

M1 Project: Railway Delay Statistics and Modelling

Lusine Shirvanyan

Erasmus Mundus Master in Complex Systems Science,
Centre for Complexity Science, University of Warwick

We consider the problem of extracting statistical information regarding the distribution of train delays in the UK. Our work seeks to extend the work of Briggs and Beck [1], who showed that the q -exponential distribution accurately fits the distribution of observed delays for UK rail network. We studied the statistics of consecutive delays to understand the spatial and temporal correlation between the delays. The intended application of this research is to implement an accurate journey planning algorithm that can incorporate real-time data.

1 Introduction

The rail network in the UK consists of over 15,000 km of track. Planning journeys with a desired arrival time is a challenging problem, as not only are there often multiple routes between two points, to accurately predict journey time an understanding is needed of the reliability of the different parts of the network. Moreover, the likelihood of a delay can vary with the time, day of the week, as well as factors such as weather and line obstructions that can affect large parts of the network, creating a complicated dependency structure.

To help people plan their journeys, there are a number of route planning algorithms made available online. However, these are mostly limited—they plan journey under the assumption that no train is ever delayed! This is particularly problematic when a journey requires a change of trains, as even a short delay can result in a long delay in the event of a missed connection.

Although, statistics regarding the reliability of the different train routes are routinely gathered, they are normally too simple to be very helpful for route planning. For example, if a train route has 20 trains per day, and is 95% reliable, this mean that one train is cancelled, or it could mean that once every five days, five trains in a row are cancelled. The reliability is the same, but it makes a great difference to the distribution of journey times experienced by passengers.

This is why it is important to have an understanding of spatial and temporal correlations between delays. For example, if we are travelling from York to King Cross and the previous

train was delayed by three minutes, what is the reliability of trains departing over the next two hours?

Because of engineering works last for few weeks, it is desirable to be able to observe changes in reliability in the network over quite short time scales. It is also useful in the case of changes in timetabling to be able to adapt in a short period of time. Even though there is potentially many years of data is available, it important to able to extract information from short amount of data. It is useful to be able to fit parametric models of data.

In their work, Keith and Beck showed that collected delays from 23 major stations in UK for the period September 2005-October 2006 can be accurately modelled by a two-parameter q -exponential distribution [1]. Their preliminary investigation implied that the following model (1) would fit well.

$$e_{q,b,c}(t) = c(1 + b(q - 1)t)^{\frac{1}{1-q}} \quad (1)$$

Here t is the delay, $0 < q < 2$ and $b > 0$ are shape parameters, and c is a normalization parameter. Levenberg-Marquardt method of solving nonlinear least squares problems has been used for finding the best fit parameters for the distribution of delays. Figure 1 illustrates some of the results from Keiths and Becks work [1] for fitting model for all delays.

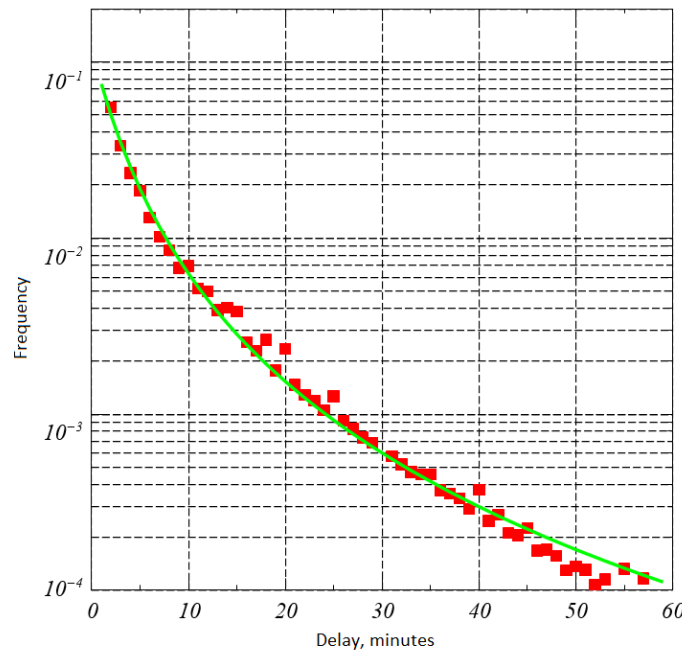


Figure 1: All train data and best-fit q -exponential: $q = 1.355 \pm 8.8 \times 10^{-5}$,
 $b = 0.524 \pm 2.5 \times 10^{-8}$ (by Keith and Beck[1])

The goal of this work is to extend understanding of the statistics of multiple trains. First of all we checked that new data for the period of time from 2014 to 2015 still fits weel the

q -exponential distribution. Then we focus on new types of data analysis, such as conditional probabilities and correlations between delays, because they can answer when a train is delayed, what is the probability that it is still delayed later in its journey? Or what is the probability that later trains on the same line are delayed? This will allow accurate simulation of passenger movement through a rail system with delays, the prediction of arrival times and maximization of departure time in a trip planning algorithm.

2 Data

2.1 The structure of the collected data

The advent of real-time train information available on the internet for the British network (<http://www.nationalrail.co.uk/ldb/livedepartures.asp>) has made it possible to gather a huge amount of data. We collected data on departure times for 10 major stations for the period January 2014 to May 2015 using software which downloads the real-time information from webpage every minute for each station. Data are collected in YAML(Yet Another Markup Language) files and each file contains information about the departures from one particular station to the other stations for one day. Each row represents a list of pairs of departure time and the corresponding delay for the final destination. As each train eventually departs, the most recent delay value is saved to a database.

Figure 2 presents an example of such yaml file. According to it, the train from YRK to KGX departing at 07:01 from YRK has been delayed for 3 minutes, while the next one at 07:37 has been delayed for 5 minutes and etc.

```

1 # Last updated 2015-02-11 01:01
2 # Trains from YRK
3
4 ABD: {'0737': 3, '1154': 0, '1532': 0, '1555': 0, '1754': 0}
5 BHM: {'1936': 0, '2035': 0, '2045': 0}
6 BFN: {'0719': 0, '0827': 0, '0918': 0, '1027': 3, '1118': 0, '1227': 5, '1318': 0, '1427': 0, '1518': 0, '1627': 0, '1718': 0, '1827': 0, '1918': 0}
7 BRI: {'1845': 0, '1945': 0}
8 BUY: {'0755': 0, '0845': 0, '0911': 0, '1011': 0, '1111': 0, '1211': 0, '1311': 0, '1411': 0, '1511': 0, '1611': 0, '1704': 0, '1811': 0, '1913': 0, '2011': 0}
9 DEE: {'1732': 0}
10 EDB: {'0855': 0, '0934': 0, '0937': 2, '0953': 4, '1054': 1, '1132': 0, '1136': 0, '1253': 0, '1332': 0, '1336': 3, '1454': 0, '1654': 0, '1836': 9, '1854': 0,
11 GLC: {'0829': 0, '1032': 0, '1232': 0, '1432': 0, '1632': 0, '1736': 2, '1833': 9}
12 HUL: {'0730': 2, '0843': 0, '1019': 0, '1145': 0, '1247': 0, '1344': 0, '1502': 5, '1611': 0, '1725': 0, '1840': 2, '2229': 0}
13 INV: {'1355': 0}
14 KGX: {'0701': 3, '0737': 5, '0802': 0, '0820': 0, '0831': 0, '0856': 2, '0931': 28, '0957': 10, '1003': 6, '1027': 0, '1031': 5, '1059': 0, '1131': 0, '1157': 6, '1181': 0, '1201': 0, '1215': 0, '1219': 0, '1223': 0, '1227': 0, '1231': 0, '1235': 0, '1239': 0, '1243': 0, '1247': 0, '1251': 0, '1255': 0, '1259': 0, '1263': 0, '1267': 0, '1271': 0, '1275': 0, '1279': 0, '1283': 0, '1287': 0, '1291': 0, '1295': 0, '1299': 0, '1303': 0, '1307': 0, '1311': 0, '1315': 0, '1319': 0, '1323': 0, '1327': 0, '1331': 0, '1335': 0, '1339': 0, '1343': 0, '1347': 0, '1351': 0, '1355': 0, '1359': 0, '1363': 0, '1367': 0, '1371': 0, '1375': 0, '1379': 0, '1383': 0, '1387': 0, '1391': 0, '1395': 0, '1399': 0, '1403': 0, '1407': 0, '1411': 0, '1415': 0, '1419': 0, '1423': 0, '1427': 0, '1431': 0, '1435': 0, '1439': 0, '1443': 0, '1447': 0, '1451': 0, '1455': 0, '1459': 0, '1463': 0, '1467': 0, '1471': 0, '1475': 0, '1479': 0, '1483': 0, '1487': 0, '1491': 0, '1495': 0, '1499': 0, '1503': 0, '1507': 0, '1511': 0, '1515': 0, '1519': 0, '1523': 0, '1527': 0, '1531': 0, '1535': 0, '1539': 0, '1543': 0, '1547': 0, '1551': 0, '1555': 0, '1559': 0, '1563': 0, '1567': 0, '1571': 0, '1575': 0, '1579': 0, '1583': 0, '1587': 0, '1591': 0, '1595': 0, '1599': 0, '1603': 0, '1607': 0, '1611': 0, '1615': 0, '1619': 0, '1623': 0, '1627': 0, '1631': 0, '1635': 0, '1639': 0, '1643': 0, '1647': 0, '1651': 0, '1655': 0, '1659': 0, '1663': 0, '1667': 0, '1671': 0, '1675': 0, '1679': 0, '1683': 0, '1687': 0, '1691': 0, '1695': 0, '1699': 0, '1703': 0, '1707': 0, '1711': 0, '1715': 0, '1719': 0, '1723': 0, '1727': 0, '1731': 0, '1735': 0, '1739': 0, '1743': 0, '1747': 0, '1751': 0, '1755': 0, '1759': 0, '1763': 0, '1767': 0, '1771': 0, '1775': 0, '1779': 0, '1783': 0, '1787': 0, '1791': 0, '1795': 0, '1799': 0, '1803': 0, '1807': 0, '1811': 0, '1815': 0, '1819': 0, '1823': 0, '1827': 0, '1831': 0, '1835': 0, '1839': 0, '1843': 0, '1847': 0, '1851': 0, '1855': 0, '1859': 0, '1863': 0, '1867': 0, '1871': 0, '1875': 0, '1879': 0, '1883': 0, '1887': 0, '1891': 0, '1895': 0, '1899': 0, '1903': 0, '1907': 0, '1911': 0, '1915': 0, '1919': 0, '1923': 0, '1927': 0, '1931': 0, '1935': 0, '1939': 0, '1943': 0, '1947': 0, '1951': 0, '1955': 0, '1959': 0, '1963': 0, '1967': 0, '1971': 0, '1975': 0, '1979': 0, '1983': 0, '1987': 0, '1991': 0, '1995': 0, '1999': 0, '2003': 0, '2007': 0, '2011': 0, '2015': 0, '2019': 0, '2023': 0, '2027': 0, '2031': 0, '2035': 0, '2039': 0, '2043': 0, '2047': 0, '2051': 0, '2055': 0, '2059': 0, '2063': 0, '2067': 0, '2071': 0, '2075': 0, '2079': 0, '2083': 0, '2087': 0, '2091': 0, '2095': 0, '2099': 0, '2103': 0, '2107': 0, '2111': 0, '2115': 0, '2119': 0, '2123': 0, '2127': 0, '2131': 0, '2135': 0, '2139': 0, '2143': 0, '2147': 0, '2151': 0, '2155': 0, '2159': 0, '2163': 0, '2167': 0, '2171': 0, '2175': 0, '2179': 0, '2183': 0, '2187': 0, '2191': 0, '2195': 0, '2199': 0, '2203': 0, '2207': 0, '2211': 0, '2215': 0, '2219': 0, '2223': 0, '2227': 0, '2231': 0, '2235': 0, '2239': 0, '2243': 0, '2247': 0, '2251': 0, '2255': 0, '2259': 0, '2263': 0, '2267': 0, '2271': 0, '2275': 0, '2279': 0, '2283': 0, '2287': 0, '2291': 0, '2295': 0, '2299': 0, '2303': 0, '2307': 0, '2311': 0, '2315': 0, '2319': 0, '2323': 0, '2327': 0, '2331': 0, '2335': 0, '2339': 0, '2343': 0, '2347': 0, '2351': 0, '2355': 0, '2359': 0, '2363': 0, '2367': 0, '2371': 0, '2375': 0, '2379': 0, '2383': 0, '2387': 0, '2391': 0, '2395': 0, '2399': 0, '2403': 0, '2407': 0, '2411': 0, '2415': 0, '2419': 0, '2423': 0, '2427': 0, '2431': 0, '2435': 0, '2439': 0, '2443': 0, '2447': 0, '2451': 0, '2455': 0, '2459': 0, '2463': 0, '2467': 0, '2471': 0, '2475': 0, '2479': 0, '2483': 0, '2487': 0, '2491': 0, '2495': 0, '2499': 0, '2503': 0, '2507': 0, '2511': 0, '2515': 0, '2519': 0, '2523': 0, '2527': 0, '2531': 0, '2535': 0, '2539': 0, '2543': 0, '2547': 0, '2551': 0, '2555': 0, '2559': 0, '2563': 0, '2567': 0, '2571': 0, '2575': 0, '2579': 0, '2583': 0, '2587': 0, '2591': 0, '2595': 0, '2599': 0, '2603': 0, '2607': 0, '2611': 0, '2615': 0, '2619': 0, '2623': 0, '2627': 0, '2631': 0, '2635': 0, '2639': 0, '2643': 0, '2647': 0, '2651': 0, '2655': 0, '2659': 0, '2663': 0, '2667': 0, '2671': 0, '2675': 0, '2679': 0, '2683': 0, '2687': 0, '2691': 0, '2695': 0, '2699': 0, '2703': 0, '2707': 0, '2711': 0, '2715': 0, '2719': 0, '2723': 0, '2727': 0, '2731': 0, '2735': 0, '2739': 0, '2743': 0, '2747': 0, '2751': 0, '2755': 0, '2759': 0, '2763': 0, '2767': 0, '2771': 0, '2775': 0, '2779': 0, '2783': 0, '2787': 0, '2791': 0, '2795': 0, '2799': 0, '2803': 0, '2807': 0, '2811': 0, '2815': 0, '2819': 0, '2823': 0, '2827': 0, '2831': 0, '2835': 0, '2839': 0, '2843': 0, '2847': 0, '2851': 0, '2855': 0, '2859': 0, '2863': 0, '2867': 0, '2871': 0, '2875': 0, '2879': 0, '2883': 0, '2887': 0, '2891': 0, '2895': 0, '2899': 0, '2903': 0, '2907': 0, '2911': 0, '2915': 0, '2919': 0, '2923': 0, '2927': 0, '2931': 0, '2935': 0, '2939': 0, '2943': 0, '2947': 0, '2951': 0, '2955': 0, '2959': 0, '2963': 0, '2967': 0, '2971': 0, '2975': 0, '2979': 0, '2983': 0, '2987': 0, '2991': 0, '2995': 0, '2999': 0, '3003': 0, '3007': 0, '3011': 0, '3015': 0, '3019': 0, '3023': 0, '3027': 0, '3031': 0, '3035': 0, '3039': 0, '3043': 0, '3047': 0, '3051': 0, '3055': 0, '3059': 0, '3063': 0, '3067': 0, '3071': 0, '3075': 0, '3079': 0, '3083': 0, '3087': 0, '3091': 0, '3095': 0, '3099': 0, '3103': 0, '3107': 0, '3111': 0, '3115': 0, '3119': 0, '3123': 0, '3127': 0, '3131': 0, '3135': 0, '3139': 0, '3143': 0, '3147': 0, '3151': 0, '3155': 0, '3159': 0, '3163': 0, '3167': 0, '3171': 0, '3175': 0, '3179': 0, '3183': 0, '3187': 0, '3191': 0, '3195': 0, '3199': 0, '3203': 0, '3207': 0, '3211': 0, '3215': 0, '3219': 0, '3223': 0, '3227': 0, '3231': 0, '3235': 0, '3239': 0, '3243': 0, '3247': 0, '3251': 0, '3255': 0, '3259': 0, '3263': 0, '3267': 0, '3271': 0, '3275': 0, '3279': 0, '3283': 0, '3287': 0, '3291': 0, '3295': 0, '3299': 0, '3303': 0, '3307': 0, '3311': 0, '3315': 0, '3319': 0, '3323': 0, '3327': 0, '3331': 0, '3335': 0, '3339': 0, '3343': 0, '3347': 0, '3351': 0, '3355': 0, '3359': 0, '3363': 0, '3367': 0, '3371': 0, '3375': 0, '3379': 0, '3383': 0, '3387': 0, '3391': 0, '3395': 0, '3399': 0, '3403': 0, '3407': 0, '3411': 0, '3415': 0, '3419': 0, '3423': 0, '3427': 0, '3431': 0, '3435': 0, '3439': 0, '3443': 0, '3447': 0, '3451': 0, '3455': 0, '3459': 0, '3463': 0, '3467': 0, '3471': 0, '3475': 0, '3479': 0, '3483': 0, '3487': 0, '3491': 0, '3495': 0, '3499': 0, '3503': 0, '3507': 0, '3511': 0, '3515': 0, '3519': 0, '3523': 0, '3527': 0, '3531': 0, '3535': 0, '3539': 0, '3543': 0, '3547': 0, '3551': 0, '3555': 0, '3559': 0, '3563': 0, '3567': 0, '3571': 0, '3575': 0, '3579': 0, '3583': 0, '3587': 0, '3591': 0, '3595': 0, '3599': 0, '3603': 0, '3607': 0, '3611': 0, '3615': 0, '3619': 0, '3623': 0, '3627': 0, '3631': 0, '3635': 0, '3639': 0, '3643': 0, '3647': 0, '3651': 0, '3655': 0, '3659': 0, '3663': 0, '3667': 0, '3671': 0, '3675': 0, '3679': 0, '3683': 0, '3687': 0, '3691': 0, '3695': 0, '3699': 0, '3703': 0, '3707': 0, '3711': 0, '3715': 0, '3719': 0, '3723': 0, '3727': 0, '3731': 0, '3735': 0, '3739': 0, '3743': 0, '3747': 0, '3751': 0, '3755': 0, '3759': 0, '3763': 0, '3767': 0, '3771': 0, '3775': 0, '3779': 0, '3783': 0, '3787': 0, '3791': 0, '3795': 0, '3799': 0, '3803': 0, '3807': 0, '3811': 0, '3815': 0, '3819': 0, '3823': 0, '3827': 0, '3831': 0, '3835': 0, '3839': 0, '3843': 0, '3847': 0, '3851': 0, '3855': 0, '3859': 0, '3863': 0, '3867': 0, '3871': 0, '3875': 0, '3879': 0, '3883': 0, '3887': 0, '3891': 0, '3895': 0, '3899': 0, '3903': 0, '3907': 0, '3911': 0, '3915': 0, '3919': 0, '3923': 0, '3927': 0, '3931': 0, '3935': 0, '3939': 0, '3943': 0, '3947': 0, '3951': 0, '3955': 0, '3959': 0, '3963': 0, '3967': 0, '3971': 0, '3975': 0, '3979': 0, '3983': 0, '3987': 0, '3991': 0, '3995': 0, '3999': 0, '4003': 0, '4007': 0, '4011': 0, '4015': 0, '4019': 0, '4023': 0, '4027': 0, '4031': 0, '4035': 0, '4039': 0, '4043': 0, '4047': 0, '4051': 0, '4055': 0, '4059': 0, '4063': 0, '4067': 0, '4071': 0, '4075': 0, '4079': 0, '4083': 0, '4087': 0, '4091': 0, '4095': 0, '4099': 0, '4103': 0, '4107': 0, '4111': 0, '4115': 0, '4119': 0, '4123': 0, '4127': 0, '4131': 0, '4135': 0, '4139': 0, '4143': 0, '4147': 0, '4151': 0, '4155': 0, '4159': 0, '4163': 0, '4167': 0, '4171': 0, '4175': 0, '4179': 0, '4183': 0, '4187': 0, '4191': 0, '4195': 0, '4199': 0, '4203': 0, '4207': 0, '4211': 0, '4215': 0, '4219': 0, '4223': 0, '4227': 0, '4231': 0, '4235': 0, '4239': 0, '4243': 0, '4247': 0, '4251': 0, '4255': 0, '4259': 0, '4263': 0, '4267': 0, '4271': 0, '4275': 0, '4279': 0, '4283': 0, '4287': 0, '4291': 0, '4295': 0, '4299': 0, '4303': 0, '4307': 0, '4311': 0, '4315': 0, '4319': 0, '4323': 0, '4327': 0, '4331': 0, '4335': 0, '4339': 0, '4343': 0, '4347': 0, '4351': 0, '4355': 0, '4359': 0, '4363': 0, '4367': 0, '4371': 0, '4375': 0, '4379': 0, '4383': 0, '4387': 0, '4391': 0, '4395': 0, '4399': 0, '4403': 0, '4407': 0, '4411': 0, '4415': 0, '4419': 0, '4423': 0, '4427': 0, '4431': 0, '4435': 0, '4439': 0, '4443': 0, '4447': 0, '4451': 0, '4455': 0, '4459': 0, '4463': 0, '4467': 0, '4471': 0, '4475': 0, '4479': 0, '4483': 0, '4487': 0, '4491': 0, '4495': 0, '4499': 0, '4503': 0, '4507': 0, '4511': 0, '4515': 0, '4519': 0, '4523': 0, '4527': 0, '4531': 0, '4535': 0, '4539': 0, '4543': 0, '4547': 0, '4551': 0, '4555': 0, '4559': 0, '4563': 0, '4567': 0, '4571': 0, '4575': 0, '4579': 0, '4583': 0, '4587': 0, '4591': 0, '4595': 0, '4599': 0, '4603': 0, '4607': 0, '4611': 0, '4615': 0, '4619': 0, '4623': 0, '4627': 0, '4631': 0, '4635': 0, '4639': 0, '4643': 0, '4647': 0, '4651': 0, '4655': 0, '4659': 0, '4663': 0, '4667': 0, '4671': 0, '4675': 0, '4679': 0, '4683': 0, '4687': 0, '4691': 0, '4695': 0, '4699': 0, '4703': 0, '4707': 0, '4711': 0, '4715': 0, '4719': 0, '4723': 0, '4727': 0, '4731': 0, '4735': 0, '4739': 0, '4743': 0, '4747': 0, '4751': 0, '4755': 0, '4759': 0, '4763': 0, '4767': 0, '4771': 0, '4775': 0, '4779': 0, '4783': 0, '4787': 0, '4791': 0, '4795': 0, '4799': 0, '4803': 0, '4807': 0, '4811': 0, '4815': 0, '4819': 0, '4823': 0, '4827': 0, '4831': 0, '4835': 0, '4839': 0, '4843': 0, '4847': 0, '4851': 0, '4855': 0, '4859': 0, '4863': 0, '4867': 0, '4871': 0, '4875': 0, '4879': 0, '4883': 0, '4887': 0, '4891': 0, '4895': 0, '4899': 0, '4903': 0, '4907': 0, '4911': 0, '4915': 0, '4919': 0, '4923': 0, '4927': 0, '4931': 0, '4935': 0, '4939': 0, '4943': 0, '4947': 0, '4951': 0, '4955': 0, '4959': 0, '4963': 0, '4967': 0, '4971': 0, '4975': 0, '4979': 0, '4983': 0, '4987': 0, '4991': 0, '4995': 0, '4999': 0, '5003': 0, '5007': 0, '5011': 0, '5015': 0, '5019': 0, '5023': 0, '5027': 0, '5031': 0, '5035': 0, '5039': 0, '5043': 0, '5047': 0, '5051': 0, '5055': 0, '5059': 0, '5063': 0, '5067': 0, '5071': 0, '5075': 0, '5079': 0, '5083': 0, '5087': 0, '5091': 0, '5095': 0, '5099': 0, '5103': 0, '5107': 0, '5111': 0, '5115': 0, '5119': 0, '5123': 0, '5127': 0, '5131': 0, '5135': 0, '5139': 0, '5143': 0, '5147': 0, '5151': 0, '5155': 0, '5159': 0, '5163': 0, '5167': 0, '5171': 0, '5175': 0, '5179': 0, '5183': 0, '5187': 0, '5191': 0, '5195': 0, '5199': 0, '5203': 0, '5207': 0, '5211': 0, '5215': 0, '5219': 0, '5223': 0, '5227': 0, '5231': 0, '5235': 0, '5239': 0, '5243': 0, '5247': 0, '5251': 0, '5255': 0, '5259': 0, '5263': 0, '5267': 0, '5271': 0, '5275': 0, '5279': 0, '5283': 0, '5287': 0, '5291': 0, '5295': 0, '5299': 0, '5303': 0, '5307': 0, '5311': 0, '5315': 0, '5319': 0, '5323': 0, '5327': 0, '5331': 0, '5335': 0, '5339': 0, '5343': 0, '5347': 0, '5351': 0, '5355': 0, '5359': 0, '5363': 0, '5367': 0, '5371': 0, '5375': 0, '5379': 0, '5383': 0, '5387': 0, '5391': 0, '5395': 0, '5399': 0, '5403': 0, '5407': 0, '5411': 0, '5415': 0, '5419': 0, '5423': 0, '5427': 0, '5431': 0, '5435': 0, '5439': 0, '5443': 0, '5447': 0, '5451': 0, '5455': 0, '5459': 0, '5463': 0, '5467': 0, '5471': 0, '5475': 0, '5479': 0, '5483': 0, '5487': 0, '5491': 0, '5495': 0, '5499': 0, '5503': 0, '5507': 0, '5511': 0, '5515': 0, '5519': 0, '5523': 0, '5527': 0, '5531': 0, '5535': 0, '5539': 0, '5543': 0, '5547': 0, '5551': 0, '5555': 0, '5559': 0, '5563': 0, '5567': 0, '5571': 0, '5575': 0, '5579': 0, '5583': 0, '5587': 0, '5591': 0, '5595': 0, '5599': 0, '5603': 0, '5607': 0, '5611': 0, '5615': 0, '5619': 0, '5623': 0, '5627': 0, '5631': 0, '5635': 0, '5639': 0, '5643': 0, '5647': 0, '5651': 0, '5655': 0, '5659': 0, '5663': 0, '5667': 0, '5671': 0, '5675': 0, '5679': 0, '5683': 0, '5687': 0, '5691': 0, '5695': 0, '5699': 0, '5703': 0, '5707': 0, '5711': 0, '5715': 0, '5719': 0, '5723': 0, '5727': 0, '5731': 0, '5735': 0, '5739': 0, '5743': 0, '5747': 0, '5751': 0, '5755': 0, '5759': 0, '5763': 0, '5767': 0, '5771': 0, '5775': 0, '5779': 0, '5783': 0, '5787': 0, '5791': 0, '5795': 0, '5799': 0, '5803': 0, '5807': 0, '5811': 0, '5815': 0, '5819': 0, '5823': 0, '5827': 0, '5831': 0, '5835': 0, '5839': 0, '5843': 0, '5847': 0, '5851': 0, '5855': 0, '5859': 0, '5863': 0, '5867': 0, '5871': 0, '5875': 0, '5879': 0, '5883': 0, '5887': 0, '5891': 0, '5895': 0, '5899': 0, '5903': 0, '5907': 0, '5911': 0, '5915': 0, '5919': 0, '5923': 0, '5927': 0, '5931': 0, '5935': 0, '5939': 0, '5943': 0, '5947': 0, '5951': 0, '5955': 0, '5959': 0, '5963': 0, '5967': 0, '5971': 0, '5975': 0, '5979': 0, '5983': 0, '5987': 0, '5991': 0, '5995': 0, '5999': 0, '6003': 0, '6007': 0, '6011': 0, '6015': 0, '6019': 0, '6023': 0, '6027': 0, '6031': 0, '6035': 0, '6039': 0, '6043': 0, '6047': 0, '6051': 0, '6055': 0, '6059': 0, '6063': 0, '6067': 0, '6071': 
```

2.2 Limitations of the collected data

There were some limitations of the representation of the collected data, which turned out to be an obstacle for some of the analysis. First of all, the fact that data are spread between many files makes it difficult to collect all the departure data for a particular station for a specific period of time. In order to obtain this data, we would have to scan all the files. This would be time consuming since only 10% or even less of data is relevant to a given station as the final destination. Another problem was the format of departure times in the file, which, as we can see in Figure 2, is not the standard format for date/time. This was not allowing us to make queries for collecting departures for some particular times of the day, for example, looking at delays in off-peak times only or comparing delays of consecutive trains. We managed to overcome these limitations by transferring all data from yaml files to SQLite database and keeping all data in an easy to use format. Before giving more details about the SQLite generated database, we should focus our attention to some limitations of data.

As mentioned before, each file contains information about departures from one source station to the final destinations only and there is no information about intermediate stations for each journey. This lack of information did not allow us to analyze the correlations between delays in different stations from the same journey. This resulted in setting the additional task of train identification, which will be discussed in more detail in section 3.1. And the last limitation is that only departure times are provided in the data, with no information about arrival times, which had some impact on the train identification task.

2.3 Generated SQLite database

Because of the limitations of data representation mentioned in subsection 2.2 and its implementation advantages, SQLite database was chosen to store the data. As we can see in Figure 3, three tables were constructed for storing the data from the yaml files. Stations table contains names of all stations collected from all files and additionally assigns unique StationID to each of them, which makes querying of data much faster, by avoiding string comparisons every time. The Departures table keeps all departure records from the files, by assigning unique DepID to each of them. Corresponding ID-s from the Stations table are recorded for the source and the destination stations (marked in green in figure). Additional Weekday column was added to table, which contains information about the day of the week extracted from the date, which makes possible the analysis considering the day of the week. The Trains table was constructed to store the results of the train identification task. TrainID and JourneyID in the table are used for grouping different records from the Departures table into the same journey or same type of train.

StationID	StationName	
1	0	Station is not recogni...
2	1	YRK
3	2	MAN
4	3	LIV
5	4	IPS
6	5	DON
7	6	COL
8	7	CHM
9	8	PBO
10	9	LDS
11	10	BHM
12	11	KGX
13	12	ABD
14	13	ALD
15	14	APP
16	15	AST
17	16	AWK
18	17	AYW

DepID	SrcStationID	DestStationID	DepTime	Delay	Date	WeekDay
1	4	6	16:52:00	3	2013-12-23	1
2	4	6	17:33:00	16	2013-12-23	1
3	4	6	19:23:00	9	2013-12-23	1
4	4	6	20:22:00	23	2013-12-23	1
5	4	6	21:28:00	15	2013-12-23	1
6	4	6	23:22:00	10	2013-12-23	1
7	4	8	13:58:00	0	2013-12-23	1
8	4	8	16:00:00	2	2013-12-23	1
9	4	8	17:49:00	0	2013-12-23	1
10	4	8	20:00:00	2	2013-12-23	1
11	4	8	11:58:00	0	2014-01-09	4
12	4	6	17:33:00	9	2014-01-14	2
13	4	6	19:23:00	2	2014-01-14	2
14	4	6	20:22:00	0	2014-01-14	2
15	4	6	21:28:00	0	2014-01-14	2
16	4	6	23:22:00	0	2014-01-14	2
17	4	8	11:58:00	0	2014-01-14	2
18	4	8	17:49:00	0	2014-01-14	2

TrainID	JourneyID	DepID
1	1	20955
1	1	68769
2	1	96429
3	1	20950
4	1	68773
5	1	96422
6	1	20959
7	1	68766
8	1	96433
9	1	20960
10	1	68776
11	1	96443
12	1	20963
13	1	68780
14	1	96446
15	1	20929
16	1	68782
17	1	96450
18	1	20955

Figure 3: Database generated from Yaml files.

3 Results

3.1 Train identification

We grouped departures for the same final destination in such order, that with high probability its the same train that goes through all that stations to the same destination. Obviously, this could not be done using only the data we have, since there is no information such as distance between stations and the route trains can take to the final station. Therefore we gathered information about the different types of trains for one particular journey and we then tried to match these to the departures from the database, by providing some basic properties about each type of train, such as in which stations the train stops and what the difference between departure times for two consecutive stations. Figure 3 illustrates the result of train identification for journeys from YRK station to KGX station. We identified 3 different types of trains going from YRK to KGX, 2 of which have stops in DON and PBO stations (red and black), whereas third type(blue) stops in PBO only. On the left side is a plot for data collected manually from the website for 01-06-2015, and on the right side are plotted trains identified from data set for some 07-05-2015. Most of the trains were identified correctly, however, there were some misidentified trains too. In addition, we can notice that we do not have last, final station in our second graph, which is because there is no information provided about arrival times. However, the results were good enough for using them for correlation analyses.

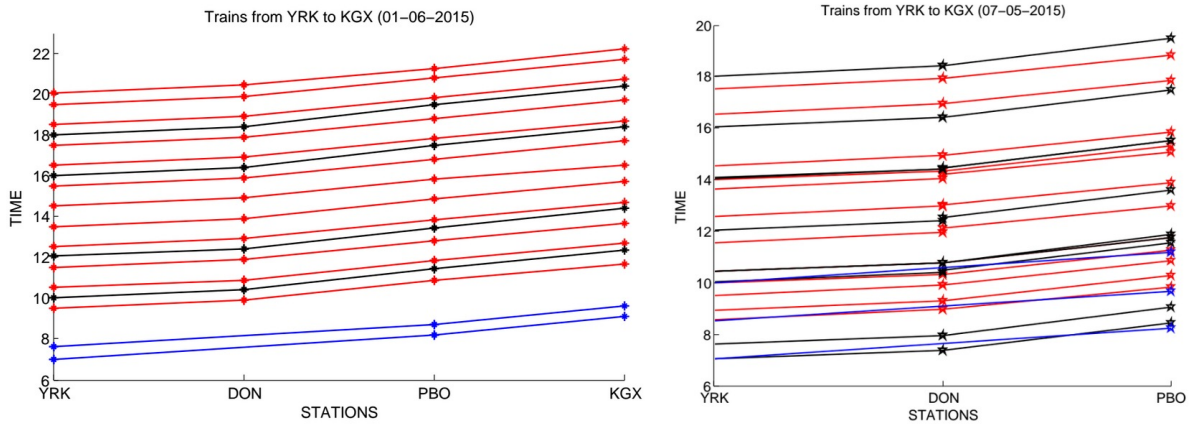


Figure 4: Train identification. On the left different trains and journeys collected manually. On the right side trains/journeys identified from our data.

3.2 Correlation between delays in different stations

Now, when we have some information about departures from same journey, we can look at correlation of delays between consecutive stations. Scatter plot in Figure 5 illustrates correlation of delays for each pair of stations from YRK to KGX. As we can see, there is mainly a positive correlation between all stations, which is quite natural, since when a train is delayed in one of the stations, it is with a high probability that it was either delayed previously on its journey, or it will be delayed later. However, we see some big delays on the plot, which did not lead to delay in next station, which is very unlikely in reality. This probably means that those delays are from misidentified trains, i.e. departures grouped together as one part of one journey are independent in real.

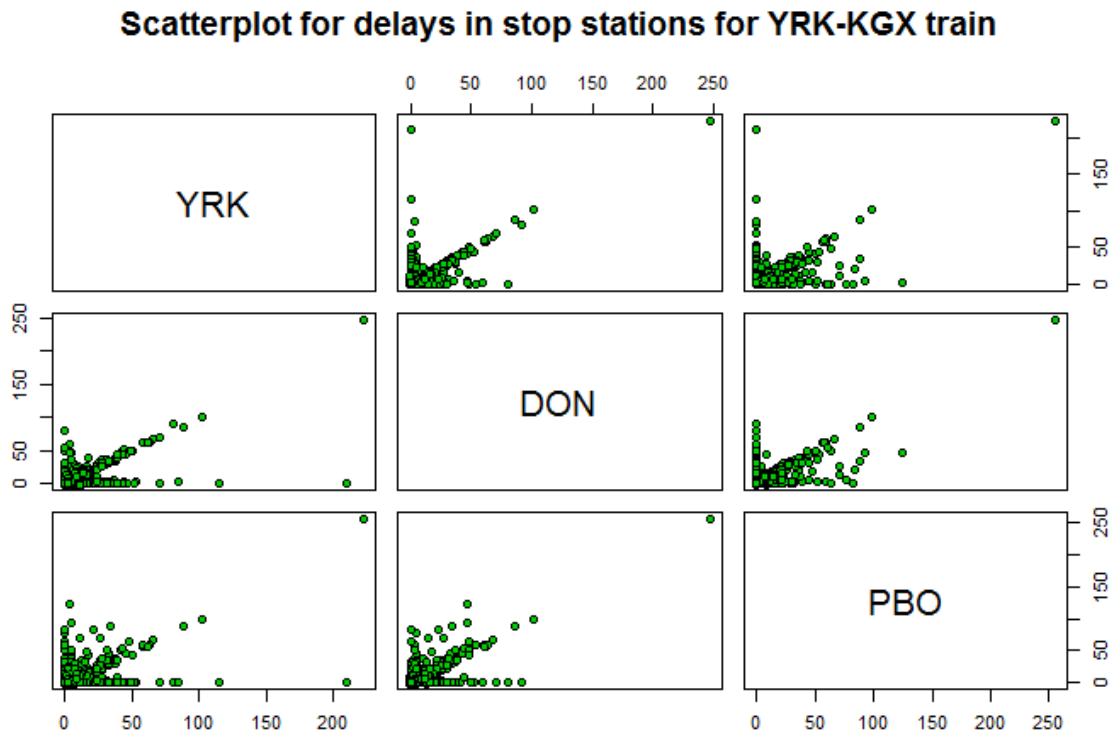


Figure 5: Correlation between delays in different stations of same journey.

3.3 Distribution of delays for the days of the week

Another thing we investigated is the distribution of train delays for the different days of the week. The violin plot in Figure 6 presents clearly the differences and similarities between delays for each day of the week. For all days the mean delay is approximately 4 minutes, whereas it's slightly more (about 2 minutes) for Monday. Another noteworthy fact is that the probability of big delays on Sunday is much higher comparing to all other days, while there is much smaller range of probable delays on Wednesday. Overall, the range of delays decreases at the beginning of the week and then increases after Wednesday. The reasons of this phenomenon can be various, however, the important thing here is that taking into account the day of the week may help doing more accurate predictions about delays.

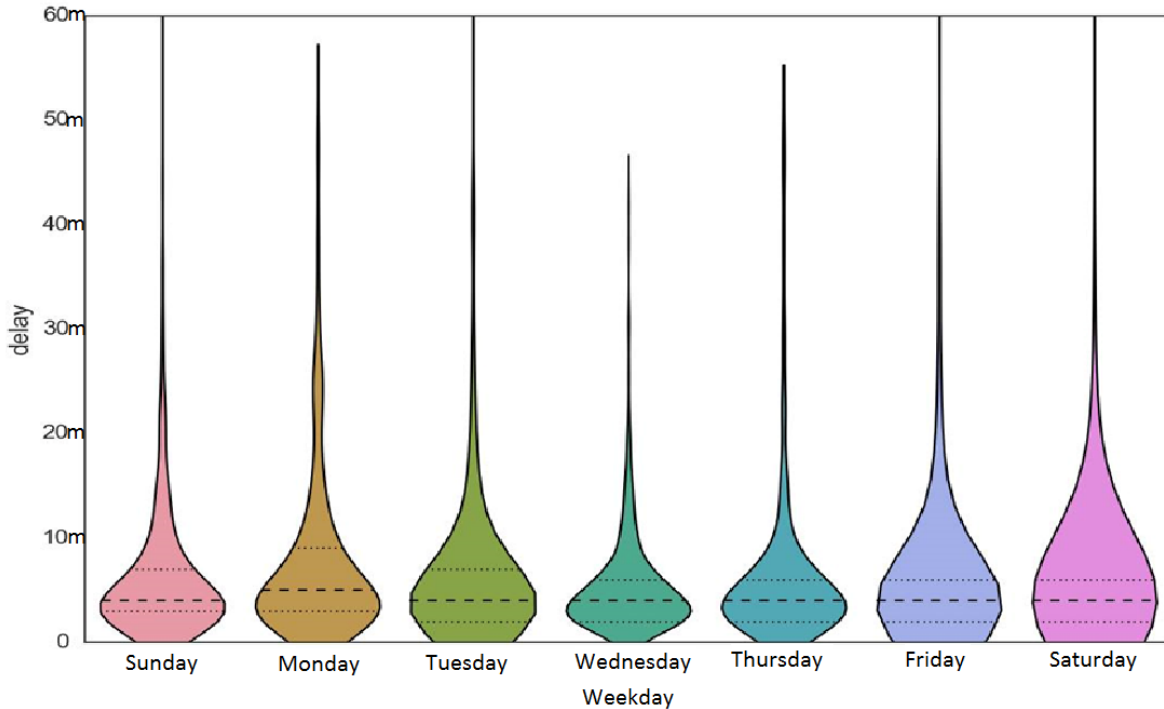


Figure 6: Distribution of delays for the days of the week.

3.4 Conditional probabilities

The last thing we analyzed in the project was the conditional probabilities of delays. More specifically, we looked at the probability distribution of the delay of next train conditional on the previous train being late by n minutes.

Because of small amount of data this distributions are fairly noisy, and it is difficult to visualise using a standard plot. To show how the distribution is changing with n , we have plotted the differences between distributions and the $n = 0$ distribution, see Figure 7. The colour, from blue to red, indicates increasing n . You can see that as n increases, the probability of having no delay for the subsequent train decreases, with a corresponding increase in the probability of having a positive delay. Here we are focusing on the probability of small delays, to understand longer delays, one could model the conditional distributions as q -exponential distributions.

We tried using singular value decomposition (SVD) for smoothing the distributions. Formally, the singular value decomposition of an $m \times n$ matrix M is a factorization of the form $M = U\Sigma V^*$, where U is an $m \times m$, Σ is an $m \times n$ rectangular diagonal matrix with non-negative real numbers on the diagonal, and V^* (the conjugate transpose of V , or simply the transpose of V if V is real) is an $n \times n$ real or complex unitary matrix. We get rid of noise which is high rank, by just leaving meaningful information. The method of validation have been applied for minimization of negative log likelihood, by using 80% of data for training, and

the rest for testing. The plot in Figure 8 demonstrates that the rank of 7 minimises negative log likelihood the most, thus this value is the best for smoothing data. We also tried to smooth data by splitting it to different times of the day (e.g. morning, afternoon, evening), but it did not improve the result that we had.

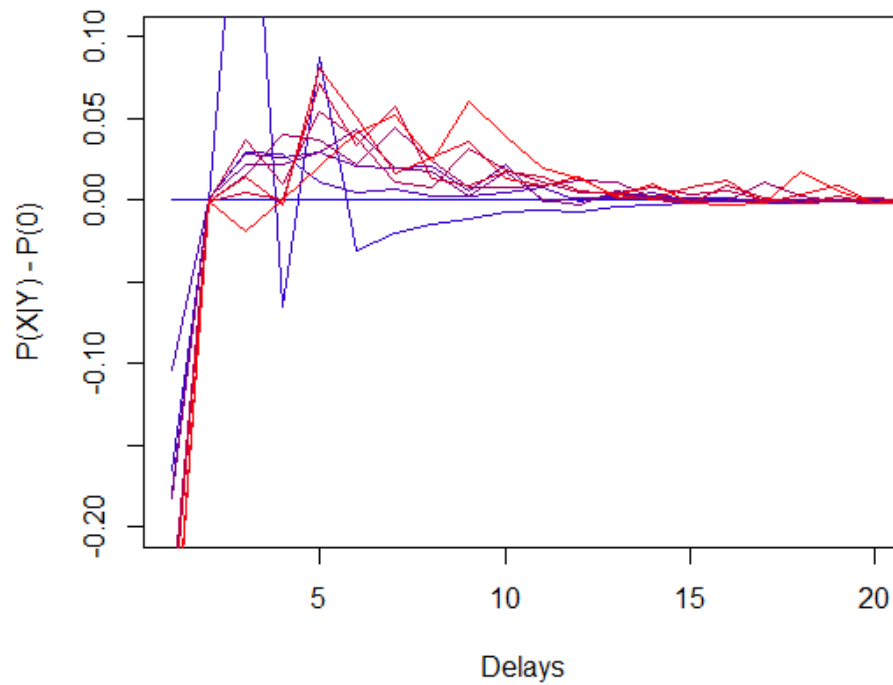


Figure 7: Conditional probabilities of delays for consecutive trains from IPS to COL. Subtracting the zero delays distributions normalizes the conditional probabilities by allowing us to compare different conditional distributions.

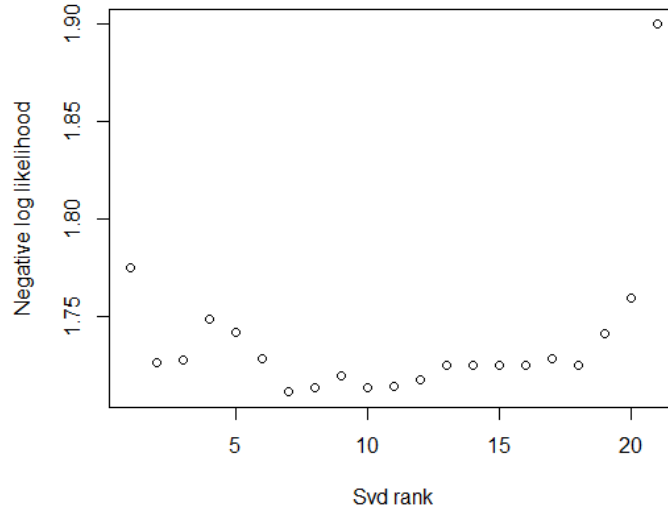


Figure 8: Negative log likelihoods after smoothing data using different ranks in SVD decomposition

4 Conclusions

We have presented analysis of train delays distribution in the UK rails system. First of all, although we were able to identify the majority of the trains for a particular journey by providing basic properties of the journey, there were misidentified trains. Because of that, the data analysis can not be very accurate. Possible solution to this could be providing more details about each journey as an input, which would leave implications on the computations. Alternatively, the possibility of collecting more information about each departure (such as arrival times and stop stations) should be considered. Concerning the other statistics of the data set, we were able to find high correlation (0.76) between delays of consecutive stations during one journey despite the presence of misidentified trains. In other words, if a train is delayed, then with high probability it will be still delayed later in its journey. Also, it was found that distribution of delays is different for the different days of week, which means that it may be reasonable considering the day of the week in journey planning. These results may allow more accurate simulation of passenger movement through a rail system with delays, which was one of the main purposes for this project.

Acknowledgement

My gratitude goes to my supervisors Dr. Keith Briggs from BT and Dr. Ben Graham from Warwick University for their guidance and help provided during the project.

Reference

1. Keith Briggs and Christian Beck. Modelling train delays with q -exponential functions (2007)
2. Keith Briggs and Peter Kin Po Tam. Optimal trip planning in timetabled transport systems possessing random delays(2011)
3. Nonlinear-Least-Squares Analysis of Slow-Motion EPR Spectra in One and Two Dimensions Using a Modified LevenbergMarquardt Algorithm
4. Shmuel Friedland University of Illinois at Chicago. The Role of Singular Value Decomposition in Data Analysis
5. Kirk Baker March 29, 2005 (Revised January 14,2013).Singular Value Decomposition Tutorial
6. J.-F. Bercher and C. Vignat.A new look at q -exponential distributions via excess statistics