

Homework #1 – SOLUTIONS

1. Prepare a one paragraph summary of what we talked about in class so far. I do not want just a bulleted list of topics, I want you to use complete sentences and establish context (Why is what we have learned relevant? How does it connect with other classes?). The more insight you give, the better.

Solution: Make sure that what is written is in coherent, complete sentences.

2. Familiarize yourself with Python.
 - (a) Install (if necessary) the Anaconda distribution of Python at <https://www.anaconda.com/products/individual>.
 - (b) Review the Python tutorial at <https://cs4540-f18.github.io/notes/python-basics>.
 - (c) Check out and consider bookmarking this site: <https://numpy.org/doc/stable/user/numpy-for-matlab-users.html>
 - (d) If you feel it would be helpful, especially if you are new to Python, consider taking the time to go through one (or both) of these more basic Python tutorials:
 - <https://developers.google.com/edu/python/>
 - <https://www.learnpython.org/>
3. In the problems below and on subsequent homeworks you can/will make frequent use of the CVXPY package. You will probably need to install this package as it does not come standard with most python distributions. This should be easy – simply open up a Python prompt and type `pip install cvxpy`. If you are on Windows you may also need to install the Visual Studio build tools if you don't already have them. If you are having any trouble, see <https://www.cvxpy.org/install/> for details. Install CVXPY and poke around on the CVXPY website to get a feeling for what it can do.
4. In our first class we mentioned a nonconvex optimization problem that nevertheless has a simple solution:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{maximize}} \quad \mathbf{x}^T \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \|\mathbf{x}\|_2 = 1.$$

Here we consider a special case of this problem (that is not hard to extend to the general case). Suppose that \mathbf{D} is a *diagonal* $N \times N$ matrix whose diagonal elements are positive. Show that one possible maximizer to

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{maximize}} \quad \mathbf{x}^T \mathbf{D} \mathbf{x} \quad \text{subject to} \quad \|\mathbf{x}\|_2 = 1$$

is given by the $\hat{\mathbf{x}}$ with a 1 in the entry corresponding to the largest diagonal element of \mathbf{D} , and is 0 elsewhere. You should approach this by rigorously arguing that any feasible $\mathbf{x}' \neq \hat{\mathbf{x}}$ cannot achieve a higher value for the objective function.

Solution: First, note that in general, we can also write our objective function as

$$\mathbf{x}^T \mathbf{D} \mathbf{x} = \sum_{n=1}^N x_n (d_n x_n) = \sum_{n=1}^N d_n x_n^2.$$

Now, suppose without loss of generality that the first entry d_1 on the diagonal is the largest, so that $d_1 \geq d_n$ for all $n \neq 1$. In this case the objective function for our proposed $\hat{\mathbf{x}}$ takes a value of $\hat{\mathbf{x}}^T \mathbf{D} \hat{\mathbf{x}} = d_1$. Now, consider any $\mathbf{x}' \neq \hat{\mathbf{x}}$ satisfying $\|\mathbf{x}'\|_2 = 1$. Note that we can write the objective function evaluated at \mathbf{x}' as

$$\begin{aligned} \mathbf{x}'^T \mathbf{D} \mathbf{x}' &= \sum_{n=1}^N d_n (x'_n)^2 \\ &\leq \sum_{n=1}^N d_1 (x'_n)^2 \\ &= d_1 \sum_{n=1}^N (x'_n)^2 \\ &= d_1, \end{aligned}$$

where the inequality follows from the assumption that $d_1 \geq d_n$ and the final equality from the fact that $\|\mathbf{x}'\|_2 = 1$. But this shows that $\mathbf{x}'^T \mathbf{D} \mathbf{x}' \leq d_1$ for any possible \mathbf{x}' , and thus $\hat{\mathbf{x}}$ is a maximizer.

5. A significant motivation in the early development of optimization algorithms involved applications in logistics and supply chain management. This is a somewhat contrived example. Suppose you are the manager of a factory that produces two products: toilet paper and paper towels. The profit for paper towels is \$10 per ton of pulp, but the profit for toilet paper is \$100 per ton of pulp. Your factory has a supply of 8 000 tons of pulp per month. In addition, you face the following conditions:

- You have a government contract that requires you to produce 200 tons of paper towels each month.¹
- You must ship all of the products that you produce using a fleet of trucks that has a capacity of 100 000 ton-miles (i.e., you could transport 100 000 tons 1 mile, or 1 ton 100 000 miles, or somewhere in between).
- Paper towels are delivered to a warehouse 10 miles from your factory, while toilet paper is delivered to a warehouse 35 miles from your factory.

You would like to decide how much of each product to produce to maximize your profit. Write down a mathematical formulation of this as an optimization problem. What kind of optimization problem is this (linear, quadratic, semidefinite, ...)?

¹Assume that 1 ton of paper towels or 1 ton of toilet paper each require 1 ton of raw pulp.

Solution: If x_1 denotes the tons of pulp devoted to paper towels and x_2 denotes the tons of pulp devoted to toilet paper, then this can be represented as

$$\begin{aligned} & \underset{x_1, x_2}{\text{minimize}} && 10x_1 + 100x_2 \\ & \text{subject to} && x_1 \geq 200, \quad x_2 \geq 0 \\ & && x_1 + x_2 \leq 8000 \\ & && 10x_1 + 35x_2 \leq 100000. \end{aligned}$$

6. Consider the optimization program

$$\min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_1,$$

where \mathbf{A} is a given $M \times N$ matrix and $\mathbf{y} \in \mathbb{R}^M$. Show how this can be converted into an equivalent linear program.

Solution: Let $\mathbf{a}_m \in \mathbb{R}^N$ be the (transpose of the) m^{th} row in \mathbf{A} . Then the solution $(\hat{\mathbf{x}}, \hat{\mathbf{u}})$ to the program

$$\min_{\mathbf{x}, \mathbf{u}} \sum_{m=1}^M u_m, \quad \text{subject to} \quad -u_m \leq y_m - \mathbf{a}_m^T \mathbf{x} \leq u_m, \quad m = 1, \dots, M,$$

will obey $|y_m - \mathbf{a}_m^T \hat{\mathbf{x}}| = \hat{u}_m$. This is clear, as if $|y_m - \mathbf{a}_m^T \hat{\mathbf{x}}| > \hat{u}_m$, the constraints above are violated, and if $|y_m - \mathbf{a}_m^T \hat{\mathbf{x}}| < \hat{u}_m$, the function is not as small as it could be while still meeting the constraints (we could simply lower \hat{u}_m).

We can rewrite the above as a linear program in \mathbb{R}^{N+M} :

$$\min_{\mathbf{z} \in \mathbb{R}^{N+M}} \mathbf{c}^T \mathbf{z} \quad \text{subject to} \quad \mathbf{C}\mathbf{z} \leq \mathbf{b},$$

with $\mathbf{c} = \begin{bmatrix} \mathbf{0}_N \\ \mathbf{1}_M \end{bmatrix}$, and

$$\mathbf{C} = \begin{bmatrix} -\mathbf{A} & -\mathbf{I} \\ \mathbf{A} & -\mathbf{I} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -\mathbf{y} \\ \mathbf{y} \end{bmatrix}.$$

7. One important optimization problem that arises in the context of machine learning has the form

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \sum_{m=1}^M \max(\mathbf{a}_m^T \mathbf{x}, z_m) + \tau \|\mathbf{x}\|_2^2.$$

Show how to express this program as a quadratic program with linear constraints. (Hint: $\max(\mathbf{a}^T \mathbf{x}, z) \leq u$ is the same as saying $\mathbf{a}^T \mathbf{x} \leq u$ and $z \leq u$.)

Solution: We introduce new optimization variables $\mathbf{u} \in \mathbb{R}^M$, replace $\sum_m \max(\mathbf{a}_m^T \mathbf{x}, z_m)$ with $\sum_m u_m$ and add the constraints that $\mathbf{a}_m^T \mathbf{x} \leq u_m$ and $z_m \leq u_m$. The final

quadratic program with linear constraints is

$$\begin{aligned} \underset{\mathbf{x} \in \mathbb{R}^N, \mathbf{u} \in \mathbb{R}^M}{\text{minimize}} \quad & \sum_{m=1}^M u_m + \tau \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad \mathbf{a}_m^T \mathbf{x} - u_m \leq 0 \\ & z_m - u_m \leq 0 \\ & m = 1, \dots, M. \end{aligned}$$

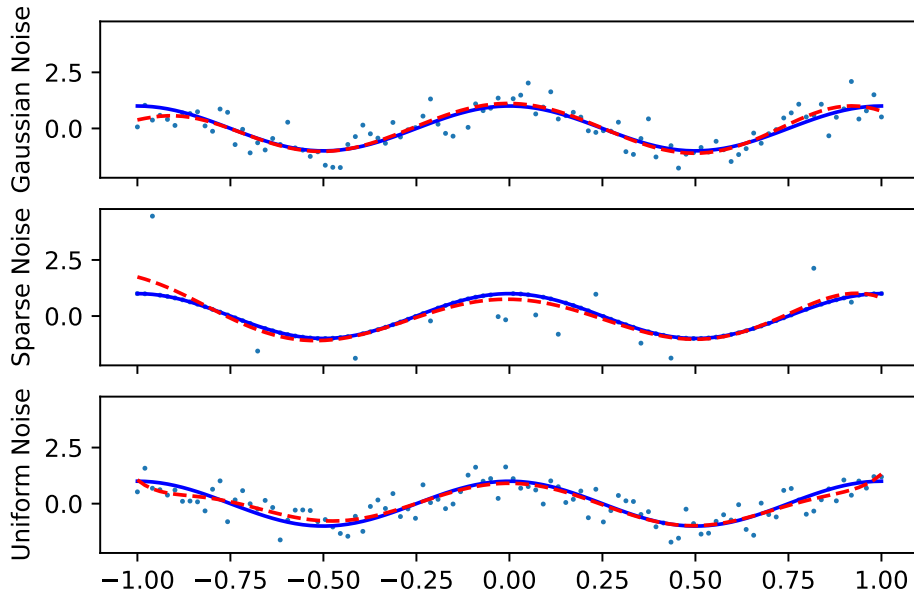
To convince yourself that these problems are equivalent, let $\hat{\mathbf{x}}, \hat{\mathbf{u}}$ be minimizers of the program above. Then it is true that $\hat{u}_m = \max(\mathbf{a}_m^T \hat{\mathbf{x}}, z_m)$ for all $m = 1, \dots, M$, and so $\hat{\mathbf{x}}$ is also a minimizer for the original program. This follows from two easy observations. First, if $\hat{u}_m < \max(\mathbf{a}_m^T \hat{\mathbf{x}}, z_m)$, then $\hat{\mathbf{x}}, \hat{\mathbf{u}}$ would not be feasible, since at least one of the constraints would be broken. Second, if $\hat{u}_m > \max(\mathbf{a}_m^T \hat{\mathbf{x}}, z_m)$ for any m , then $\hat{\mathbf{u}}$ is not optimal, as we could simply decrease \hat{u}_m , while simultaneously maintaining feasibility and lowering the functional value.

8. You hopefully have some prior experience with least squares problems in which we aim to find an \mathbf{x} that minimizes $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$. While minimizing the ℓ_2 norm of the error is *often* a good idea, a big part of why it is so popular is just that it is so easy to solve the underlying optimization problem. However, with modern tools we can easily explore situations where other ℓ_p norms might be more appropriate. In particular, we will consider minimizing $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_1$ and $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_\infty$.

- (a) Download the file `hw01_prob8.py`. This file sets up a simple regression problem in which \mathbf{y} consists of noisy observations of a smooth function $f(t)$ and considers three noise models: Gaussian noise, sparse Gaussian noise, and uniform noise. Compute and plot the least squares solution for each case.

Solution: See the python code in Canvas.

Least squares (ℓ_2) estimates

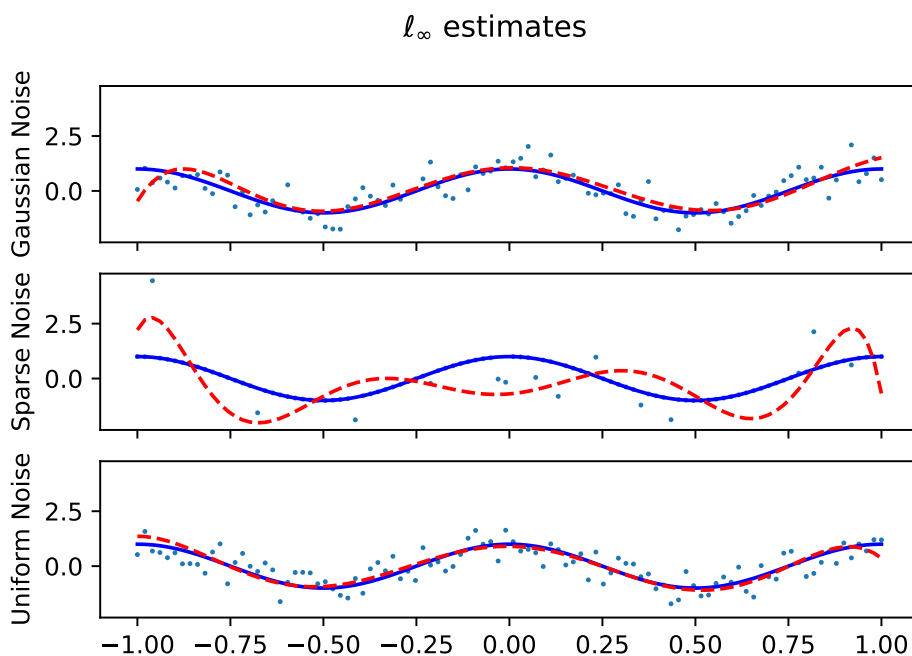


(b) Now estimate the solution for the Chebyshev approximation problem:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_\infty.$$

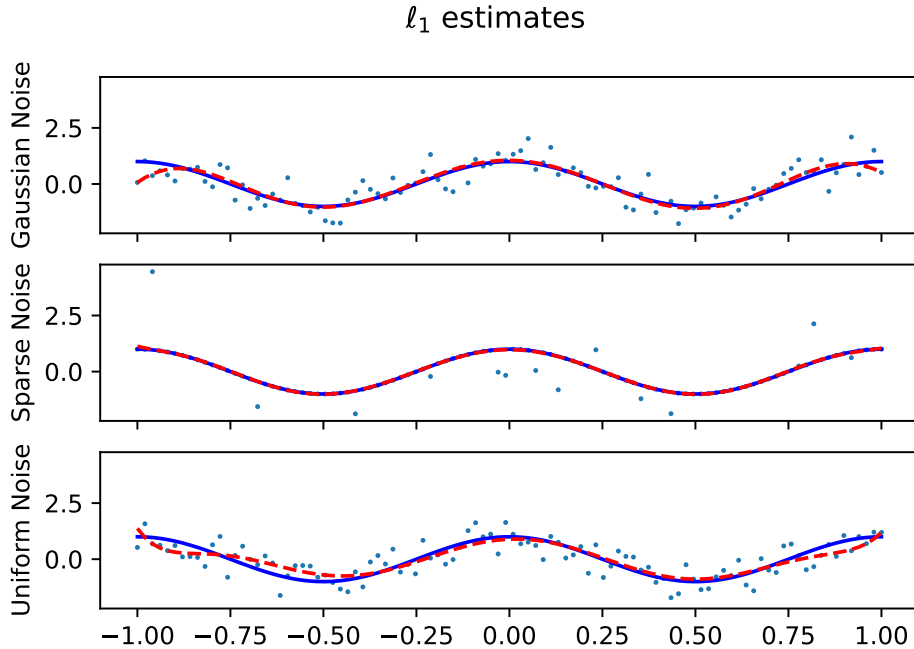
You do not have to convert this to a linear program first; CVXPY has built-in tools for handling the ℓ_∞ norm. Compute and plot the resulting estimate of \mathbf{x} for each of the three noise cases.

Solution: [See the python code in Canvas.](#)



(c) Next repeat part (b) for $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_1$. Again, compute and plot the results for each of the three noise cases.

Solution: [See the python code in Canvas.](#)



- (d) Compare the results that you obtained in the previous parts. Which method seems best suited for each case?

Solution: In the case of Gaussian noise, least squares clearly performs the best. In the case of sparse noise, minimizing the ℓ_1 norm of the error does dramatically better than the others. In the case of uniform noise the performance is more mixed, although you might give the ℓ_∞ norm a slight edge over the others. In this case, due to some large errors near the boundary, none of the methods perform particularly well at the edge of the interval (a quite common phenomenon when fitting a polynomial.)

9. In this problem we will explore the idea of *group testing* as a strategy for testing a large population for a rare disease by pooling samples together. Suppose that we have a population of N people and we collect saliva samples from each of them. We are looking for genetic signatures of a particular virus in these samples. Let x_n denote the concentration of this material in the sample for the n^{th} person. We will assume that for healthy people $x_n = 0$, but for infected people, $x_n > 0$. Our testing procedure will be to form a series of M mixtures of samples from different subsets of people, and then only run tests on these mixtures. The goal here is to set $M < N$, and the question is then whether we can identify the infected people from the results of these tests.

We will consider the following approach:² we will form mixtures by constructing *random* combinations of samples, and we will attempt to recover the original \mathbf{x} using a simple convex optimization problem.

To mathematically represent the sampling/testing process, assume that we will ultimately run M tests, each of which will tell us the concentration of viral material

²Just to be clear, this is not exactly what Georgia Tech is doing, but it is similar in spirit. The approach that Georgia Tech is somewhat simpler and builds on ideas that date back to World War II, when they were used to test soldiers for syphilis.

in the combined sample being tested. For each person, their sample will be divided into $K < M$ equal portions, which will be assigned at random to the M tests. We will do this independently for each of the N people. We can ultimately represent the concentration of viral material in each of the mixed samples that we will ultimately test as a vector $\mathbf{y} \in \mathbb{R}^M$. We can write \mathbf{y} as

$$\mathbf{y} = \mathbf{A}\mathbf{x},$$

where \mathbf{A} is a matrix that represents the assignment of people to mixed samples/tests. Specifically, \mathbf{A} is a $M \times N$ matrix where each column is constructed independently by picking K entries at random, setting them to 1, and setting the remaining entries to 0.

Suppose there is no noise in our tests, so that we can estimate \mathbf{y} perfectly. Our inference problem is now to estimate \mathbf{x} given knowledge of \mathbf{y} and \mathbf{A} . In general, since $M < N$, recovering \mathbf{x} is impossible. However, when \mathbf{x} is sparse, meaning that it has only a few nonzeros (in this case meaning that most of the population is negative), then recovering \mathbf{x} is possible, although this fact was only broadly appreciated within the last 15 years or so.

We will try to estimate \mathbf{x} by solving the following optimization problem:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{y}, \quad \mathbf{x} \geq \mathbf{0}.$$

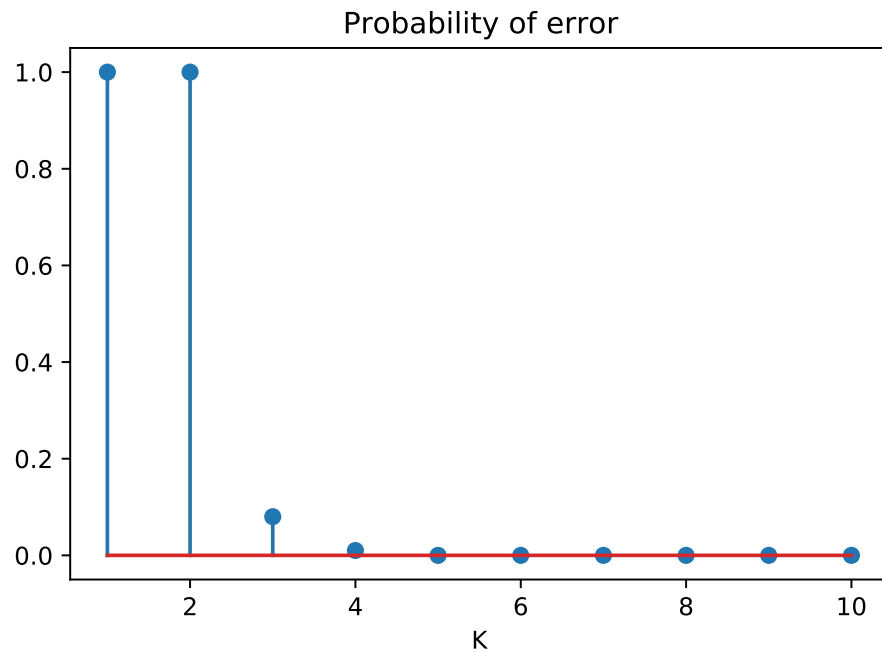
Below we will explore when and how well this works.

- (a) Suppose that you are testing a population of size $N = 1000$, but you can only process $M = 100$ tests. Assume that only 1% of the population is positive (meaning that there will be 10 infected individuals). Each person's sample will be split and added to $K = 10$ different batches. The file `hw01_prob09.py` contains code that sets up this problem. Use CVXPY to solve the optimization problem above and verify that this approach correctly identifies the 10 infected individuals.

Solution: See the python code in Canvas.

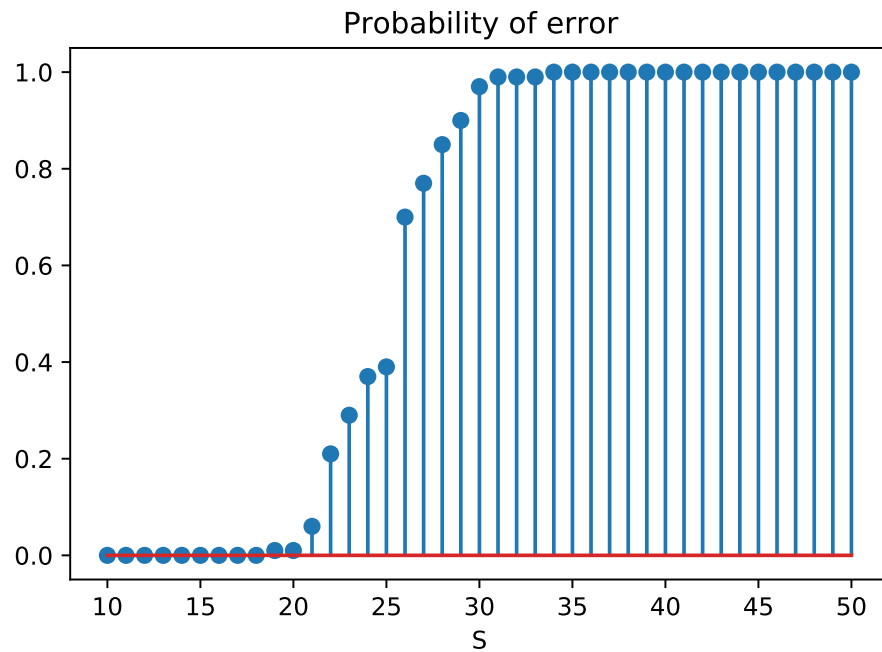
- (b) Experiment with K . In practice, you might not want to divide a person's sample into too many tests. How low can you set K before things begin to fail?

Solution: See the python code in Canvas. In this part I did 100 random trials with different random assignment to pools. We can take K as low as 5 pools without there being any impact on performance. However, for $K = 4$ there becomes a small (1%) chance of error. When $K = 1$ or $K = 2$, the algorithm failed on every trial. This should make some intuitive sense – when $K = 1$ you need to be *extremely* lucky to be able to unambiguously identify the locations of the nonzeros. The only way this would work (for $K = 1$) is if each nonzero happened to coincide with a column of \mathbf{A} where the one nonzero in that column also happened to be the only nonzero in that row. The situation is not too different for $K = 2$, but there is a pretty sharp improvement as K grows.



- (c) Suppose that the prevalence of the disease begins to grow beyond 1%. How widespread can the disease become before the approach begins to fail (holding M and N fixed). As the disease becomes more widespread, you may need to adjust K . What value of K seems to work best when the spread of the disease is just below the threshold where the approach begins to fail?

Solution: See the python code in Canvas. Using a similar approach as in the previous part, we can look at the probability of failure as S (the number of infected individuals) grows. We see that we still get reliable recovery for $S \leq 18$, but after that if S continues to grow the performance degrades fairly rapidly:



For $S = 18$ if we examine the impact of K we find that as we decrease K the algorithm fails sooner (when $S = 10$ the algorithm still had OK performance for $K = 3, 4$, but as S increases the performance degrades for these particularly small values of K).

