

# Optimal City Pairs for High Speed Rail in the United States

Ellie Jensen, Niels Vanderloo

December 21, 2022

## Abstract

We provide a preliminary study to identify the optimal city pairs to include in connections between future high speed rail networks in the United States. We introduce a linear programming model to accomplish this task. Making use of Census Bureau datasets and the Gurobi optimization software, we find the optimal city pairs headlined by New York City  $\iff$  Washington D.C.

## Introduction

While high speed rail (HSR) has been a staple of transportation in many countries for decades, the U.S. does not currently have a fully high speed train line. HSR is generally categorized by speeds greater than 125 miles per hour; the Acela Express, which runs along the Northeast Corridor between Boston, New York, and Washington D.C. reaches 150 miles per hour, but averages only 66 miles per hour. Comparatively, the Tokaido Shinkansen, which was built in 1964, runs between Tokyo and Osaka at an average speed of 177 miles per hour [1].

The implementation of HSR in the U.S. has many benefits. HSR is fully electric and more energy efficient than airplanes and cars, meaning that it would decrease greenhouse gas emissions, improve air quality, and decrease the U.S.'s dependence on foreign oil [2]. With fewer drivers on the road also comes reduced traffic fatalities and alleviated road congestion. Furthermore, for many distances HSR is a faster and more efficient form of travel.

If HSR is so efficient and has been around since the 1960s, the question that remains is why the U.S. does not have it. There are many political as well as logistical issues that arise when considering HSR. Compared to Europe and Asia, U.S. cities generally have a lower population density, making it less profitable. The U.S. also has higher car usage than most countries with HSR systems, meaning that people typically drive distances that could be traveled via HSR. There is also strong political opposition to and inaction on funding HSR.

The difficulties posed by implementing HSR have not stopped speculation of how we can build an HSR system in the U.S. Creating an extensive HSR network in the U.S. is possible, but the process of building rail across the the country is an expensive, years-long

project. But a network has to start somewhere, and if we start with the most efficient, cost-effective lines, we will see the most benefits, which would presumably expedite the process of adopting HSR in the U.S.

We intend to create a linear programming model that creates a list of optimal city pairs between which to build HSR subject to certain cost constraints. For the sake of this problem, we assume that there are no networks of rail; for someone to get from one city to another they must take a direct route.

## Initial Model

The problem we are going to solve is finding the optimal city pairs to include in an HSR system. To formulate our problem as a linear program, we need to propose some quantity to maximize or minimize. To find the “optimal” city pairs to include in our system we are really saying that we aim to find the routes that will serve the most people. Since we can’t know the actual ridership of a proposed system, we will reframe this as “travel demand.” In other words, we want our HSR system to maximize the total travel demand. We can observe potential factors influencing the travel demand between two cities. First, cities that are close together will have higher travel demand than cities that are far apart. Next, cities with higher population will have higher travel demand. In linear programming we also need to constrain the feasible region. In our problem we can impose a budget constraint where we only have the budget to build a certain length of total track.

## Data

Since we intend to use cities’ populations to calculate travel demand between cities, we need data on the population of cities in the U.S. Naively, we may choose to scrape data on population of cities proper. However, official city populations depend on arbitrary and inconsistent city limits defined by local governments, so this data is often unreliable. Instead, we will use population data for combined statistical areas (CSA). CSAs are defined by the U.S. Census Bureau and are comprised of geographically close metropolitan areas that have economic and social ties. There are currently 172 CSAs in the contiguous U.S. [3]. When we refer to a city  $i$  going forward we are actually referring to a CSA with population  $P_i$ .

To calculate travel demand we also need the distances between CSAs. Within the contiguous U.S. it is often fair to assume that a train between two cities may be built in an approximately straight path. To calculate the distance between CSAs we identify the latitude and longitude of the most populous city in each CSA and use the Haversine formula to find the geodesic distance between them [4]. Then, we can create a  $172 \times 172$  distance matrix with the distance between any two cities  $i$  and  $j$  being  $d_{i,j}$ .

A common method to compute travel demand is based on a simple gravity model [5]. In the same way that the strength of gravity increases as two objects become larger and closer to each other, there will be more demand to travel between larger cities in close proximity. The travel demand  $T_{i,j}$  between two cities  $i$  and  $j$  is often computed by dividing the product of their populations by the square of the distance between them. That is,

$$T_{i,j} = \frac{P_i \times P_j}{d_{i,j}^2}.$$

However, looking at the map of CSAs in Figure 2 we can see that some CSAs are close enough to each other such that implementing HSR between them would be impractical. The ideal distance for HSR is approximately 300 miles; for distances shorter than 100 miles, it makes more sense to travel by car, and it is more efficient to travel by plane once the distance between cities grows beyond 500 miles [6]. In our model, cities that are close together will benefit disproportionately unless we account for the practical minimum distance between cities. We therefore make the denominator of the travel demand function the maximum of  $d_{i,j}$  and 300 miles, the optimal distance for travel by HSR [5]. This way, cities that are less than 300 miles apart will not benefit from being closer to each other.

Additionally, with larger cities, we assume that the travel time from the train station to a person's final destination is longer than it would be in a smaller city. We account for these diseconomies of scale by adding an exponent of 0.8 to  $P_i$  and  $P_j$  [5]. The travel demand is now expressed as

$$T_{i,j} = \frac{P_i^{0.8} \times P_j^{0.8}}{\max(d_{i,j}, 300)^2}.$$

## Variables & Objective Function

In accordance to how we've defined  $T_{i,j}$ , our basic variables will be binary with

$$x_{i,j} = \begin{cases} 1 & \text{if we build HSR between cities } i \text{ and } j \\ 0 & \text{otherwise.} \end{cases}$$

We now have the framework to explicitly define our objective function. Recall that we wish to maximize the travel demand of our whole HSR system. The objective function becomes

$$\text{maximize } z = \sum_{\forall i} \sum_{\forall j} T_{i,j} x_{i,j}.$$

Note that using appropriate units lets us use  $z$  as a proxy for the number of riders (in millions) served by our system each year [5]. Similarly,  $T_{i,j}$  lets us estimate the number of riders (in millions) per year between any two cities in our system.

## Constraints

In the simplest form of this linear program, we have only two constraints. First we want to ensure that only one route is built between two cities, as there only needs to be one set of rails for trains to go both from  $i$  to  $j$  and  $j$  to  $i$ . The corresponding constraint is

$$x_{i,j} + x_{j,i} \leq 1 \quad \forall i, j.$$

As mentioned earlier, building HSR costs money, and government funding is hard to come by, so we may want to impose a budget constraint. For our purposes we set the limit to the amount of rail we build at 2,000 miles. For reference, New York to Los Angeles, an ideal airplane route, is 2,500 miles. In our linear program this constraint looks like

$$\sum_{\forall i} \sum_{\forall j} d_{i,j} x_{i,j} \leq 2000.$$

Our full model becomes

$$\begin{aligned} \text{maximize} \quad & z = \sum_{\forall i} \sum_{\forall j} T_{i,j} x_{i,j} \\ \text{subject to} \quad & x_{i,j} + x_{j,i} \leq 1 \\ & \sum_{\forall i} \sum_{\forall j} d_{i,j} x_{i,j} \leq 2000. \end{aligned}$$

We find the optimal solution to this linear program using the Gurobi optimization software. The optimal solution is written in Figure 1 and illustrated in Figure 2 [7]. The results of our initial model gives us twelve city pairs to include in a HSR network. Many of these cities, namely the Northeast Corridor, are approximately collinear which is good news for building a coherent network. In reality, we would only need to build one set of tracks connecting those cities; in our model, we treat each of them separately and neglect some potential efficiency. Introducing a network system to our linear program is beyond the scope of this paper, but we can introduce other improvements to our model.

City Pair ( $i, j$ )	Travel Demand ( $T_{i,j}$ )	Distance in mi ( $d_{i,j}$ )
Washington, D.C. $\iff$ New York	23.673873	203.520986
New York $\iff$ Boston	20.765652	189.984500
New York $\iff$ Philadelphia	18.604855	80.681696
San Francisco $\iff$ Los Angeles	18.493843	305.738353
Philadelphia $\iff$ Washington, D.C.	9.346010	123.101872
Boston $\iff$ Philadelphia	8.197897	270.602821
Houston $\iff$ Dallas	7.871459	224.299357
Chicago $\iff$ Detroit	7.280050	237.032099
Hartford $\iff$ New York	5.151030	100.470278
Sacramento $\iff$ San Francisco	4.070578	87.663661
Chicago $\iff$ Milwaukee	3.362160	82.232068
Detroit $\iff$ Cleveland	3.242590	92.070733

Figure 1: The results obtained using a “budget” of 2000 track miles and no further constraints (besides  $x_{i,j} + x_{j,i} \leq 1$ ).  $z = 130.06$  million people per year served.

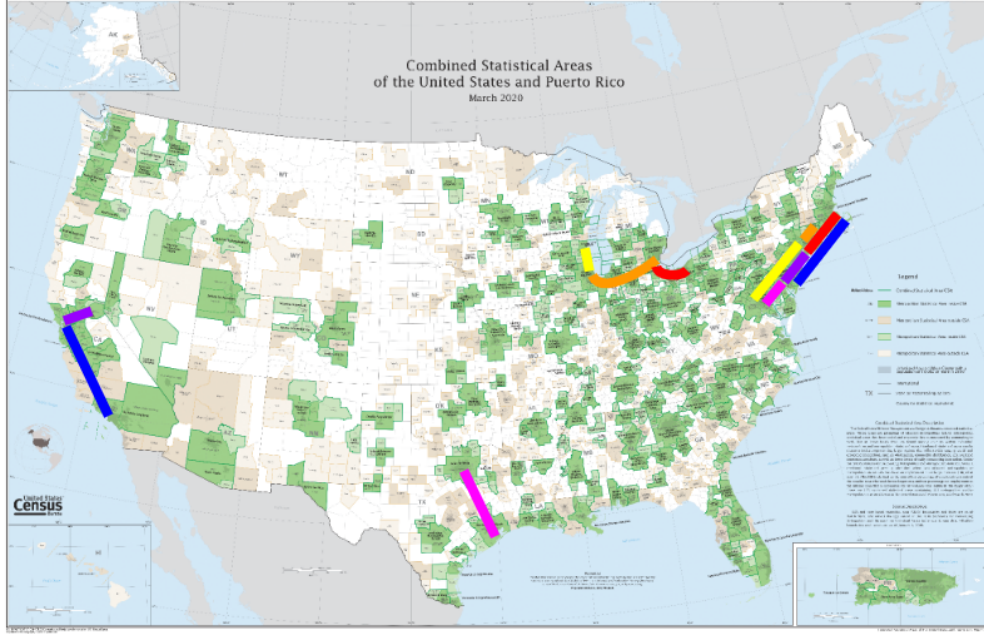


Figure 2: A map of the optimal city pairs to include in a U.S. HSR system solved using our initial linear program. [8]

## Improved Model

To improve our model we can think of additional constraints that would enhance the realism. First, we consider the possibility that there are fixed costs associated with building new rail lines regardless of the rail line's length. This includes station building or modifying and upgrading costs, among other factors. In effect, this adjustment would weight the model so that it favors longer lines over shorter lines. We do not actually add a new constraint to account for this, but instead modify our budget constraint to include a penalty for each new route built. Specifically, for each new route, we take a 200 track-mile cut to our budget. That is, for each new route we can build 200 fewer miles of track in our system. The constraint now looks like

$$\sum_{\forall i} \sum_{\forall j} d_{i,j} x_{i,j} + 200x_{i,j} \leq 4000.$$

Notice that we also have increased the budget to 4000 track-miles to account for the added costs while still building a similar number of routes.

As evidenced by our original results, our model disproportionately favors cities in the Northeast, especially New York. This makes sense given New York's large population and proximity to other major cities, but we want to ensure that no single city is overrepresented. To include a wider variety of cities we can declare that only three lines may go to any given city. The constraint looks like

$$\sum_{\forall i} \sum_{\forall j} x_{i,j} + x_{j,i} \leq 3 \quad \forall i.$$

The improved model, in full, is now,

$$\begin{aligned}
& \text{maximize} && z = \sum_{\forall i} \sum_{\forall j} T_{i,j} x_{i,j} \\
& \text{subject to} && x_{i,j} + x_{j,i} \leq 1 \\
& && \sum_{\forall i} \sum_{\forall j} d_{i,j} x_{i,j} + 200 x_{i,j} \leq 4000 \\
& && \sum_{\forall i} \sum_{\forall j} x_{i,j} + x_{j,i} \leq 3 \quad \forall i.
\end{aligned}$$

Using the Gurobi optimization software we find the optimal solution which is tabulated in Figure 3 and illustrated in Figure 4. Notice that in our improved model  $z = 124.41$  million people are served by our system each year. But in our original model we had  $z = 130.6$  million people. This makes sense because we have only added to our constraints meaning our solution will be less “optimal,” but more realistic.

City Pair $(i, j)$	Travel Demand $(T_{i,j})$	Distance in mi $(d_{i,j})$
New York $\iff$ Washington, D.C.	23.673873	203.520986
Boston $\iff$ New York	20.765652	189.984500
Philadelphia $\iff$ New York	18.604855	80.681696
San Francisco $\iff$ Los Angeles	18.493843	305.738353
Washington, D.C. $\iff$ Philadelphia	9.346010	123.101872
Philadelphia $\iff$ Boston	8.197897	270.602821
Houston $\iff$ Dallas	7.871459	224.299357
Chicago $\iff$ Detroit	7.280050	237.032099
Los Angeles $\iff$ Las Vegas	6.102585	228.347259
San Francisco $\iff$ Sacramento	4.070578	87.663661

Figure 3: The results obtained using a “budget” of 4000 track miles, a penalty of 200 track miles per new line, and a maximum of three lines per city.  $z = 124.41$  million people per year served.

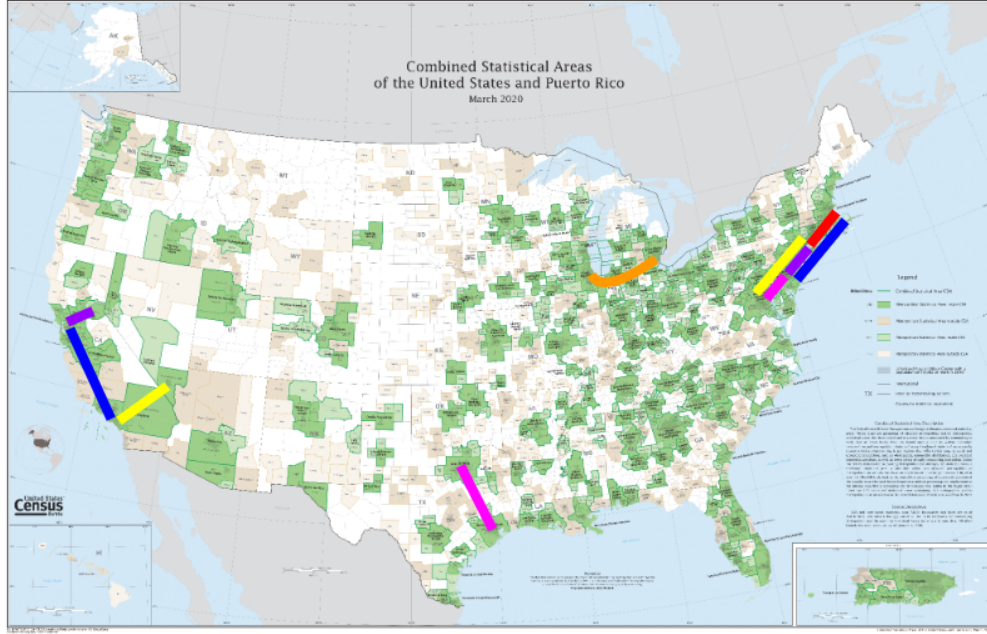


Figure 4: A map of the optimal city pairs to include in a U.S. HSR system solved using our improved linear program. [8]

The second model gives us ten optimal city pairs as opposed to the first model's twelve. Three city pairs, Hartford  $\longleftrightarrow$  New York, Chicago  $\longleftrightarrow$  Milwaukee, and Detroit  $\longleftrightarrow$  Cleveland, are no longer optimal under the second model's constraints. Only one optimal city pair from the second model, Los Angeles  $\longleftrightarrow$  Las Vegas, is not included in the optimal pairs of the original model. This shows that our fixed cost constraint is effective in prioritizing city pairs that have a greater distance between them. Since there are ten cities here, the fixed cost of building each line summed to 2000 track-miles. Since our total budget was 4000 track-miles, we know that less than 2000 miles of track were actually built in this system, making these results comparable to that of our original model.

## Discussion

While our model successfully provides us with potential city pairs for HSR implementation, it is only a rudimentary model. In reality, lines connect with each other to create a network of rail. However, modeling a rail network that spans the entirety of the U.S. is beyond the scope of this project, so the model operates under the assumption of direct routes only. This creates redundancy in our results: if we are already building the rail from Washington D.C.  $\longleftrightarrow$  Philadelphia and Philadelphia  $\longleftrightarrow$  New York, there is no additional cost for creating a route from Washington D.C.  $\longleftrightarrow$  New York.

Theoretically, we could address this issue in the linear program. However, our current model does not take into account the geographic location of cities in relation to each other, only the distance between them. So, while Washington D.C., Philadelphia, and New York are collinear, making it reasonable to travel from Washington D.C. to New York

via Philadelphia, this is not always the case. If three cities are arranged in a triangular formation, like San Antonio, Houston, and Dallas, it makes less sense to expect people to travel to one city via another. For example, it is inefficient to travel from San Antonio to Houston via Dallas, and a network model would need to account for this.

Another area that can be further explored is adding some sort of regional constraint. Even after adding constraints to our original model, the majority of our city pairs are located on the coasts, especially in the Northeast. This might be the most cost-efficient solution given the current constraints, but it underserves large areas of the U.S., like the Southeast and Plains. If HSR is to be implemented in the U.S., there would likely need to be lines in these areas to avoid political animosity.

Another potential change to the model that would adjust the results is our use of CSAs. We use CSAs for cities as opposed to metropolitan statistical areas so that the city pairs our model produced would be at least reasonably far apart (e.g. no Baltimore  $\longleftrightarrow$  D.C. route). HSR is intended to traverse mid to long distances, so we want to avoid putting HSR where it is more efficient to have regional rail. However, CSAs are not necessarily representative of the largest cities. For example, Sioux Falls, SD is not a CSA as there are no substantial bordering cities, but its population is approximately twice that of Rapid City-Spearfish, SD, which is a CSA. Furthermore, population of a CSA is not always indicative of travel demand to that city, which would likely be the case for many popular vacation areas.

Our model is useful as a preliminary study on how to approach the implementation of HSR and identifying optimal city pairs to include in future HSR networks in the U.S.

Our model is useful as a preliminary study for identifying optimal city pairs to include in future high speed rail networks in the U.S.

Our model identifies the optimal city pairs to include in an HSR system in the U.S. based on distances between cities and their population, and is useful as a preliminary study for future HSR networks.

## References

- [1] Richard Nunno. “High Speed Rail Development Worldwide” July 19, 2018.  
<https://www.eesi.org/papers/view/fact-sheet-high-speed-rail-development-worldwide>
- [2] American Public Transportation Association. “Benefits of High-Speed Rail for the United States” March 2021. <https://www.apta.com/research-technical-resources/high-speed-passenger-rail/benefits-of-high-speed-rail-for-the-united-states/>
- [3] United States Census Bureau. “Metropolitan and Micropolitan Statistical Areas Population Totals and Components of Change: 2020-2021.” July 2021.  
<https://www.census.gov/data/tables/time-series/demo/popest/2020s-total-metro-and-micro-statistical-areas.html>
- [4] “Haversine Formula.” In Wikipedia, June 1, 2022. [https://en.wikipedia.org/w/index.php?title=Haversine\\_formula&oldid=1090892412](https://en.wikipedia.org/w/index.php?title=Haversine_formula&oldid=1090892412).



- [5] Pedestrian Observations. “Metcalf’s Law for High-Speed Rail,” February 13, 2020.  
<https://pedestrianobservations.com/2020/02/13/metcalfes-law-for-high-speed-rail/>.
- [6] Martha Lawrence. “Should Countries Invest in High-Speed Rail?” Accessed December 21, 2022. <https://blogs.worldbank.org/transport/should-countries-invest-high-speed-rail>.
- [7] Gurobi Optimization, LLC. “Gurobi Optimizer Reference Manual.” 2022. <https://gurobi.com>.
- [8] United States Census Bureau. “Combined Statistical Areas Map.” March 2020.  
[https://www2.census.gov/geo/maps/metroarea/us\\_wall/Mar2020/CSA\\_WallMap\\_Mar2020.pdf](https://www2.census.gov/geo/maps/metroarea/us_wall/Mar2020/CSA_WallMap_Mar2020.pdf)

## Code Availability

The code used to perform our analysis is available on GitHub at <https://github.com/n-loo/HSR-Optimization>

We have adhered to the honor code on this assignment. *Ellie Jensen. Niels Vanderloo.*