

Hazardous asteroids forecast via Markov random fields

Project for the exam: Probabilistic Modelling (DSE)

Marzio De Corato

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 [12]: equations of motion

Lets Consider 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

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

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 [12]: equations of motion

With polar coordinates $\hat{\mathbf{r}}$ and $\hat{\boldsymbol{\theta}}$

$$\mathbf{r} = r\hat{\mathbf{r}} \quad (4)$$

$$\dot{\mathbf{r}} = \dot{r}\hat{\mathbf{r}} + r\dot{\theta}\hat{\boldsymbol{\theta}} \quad (5)$$

$$\ddot{\mathbf{r}} = \left(\ddot{r} - r\dot{\theta}^2\right)\hat{\mathbf{r}} + \left[\frac{1}{r}\frac{d}{dt}\left(r^2\dot{\theta}\right)\right]\hat{\boldsymbol{\theta}} \quad (6)$$

$$\mathbf{h} = r^2\dot{\theta}\hat{\mathbf{z}} \quad (7)$$

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

Celestial mechanics [12]: 2th Kepler law

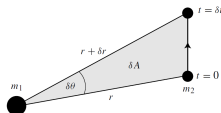


Figure 1: [12]

$$\delta A \approx \frac{1}{2} r(r + dr) \sin(\delta\theta) \approx \frac{1}{2} r^2 \delta\theta \quad (9)$$

$$\frac{dA}{dt} = \frac{1}{2} r^2 \frac{d\theta}{dt} = \frac{1}{2} h \quad (10)$$

h is constant \implies 2th Kepler law

Celestial mechanics [12]: 1th Kepler law

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

$$\dot{r} = -\frac{1}{u} \frac{du}{d\theta} \dot{\theta} = -h \frac{du}{d\theta} \quad (11)$$

$$\ddot{r} = -h \frac{d^2 u}{d\theta^2} \dot{\theta} = -h^2 u^2 \frac{d^2 u}{d\theta^2} \quad (12)$$

$$\frac{d^2 u}{d\theta^2} + u = \frac{\mu}{h^2} \quad (13)$$

$$u = \frac{\mu}{h^2} [1 + e \cos(\theta - \phi)] \quad (14)$$

Celestial mechanics [12]: 1th Kepler law

$$r = \frac{p}{1 + e \cos(\theta - \phi)} \quad (15)$$

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

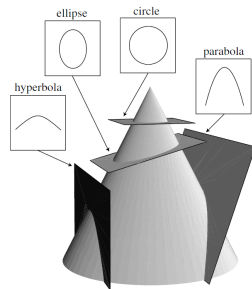


Figure 2: [12]

Celestial mechanics [12]: 3th Kepler law

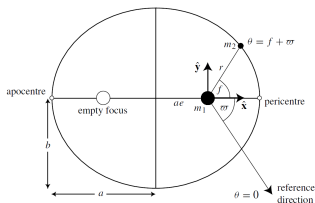


Figure 3: [12]

$$b^2 = a^2(1 - e^2) \quad (16)$$

$$r = \frac{a(1 - e^2)}{1 + e \cdot \cos(\theta - \phi)} \quad (17)$$

Area swept in one **orbital period** T

$$A = \pi ab$$

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

Therefore

$$T^2 = \frac{4\pi^2}{\mu} a^3 \quad (18)$$

Celestial mechanics [12]: 3th Kepler law

$$\frac{m_c + m}{m_c + m'} = \left(\frac{a}{a'}\right)^3 \left(\frac{T'}{T}\right)^2 \quad (19)$$

But since $m, m' \ll m_c$

$$\left(\frac{a}{a'}\right)^3 \approx \left(\frac{T'}{T}\right)^2 \quad (20)$$

And therefore

$$T' \approx a'^{3/2} \quad (21)$$

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

Celestial mechanics [12]: 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 [12]: Orbital parameters

Mean anomaly

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

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

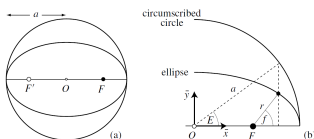


Figure 4: [12]

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

Jupiter Tisserard invariant

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

Celestial mechanics [12]: Orbital parameters

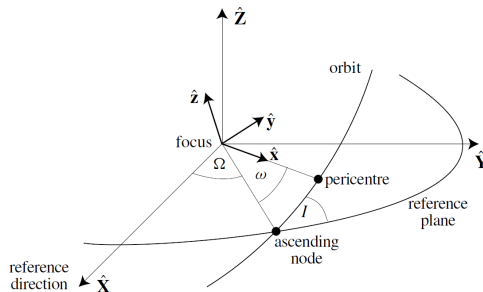


Figure 5: [12]

I : inclination of the orbit

Ω : longitude of the ascending node

Celestial mechanics [12]: Magnitude

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

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

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

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

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

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

Celestial mechanics [12]: Magnitude

$$\Phi = \frac{L}{4\pi r^2} \quad (32)$$

$$m = -2.5 \log_{10} \Phi + C \quad (33)$$

$$m_1 - m_2 = -2.5 \log_{10} \frac{\Phi_1}{\Phi_2} \quad (34)$$

$$M - m = -2.5 \log_{10} \frac{\Phi \cdot d^2}{\Phi \cdot 10^2} \quad (35)$$

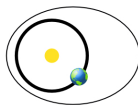
$$M = m + 5 - 5 \log_{10} d \quad (36)$$

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

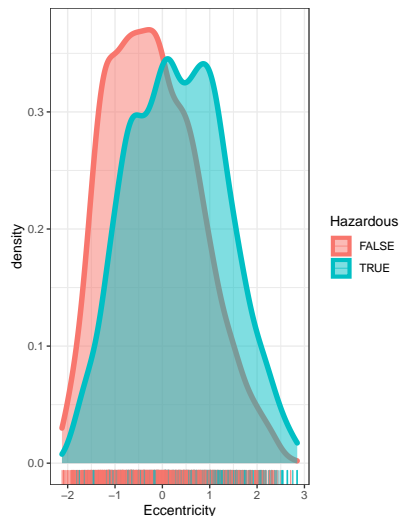
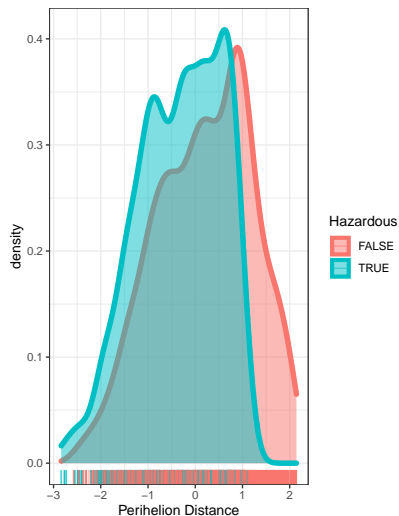
Features

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

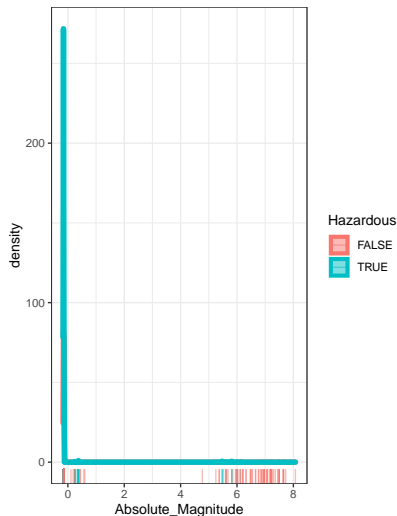
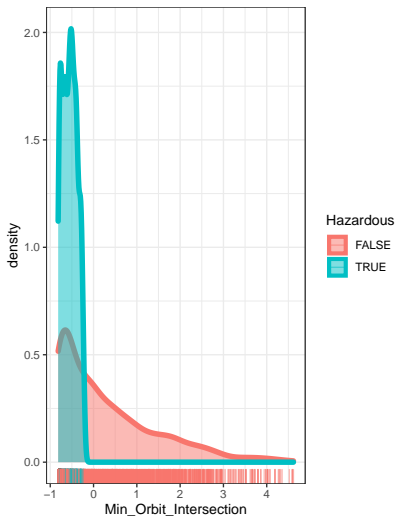
Features

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)

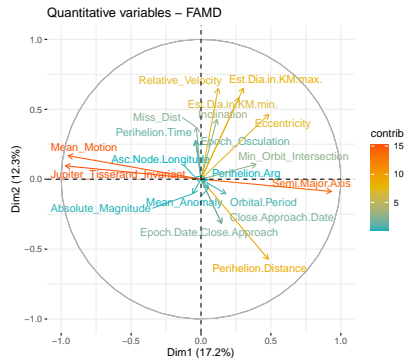
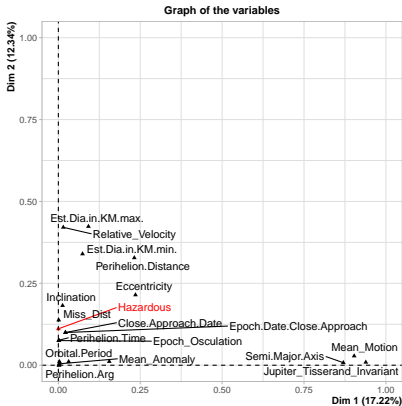
Density Plot



Density Plot

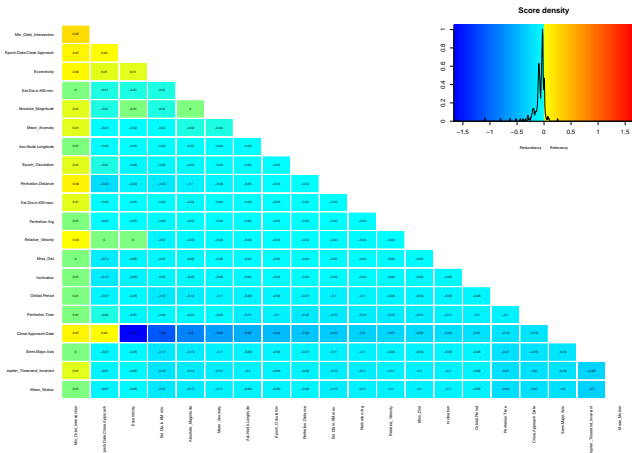


FAMD



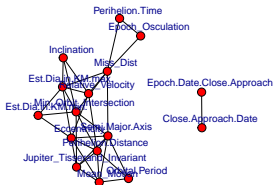
Performed with the FactoMineR package [10]

Mutual information analysis

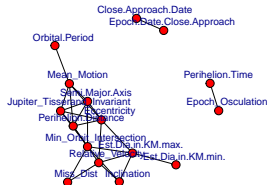


Performed with the varrank package [9]

GLASSO



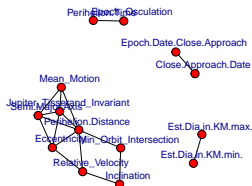
$\rho=0.2$



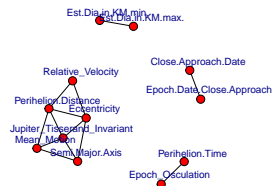
$\rho=0.2$

Performed with the GLASSO package [4]

GLASSO



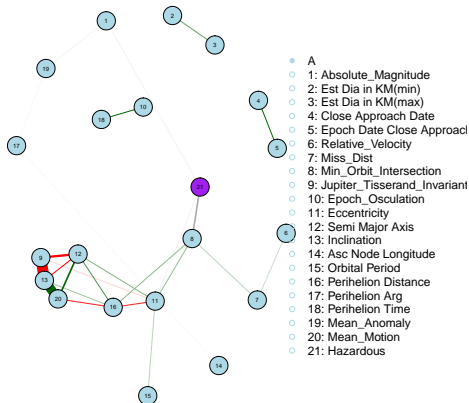
$$\rho=0.3$$



$$\rho=0.4$$

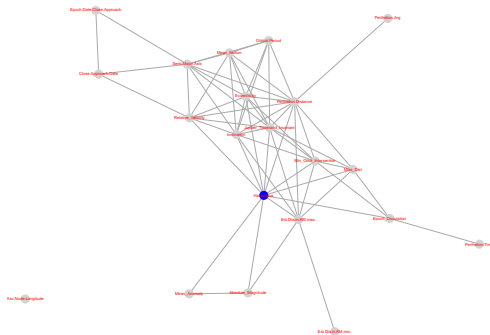
Performed with the GLASSO package [4]

Mixed interactions: mgm



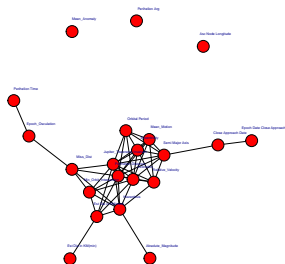
Performed with the mgm package [7]

Mixed interactions: minforest



Performed with the gRapHD package [6]

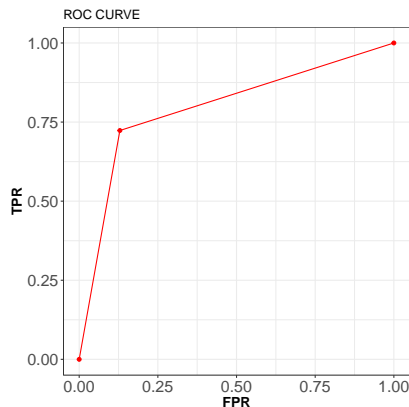
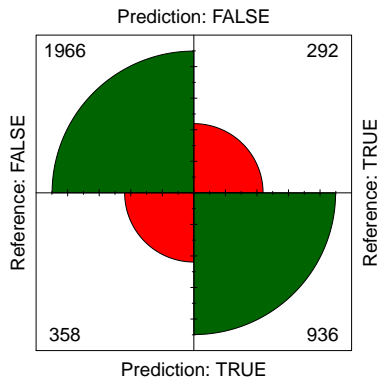
Mixed interactions: mmod



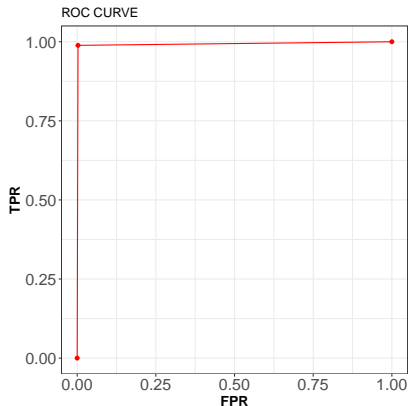
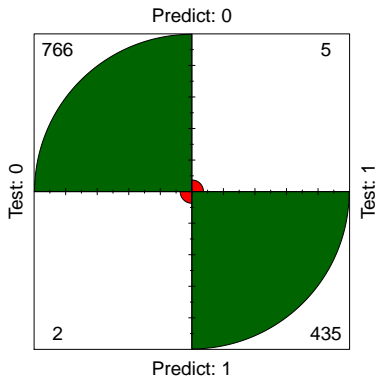
Performed with the gRim package [8]

Mixed interactions

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

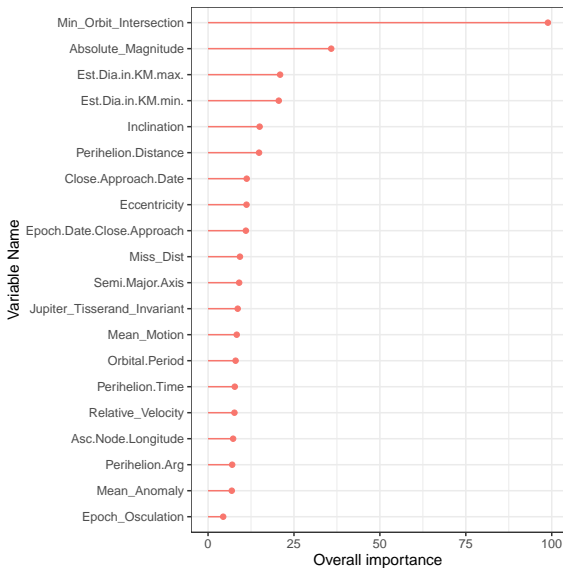


Random Forest

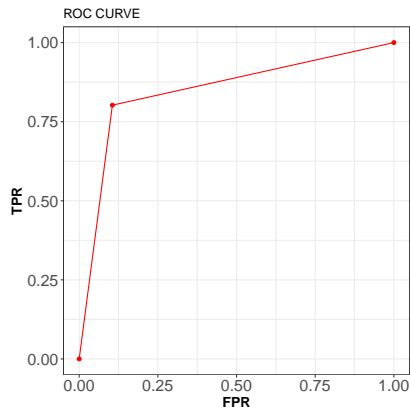
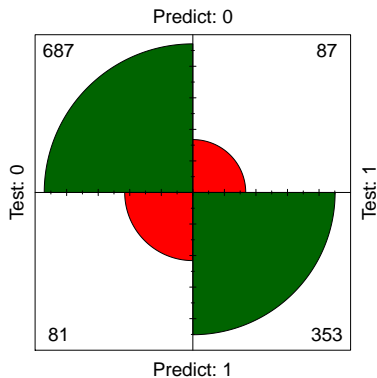


Performed with the rfor package [11]

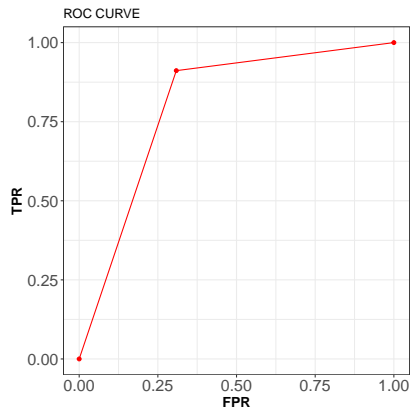
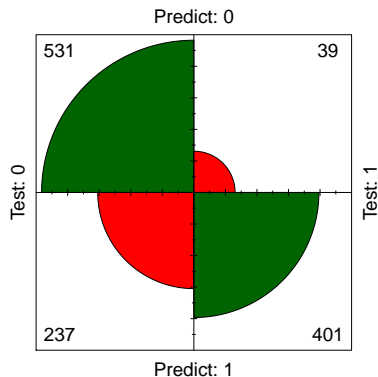
Random Forest



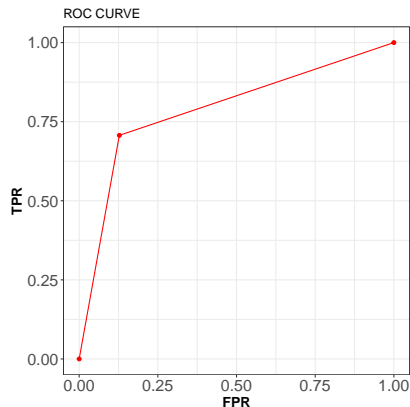
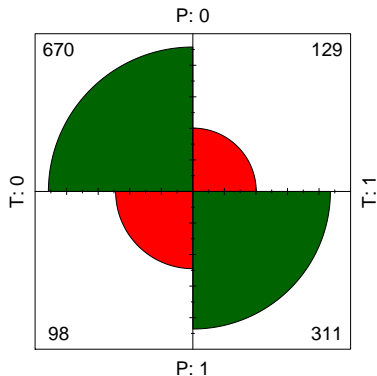
Support Vector Machines



Quadratic Discriminant Analysis (QDA)



Logistic regression



ϕ coefficient

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: forecast performances vs intepretability

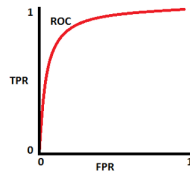
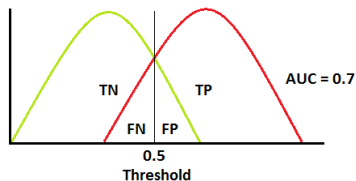
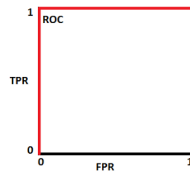
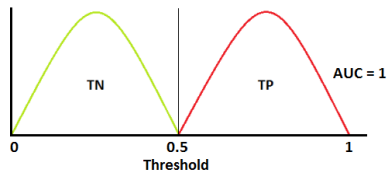
- The mgm algorithm is not the best one in term of performances, but it provides the connections between the features. On the other side, except for the variable importance in RF, the other are black box one
- The mgm model, as the other graphical model is open to a true scientific validation, the other not.
- The probabilistic models lack in the forecast is definitely compensated by their interetability
- This is meaningful since this two features are in conflict
- The probabilistic models provide a good trade-off between intepretability and forecast performances, as long as one is interest to produce a really scientific result (e.g if the only aim is the forecast the RF is definitely better. However how long one can trust to the RF result ?)

ϕ coefficient

	Actual - N	Actual - P
Predicted - N	#TP	#FN
Predicted - P	#FP	#TP

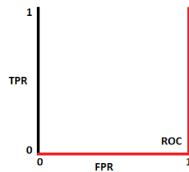
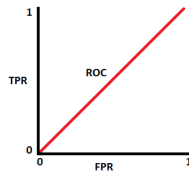
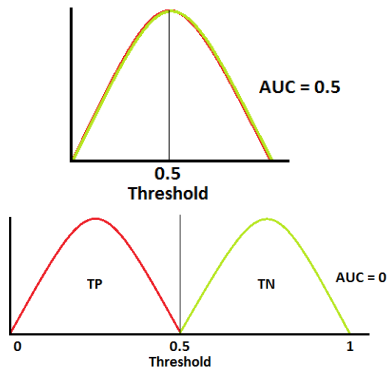
$$\phi = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Receiver operating characteristic



Images taken from [5]

Receiver operating characteristic



Images taken from [5]

Factor analysis of mixed data - FAMD

$r(z,k)$ correlation coefficient (z and k quantitative)

$\eta^2(z,q)$ correlation ratio (z quantitative and q qualitative)

$$\text{PCA} \rightarrow \max \sum_k r^2(z, k)$$

$$\text{MCA} \rightarrow \max \sum_q \eta^2(z, q)$$

$$\text{FAMD} \rightarrow \max \sum_k r^2(z, k) + \max \sum_q \eta^2(z, q)$$

mgm algorithm

$$P(X_s|X_{\setminus s}) = \exp \{ E_s(X_{\setminus s})\phi_s(X_s) + B_s(X_s) - \Phi(X_{\setminus s}) \} \quad (37)$$

ϕ_s function of sufficient statistics B_s base measure

$$P(X) = \exp \left\{ \sum_{s \in V} \theta_s \phi_s(X_s) + \sum_{s \in V} \sum_{r \in N(s)} \theta_{s,r} \phi_s(X_s) \phi_r(X_r) + \dots + \right. \\ \left. \sum_{r_1, \dots, r_k \in C} \theta_{r_1, \dots, r_k} \prod_{j=1}^k \phi_{r_j}(X_{r_j}) + \sum_{s \in V} B_s(X_s) - \Phi(\theta) \right\} \quad (38)$$

$$\hat{\theta} = \arg \min_{\theta} \{ -\mathcal{L}(\theta, X) + \lambda \|\theta\|_1 \} \quad \|\theta\|_1 = \sum_{j=1}^J |\theta_j| \quad (39)$$

GLASSO

$$L_{pen}(K, \hat{\mu}) = \log \det(K) - \text{tr}(K S) - \rho \|K\| \quad (40)$$

$$K = \Sigma^{-1}$$

S: empirical covariance matrix

Parameters tuning

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left(\sum_c \boldsymbol{\theta}_c^T \phi_c(\mathbf{y}) \right) \quad (41)$$

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{N} \sum_i \log p(\mathbf{y}_i|\boldsymbol{\theta}) = \frac{1}{N} \sum_i \left[\sum_c \boldsymbol{\theta}_c^T \phi_c(y_i) - \log Z(\boldsymbol{\theta}) \right] \quad (42)$$

Parameters tuning

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}_c} = \frac{1}{N} \sum_i \left[\phi_c(y_i) - \frac{\partial}{\partial \boldsymbol{\theta}_c} \log Z(\boldsymbol{\theta}) \right] \quad (43)$$

$$\frac{\partial \log Z(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E} [\phi_c(\mathbf{y}) | \boldsymbol{\theta}] = \sum_{\mathbf{y}} \phi_c(\mathbf{y}) p(\mathbf{y} | \boldsymbol{\theta}) \quad (44)$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}_c} = \left[\frac{1}{N} \sum_i \phi_c(y_i) \right] - \mathbb{E} [\phi_c(\mathbf{y})] \quad (45)$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}_c} = \mathbb{E}_{p_{emp}} [\phi_c(\mathbf{y})] - \mathbb{E}_{p_{(\cdot | \boldsymbol{\theta})}} [\phi_c(\mathbf{y})] \quad (46)$$

$$\mathbb{E}_{p_{emp}} [\phi_c(\mathbf{y})] = \mathbb{E}_{p_{(\cdot | \boldsymbol{\theta})}} [\phi_c(\mathbf{y})] \quad (47)$$

mmod

$$f(i, y) = p(i) (2\pi)^{-q/2} \det(\Sigma)^{-1/2} \exp \left[-\frac{1}{2} (y - \mu(i))^T \Sigma^{-1} (y - \mu(i)) \right] \quad (48)$$

$$\begin{aligned} f(i, y) &= \exp \left\{ g(i) + \sum_u h^u(i) y_u - \frac{1}{2} \sum_{uv} y_u y_v k_{uv} \right\} \\ &= \exp \left\{ g(i) + h(i)^T y - \frac{1}{2} y^T K y \right\} \end{aligned} \quad (49)$$

where $g(i)$, $h(i)$ and K are the canonical parameters

mmod

$$\begin{aligned}
 K &= \Sigma^{-1} \\
 h(i) &= \Sigma^{-1} \mu(i) \\
 g(i) &= \log p(i) - \frac{1}{2} \log \det(\Sigma) \\
 &\quad - \frac{1}{2} \mu(i)^T \Sigma^{-1} \mu(i) - \frac{q}{2} \log 2\pi
 \end{aligned} \tag{50}$$

Graphical models

$$p(x_1, x_2, \dots, x_n) \quad (51)$$

$$p(x_{1:v}) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)\dots p(x_v|x_{1:v-1}) \quad (52)$$

$$X \perp Y|Z \iff p(X, Y|Z) = p(X|Z)p(Y|Z) \quad (53)$$

$$p(\mathbf{x}_{1:v}) = p(x_1) \prod_{t=1}^v p(x_t|x_{t-1}) \quad (54)$$

Graphical models

Theorem (Hammersley-Clifford)

A positive distribution $p(\mathbf{y}) \geq 0$ satisfies the CI properties of an indirect graph G iif p can be represented as a product of factor, one per maximal clique, i.e.

$$p(\mathbf{y}|\theta) = \frac{1}{Z(\theta)} \prod_{c \in C} \psi_c(\mathbf{y}_c|\theta_c) \quad (55)$$

where C is the set of all the (maximal) cliques of G , and $Z(\theta)$ is the partition function given by

$$Z(\theta) := \sum_{\mathbf{y}} \prod_{c \in C} \psi_c(\mathbf{y}_c|\theta_c) \quad (56)$$

Note that this partition function is what ensures the overall distribution sums to 1

Graphical models

$$p(y|\theta) = \frac{1}{Z(\theta)} \exp \left(- \sum_c E(y_c|\theta_c) \right) \quad (57)$$

$$\psi_c(y_c|\theta_c) = \exp(-E(y_c|\theta_c)) \quad (58)$$

Information theory

$$H(X) = - \sum_{x \in X} p(x) \log p(x) \quad (59)$$

$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y) \quad (60)$$

$$\begin{aligned} H(X|Y) &= \sum_{x \in \mathcal{X}} p(x) H(Y|X = x) \\ &= - \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x) \\ &= - E \log p(Y|X) \end{aligned} \quad (61)$$

$$H(X, Y) = H(X) + H(Y|X) \quad (62)$$

$$D(p||q) = \sum p(x) \log \frac{p(x)}{q(x)} \quad D(p||q) \geq 0 \quad (63)$$

Information theory

$$I(X; Y) = \sum (x, y) \log \frac{p(x, y)}{p(x)p(y)} = D(p(x, y) || p(x)p(y)) \quad (64)$$

$$= H(X) - H(X|Y) = H(Y) - H(Y|X)$$

$$I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, X_{i-2}, \dots, X_1) \quad (65)$$

$$g(\alpha, \mathbf{C}, \mathbf{S}, f_i) = MI(f_i; \mathbf{C}) - \sum_{f_s \in S} \alpha(f_i, f_s, \mathbf{C}, \mathbf{S}) MI(f_i; f_s) \quad (66)$$

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://towardsdatascience.com/understanding-auc-roc-curve-68b2303cc9c5>.
- [6] 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

- [7] Jonas Haslbeck and Lourens J Waldorp. “mgm: Estimating time-varying mixed graphical models in high-dimensional data”. In: *arXiv preprint arXiv:1510.06871* (2015).
- [8] 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.
- [9] 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).
- [10] 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.

Bibliography III

- [11] 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/>.
- [12] Carl D Murray and Stanley F Dermott. *Solar system dynamics*. Cambridge university press, 1999.