equal N , and the formula simplifies to $\widehat{N} = N$.

We break the proof into two parts. While we are fortunate in that we are able to obtain a closed-form expression, if we have a good guess as to the relationship we can use statistics to test its reasonableness; we do that in Appendix A. For the proof we first determine the probability that the observed largest value is m . Next we compute the expected value, and show how to pass from that to an estimate for N . We need two combinatorial results.

The first is Pascal's identity:

$$\binom{n+1}{r} = \binom{n}{r} + \binom{n}{r-1}. \tag{2.1}$$

There are many approaches to proving this; the easiest is to interpret both sides as two different ways of counting how many ways we can choose r people from a group of $n+1$ people, where exactly n of these people are in one set and exactly one person is in another set. It is easier to see if we rewrite it as

$$\binom{n+1}{r} = \binom{1}{0} \binom{n}{r} + \binom{1}{1} \binom{n}{r-1};$$

this is permissible because $\binom{1}{0} = \binom{1}{1} = 1$. Note the left side is choosing r people from the combined group of $n+1$ people, while the right is choosing r people with the first summand corresponds to not choosing the person from the group with just one person, and the second summand to requiring we chose that person. □

The second identity involves sums of binomial coefficients:

$$\sum_{m=k}^N \binom{m}{k} = \binom{N+1}{k+1}. \tag{2.2}$$

We can prove this by induction on N , noting that k is fixed. The base case is readily established. Letting $N = k$, we find

$$\sum_{m=k}^N \binom{m}{k} = \sum_{m=k}^k \binom{m}{k} = \binom{k}{k} = 1 = \binom{k+1}{k+1}.$$

For the inductive step, we assume

$$\sum_{m=k}^N \binom{m}{k} = \binom{N+1}{k+1}.$$

Then

$$\begin{aligned} \sum_{m=k}^{N+1} \binom{m}{k} &= \left(\sum_{m=k}^N \binom{m}{k} \right) + \binom{N+1}{k} \\ &= \binom{N+1}{k+1} + \binom{N+1}{k} \\ &= \binom{N+2}{k+1}, \end{aligned}$$

where the last equality follows from Pascal's identity, (2.1). This completes the proof. \square

While this identity suffices for the original formulation of the German Tank Problem, when we do not know the starting serial number the combinatorics become slightly more involved, and we need a straightforward generalization:

$$\sum_{\ell=a}^b \binom{\ell}{a} = \binom{b+1}{a+1}; \quad (2.3)$$

the proof follows similarly.

2.1. The probability that the sample maximum is m . Let M be the random variable for the maximum number observed, and let m be the value we see. Note that there is zero probability of observing a value smaller than k or larger than N . We claim for $k \leq m \leq N$ that

$$\text{Prob}(M = m) = \frac{\binom{m}{k} - \binom{m-1}{k}}{\binom{N}{k}} = \frac{\binom{m-1}{k-1}}{\binom{N}{k}}.$$

We give two proofs. The first is to note there are $\binom{N}{k}$ ways to choose k numbers from N when order does not matter. The probability that the largest observed is exactly m equals the probability the largest is at most m minus the probability the largest is at most $m-1$. The first probability is just $\binom{m}{k} / \binom{N}{k}$, as if the largest value is at most m then all k observed numbers must be taken from $\{1, 2, \dots, m\}$. A similar argument gives the second probability is $\binom{m-1}{k} / \binom{N}{k}$, and the claim now follows by using Pascal's identity to simplify the difference of the binomial coefficients.

We could also argue as follows. If the largest is m then we have to choose that serial number, and now we must choose $k-1$ tanks from the $m-1$ smaller values; thus we find the probability is just $\binom{m-1}{k-1} / \binom{N}{k}$. \square

Remark 2.1. *Interestingly, we can use the two equivalent arguments above as yet another way to prove the Pascal identity.*

2.2. **The best guess for \hat{N} .** We now compute the best guess for N by first finding the expected value of M . Recall the expected value of a random variable M is the sum of all the possible values of M times the probability of observing that value. We write $\mathbb{E}[M]$ for this quantity, and thus we must compute

$$\mathbb{E}[M] := \sum_{m=k}^N m \text{Prob}(M = m)$$

(note we only need to worry about m in this range, as for all other m the probability is zero and thus does not contribute). Once we find a formula for $\mathbb{E}[M]$ we will convert that to one for the expected number of tanks.

Our first step is to substitute in the probability that M equals m , obtaining

$$\mathbb{E}[M] = \sum_{m=k}^N m \frac{\binom{m-1}{k-1}}{\binom{N}{k}}.$$

Fortunately this sum can be simplified into a nice closed-form expression; it is this simplification that allows us to obtain a simple formula for \hat{N} . We expand the binomial coefficients in the expression for $\mathbb{E}[M]$ and then use our second combinatorial identity, (2.2), to simplify the sum of $\binom{m}{k}$ which emerges as we manipulate the quantities below. We find

$$\begin{aligned} \mathbb{E}[M] &= \sum_{m=k}^N m \frac{\binom{m-1}{k-1}}{\binom{N}{k}} \\ &= \sum_{m=k}^N m \frac{(m-1)!}{(k-1)!(m-k)!} \frac{k!(N-k)!}{N!} \\ &= \sum_{m=k}^N \frac{m!k}{k!(m-k)!} \frac{k!(N-k)!}{N!} \\ &= \frac{k \cdot k!(N-k)!}{N!} \sum_{m=k}^N \binom{m}{k} \\ &= \frac{k \cdot k!(N-k)!}{N!} \binom{N+1}{k+1} \\ &= \frac{k \cdot k!(N-k)!}{N!} \frac{(N+1)!}{(k+1)!(N-k)!} \\ &= \frac{k(N+1)}{k+1}. \end{aligned}$$

As we have such a clean expression, it's trivial to solve for N in terms of k and $\mathbb{E}[M]$:

$$N = \mathbb{E}[M] \left(1 + \frac{1}{k}\right) - 1.$$

Thus if we substitute in m (our observed value for M) as our best guess for $\mathbb{E}[M]$, we obtain our estimate for the number of tanks produced:

$$\hat{N} = m \left(1 + \frac{1}{k}\right) - 1,$$

completing the proof. □

Remark 2.2. *A more advanced analysis can prove additional results about our estimator, for example, whether or not it is unbiased.*

Remark 2.3. *There are many ways to see this formula is reasonable. The first is to try extreme cases, such as $k = N$ (which forces m to be N and gives N as the answer), or to try $k = 1$. In that case we expect our one observation to be around $N/2$, and thus a formula that has the best guess being doubling the observation is logical. We can also get close to this formula from by trying to guess the functional form (for more details see Appendix A). We know our best guess must be at least m , so let's write it as $m + f(m, k)$. For a fixed k as m increases we might expect our guess to increase, while for fixed m as k increases we would expect a smaller boost. These heuristics suggest $f(m, k)$ increases with m and decreases with k ; the simplest such function is bm/k for some constant b . This leads to a guess of $m + bm/k$, and again looking at extreme cases we get very close to the correct formula.*

3. DERIVATION WITH AN UNKNOWN MINIMUM

Not surprisingly, when we do not know the lowest serial number the resulting algebra becomes more involved; fortunately, though, with a bit of work we are still able to get nice closed-form expressions for the needed sums and obtain again a clean answer for the estimated number of tanks. We still assume the tanks are numbered sequentially, and focus on the spread (the difference between the largest and smallest observed values). Similar to the previous section, we derive a formula to inflate the observed spread to be a good estimate of the number of total tanks.

We first set some notation:

- the minimum tank serial value, N_1 ,
- the maximum tank serial value, N_2 ,
- the total number of tanks, N ($N = N_2 - N_1 + 1$),
- the observed minimum value, m_1 (with corresponding random variable M_1),
- the observed maximum value, m_2 (with corresponding random variable M_2),
- the observed spread s (with corresponding random variable S).

As $s = m_2 - m_1$, in the arguments below we can focus on just s and S . We will prove the best guess is $s \left(1 + \frac{2}{k-1}\right) - 1$.

Remark 3.1. *There are two differences between this formula and the case when the smallest serial number is known. The first is we divide by $k - 1$ and not k ; however, as we cannot estimate a spread with one observation this is reasonable. Note the similarity here with the sample standard deviation, where we divide by one less than the number of observations; while one point suffices to estimate a mean, we need at least two for the variance. The second difference is that we have a factor of 2, which can be interpreted as movement in both directions.*

3.1. The probability that the spread is s . We claim that if we observe k tanks then for $k - 1 \leq s \leq N_2 - N_1$ we have

$$\text{Prob}(S = s) = \frac{\sum_{m=N_1}^{N_2-s} \binom{s-1}{k-2}}{\binom{N_2-N_1+1}{k}} = \frac{(N_2 - N_1 + 1 - s) \binom{s-1}{k-2}}{\binom{N_2-N_1+1}{k}} = \frac{(N - s) \binom{s-1}{k-2}}{\binom{N}{k}},$$

and for all other s the probability is zero.

To see this, note that the spread s must be at least $k - 1$ (as we have k observations), and cannot be larger than $N_2 - N_1$. If we want a spread of s , if the smallest observed value is m then the largest is $m + s$. We must choose exactly $k - 2$ of the $s - 1$ numbers in $\{m + 1, m + 2, \dots, m + s - 1\}$; there are $\binom{s-1}{k-2}$ ways to do so. This proves the first equality, the sum over m . As all the summands are the same we get the second equality, and the third follows from our definition of N . \square

3.2. The best guess for \widehat{N} . We argue similarly as in the previous section. In the algebra below we will use our second binomial identity, (2.2); relabeling the parameters it is

$$\sum_{\ell=a}^b \binom{\ell}{a} = \binom{b+1}{a+1}. \quad (3.1)$$

We begin by computing the expected value of the spread. We include all the details of the algebra; the idea is to manipulate the expressions and pull out terms that are independent of the summation variable, and rewrite expressions so that we can identify binomial coefficients and then apply our combinatorial results. We have

$$\begin{aligned} \mathbb{E}[S] &= \sum_{s=k-1}^{N-1} s \text{Prob}(S = s) \\ &= \sum_{s=k-1}^{N-1} s \frac{\binom{N-s}{k-2}}{\binom{N}{k}} \\ &= \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} s (N-s) \binom{s-1}{k-2} \\ &= \binom{N}{k}^{-1} N \sum_{s=k-1}^{N-1} \frac{s(s-1)!}{(s-k+1)!(k-2)!} - \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{s^2(s-1)!}{(s-k+1)!(k-2)!} \\ &= \binom{N}{k}^{-1} N \sum_{s=k-1}^{N-1} \frac{s!(k-1)}{(s-k+1)!(k-1)!} - \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{ss!(k-1)}{(s-k+1)!(k-1)!} = T_1 - T_2. \end{aligned}$$

We first simplify T_1 ; below we always try to multiply by 1 in such a way that we can combine ratios of factorials into binomial coefficients:

$$\begin{aligned} T_1 &= \binom{N}{k}^{-1} N \sum_{s=k-1}^{N-1} \frac{s!(k-1)}{(s-k+1)!(k-1)!} \\ &= \binom{N}{k}^{-1} N(k-1) \sum_{s=k-1}^{N-1} \binom{s}{k-1} \\ &= \binom{N}{k}^{-1} N(k-1) \binom{N}{k} = N(k-1), \end{aligned}$$

where we used (3.1) with $a = k - 1$ and $b = N - 1$.

Turning to T_2 we argue similarly, at one point replacing s with $(s - 1) + 1$ to assist in collecting factors into a binomial coefficient:

$$\begin{aligned}
T_2 &= \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{s s!(k-1)}{(s-k+1)!(k-1)!} \\
&= \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{(s+1-1)s!(k-1)}{(s-(k-1))!(k-1)!} \\
&= \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{(s+1)!(k-1)k}{(s+1-k)!(k-1)!k} - \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{s!(k-1)}{(s-(k-1))!(k-1)!} \\
&= \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \frac{(s+1)!k(k-1)}{(s+1-k)!k!} - \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} (k-1) \binom{s}{k-1} \\
&= \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} k(k-1) \binom{s+1}{k} - \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} (k-1) \binom{s}{k-1} = T_{21} + T_{22}.
\end{aligned}$$

We can immediately evaluate T_{22} by using (2.3) with $a = k - 1$ and $b = N - 1$, and find

$$T_{22} = \binom{N}{k}^{-1} (k-1) \binom{N}{k} = k-1.$$

Thus all that remains is analyzing T_{21} :

$$T_{21} = \binom{N}{k}^{-1} \sum_{s=k-1}^{N-1} \binom{s+1}{k} k(k-1).$$

We pull $k(k-1)$ outside the sum, and letting $w = s + 1$ we see that

$$T_{21} = \binom{N}{k}^{-1} k(k-1) \sum_{w=k}^N \binom{w}{k},$$

and then from (2.3) with $a = k$ and $b = N$ we obtain

$$T_{21} = \binom{N}{k}^{-1} k(k-1) \sum_{w=k}^N \binom{w}{k} = \binom{N}{k}^{-1} k(k-1) \binom{N+1}{k+1}.$$

Thus substituting everything back yields

$$\mathbb{E}[S] = N(k-1) + (k-1) - \binom{N}{k}^{-1} k(k-1) \binom{N+1}{k+1}.$$

We can simplify the right hand side:

$$\begin{aligned}
(N+1)(k-1) - k(k-1) \frac{\frac{(N+1)!}{(N-k)!(k+1)!}}{\frac{N!}{(N-k)!k!}} &= (N+1)(k-1) - k(k-1) \frac{(N+1)!(N-k)!k!}{N!(N-k)!(k+1)!} \\
&= (N+1)(k-1) - k(k-1) \frac{N+1}{k+1} \\
&= (N+1)(k-1) - \frac{k(k-1)(N+1)}{k+1} \\
&= (N+1)(k-1) \left(1 - \frac{k}{k+1}\right) \\
&= (N+1) \frac{k-1}{k+1},
\end{aligned}$$

and thus obtain

$$\mathbb{E}[S] = (N+1) \frac{k-1}{k+1}.$$

The analysis is completed as before, where we pass from our observation of s for S to a prediction \hat{N} for N :

$$\hat{N} = \frac{k+1}{k-1}s - 1 = s \left(1 + \frac{2}{k-1}\right) - 1,$$

where the final equality is due to rewriting the algebra to mirror more closely the formula from the case where the first tank is numbered 1. Note that this formula passes the same smell checks the other did; for example $s \frac{2}{k-1} - 1$ is always at least 1 (remember k is at least 2), and thus the lowest estimate we can get for the number of tanks is $s + 1$.

4. COMPARISON OF APPROACHES

So, which did better: statistics or spies? Once the Allies won the war, they could look into Albert Speer's, the Nazi Minister of Armaments, records to see the exact number of tanks produced each month; see Figure 2.

| Month | Statistical estimate | Intelligence estimate | German records |
|-------------|----------------------|-----------------------|----------------|
| June 1940 | 169 | 1,000 | 122 |
| June 1941 | 244 | 1,550 | 271 |
| August 1942 | 327 | 1,550 | 342 |

FIGURE 2. Comparison of estimates from statistics and spies to the true values. Table from [6].

The meticulous German record keeping comes in handy for the vindication of the statisticians; these estimates were astoundingly more accurate. While certainly not perfect (an underestimation

of 30 tanks could have pretty dire consequences when high command is allocating resources), the statistical analysis was tremendously superior to the intelligence estimates, which were off by factors of 5 or more. We mentioned earlier the lessons to be learned from McClellan's caution. He was the first General of the Army of the Potomac (which was the Union army headquartered near Washington), and he repeatedly missed opportunities to deliver a debilitating blow to General Robert E. Lee's army of Northern Virginia, most famously during Lee's retreat from Antietam. Despite vastly outnumbering Lee in men and supplies, McClellan chronically overestimated Lee's forces, causing him to be overly cautious and far too timid a commander. Ultimately the Civil War would drag on for four years and costing over 650,000 American lives, and one wonders how the outcome would have been different if McClellan had been more willing to take the field.

We encourage the reader to write some simple code to simulate both problems discussed here (or see [5]), namely when we know and when we don't know the number of the lowest tank. These problems provide a valuable warning on how easy it is to accidentally convey information. In many situations today numbers are randomly generated to prevent such an analysis. Alternatively, sometimes numbers are deliberately started higher to fool an observer into thinking that more is present than actually is (examples frequently seen are the counting done during a workout, putting money in the tip jar at the start of the shift to encourage future patrons to be generous, or checkbooks starting with the first check as 100 or higher so the recipient does not believe it is from a new account).

APPENDIX A. THE GERMAN TANK PROBLEM AND LINEAR REGRESSION

The German Tank Problem is frequently used in probability or discrete math classes, as it illustrates the power of those two disciplines to use binomial identities to great advantage. It's also seen in statistics classes in discussing how to find good estimators of population values. Focusing on these examples, however, neglects another great setting where it may be used effectively: as an application of the power of Linear Regression (or the Method of Least Squares). We quickly review how these methods yield the best-fit line or hyperplane, and then generalize to certain non-linear relationships. We show how simulations can be used to provide support for formulas. This is extremely important, as often we are unable to prove conjectured relationships. Returning to World War II, the Allies could run trials (say drawing numbered pieces of paper from a bag) to model the real world problem, and use the gathered data to sniff out the relationship m , k and N .

Additionally, we use this appendix as an opportunity to discuss some of the issues that can arise when implementing the Method of Least Squares to find the best fit line. While these do not occur in most applications, it is worth knowing that they can happen and seeing solutions.

A.1. Theory of Regression. Suppose we believe there are choices of a and b such that given an input x we should observe $y = ax + b$, but we don't know what these values are. We could observe a large number of pairs of data $\{x_i, y_i\}_{i=1}^I$, and use these to find the values of a and b that minimize the sum of the squares of the errors⁵ between the observed and predicted values:

$$E(a, b) = \sum_{i=1}^I (y_i - (ax_i + b))^2.$$

⁵We cannot just add the errors, as then a positive error could cancel with a negative error. We could take the sum of the absolute values, but the absolute value function is not differentiable; it is to have calculus available that we measure errors by sums of squares.

By setting

$$\frac{\partial E}{\partial a} = \frac{\partial E}{\partial b} = 0,$$

after some algebra⁶ we find the best fit values are

$$\begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^I x_i^2 & \sum_{i=1}^I x_i \\ \sum_{i=1}^I x_i & \sum_{i=1}^I 1 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^I x_i y_i \\ \sum_{i=1}^I y_i \end{pmatrix};$$

see for example the supplemental material online for [2]. What matters is that the relation is linear in the unknown parameters a and b (or more generally a_1, \dots, a_ℓ); similar formulas hold for

$$y = a_1 f_1(x) + \dots + a_\ell f_\ell(x).$$

For a linearly algebraic approach to regression see for example [3].

Regression is a rich subject; we wish to try to find the best fit parameters to relate N to m and k ; however, we'll shortly see that our initial guess at a relationship is non-linear. Fortunately, by taking logarithms, we can convert many non-linear relations to linear ones, and thus the formulas above are available again. The idea is that by doing extensive simulations we can gather enough data to make a good conjecture on the relationship. Sometimes, as will be the case with the German Tank Problem, we are able to do a phenomenal job in predicting the functional form and coefficients, while other times we can only get some values with confidence.

To highlight these features we first quickly review a well-known problem: The Birthday Paradox (see for example [2]). The standard formulation assumes we have a year with D days, and asks how many people do we need in a room to have a 50% chance that at least two share a birthday, under the assumption that the birthdays are independent and uniformly distributed from 1 to D . A straightforward analysis shows the answer is approximately $D^{1/2} \sqrt{\log 4}$. We now consider the closely related but less well-known problem of what is the expected number of people P we need in a room before there is a match.⁷ Based on the first problem it is reasonable to expect the answer to also be on the order of $D^{1/2}$, but what is the constant factor? We can try a relation of the form $P = BD^a$, and then taking logs (and setting $b = \log B$) we would get $\log P = a \log D + b$. See Figure 3.

The two simulations both have similar values for a , with both of them consistent with an exponent of $1/2$. Unfortunately the values for b wildly differ, though of the two parameters we care more about a as it tells us how our answer changes with the number of days. There is an important lesson here: data analysis can often suggest much of the answer, but it is not always the full story and there is a role for theory in supplementing such analysis.

A.2. Issues in Applying to the German Tank Problem. Building on this lesson, we return to the German Tank Problem. What is a reasonable choice for N as a function of m and k ? Clearly N is at least m , so we try $N = m + f(m, k)$, which transfers the problem to estimating $f(m, k)$. We expect that as m increases this should increase, and as k increases it should decrease. Looking at extreme cases is useful; if $k = N$ then $f(N, N)$ should vanish, as then m must equal N . The

⁶The resulting matrix is invertible, and hence there is a unique solution, so long as at least two of the x_i 's differ. One can see this through some algebra, where the determinant of the matrix is essentially the variance of the x_i 's; if they are not all equal then the variance is positive. If the x_i 's are all equal the inverse does not exist, but in such a case we should not be able to predict how y varies with x as we are not varying x !

⁷As a nice exercise, use linearity of expectation to show that we expect at least two people to share a birthday when $P = D^{1/2} \sqrt{2} + 1$.

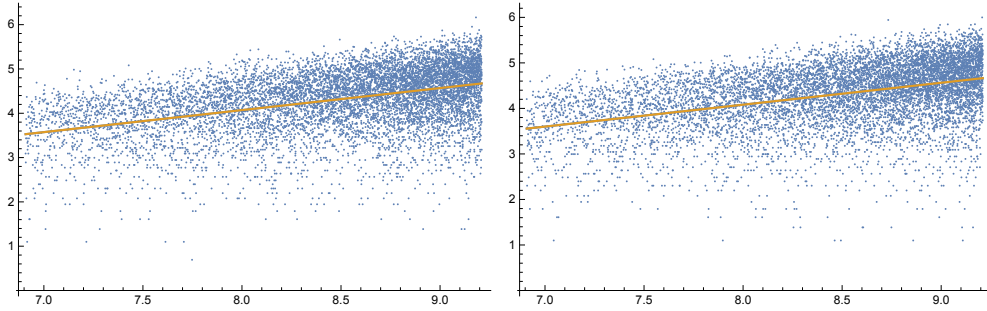


FIGURE 3. Plot of best fit line for P as a function of D . We twice ran 10,000 simulations with D chosen from 10,000 to 100,000. Best fit values were $a \approx 0.506167, b \approx -0.0110081$ (left) and $a \approx 0.48141, b \approx 0.230735$ (right).

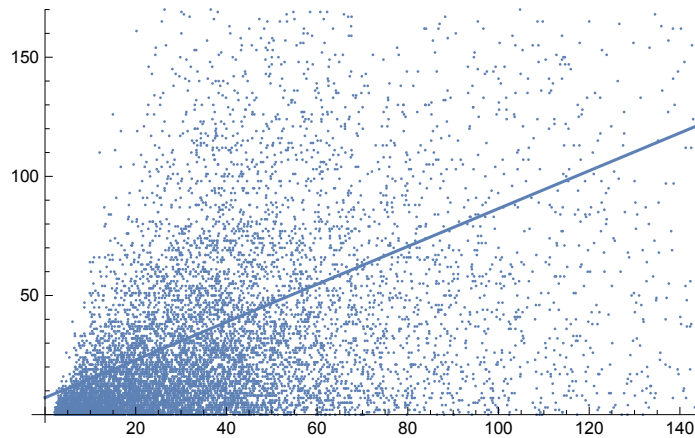


FIGURE 4. Plot of best fit line for $N - m$ as a function of m/k . We ran 10,000 simulations with N chosen from $[100, 2000]$ and k from $[10, 50]$. Best fit values for $N - m = a(m/k) + b$ for this simulation were $a \approx 0.793716, b \approx 7.10602$.

simplest function that fits this is $f(m, k) = b \cdot m/k$ with b as our free parameter, and we are led to conjecturing a relationship of the form

$$N = m + b \frac{m}{k} = m \left(1 + \frac{b}{k} \right).$$

Note that this guess is quite close to the true answer, but because the observed quantities m and k appear as they do, it is not a standard regression problem. We could try to fix this by looking at $N - m$, the number of tanks we need to add to our observed largest value to get the true number. We could then try to write this as a linear function of the ratio m/k :

$$N - m = a \frac{m}{k} + b,$$

where we allowed ourselves a constant term to increase our flexibility of what we can model. While for $a = -b = 1$ this reproduces the correct formula, finding the best fit values leads to a terrible fit, as evidenced in Figure 4.

Why is the agreement so poor, given that proper choices exist? The problem is the way m and k interact, and in the set-up above we have the observed quantity m both as an input variable and

as an output in the relation. We thus need a way to separate m and k , keeping both on the input side. As remarked, we can do this through logarithms; we discuss another approach in the next subsection.

A.3. Resolving Implementation Issues. We look at our best fit line for two choices of k ; The left side of Figure 5 does $k = 1$ while Figure 6 is $k = 5$. Both of these show a terrible fit of N as a linear function of m (for a fixed k). In particular, when $k = 1$ we expect N to be $2m - 1$ but our best fit line is about $.784m + 2875$; this is absurd as for large m we predict N to be less than m ! Note, however, the situation is completely different if instead we plot m against N (the right hand side of those figures). Clearly if N linearly depends on m then m linearly depends on N . When we do the fits this way, the results are excellent.

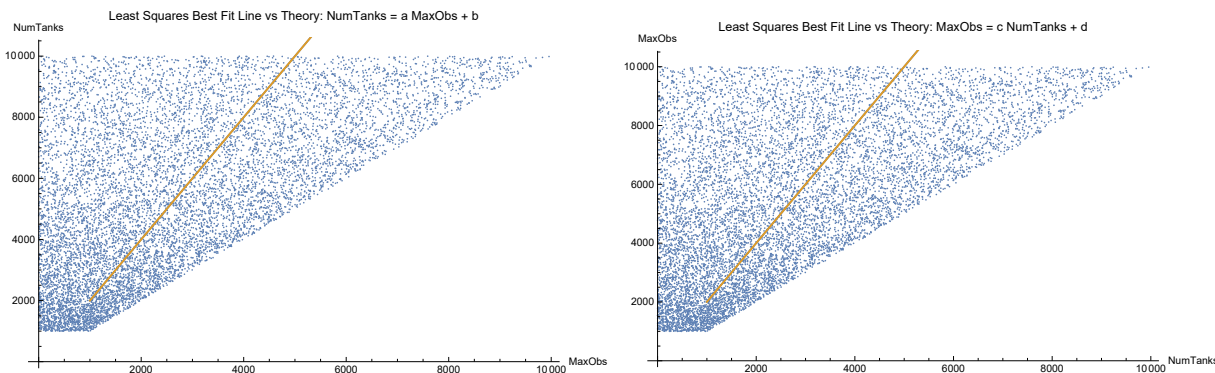


FIGURE 5. Left: Plot of N vs maximum observed tank m for fixed $k = 1$. Theory: $N = 2m - 1$, best fit $N = .784m + 2875$. Right: Plot of maximum observed tank m vs N for fixed $k = 1$. Theory: $m = .5N + .5$, best fit $m = .496N + 10.5$.

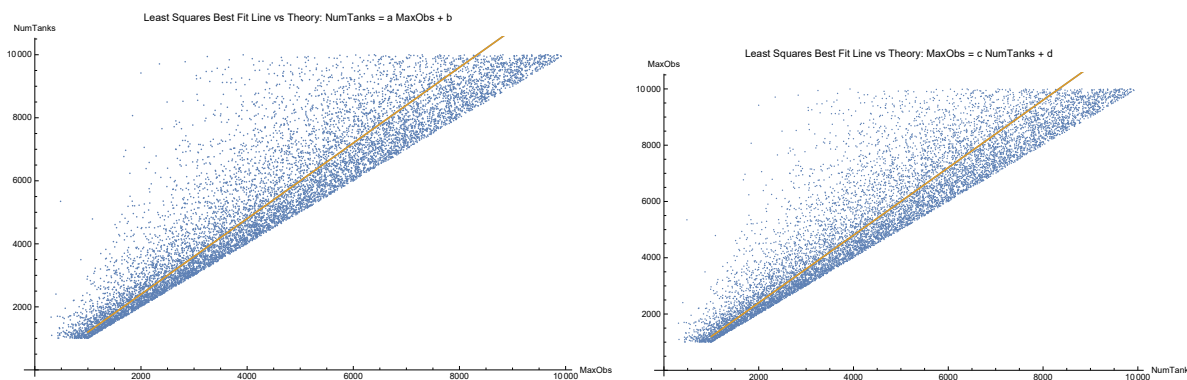


FIGURE 6. Left: Plot of N vs maximum observed tank m for fixed $k = 5$. Theory: $N = 1.2m - 1$, best fit $N = 1.037m + 749$. Right: Plot of maximum observed tank m vs N for fixed $k = 5$. Theory: $m = .883N + .883$, best fit $m = .828N + 25.8$.

Note that from the point of view of an experiment, it makes more sense to plot m as the dependent variable and N as the independent, input variable. The reason is the way we simulate; we fix a k and an N and then choose k distinct numbers uniformly from $\{1, \dots, N\}$.

We end with another approach which works well, and allows us to view N as a function of m . Instead of plotting each pair (m, N) for a fixed k , we instead fix k , choose an N , and then do 100 trials. For each trial we record the largest serial number m , and then we average these, and plot (\bar{m}, N) where \bar{m} is the average. This greatly decreases the variability, and we now obtain a nearly perfect straight line and fit; see Figure 7.

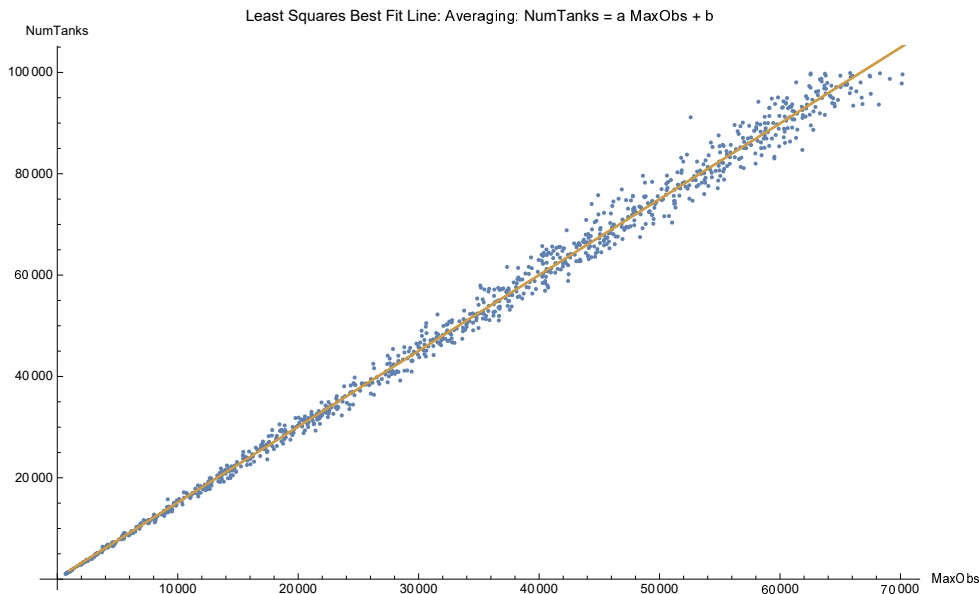


FIGURE 7. Plot of N vs maximum observed tank m for fixed $k = 1$. Theory: $N = 1.5m - 1$, best fit $N = 1.496m + 171.2$.

A.4. Determining the Functional Form. We consider the more general relation

$$N = Cm^a \left(1 + \frac{b}{k}\right),$$

where we expect $C = a = b = 1$; note this won't have the -1 summand we know should be there, but for large m that should have negligible impact. Letting $C = e^c$ for notational convenience, we find

$$\log N = c + a \log(m) + \log \left(1 + \frac{b}{k}\right).$$

If x is large then $\log(1 + 1/x) \approx 1/x$, so we try the approximation

$$\log N \approx c + a \log(m) + b \frac{1}{k}.$$

Figure 8 shows the analysis when $C = 1$ (so $c = 0$), as the analysis then reduces to the usual case with two unknown parameters. We chose to take $C = 1$ from the lesson we learned in the analysis of the Birthday Problem.

The best fit values of the parameters are $a = 0.999911$ and $b = 0.961167$, which are reasonably close to $a = b = 1$. Thus these numerics strongly support our conjectured relation $N = m(1 + 1/k)$, and shows the power of statistics. While we were able to see the arguments needed to prove this relation essentially holds, imagine we could not prove it but still have our heuristic

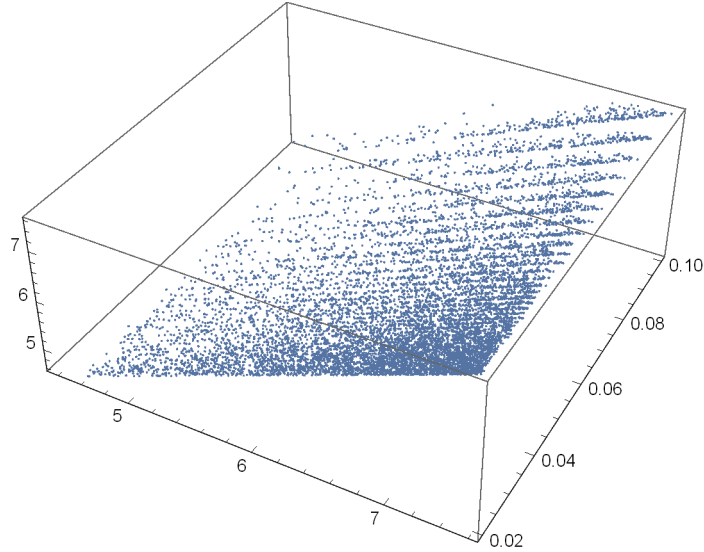


FIGURE 8. Plot of $\log N$ against $\log m$ and $1/k$. We ran 10,000 simulations with N chosen from $[100, 2000]$ and k from $[10, 50]$. The data is well-approximated by a plane (we do not draw it in order to prevent our image from being too cluttered).

arguments and analysis of extreme cases which suggest it is true. By simulating data and running the regression, we see that our formula does a stupendous job explaining our observations, and thus gain confidence to use it in the field.

We end with one last approach. Let us guess a relationship of the form $N = a(k)m + b(k)$, where $a(k) = 1 + f(k)$ (we write $a(k)$ as $1 + f(k)$ as we know there have to be at least m tanks). We can fix k , and find the best fit values of $a(k)$ and $b(k)$. In Figure 9 we plot the best fit slope $a(k)$ versus k , as well as a log-log plot. For the log-log plot we look at $a(k) - 1$, subtracting the known component. We see a beautiful linear relation, and thus even if we did not know it should be m plus a constant times m/k , the data suggests that beautifully! Specifically, we found the best fit line was $\log(a(k) - 1) = -.999 \log(k) - .007$, suggesting that $a(k) = 1 + 1/k$; we obtain the correction functional form just by running simulations!

REFERENCES

- [1] M. Cozzens and S. J. Miller, *The Mathematics of Encryption: An Elementary Introduction*, AMS Mathematical World series **29**, Providence, RI, 2013.
- [2] S. J. Miller, *The Probability Lifesaver*, Princeton University Press, Princeton, NJ, 2018. https://web.williams.edu/Mathematics/sjmiller/public_html/probabilitylifesaver/index.htm.
- [3] G. Strang, *Introduction to Linear Algebra*, Fifth Edition, Wellesley-Cambridge Press, 2016.
- [4] Probability and Statistics Blog, *How many tanks? MC testing the GTP*, <https://statisticsblog.com/2010/05/25/how-many-tanks-gtp-gets-put-to-the-test/>.
- [5] Statistical Consultants Ltd, *The German Tank Problem*, <https://www.statisticalconsultants.co.nz/blog/the-german-tank-problem.html>.
- [6] WikiEducator, *Point Estimation - German Tank Problem*, https://wikieducator.org/Point_estimation_-_German_tank_problem.
- [7] Wikipedia, *German Tank Problem*, Wikimedia Foundation, https://en.wikipedia.org/wiki/German_tank_problem.

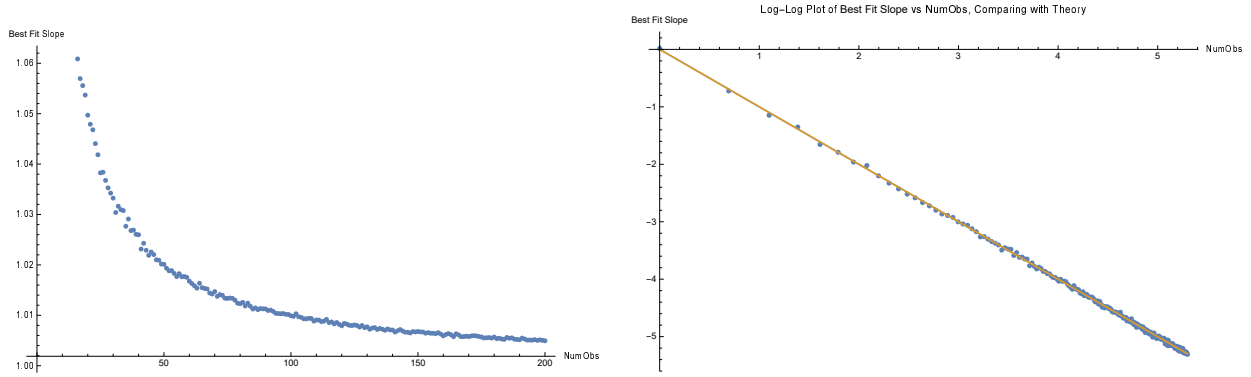


FIGURE 9. Left: Plot of $a(k)$, the slope in $N = a(k)m + b$, versus k . Right: Log-Log Plot of $a(k) - 1$ versus k . In $\log(a(k) - 1)$ versus $\log k$, theory is $\log(a(k) - 1) = -1 \log(k)$, best fit line is $\log(a(k) - 1) = -.999 \log(k) - .007$.

Email address: gpclark96@gmail.com

DEPARTMENT OF MATHEMATICS AND STATISTICS, WILLIAMS COLLEGE, WILLIAMSTOWN, MA 01267

Email address: aag2@williams.edu

DEPARTMENT OF MATHEMATICS AND STATISTICS, WILLIAMS COLLEGE, WILLIAMSTOWN, MA 01267

Email address: sjm1@williams.edu, Steven.Miller.MC.96@aya.yale.edu

DEPARTMENT OF MATHEMATICS AND STATISTICS, WILLIAMS COLLEGE, WILLIAMSTOWN, MA 01267

Current address: Department of Mathematical Sciences, Carnegie Mellon University, Pittsburgh, PA 15213