# EY NextWave

# Data Science Competition

**Marcus Gawronsky**

University of Cape Town
South Africa
May 20, 2019

# Introduction

This report forms part of an entry in the EY NextWave Data Science Competition 2019. It outlines the analysis and applications of predictive models trained on a dataset of user GPS data. The report provides details on how to replicate the analysis performed, as well as a criteria for evaluating the chosen model.
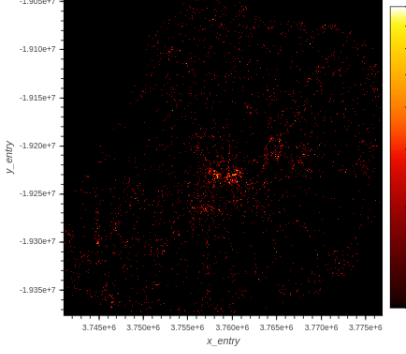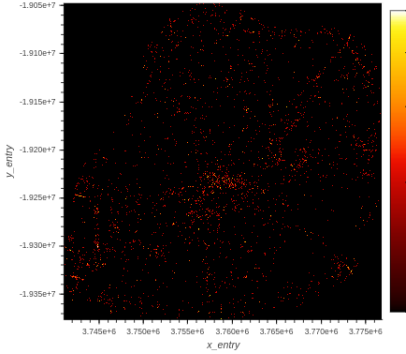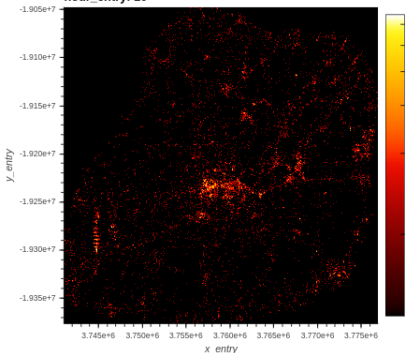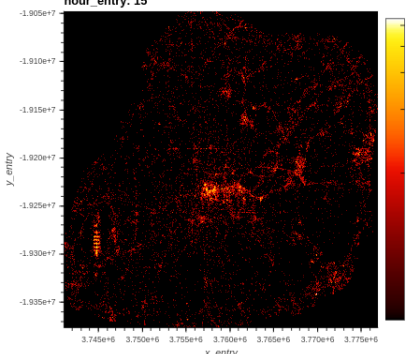
# Reproducibility

The analysis in the report was performed using Python 3.7.3 on Ubuntu 19.04. A full list of package dependencies along with the code used to run the analysis is provided in the appendix and can be installed and executed using Anaconda 4.6.11 or using the latest Miniconda Docker Container.

# Data

The datasets provided for the competition came split into a training and test-set with 814262 and 202937 entries respectively. The table below provides a sample of the data containing datapoints trips by individuals. Each individual is given a *hash*, which is reset each day, and each trip is given a unique *trajectory_id*.

| hash | trajectory_id | time_entry | time_exit | vmax | vmin | vmean |
|------|---------------|------------|-----------|------|------|-------|
| 005... | traj_005a... | 10:20:00 | 10:28:17 | 0.0 | 0.0 | 0.0 |
| 005... | traj_005a... | 10:40:51 | 10:52:47 | NaN | NaN | NaN |
| 005... | traj_005a... | 11:22:09 | 11:45:23 | 0.0 | 0.0 | 0.0 |

| x_entry | y_entry | x_exit | y_exit |
|---------|---------|--------|--------|
| 3.758e+06 | -1.904e+07 | 3.75e+06 | -1.905e+07 |
| 3.759e+06 | -1.911e+07 | 3.75e+06 | -1.913e+07 |
| 3.767e+06 | -1.920e+07 | 3.77e+06 | -1.912e+07 |

In the table below, one can see imaged which show the distribution of individuals, based on their GPS locations at different points accross the day. One can see that as the day progresses, we see a greater density of individuals in the city centre.

| Time | Image |
|---|---|
| 00:00 |  |
| 05:00 |  |
| 10:00 |  |
| 15:00 |  |

Given this data the aim of the challenge was two-fold:

1. To predict the final location of an individual at 15:00

2. To determine a decision-boundary for what constituted the city-centre in order to determine if individuals were in the city centre or not

# Applications

There exists a large miriad of application for this data and the application of predictive modelling, as shown in the table below. While these applications have vastly different requirements in terms of accuracy, explainability and security, they can be broadly split into simulation-based application and non-simulation-based applications. The primary distinction between these two applications is that in simulation-based applications we require a global model, which given a starting point on a map can continuously predict the next location at which the person is observed. Non-simulation-based models may have different requirements and allow for models which only apply to an individual as so do not generalize to arbitary groups of users.

| Simulation | Non-simulation |
|---|---|
| Disaster Relief | **Fraud detection (Credible credit card use)** |
| Infectious Disease Modelling | Policing (Testimony, finding suspects) |
| Urban Planning | Marketing/ Product Placement |

While the simulation distinction is import, other applications may have different requirements such as explainability or a requirement in quantifying uncertainty.

Applications such as Fraud detection would use GPS coordinates of an individual given through a cellphone or records of where last a transaction was made in order to try determine if it was likely the person made the next transaction. In this application, location is not the only priority, security and the quantification of risk or uncertainty is vitally important is we may set a threshold for the the probability of a person being at an store in order to determine if the card has been duplicated. Similarly we may prefer models which have the ability to be better secured along with the users data and which may be stored and used on edge devices, such as credit card machine or EPOS devices.

# Priorities

In order to best determine our ideal application and methodology, we must first determine a framework by which to consider and evaluate these use-cases. These can be broken up into four main subsections: Sustainability, Security, Resources and Predictive Requirements.
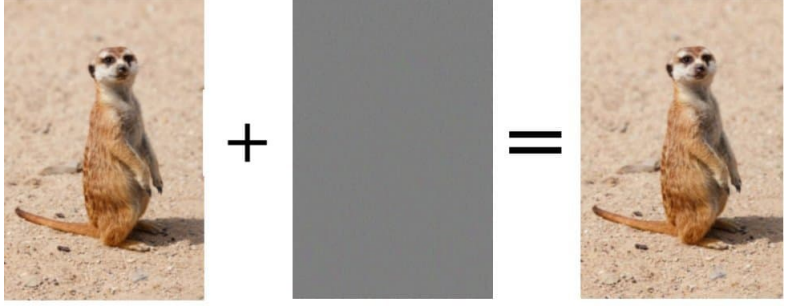
## Sustainability

A major consideration in decising upon our applications is around sustainability. In order to identify if our application is sustainable we need to determine if there a market or need for such as product which can cover even server costs and development. We also need to determine if certain applications have requirements for greater or continuous data aquisition and whether the market allows for that.

Arguable, all the application present some oppotunity for sustainability and have different markets in which they may compete. However, while simulation-based methods require little in terms of further data aquisition, our non-simulation-based methods may require continuous data aquisition of where individuals are or have been in order to provide utility.

While this is not a limiting factor applications such as Marketing/ Product Placement may lack the sufficient incentives on the part of users to provide such data in order to optimally market to them. In the case of Google Maps, the incentives do work by providing users with a facility by which to navigate and valuable recommendation in exchange for marketing data.

## Security

Security is a major consideration for Data Science and IT applications and is a concern which has been growing by legislators and individuals. In the table below a list of security concerns and examples or explanations is provided:

| Type | Example |
|---|---|
| Fooled by Data | <br><br>meerkat, mierkat (score = 0.90021)<br>mongoose (score = 0.02666)<br>Windsor tie (score = 0.00072)<br>otter (score = 0.00069)<br>doormat, welcome mat (score = 0.00055)<br><br>kite (score = 0.07896)<br>bald eagle, American eagle, Haliaeetus leucocephalus (score = 0.04153)<br>bee eater (score = 0.03940)<br>parachute, chute (score = 0.02724)<br>hummingbird (score = 0.02334)<br><br>doormat, welcome mat (score = 1.00000)<br>prayer rug, prayer mat (score = 0.00000)<br>manhole cover (score = 0.00000)<br>miniature poodle (score = 0.00000)<br>palace (score = 0.00000) |
| Adverserial Model Drift | Agents feed false data |
| Model Leaks | 1. Access to accurate global model (GPS spoof me to learn about you) 2. Parameterizes personal data (Leaks means I know where you live) |
| Data Leaks | Centralized data is accessed |

A major trade-off in our application is the need for a centralized vs a distributed model. A centralized model is one in which all individuals are captured by a single model which must be centrally trained and queried by all users. While such models can scale, they necessitate data being centralized, which creates a risk or security leaks, and provide a tools which malicous agents can use in order to perform queries about the movements of other individuals. For example, if I have a model which only works for me and which I only have access to then if I put in your location it will not tell me where you will likely be next, however if I have a model which can make predictions about all individuals, I can spoof the model with your coordinates and determine where you will be. Given the concerns over our the privacy of our data and application this renders simulation-based applications less desireable, as they rely on a centralized model.

Similarly models such as random forests, which paramterize data through decision-rules (such as *if the person is at this location then they will be at this location* given us information about a persons location and) may be less preferable as they raise the security concerns about malacious agents gaining access to either our model or our data, not just the data. While this concern may be handles in some may through data obfuscation, this still remain a criterea for evaluation.

Lastly, as with many Deep Learning, concerns must be raised in certain applications about our ability to have assurances about our model, as with many cases our inability to properly reason about the model means, as with images, they adding noise to images can lead to castly inaccurate results, which may allow persons to trick models or feed it data which renders it less accurate over time without us knowing or being easily able to identify.

## Resources

Across many applications an unavoidable limitation is resource constraints. Algorithms which rely on matrix factorization do no scale to large datasets easily, as they require all the data being in memory and are thus heabily io-bound. Similarly extremely Deep Neural Networks can be gigabytes in size and can require specialized and expensive hardware in order to use or may require communication with a server, limiting their application.

| Scenario | Method | Cheap | Scallable |
|---|---|---|---|
| Embedded System | Least-square Regression | Yes | No |
| Large cluster | Alex-net (Deep CNN's) | No | Yes |

## Predictive Requirements

Different models and applications have different predictive requirements. In the table below, we detail three main requirements: Correctness, Confidence and Explainability. Correctness, here, is the accuracy of a given model which can mean common loss functions like Mean Squared-Error or Binary Cross-entropy, which may be evaluated using some form of validation, such ass using a Test-Train-Split or using Stratified K-fold Cross-validation.

Confidence asks whether our prediction gives some measure of uncertaintly of its estimate. For many classification tasts, using Binary cross-entropy a pseudo-probability is specified using a link function and using our loss function a probaility distribution is assumed, meaning we have some assumptions of the distribution of our uncertainty and some quantization in the form of a pseudo-probability. In many bayesian methods the uncertainty of our predicted output can be extensive, in which we are given a probability distribution of an output and can use a range of methods for point estimation based on the risk tolerance of the application.
Lastly, in explainability, some models allow for various levels of explainability in determining what variables in our data and what parameters led to such a prediction. Here linear models and simple decision tree provide high levels of explainabilitu and

assurance, many deep learning methods provide limited methods by which to explain a given prediction.

| Example | Correctness | Confidence | Explainability |
|---|---|---|---|
| Siri | High | Low | Low |
| Stock Prices | High | High | Low |
| Macroeconomic Models | Low | High | High |
| Medical Imaging | High | High | High |

Looking at the table above, applications such as Siri demand the model used in that application to be correct or users become frustrated, but the certainty of a particular prediction or the ability to explain why a particular response was given is not particularly important. Alternatively in stock prediction, quantifying uncertainty is crucial in deciding optimal asset allocation.

In our applications discussed, applications such as Fraud Detection and Disaster Relief require high levels of correctness and confidence, as we need to know crucially the confidence of a prediction in order to make an optimal decision.

## Methodology

From the framework above, it is clear there are obvious downsides to the large centralized models required for simulation-based application, both in terms of the resources requirements and the security risks.

Given these constraints, we must then tackle the issue of sustainability. Looking at Fraud Detection, much of the infrastructure is currently in place for businesses with a wide distribution of customer app installs, EPOS devices and credit card machines for data collection. A clear metric exists to evaluate the value to the business and models can easily be used along with other datasets and models already present in the business.

With this application it is clear that correctness, confidence and security are the primary concerns. Using Bayesian Modelling we can specify a seperate model unique to each person, which provides us a probability distribution for their location at a point in time. Here, we can make use of some basis explansion to create polynomial and piecewise functions of time and use Bayesian Ridge Regression, with a differing prior, in order to find an optimal basis in determining the location of a user at a point in time.

For this analysis, rather than looking at the data as a series of entries and exits, we just look at the locations of individuals at points in time, using the following model:

$$y_{p,i} = \beta_{0,p} + \beta_{1,p}t_{p,i} + \beta_{2,p}t_{p,i}^2 + \beta_{3,p}t_{p,i}^3 + \beta_{4,p}B(||t_{p,i} - 10\text{:}00\text{am}||_+)$$

$$x_{p,i} = \beta_{0,p} + \beta_{1,p}t_{p,i} + \beta_{2,p}t_{p,i}^2 + \beta_{3,p}t_{p,i}^3 + \beta_{4,p}B(||t_{p,i} - 10\text{:}00\text{am}||_+)$$

$$B(t) = x^2 log(x)$$

$p$ is a person with parameter $(\beta)$ and times $(t)$
$i$ is a point in time $(t)$
$y_{p,i}$, $y_{p,i}$ are the longitude and latitude of an individual
$\beta_{k,p}$ are the parameters of each model used for each person

In order to fit our model, we used the specification foundin Tipping (2001) and MacKay (1992). This defines the distribution over our weights of:

$$Pr(\beta|\lambda) = N(\beta|0, I\lambda^{-1})$$

The priors for $\lambda$ are chosen to be gamma distributions, using the conjugate prior for the precision of the Gaussian Distribution. Non-informative priors were chosen for our hyperparamters of $10^{-6}$.

This model was tested over a range of Basis Functions, an example of a predicted path is shown in red below against the path of the user:

**Means**

A final decision boundary for the city was chosen to be the 75th and 25th quantiles over our GPS coordinates throughout the day over the longitude and lattitude. This gave an f1 score of 0.72.

# Conclusions

While model does not provide state of the art accuracy compared to other approached, this approach offers a miriad of advantages in security, scalability, explainability, compute and confidence. The abilty to analyze a moving probability distribution of the possible location of an individual over time as they move over a map, given their past movements provides a range of advantages, aiding decision-making across a range of applications.

# Appendix

## Anaconda Python Dependencies

```
 1  name: EY
 2  channels:
 3    - pytorch
 4    - conda-forge
 5    - intel
 6    - defaults
 7  dependencies:
 8    - _tflow_select=2.3.0=mkl
 9    - absl-py=0.7.1=py37_0
10    - arrow-cpp=0.13.0=py37hfe0260e_2
11    - asn1crypto=0.24.0=py37_0
12    - astor=0.7.1=py37_0
13    - attrs=19.1.0=py37_1
14    - backcall=0.1.0=py37_0
15    - blas=1.0=mkl
16    - bleach=3.1.0=py37_0
17    - bokeh=1.0.4=py37_0
18    - boost-cpp=1.68.0=h11c811c_1000
19    - brotli=1.0.7=hf484d3e_1000
20    - bzip2=1.0.6=h14c3975_5
21    - c-ares=1.15.0=h7b6447c_1
22    - ca-certificates=2019.3.9=hecc5488_0
23    - cairo=1.14.12=h8948797_3
24    - cartopy=0.17.0=py37hbb7e04d_1
25    - certifi=2019.3.9=py37_0
26    - cffi=1.12.2=py37h2e261b9_1
27    - cftime=1.0.3.4=py37hdd07704_0
28    - chardet=3.0.4=py37_1
29    - click=7.0=py37_0
30    - click-plugins=1.0.4=py37_0
31    - cligj=0.5.0=py37_0
32    - cloudpickle=0.8.1=py_0
33    - colorcet=1.0.0=py37_0
34    - cryptography=2.6.1=py37h1ba5d50_0
35    - cudatoolkit=10.0.130=0
36    - cudnn=7.3.1=cuda10.0_0
37    - curl=7.64.0=hbc83047_2
38    - cycler=0.10.0=py37_0
39    - cytoolz=0.9.0.1=py37h14c3975_1
40    - dask-core=1.2.0=py_0
41    - datashader=0.6.9=py_0
42    - datashape=0.5.4=py37_1
43    - dbus=1.13.6=h746ee38_0
44    - decorator=4.4.0=py37_1
45    - defusedxml=0.5.0=py37_1
```

```
46    − descartes=1.1.0=py37_0
47    − entrypoints=0.3=py37_0
48    − expat=2.2.6=he6710b0_0
49    − fiona=1.8.4=py37hc38cc03_0
50    − fontconfig=2.13.0=h9420a91_0
51    − freetype=2.9.1=h8a8886c_1
52    − freexl=1.0.5=h14c3975_0
53    − gast=0.2.2=py37_0
54    − gdal=2.3.3=py37hbb2a789_0
55    − geopandas=0.4.1=py_0
56    − geos=3.7.1=he6710b0_0
57    − geoviews=1.6.2=py_0
58    − geoviews−core=1.6.2=py_0
59    − gflags=2.2.2=hf484d3e_1001
60    − giflib=5.1.4=h14c3975_1
61    − glib=2.56.2=hd408876_0
62    − glog=0.3.5=hf484d3e_1001
63    − gmp=6.1.2=h6c8ec71_1
64    − grpcio=1.16.1=py37hf8bcb03_1
65    − gst−plugins−base=1.14.0=hbbd80ab_1
66    − gstreamer=1.14.0=hb453b48_1
67    − h5py=2.9.0=py37h7918eee_0
68    − hdf4=4.2.13=h3ca952b_2
69    − hdf5=1.10.4=hb1b8bf9_0
70    − heapdict=1.0.0=py37_2
71    − holoviews=1.11.3=py_0
72    − hvplot=0.4.0=py_0
73    − icu=58.2=h9c2bf20_1
74    − idna=2.8=py37_0
75    − imageio=2.5.0=py37_0
76    − intel−openmp=2019.3=199
77    − ipykernel=5.1.0=py37h39e3cac_0
78    − ipython=7.1.1=py37h39e3cac_0
79    − ipython_genutils=0.2.0=py37_0
80    − ipywidgets=7.4.2=py37_0
81    − jedi=0.13.3=py37_0
82    − jinja2=2.10.1=py37_0
83    − jpeg=9c=h14c3975_1001
84    − json−c=0.13.1=h1bed415_0
85    − jsonschema=3.0.1=py37_0
86    − jupyter=1.0.0=py37_7
87    − jupyter_client=5.2.4=py37_0
88    − jupyter_console=6.0.0=py37_0
89    − jupyter_contrib_core=0.3.3=py_2
90    − jupyter_contrib_nbextensions=0.5.1=py37_0
91    − jupyter_core=4.4.0=py37_0
92    − jupyter_highlight_selected_word=0.2.0=py37_1000
93    − jupyter_latex_envs=1.4.4=py37_1000
94    − jupyter_nbextensions_configurator=0.4.1=py37_0
```

```
95    - jupyterlab=0.35.4=py37hf63ae98_0
96    - jupyterlab_server=0.2.0=py37_0
97    - kealib=1.4.7=hd0c454d_6
98    - keras=2.2.4=0
99    - keras-applications=1.0.7=py_0
100   - keras-base=2.2.4=py37_0
101   - keras-preprocessing=1.0.9=py_0
102   - kiwisolver=1.0.1=py37hf484d3e_0
103   - krb5=1.16.1=h173b8e3_7
104   - libboost=1.67.0=h46d08c1_4
105   - libcurl=7.64.0=h20c2e04_2
106   - libdap4=3.19.1=h6ec2957_0
107   - libedit=3.1.20181209=hc058e9b_0
108   - libffi=3.2.1=hd88cf55_4
109   - libgcc-ng=8.2.0=hdf63c60_1
110   - libgdal=2.3.3=h2e7e64b_0
111   - libgfortran-ng=7.3.0=hdf63c60_0
112   - libkml=1.3.0=h590aaf7_4
113   - libnetcdf=4.6.1=h11d0813_2
114   - libpng=1.6.37=hbc83047_0
115   - libpq=11.2=h20c2e04_0
116   - libprotobuf=3.7.1=h8b12597_0
117   - libsodium=1.0.16=h1bed415_0
118   - libspatialindex=1.8.5=h20b78c2_2
119   - libspatialite=4.3.0a=hb08deb6_19
120   - libssh2=1.8.0=h1ba5d50_4
121   - libstdcxx-ng=8.2.0=hdf63c60_1
122   - libtiff=4.0.10=h9022e91_1002
123   - libuuid=1.0.3=h1bed415_2
124   - libxcb=1.13=h1bed415_1
125   - libxml2=2.9.9=he19cac6_0
126   - libxslt=1.1.33=h7d1a2b0_0
127   - llvmlite=0.28.0=py37hd408876_0
128   - locket=0.2.0=py37_1
129   - lxml=4.3.2=py37hefd8a0e_0
130   - lz4-c=1.8.3=he1b5a44_1001
131   - mapclassify=2.0.1=py_0
132   - markdown=3.0.1=py37_0
133   - markupsafe=1.1.1=py37h7b6447c_0
134   - matplotlib=3.0.3=py37h5429711_0
135   - mistune=0.8.4=py37h7b6447c_0
136   - mkl=2019.3=199
137   - mkl_fft=1.0.12=py37ha843d7b_0
138   - mkl_random=1.0.2=py37hd81dba3_0
139   - mock=2.0.0=py37_0
140   - msgpack-python=0.6.1=py37hfd86e86_1
141   - multipledispatch=0.6.0=py37_0
142   - munch=2.3.2=py37_0
143   - nbconvert=5.4.1=py37_3
```

```
144    − nbformat =4.4.0=py37_0
145    − ncurses =6.1=he6710b0_1
146    − netcdf4 =1.4.2=py37h808af73_0
147    − networkx =2.2=py37_1
148    − ninja =1.9.0=py37hfd86e86_0
149    − nodejs =10.13.0=he6710b0_0
150    − notebook =5.7.8=py37_0
151    − numba =0.43.1=py37h962f231_0
152    − numpy =1.16.3=py37h7e9f1db_0
153    − numpy−base =1.16.3=py37hde5b4d6_0
154    − olefile =0.46=py37_0
155    − openjpeg =2.3.0=h05c96fa_1
156    − openssl =1.1.1b=h14c3975_1
157    − owslib =0.17.1=py_0
158    − packaging =19.0=py37_0
159    − pandas =0.24.2=py37he6710b0_0
160    − pandoc =2.2.3.2=0
161    − pandocfilters =1.4.2=py37_1
162    − panel =0.4.0=0
163    − param =1.8.2=py_0
164    − parquet−cpp =1.5.1=2
165    − parso =0.3.4=py37_0
166    − partd =0.3.10=py37_1
167    − pbr =5.1.3=py_0
168    − pcre =8.43=he6710b0_0
169    − pexpect =4.6.0=py37_0
170    − pickleshare =0.7.5=py37_0
171    − pillow =6.0.0=py37h34e0f95_0
172    − pip =19.1=py37_0
173    − pixman =0.38.0=h7b6447c_0
174    − poppler =0.65.0=h581218d_1
175    − poppler−data =0.4.9=0
176    − proj4 =5.2.0=he6710b0_1
177    − prometheus_client =0.6.0=py37_0
178    − prompt_toolkit =2.0.9=py37_0
179    − protobuf =3.7.1=py37he6710b0_0
180    − psutil =5.6.2=py37h7b6447c_0
181    − psycopg2 =2.7.6.1=py37h1ba5d50_0
182    − ptyprocess =0.6.0=py37_0
183    − pyarrow =0.13.0=py37h0978efd_0
184    − pycparser =2.19=py37_0
185    − pyct =0.4.6=py37_0
186    − pyepsg =0.4.0=py37_0
187    − pygments =2.3.1=py37_0
188    − pykdtree =1.3.1=py37hdd07704_2
189    − pyopenssl =19.0.0=py37_0
190    − pyparsing =2.4.0=py_0
191    − pyproj =1.9.6=py37h14380d9_0
192    − pyqt =5.9.2=py37h05f1152_2
```

```
193    −  pyrsistent=0.14.11=py37h7b6447c_0
194    −  pyshp=2.0.1=py37_0
195    −  pysocks=1.6.8=py37_0
196    −  python=3.7.3=h0371630_0
197    −  python−dateutil=2.8.0=py37_0
198    −  pytorch−cpu=1.1.0=py3.7_cpu_0
199    −  pytz=2019.1=py_0
200    −  pyviz_comms=0.7.0=py37_0
201    −  pywavelets=1.0.2=py37hdd07704_0
202    −  pyyaml=5.1=py37h7b6447c_0
203    −  pyzmq=18.0.0=py37he6710b0_0
204    −  qt=5.9.7=h5867ecd_1
205    −  qtconsole=4.4.3=py37_0
206    −  re2=2019.04.01=he1b5a44_0
207    −  readline=7.0=h7b6447c_5
208    −  requests=2.21.0=py37_0
209    −  rise=5.5.0=py37_0
210    −  rtree=0.8.3=py37_0
211    −  scikit−image=0.14.2=py37he6710b0_0
212    −  scikit−learn=0.20.3=py37hd81dba3_0
213    −  scipy=1.2.1=py37h7c811a0_0
214    −  send2trash=1.5.0=py37_0
215    −  setuptools=41.0.1=py37_0
216    −  shapely=1.6.4=py37h86c5351_0
217    −  sip=4.19.8=py37hf484d3e_0
218    −  six=1.12.0=py37_0
219    −  snappy=1.1.7=hf484d3e_1002
220    −  sortedcontainers=2.1.0=py37_0
221    −  sqlalchemy=1.3.1=py37h7b6447c_0
222    −  sqlite=3.28.0=h7b6447c_0
223    −  tblib=1.3.2=py37_0
224    −  tensorboard=1.13.1=py37hf484d3e_0
225    −  tensorflow=1.13.1=mkl_py37h54b294f_0
226    −  tensorflow−base=1.13.1=mkl_py37h7ce6ba3_0
227    −  tensorflow−estimator=1.13.0=py_0
228    −  termcolor=1.1.0=py37_1
229    −  terminado=0.8.1=py37_1
230    −  testpath=0.3.1=py37_0
231    −  thrift−cpp=0.12.0=h0a07b25_1002
232    −  tk=8.6.8=hbc83047_0
233    −  toolz=0.9.0=py37_0
234    −  tornado=6.0.2=py37h7b6447c_0
235    −  traitlets=4.3.2=py37_0
236    −  urllib3=1.24.1=py37_0
237    −  wcwidth=0.1.7=py37_0
238    −  webencodings=0.5.1=py37_1
239    −  werkzeug=0.15.2=py_0
240    −  wheel=0.33.1=py37_0
241    −  widgetsnbextension=3.4.2=py37_0
```

```
242    - xarray=0.11.3=py37_0
243    - xerces-c=3.2.2=h780794e_0
244    - xz=5.2.4=h14c3975_4
245    - yaml=0.1.7=had09818_2
246    - zeromq=4.3.1=he6710b0_3
247    - zict=0.1.4=py37_0
248    - zlib=1.2.11=h7b6447c_3
249    - zstd=1.3.3=1
250    - pip:
251      - atomicwrites==1.3.0
252      - colorama==0.4.1
253      - dask==1.1.0
254      - distributed==1.25.0
255      - filelock==3.0.10
256      - flatbuffers==1.11
257      - funcsigs==1.0.2
258      - jupyterlab-latex==0.4.1
259      - more-itertools==7.0.0
260      - pluggy==0.11.0
261      - py==1.8.0
262      - pytest==4.4.2
263      - redis==3.2.1
264      - setproctitle==1.1.10
265      - typing==3.6.6
266  prefix: /home/marcusskky/.conda/envs/EY
```

## Model

```python
1  import os
2  import abc
3  import datetime
4  from toolz.curried import *
5  import pandas as pd
6  import holoviews as hv
7  from sklearn import linear_model
8  from sklearn.preprocessing import PolynomialFeatures, \
9                                    StandardScaler, \
10                                   FunctionTransformer
11 from sklearn.impute import SimpleImputer
12 from sklearn.pipeline import make_pipeline, make_union
13
14
15 def Imputer(model_: abc.ABCMeta,
16             data_: pd.DataFrame,
17             X_: list, y_: list)-> pd.DataFrame:
18
19     def predict(model___: abc.ABCMeta,
20                 data___: pd.DataFrame,
21                 X__: list, y__: list) -> pd.np.array:
```

```python
22            not_nul = data_.loc [:,X_+ y__].notnull().all(1)
23
24            model_instance = model__.fit(X = data__.loc[not_nul,X__],
25                                          y = data__.loc[not_nul,y__].values
      .ravel())
26
27            return model_instance.predict(X=data_.loc[:,X__])
28
29
30       return data_.assign(**{ys+'_predicted': predict(model__ = model_,
31                                                         data__ = data_,
32                                                         X__ = X_,
33                                                         y__ = list(ys))
34                              for ys in y_})
35
36 def basis(x:pd.np.array, knot:float = 12*60*60)-> pd.np.array:
37       x = x − knot
38       x = x * (x>0.0)
39       x = x**2 * pd.np.log(x, where= x > 0.0)
40
41       return x.astype(pd.np.float64)
42
43
44 if __name__ == '__main__':
45       path = os.path.join('Data','data_test.csv')
46       data = pd.read_csv(path, index_col = 'Unnamed: 0')
47
48       id_ = ['hash','trajectory_id','vmax','vmin','vmean']
49       midnight = pd.to_datetime('00:00:00')
50
51       views = data.pipe(partial(pd.wide_to_long,
52                                 stubnames = ['x','y','time'],
53                                 i = id_,
54                                 j = 'pop',
55                                 sep = '_',
56                                 suffix = '(!?entry|exit)'))\
57                  .reset_index()\
58                  .replace(['entry','exit'],[0,1])\
59                  .assign(time = lambda d: d.time\
60                                            .pipe(pd.to_datetime)\
61                                            .subtract(midnight)\
62                                            .dt.total_seconds(),
63                          missing = lambda d: d.x.isna().astype(pd.np.int
      ))
64
65        pipeline = make_pipeline(make_union(PolynomialFeatures(degree=3,
      include_bias=False),
66                                            FunctionTransformer(basis,
67                                                                validate=False)
```

16

```python
      ),
68                            linear_model.BayesianRidge(fit_intercept=True)
      )

69
70    imputed = views.groupby(['hash'])\
71                    .apply(lambda g: Imputer(model_=pipeline, data_=g,
      X_=['time'], y_=['y','x']))\
72                    .reset_index(drop=True)

73
74    imputed_path = pipe(datetime.datetime.now().timestamp(),
75                        partial(round, ndigits=0),
76                        int,
77                        lambda x: os.path.join('Data',f'imputed_{x}.
      csv'))
78    imputed.to_csv(imputed_path, index=False, sep=',')

79
80    quantile = 0.25
81    output_path = pipe(datetime.datetime.now().timestamp(),
82                        partial(round, ndigits=0),
83                        int,
84                        lambda x: os.path.join('Data',f'submission_{
      x}.csv'))

85
86    imputed.assign(city = lambda d: ((d.x_predicted <d.x.quantile(0.5+
      quantile)) &\
87                                     (d.x_predicted >d.x.quantile(0.5-
      quantile)) &\
88                                     (d.y_predicted <d.y.quantile(0.5+
      quantile)) &\
89                                     (d.y_predicted >d.y.quantile(0.5-
      quantile))).astype(pd.np.int))\
90            .where(lambda x: x.missing==1)\
91            .dropna(axis=0, how='all')\
92            .loc[:,['trajectory_id','city']]\
93            .rename(columns={'trajectory_id':'id','city':'target'})\
94            .assign(target = lambda x: x.target.astype(int))\
95            .to_csv(output_path, index=False, sep=',')
```

# Bibliography

Section 3.3 in Christopher M. Bishop: Pattern Recognition and Machine Learning, 2006

David J. C. MacKay, Bayesian Interpolation, 1992.

Michael E. Tipping, Sparse Bayesian Learning and the Relevance Vector Machine,

2001.