

Hazardous asteroids forecast via Markov random fields

Project for the exam: Probabilistic Modelling (DSE)

Marzio De Corato

December 16, 2021

Introduction

- **Final Goal** Assessment of forecasts and interpretability for different machine learning algorithms, including the probabilistic models
- **Method** Use a dataset for which the laws that interconnect the different features are known from general principles
- **Dataset** CNEOS asteroids dataset for more than 3500 asteroids
- **Theoretical laws** Celestial mechanics
- **Algorithms involved - probabilistic models** GLASSO, mgm, minforest, mmod
- **Algorithms involved - others** Random forest, Support Vector Machines, Quadratic Discriminant Analysis, Logistic Regression

Celestial mechanics

Celestial mechanics [14]: equations of motion

The interaction between a planet of mass m_1 at the position r_1 (inertial frame) and an asteroid of mass m_2 at the position r_2 is given by:

$$\mathbf{F}_1 = \mathcal{G} \cdot \frac{m_1 m_2}{r^3} \mathbf{r} = m_1 \ddot{\mathbf{r}}_1 \quad \mathbf{F}_2 = -\mathcal{G} \cdot \frac{m_1 m_2}{r^3} \mathbf{r} = m_2 \ddot{\mathbf{r}}_2 \quad (1)$$

Where \mathcal{G} is the universal gravitational constant. If we consider the motion of the second item with respect to the first one

$$\ddot{\mathbf{r}} = \ddot{\mathbf{r}}_2 - \ddot{\mathbf{r}}_1 \quad \mu = \mathcal{G}(m_1 + m_2) \quad (2)$$

$$\frac{d^2 \mathbf{r}}{dt^2} + \mu \frac{\mathbf{r}}{r^3} = 0 \quad (3)$$

$\mathbf{r} \times \ddot{\mathbf{r}} = 0 \implies \mathbf{r}$ and $\dot{\mathbf{r}}$ lies in the same plane

Celestial mechanics [14]: equations of motion

Integrating $\mathbf{r} \times \ddot{\mathbf{r}} = 0$

$$\mathbf{r} \times \dot{\mathbf{r}} = \mathbf{h} \quad (4)$$

Where \mathbf{h} is a constant of Integration. Using the polar coordinates $\hat{\mathbf{r}}$ and $\hat{\boldsymbol{\theta}}$

$$\mathbf{r} = r\hat{\mathbf{r}} \quad (5)$$

$$\dot{\mathbf{r}} = \dot{r}\hat{\mathbf{r}} + r\dot{\theta}\hat{\boldsymbol{\theta}} \quad (6)$$

$$\ddot{\mathbf{r}} = \left(\ddot{r} - r\dot{\theta}^2\right)\hat{\mathbf{r}} + \left[\frac{1}{r}\frac{d}{dt}\left(r^2\dot{\theta}\right)\right]\hat{\boldsymbol{\theta}} \quad (7)$$

$$\mathbf{h} = r^2\dot{\theta}\hat{\mathbf{z}} \quad (8)$$

$$h = r^2\dot{\theta} \quad (9)$$

Celestial mechanics [14]: 2th Kepler law

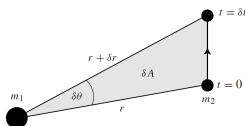


Figure 1: [14]

$$\delta A \approx \frac{1}{2} r(r + dr) \sin(\delta\theta) \approx \frac{1}{2} r^2 \delta\theta \quad (10)$$

$$\frac{dA}{dt} = \frac{1}{2} r^2 \frac{d\theta}{dt} = \frac{1}{2} h \quad (11)$$

h is constant \implies 2th Kepler law

Celestial mechanics [14]: 1th Kepler law

Using the substitution $u = \frac{1}{r}$ $h = r^2 \dot{\theta}$

$$\dot{r} = -\frac{1}{u} \frac{du}{d\theta} \dot{\theta} = -h \frac{du}{d\theta} \quad (12)$$

$$\ddot{r} = -h \frac{d^2 u}{d\theta^2} \dot{\theta} = -h^2 u^2 \frac{d^2 u}{d\theta^2} \quad (13)$$

$$\frac{d^2 u}{d\theta^2} + u = \frac{\mu}{h^2} \quad (14)$$

$$u = \frac{\mu}{h^2} [1 + e \cos(\theta - \phi)] \quad (15)$$

Celestial mechanics [14]: 1th Kepler law

$$r = \frac{p}{1 + e \cos(\theta - \phi)} \quad (16)$$

e is **eccentricity**

- circle: $e = 0$ $p = a$
- ellipse: $0 < e < 1$
 $p = a(1 - e^2)$
- parabola: $e = 1$ $p = 2q$
- hyperbola: $e > 1$
 $p = a(e^2 - 1)$

a is the **semi-major axis** of the conic

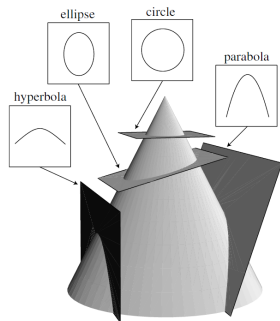


Figure 2: [14]

Celestial mechanics [14]: 3th Kepler law

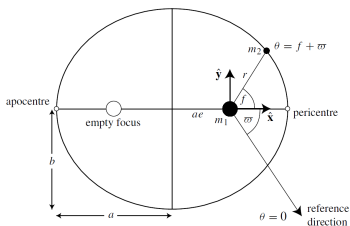


Figure 3: [14]

$$b^2 = a^2(1 - e^2) \quad (17)$$

$$r = \frac{a(1 - e^2)}{1 + e \cdot \cos(\theta - \phi)} \quad (18)$$

Area swept in one **orbital period** T

$$A = \pi ab$$

We know that: $hT/2 \quad h^2 = \mu a(1 - e^2)$

Therefore

$$T^2 = \frac{4\pi^2}{\mu} a^3 \quad (19)$$

Celestial mechanics [14]: 3th Kepler law

Consider two asteroids of mass m and m' orbiting the Earth m_c , with semi-major axes a and a' and orbital periods T and T'

$$\frac{m_c + m}{m_c + m'} = \left(\frac{a}{a'}\right)^3 \left(\frac{T'}{T}\right)^2 \quad (20)$$

But since $m, m' \ll m_c$

$$\left(\frac{a}{a'}\right)^3 \approx \left(\frac{T}{T'}\right)^2 \quad (21)$$

Remark: The mass of the asteroids is **not** involved

Celestial mechanics [14]: Orbital parameters

Mean motion $n = \frac{2\pi}{T}$

$$v_{perihelion} = na\sqrt{\frac{1+e}{1-e}} \quad (22)$$

$$v_{aphelion} = na\sqrt{\frac{1-e}{1+e}} \quad (23)$$

Remark: The mean motion of an asteroid is different with respect to the the asteroid relative velocity (measured from Earth), since the latter is different at the perihelion an at the aphelion

Celestial mechanics [14]: Orbital parameters

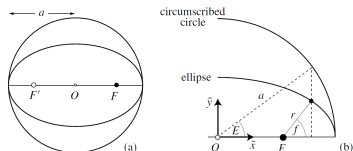


Figure 4: [14]

Mean anomaly

$$M = n(t - \tau) \quad (24)$$

- $M = f = 0 \quad t = \tau$ Perihelion
- $M = f = \pi \quad t = \tau + T/2$ Aphelion

$$M = E - e \sin E \quad (25)$$

Jupiter Tisserard invariant

$$T_P = \frac{a_p}{a} + 2 \cos I \sqrt{\frac{a}{a_p} (1 - e^2)} \quad (26)$$

Celestial mechanics [14]: Orbital parameters

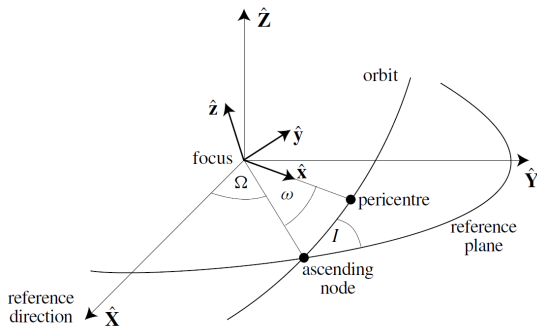


Figure 5: [14]

i : inclination of the orbit

Ω : longitude of the ascending node

Celestial mechanics [14]: Magnitude

$$\Phi = \frac{L}{4\pi d^2} \quad (27)$$

$$m = -2.5 \log_{10} \Phi + C \quad (28)$$

$$m_1 - m_2 = -2.5 \log_{10} \frac{\Phi_1}{\Phi_2} \quad (29)$$

$$M - m = -2.5 \log_{10} \frac{\Phi \cdot d^2}{\Phi \cdot 10^2} \quad (30)$$

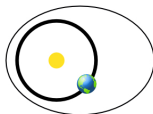
$$M = m + 5 - 5 \log_{10} d \quad (31)$$

Where Φ is the flux for a sphere of radius r , m the relative magnitude and M the **Absolute magnitude**

Celestial mechanics [1]: Classification

Amors

Earth-approaching NEAs with orbits exterior to Earth's but interior to Mars' (named after asteroid (1221) Amor)



$$a > 1.0 \text{ AU} \\ 1.017 \text{ AU} < q < 1.3 \text{ AU}$$

Apollos

Earth-crossing NEAs with semi-major axes larger than Earth's (named after asteroid (1862) Apollo)



$$a > 1.0 \text{ AU} \\ q < 1.017 \text{ AU}$$

Atens

Earth-crossing NEAs with semi-major axes smaller than Earth's (named after asteroid (2062) Aten)



$$a < 1.0 \text{ AU} \\ Q > 0.983 \text{ AU}$$

Atiras

NEAs whose orbits are contained entirely within the orbit of the Earth (named after asteroid (163693) Atira)



$$a < 1.0 \text{ AU} \\ Q < 0.983 \text{ AU}$$

(q = perihelion distance, Q = aphelion distance, a = semi-major axis)

Celestial mechanics [1]: Classification

- **Potentially Hazardous Asteroids:** $\text{MOID} \leq 0.05 \text{ au}$ $M \leq 22.0$
NEAs whose Minimum Orbit Intersection Distance (MOID) with the Earth is 0.05 au or less and whose absolute magnitude (M) is 22.0 or brighter

Dataset

- The asteroid dataset was retrieved from Kaggle [2], which reports into a more machine readable form the dataset of The Center for Near-Earth Object Studies (CNEOS) [3], a NASA research centre.
- 3552 Asteroids
- Among the 40 the features, the ones connected only to the other name of the asteroid, or connected only to the name of the orbit and the one connected with the orbiting planet (since for all it was the Earth) were discarded
- The proportion hazardous/not hazardous was set 1:5
- The continuous measures were standardised and demeaned

Features

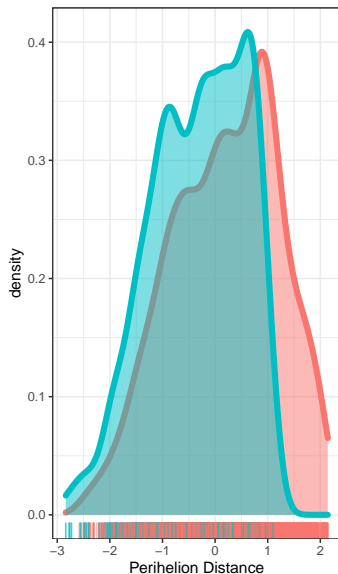
| Features | Type |
|-----------------------------|------------|
| Neo Reference ID | not used |
| Absolute Magnitude | Continuous |
| Est Dia in KM (min) | Continuous |
| Est Dia in KM (max) | Continuous |
| Close Approach Date | Continuous |
| Epoch Date Close Approach | Continuous |
| Relative_Velocity | Continuous |
| Miss_Dist | Continuous |
| Min_Orbit_Intersection | Continuous |
| Jupiter_Tisserand_Invariant | Continuous |
| Epoch_Osculation | Continuous |
| Eccentricity | Continuous |

Features

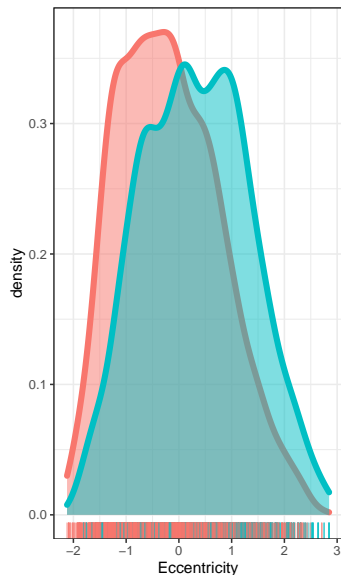
| Features | Type |
|---------------------|----------------------|
| Semi Major Axis | Continuous |
| Inclination | Continuous |
| Asc Node Longitude | Continuous |
| Orbital Period | Continuous |
| Perihelion Distance | Continuous |
| Perihelion Arg | Continuous |
| Perihelion Time | Continuous |
| Mean_Anomaly | Continuous |
| Mean_Motion | Continuous |
| Hazardous | Categorical (Binary) |

Prelim. analysis

Density Plot

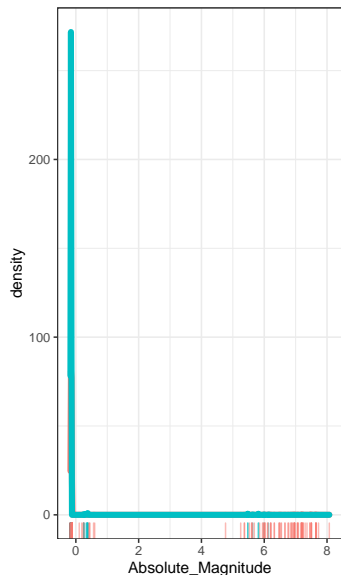
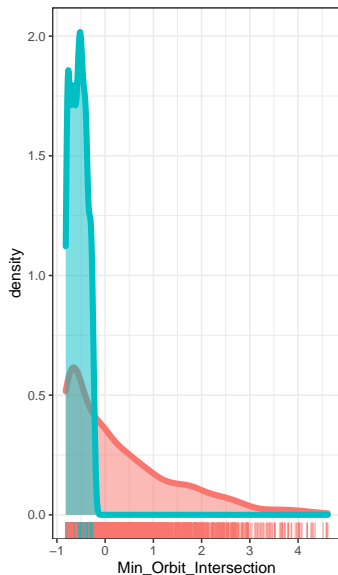


Hazardous
FALSE
TRUE

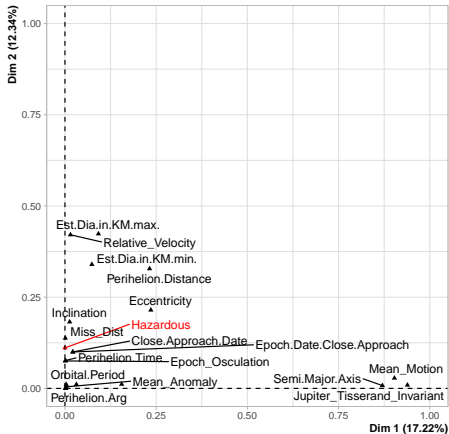


Hazardous
FALSE
TRUE

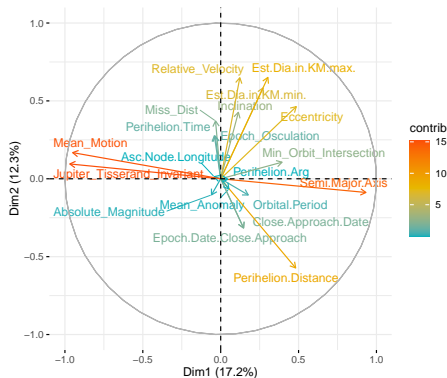
Density Plot



Graph of the variables

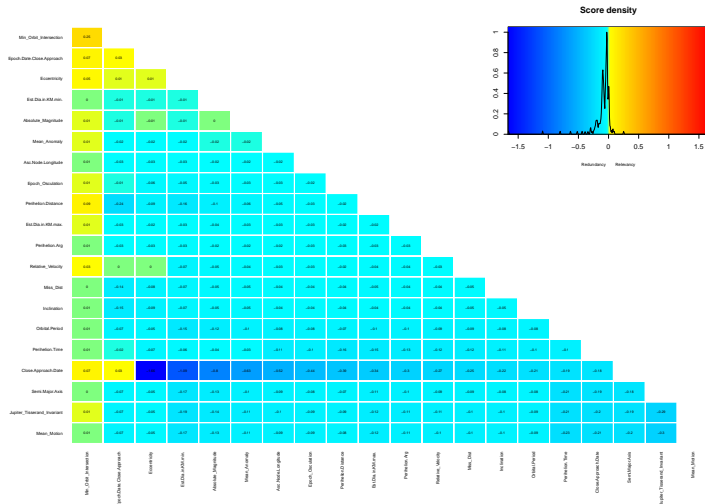


Quantitative variables – FAMD



Performed with the FactoMineR package [12]

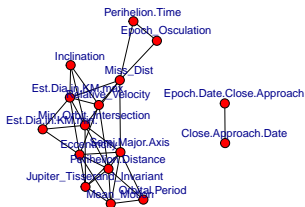
Mutual information analysis



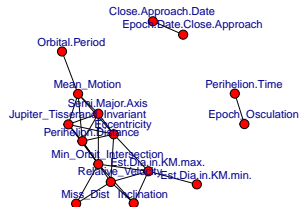
Performed with the varrank package [11]

Probabilistic modelling

GLASSO



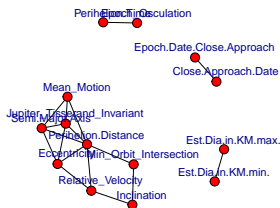
$\rho=0.1$



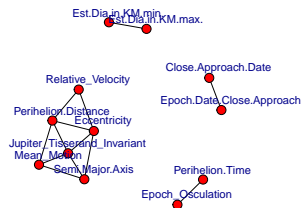
$\rho=0.2$

Performed with the GLASSO package [4]

GLASSO



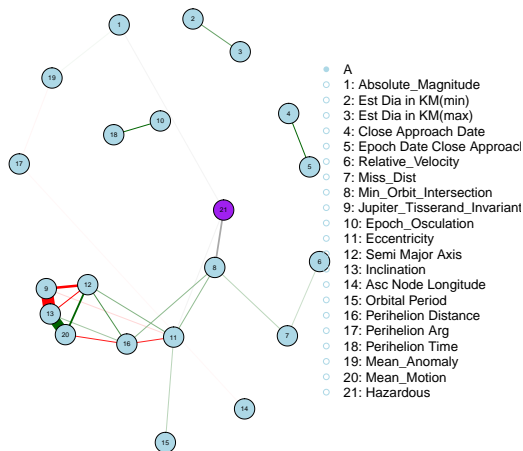
$\rho=0.3$



$\rho=0.4$

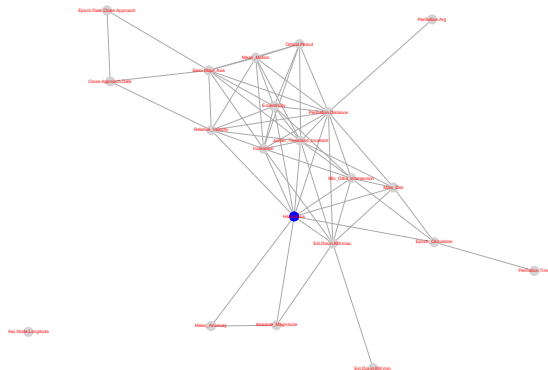
Performed with the GLASSO package [4]

Mixed interactions: mgm



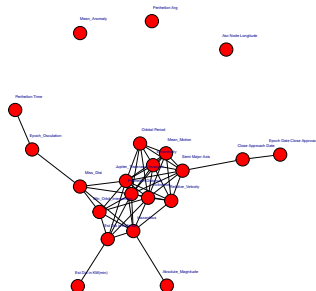
Performed with the mgm package [9]

Mixed interactions: minforest



Performed with the gRapHD package [7]

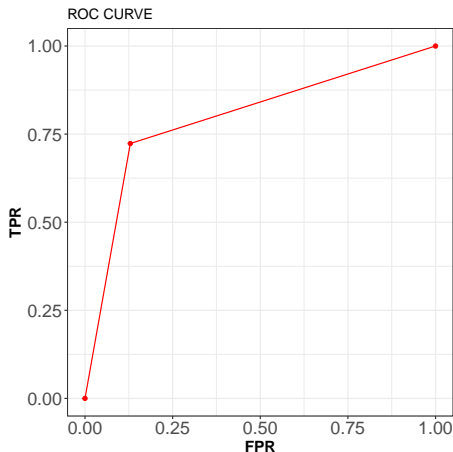
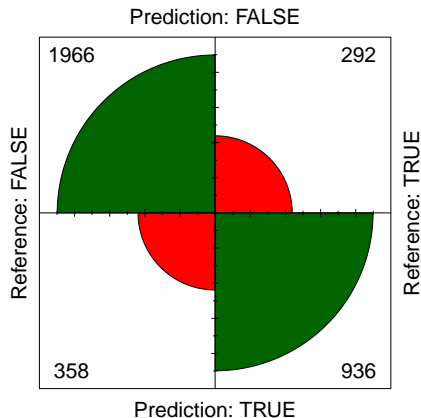
Mixed interactions: mmod



Performed with the gRim package [10]

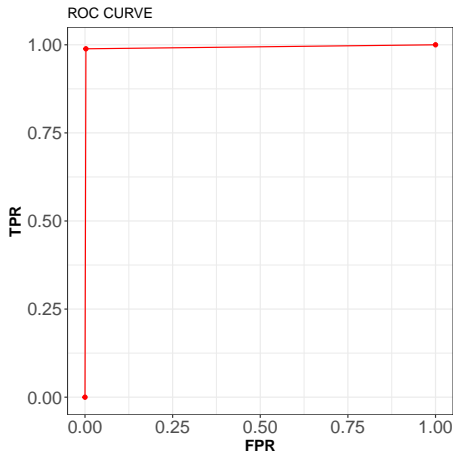
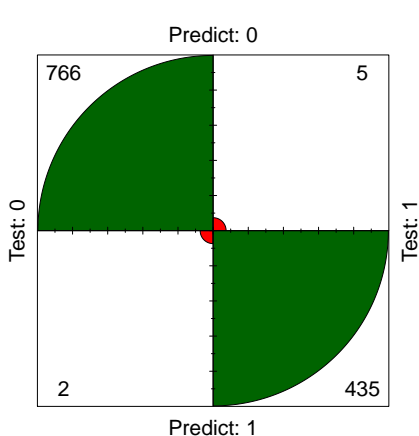
Mixed interactions

The mgm model is the one that has the list of connection more coherent with the celestial mechanics laws.



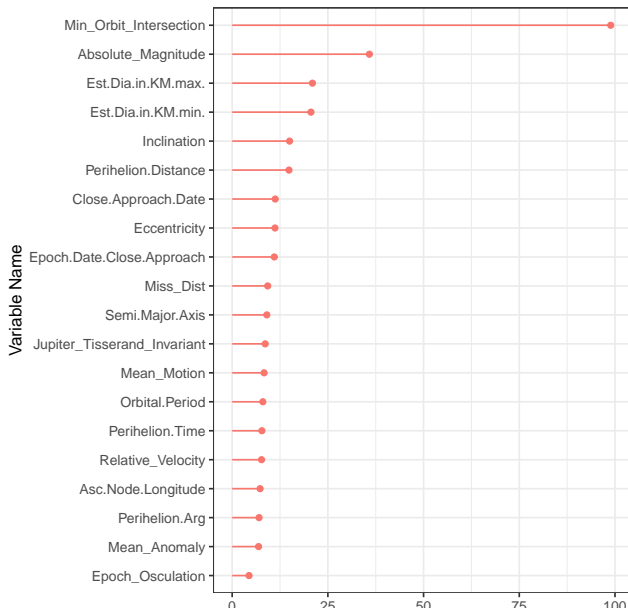
Other ML algorithms

Random Forest

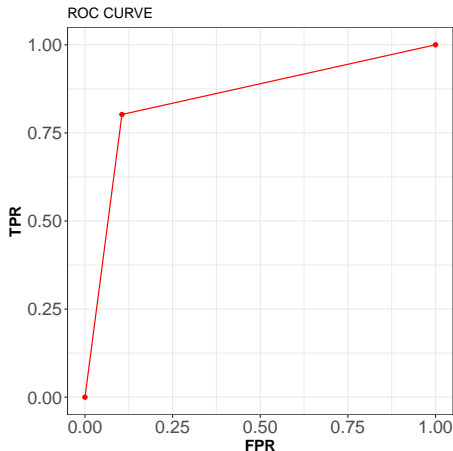
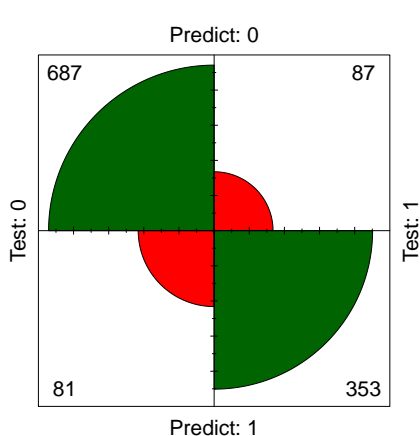


Performed with the rfor package [13]

Random Forest

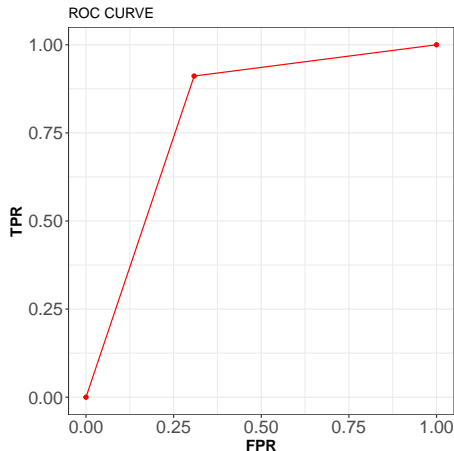
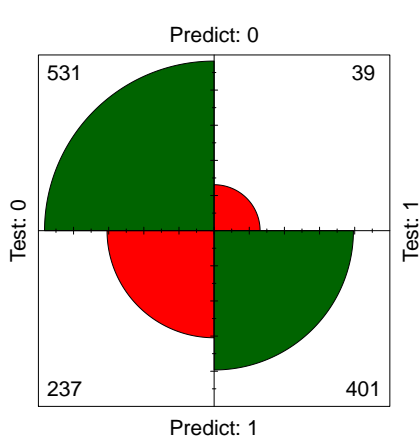


Support Vector Machines



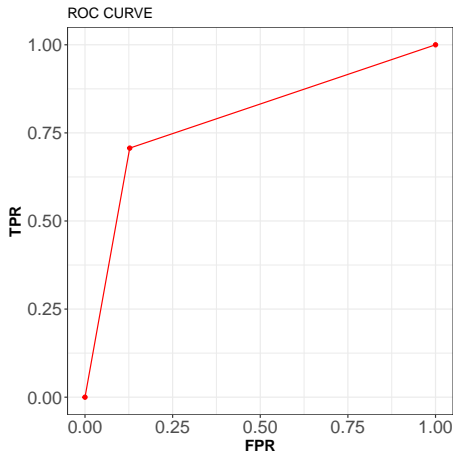
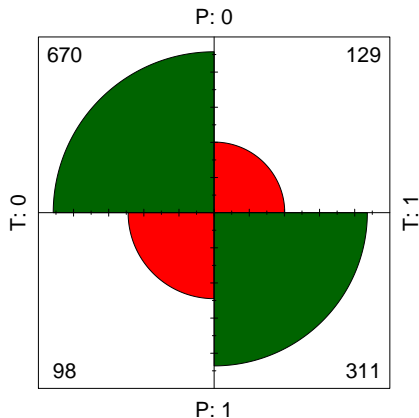
Performed with the e1071 package [8]

Quadratic Discriminant Analysis (QDA)



Performed with the MASS package [17]

Logistic regression



Performed with the stats package [15]

Table 1: ϕ coefficient (also known as Matthews correlation coefficient)

| Algorithm | ϕ |
|-----------|--------|
| RF | 0.9876 |
| SVM | 0.7111 |
| logistic | 0.6173 |
| mgm | 0.5997 |
| QDA | 0.5562 |

Conclusions

Remark (Interpretability - Tarski definition)

The formal theory T can be translated into S if and only if S can prove the theorem of T in its language [16]

Remark (Scientific method - Einstein definition)

Science uses the totality of the primary concepts, i.e., concepts directly connected with sense experiences, and propositions connecting them. Such a state of affairs cannot, however, satisfy a spirit which is really scientifically minded; because the totality of concepts and relations obtained in this manner is utterly lacking in logical unity. In order to supplement this deficiency, one invents a system poorer in concepts and relations, a system retaining the primary concepts and relations of the first layer as logically derived concepts and relations. This new secondary system pays for its higher logical unity by having elementary concepts (concepts of the second layer), which are no longer directly connected with complexes of sense experiences [5]

Conclusions: forecast performances vs interpretability

- The mgm algorithm is not the best one in terms of performance, but it provides the connections between the features. On the other side, except for the variable importance in RF, the SVM and RF are a black box one
- The mgm model, as the different graphical model is open to a proper scientific validation, the SVM and the RF not.
- The probabilistic models lack in the forecast is compensated by their interpretability
- This is meaningful since interpretability and completeness/accuracy are in conflict
- It is worth noting that the logistic model, which is highly interpretable, has performances similar to the mgm
- The probabilistic models provide a good trade-off between interpretability and forecast performances. As long as one is interested in producing a really scientific result, such algorithms as the logistic regression should be considered (e.g if the only aim is the forecast, the RF is better. However, how long one can trust the RF result ?)

In its efforts to learn as much as possible about nature, modern physics has found that certain things can never be “known” with certainty. Much of our knowledge must always remain uncertain. The most we can know is in terms of probabilities. Richard P. Feynman (1918-1988)

Bibliography I

- [1] https://cneos.jpl.nasa.gov/about/neo_groups.html.
- [2] <https://www.kaggle.com/shrutimehta/nasa-asteroids-classification>.
- [3] <https://cneos.jpl.nasa.gov/>.
- [4] <https://cran.r-project.org/web/packages/glasso/glasso.pdf>.
- [5] <https://www.amacad.org/publication/physics-reality>.
- [6] <https://towardsdatascience.com/understanding-auc-roc-curve-68b2303cc9c5>.
- [7] Gabriel CG de Abreu, Rodrigo Labouriau, and David Edwards.
“High-dimensional graphical model search with graphd R package”.
In: *arXiv preprint arXiv:0909.1234* (2009).

Bibliography II

- [8] Evgenia Dimitriadou et al. “Misc functions of the Department of Statistics (e1071), TU Wien”. In: *R package 1* (2008), pp. 5–24.
- [9] Jonas Haslbeck and Lourens J Waldorp. “mgm: Estimating time-varying mixed graphical models in high-dimensional data”. In: *arXiv preprint arXiv:1510.06871* (2015).
- [10] Søren Højsgaard, David Edwards, and Steffen Lauritzen. *Graphical Models with R*. ISBN 978-1-4614-2298-3. New York: Springer, 2012. DOI: 10.1007/978-1-4614-2299-0.
- [11] Gilles Kratzer and Reinhard Furrer. “varrank: an R package for variable ranking based on mutual information with applications to observed systemic datasets”. In: *arXiv preprint arXiv:1804.07134* (2018).

Bibliography III

- [12] Sébastien Lê, Julie Josse, and François Husson. “FactoMineR: an R package for multivariate analysis”. In: *Journal of statistical software* 25.1 (2008), pp. 1–18.
- [13] Andy Liaw and Matthew Wiener. “Classification and Regression by randomForest”. In: *R News* 2.3 (2002), pp. 18–22. URL: <https://CRAN.R-project.org/doc/Rnews/>.
- [14] Carl D Murray and Stanley F Dermott. *Solar system dynamics*. Cambridge university press, 1999.
- [15] R Core Team. *R: A Language and Environment for Statistical Computing*. ISBN 3-900051-07-0. R Foundation for Statistical Computing. Vienna, Austria, 2013. URL: <http://www.R-project.org/>.
- [16] Alfred Tarski, Andrzej Mostowski, and Raphael Mitchel Robinson. *Undecidable theories*. Vol. 13. Elsevier, 1953.

Bibliography IV

- [17] W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Fourth. ISBN 0-387-95457-0. New York: Springer, 2002. URL: <https://www.stats.ox.ac.uk/pub/MASS4/>.