# Data Modeling Project for MATH 7241

October 2020

Due date : Tuesday December 1, week after Thanksgiving.

Groups: if you wish you may collaborate in a group of at most
two people. In this case each member of the group must make
a substantial contribution to the project, and your project report
must describe which contributions were made by each member.

# Contents

Step 1: find a good data set!

Recommended source: UCI archive. Search for Time Series.

https://archive.ics.uci.edu/ml/index.php

# UCI

## Machine Learning Repository
### Center for Machine Learning and Intelligent Systems

Search
Repository  Web
Google™

**View ALL Data Sets**

**Welcome to the UC Irvine Machine Learning Repository!**

We currently maintain 394 data sets as a service to the machine learning community. You may **view all data sets** through our searchable interface. Our old web site is still available, for those who prefer the old format. For a general overview of the Repository, please visit our About page. For information about citing data sets in publications, please read our citation policy. If you wish to donate a data set, please consult our donation policy. For any other questions, feel free to contact the Repository librarians. We have also set up a mirror site for the Repository.

Supported By:    In Collaboration With:  Rexa.info

| Latest News: | Newest Data Sets: | Most Popular Data Sets (hits since 2007): |
|---|---|---|

**Latest News:**

**04-04-2013:** Welcome to the new Repository admins Kevin Bache and Moshe Lichman!
**03-01-2010:** Note from donor regarding Netflix data
**10-16-2009:** Two new data sets have been added.
**09-14-2009:** Several data sets have been added.
**07-23-2008:** Repository mirror has been set up.
**03-24-2008:** New data sets have been added!
**06-25-2007:** Two new data sets have been added: UJI Pen Characters, MAGIC Gamma Telescope

**Featured Data Set: CMU Face Images**

Task: Classification
Data Type: Image
# Instances: 640

This data consists of 640 black and white face images of people taken with varying pose (straight, left, right, up), expression (neutral, happy, sad, angry), eyes (wearing sunglasses or not), and size

**Newest Data Sets:**

**08-28-2017:** Burst Header Packet (BHP) flooding attack on Optical Burst Switching (OBS) Network
**07-23-2017:** Eco-hotel
**07-23-2017:** Las Vegas Strip
**07-20-2017:** Parkinson Disease Spiral Drawings Using Digitized Graphics Tablet
**07-18-2017:** PM2.5 Data of Five Chinese Cities
**07-16-2017:** Sales_Transactions_Dataset_Weekly
**06-29-2017:** Data for Software Engineering Teamwork Assessment in Education Setting
**05-30-2017:** chestnut – LARVIC
**05-24-2017:** Epileptic Seizure Recognition

**Most Popular Data Sets (hits since 2007):**

1503113: Iris
989310: Adult
741287: Wine
645288: Car Evaluation
579500: Breast Cancer Wisconsin (Diagnostic)
572635: Forest Fires
542310: Human Activity Recognition Using Smartphones
523473: Heart Disease
516585: Wine Quality

Search for a time series:

Select one time series:

UCI

Machine Learning Repository
Center for Machine Learning and Intelligent Systems

Search

● Repository ○ Web

View ALL Data Sets

## Air Quality Data Set
*Download:* Data Folder, Data Set Description

Abstract: Contains the responses of a gas multisensor device deployed on the field in an Italian city. Hourly responses averages are recorded along with gas concentrations references from a certified analyzer.

| Data Set Characteristics: | Multivariate, Time-Series | Number of Instances: | 9358 | Area: | Computer |
|---|---|---|---|---|---|
| Attribute Characteristics: | Real | Number of Attributes: | 15 | Date Donated | 2016-03-23 |
| Associated Tasks: | Regression | Missing Values? | Yes | Number of Web Hits: | 131942 |

**Source:**

Saverio De Vito (saverio.devito '@' enea.it), ENEA - National Agency for New Technologies, Energy and Sustainable Economic Development

**Data Set Information:**

The dataset contains 9358 instances of hourly averaged responses from an array of 5 metal oxide chemical sensors embedded in an Air Quality Chemical Multisensor Device. The device was located on the field in a significantly polluted area, at road level,within an Italian city. Data were recorded from March 2004 to February 2005 (one year)representing the longest freely available recordings of on field deployed air quality chemical sensor devices responses. Ground Truth hourly averaged concentrations for CO, Non Metanic Hydrocarbons, Benzene, Total Nitrogen Oxides (NOx) and Nitrogen Dioxide (NO2) and were provided by a co-located reference certified analyzer. Evidences of cross-sensitivities as well as both concept and sensor drifts are present as described in De Vito et al., Sens. And Act. B, Vol. 129,2,2008 (citation required) eventually affecting sensors concentration estimation capabilities. Missing values are tagged with -200 value.
This dataset can be used exclusively for research purposes. Commercial purposes are fully excluded.
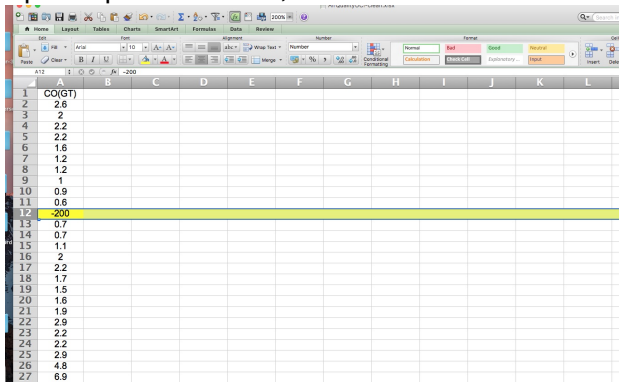
**Attribute Information:**

0 Date (DD/MM/YYYY)
1 Time (HH.MM.SS)
2 True hourly averaged concentration CO in mg/m^3 (reference analyzer)
3 PT08.S1 (tin oxide) hourly averaged sensor response (nominally CO targeted)
4 True hourly averaged overall Non Metanic HydroCarbons concentration in microg/m^3 (reference analyzer)
5 True hourly averaged Benzene concentration in microg/m^3 (reference analyzer)
6 PT08.S2 (titania) hourly averaged sensor response (nominally NMHC targeted)
7 True hourly averaged NOx concentration in ppb (reference analyzer)
8 PT08.S3 (tungsten oxide) hourly averaged sensor response (nominally NOx targeted)
9 True hourly averaged NO2 concentration in microg/m^3 (reference analyzer)
10 PT08.S4 (tungsten oxide) hourly averaged sensor response (nominally NO2 targeted)

Step 2: download and clean the data set to prepare for modeling. For example, remove unnecessary data, correct 'error' entries etc

Raw data set:

Clean it up: keep one column, remove error entries:

Step 3: map your data into a Markov chain. To do this, you must choose the states for your model. Each entry in the time series should map into a unique state.

In this example, we choose a 9-state model, with the states $\{1, 2, \ldots, 9\}$. Map each entry into a state by rounding to the nearest integer.

Step 4: transfer the time series to a platform for analysis, eg Excel, Matlab, R etc.

For example, here we import the spreadsheet into Matlab:

and here is the complete time series:

Step 5: compute the occupation frequencies for each state, and turn this into a probability distribution. This is the empirical distribution of your chain.

Empirical distribution from time series: fraction of time spent in each state

Step 6: compute the frequencies of jumps between each pair of states. Divide by the occupation frequency at each state so that the total jump probability out of each state is 1. This is your transition matrix.

Transition matrix from time series:

Step 7: find the stationary vector of your transition matrix. In case it is not unique, find all stationary vectors.

Stationary distribution of Markov chain:

```
ans =

    0.0522    0.3425    0.2879    0.1623    0.0869    0.0397    0.0180    0.0069    0.0034

fx >>
```

Step 8: compare the empirical distribution of the data set and the stationary distribution of your chain. Note any similarities!

Empirical distribution from time series compared to the stationary distribution of chain:

Step 10:

Build a simulation of the Markov chain, using the transition matrix that you computed in Step 6. Generate a typical time series using your simulation, and compare with the original time series. Does it look like a good model?

Compare simulation of the Markov chain with original time series:

Step 11: Compare your simulation with the original time series using the autocorrelation function. Given a time series $X_1, X_2, \ldots, X_N$, let $\overline{X}$ be the average, then for $k = 0, 1, 2, \ldots$ define

$$R(k) = \frac{\sum_{i=1}^{N-k} \left(X_i - \overline{X}\right) \left(X_{i+k} - \overline{X}\right)}{\sum_{i=1}^{N} \left(X_i - \overline{X}\right)^2}$$

Compute the autocorrelation $R(0), R(1), R(2), \ldots$ for the original time series and for the time series you generated in Step 10. Does it look like they describe the same series?

Step 12: Compare your simulation with the original time series using a goodness of fit test for the 2-step transition probabilities. Namely, let $\widehat{p}_{ij}$ be the transition matrix that you computed in Step 6 above. Do the same kind of calculation on the original time series to get the frequency of going from each state $i$ to each state $j$ in *two steps*. Call this frequency $N_{ij}$, and let $N_i = \sum_j N_{ij}$. Compare these frequencies with the 2-step frequencies computed using $\widehat{p}_{ij}$, that is

$$M_{ij} = N_i \, q_{ij} = N_i \sum_k \widehat{p}_{ik} \, \widehat{p}_{kj}$$

Use a goodness of fit test to compare these at the 5% significance level for each state $i$ (see notes on 'Goodness of Fit Test') and decide if the Markov chain $\{q_{ij}\}$ is a good model for the 2-step transitions of the data set.

Step 13: write a report (maximum 7 pages) explaining how you carried out the above steps, including: source and nature of raw data set, how it was cleaned, choice of states for Markov chain, empirical distribution, transition matrix, stationary distribution, comparison of empirical and stationary distributions, autocorrelation function, goodness of fit test. At the end of your report, answer this question: 'do you consider that the Markov chain method produces a good model for this time series? explain your answer'.

**Table A.3** Upper and Lower Percentiles of $\chi^2$ Distributions



$\chi^2$ distribution with $df$ degrees of freedom

Area $= 1 - p$

$\chi^2_{p, df}$

| df | $p$ | | | | | | | | |
|----|-------|-------|-------|-------|-------|-------|-------|-------|-------|
|    | 0.010 | 0.025 | 0.050 | 0.10 | 0.90 | 0.95 | 0.975 | 0.99 |
| 1 | 0.000157 | 0.000982 | 0.00393 | 0.0158 | 2.706 | 3.841 | 5.024 | 6.635 |
| 2 | 0.0201 | 0.0506 | 0.103 | 0.211 | 4.605 | 5.991 | 7.378 | 9.210 |
| 3 | 0.115 | 0.216 | 0.352 | 0.584 | 6.251 | 7.815 | 9.348 | 11.345 |
| 4 | 0.297 | 0.484 | 0.711 | 1.064 | 7.779 | 9.488 | 11.143 | 13.277 |
| 5 | 0.554 | 0.831 | 1.145 | 1.610 | 9.236 | 11.070 | 12.832 | 15.086 |
| 6 | 0.872 | 1.237 | 1.635 | 2.204 | 10.645 | 12.592 | 14.449 | 16.812 |
| 7 | 1.239 | 1.690 | 2.167 | 2.833 | 12.017 | 14.067 | 16.013 | 18.475 |
| 8 | 1.646 | 2.180 | 2.733 | 3.490 | 13.362 | 15.507 | 17.535 | 20.090 |
| 9 | 2.088 | 2.700 | 3.325 | 4.168 | 14.684 | 16.919 | 19.023 | 21.666 |
| 10 | 2.558 | 3.247 | 3.940 | 4.865 | 15.987 | 18.307 | 20.483 | 23.209 |
| 11 | 3.053 | 3.816 | 4.575 | 5.578 | 17.275 | 19.675 | 21.920 | 24.725 |
| 12 | 3.571 | 4.404 | 5.226 | 6.304 | 18.549 | 21.026 | 23.336 | 26.217 |
| 13 | 4.107 | 5.009 | 5.892 | 7.042 | 19.812 | 22.362 | 24.736 | 27.688 |
| 14 | 4.660 | 5.629 | 6.571 | 7.790 | 21.064 | 23.685 | 26.119 | 29.141 |
| 15 | 5.229 | 6.262 | 7.261 | 8.547 | 22.307 | 24.996 | 27.488 | 30.578 |
| 16 | 5.812 | 6.908 | 7.962 | 9.312 | 23.542 | 26.296 | 28.845 | 32.000 |
| 17 | 6.408 | 7.564 | 8.672 | 10.085 | 24.769 | 27.587 | 30.191 | 33.409 |
| 18 | 7.015 | 8.231 | 9.390 | 10.865 | 25.989 | 28.869 | 31.526 | 34.805 |
| 19 | 7.633 | 8.907 | 10.117 | 11.651 | 27.204 | 30.144 | 32.852 | 36.191 |
| 20 | 8.260 | 9.591 | 10.851 | 12.443 | 28.412 | 31.410 | 34.170 | 37.566 |
| 21 | 8.897 | 10.283 | 11.591 | 13.240 | 29.615 | 32.671 | 35.479 | 38.932 |
| 22 | 9.542 | 10.982 | 12.338 | 14.041 | 30.813 | 33.924 | 36.781 | 40.289 |
| 23 | 10.196 | 11.688 | 13.091 | 14.848 | 32.007 | 35.172 | 38.076 | 41.638 |
| 24 | 10.856 | 12.401 | 13.848 | 15.659 | 33.196 | 36.415 | 39.364 | 42.980 |
| 25 | 11.524 | 13.120 | 14.611 | 16.473 | 34.382 | 37.652 | 40.646 | 44.314 |
| 26 | 12.198 | 13.844 | 15.379 | 17.292 | 35.563 | 38.885 | 41.923 | 45.642 |
| 27 | 12.879 | 14.573 | 16.151 | 18.114 | 36.741 | 40.113 | 43.194 | 46.963 |
| 28 | 13.565 | 15.308 | 16.928 | 18.939 | 37.916 | 41.337 | 44.461 | 48.278 |
| 29 | 14.256 | 16.047 | 17.708 | 19.768 | 39.087 | 42.557 | 45.722 | 49.588 |
| 30 | 14.953 | 16.791 | 18.493 | 20.599 | 40.256 | 43.773 | 46.979 | 50.892 |
| 31 | 15.655 | 17.539 | 19.281 | 21.434 | 41.422 | 44.985 | 48.232 | 52.191 |
| 32 | 16.362 | 18.291 | 20.072 | 22.271 | 42.585 | 46.194 | 49.480 | 53.486 |
| 33 | 17.073 | 19.047 | 20.867 | 23.110 | 43.745 | 47.400 | 50.725 | 54.776 |
| 34 | 17.789 | 19.806 | 21.664 | 23.952 | 44.903 | 48.602 | 51.966 | 56.061 |