

Una aplicación Estadística de los métodos de clasificación en Astronomía

A Statistical application of classification methods in Astronomy

Resumen

En los últimos años, los avances en astrofísica y cosmología han sido impulsados por grandes y complejos conjuntos de datos, los cuales sólo pueden ser analizados e interpretados con el uso de métodos estadísticos muy refinados. Esto ha llevado a que estas disciplinas se complementen para formar una rama llamada la astroestadística. En este trabajo se da a conocer un método de clasificación estadístico usando modelos de mezclas de Gausianas. Este método se aplicará para encontrar estrellas que pertenecen al cúmulo de las Hyades usando una muestra de 2678 estrellas de la base de datos de Hipparcos. Se realiza una descripción breve de las características del cúmulo y se estudia la evidencia de outliers. Con este método se encuentra que la clasificación arroja tres grupos de los cuales podemos estudiar la pertenencia al cúmulo y se encuentra que la mayoría de estrellas pertenecientes al mismo están de acuerdo con la literatura. También se muestra el diagrama de Hertzsprung-Russell obtenido para el cúmulo, muy importante en estudios de evolución estelar. Finalmente, se analiza un tercer grupo obtenido por el método el cual fue analizado a través de filtros y otros métodos estadísticos para el manejo de outliers y determinar con más precisión la pertenencia de las estrellas en el cúmulo de las Hyades.

Palabras clave: Cúmulos abiertos, Diagrama Hertzsprung-Russell, Clasificación basada en modelos.

Abstract

In recent years, advances in astrophysics and cosmology have been guided by large and complex data sets, which can only be analyzed and interpreted with the use of highly refined statistical methods. This has caused these disciplines complement each other forming a research field known as astrostatistics. In this paper we provide a classification method called classification based on Gaussian mixture models. This method was used to find stars that belong to the Hyades cluster using 2678 stars sampling from the Hipparcos database. We make a brief description of characteristics of the cluster and we explore the evidence of outliers. With this method it was found that classification yields to three groups of which we can study the membership, and we show the agreement with literature. We also show the Hertzsprung-Russell diagram obtained for the cluster, extremely important for

studies of stellar evolution. Finally, the third group found was analyzed through filters and other statistical methods, for determining the membership of the stars in the Hyades cluster.

Key words: Open Cluster, Hertzsprung-Russell diagram, Model-Based Classification.

1. Introducción

El desarrollo y la aplicación de métodos estadísticos a los problemas de la astronomía viene de hace mucho tiempo. Se tiene evidencia de que Hiparco filósofo natural Griego, hizo una de las primeras aplicaciones de los principios estadísticos en el ámbito de la astronomía, al hacer mediciones de las duraciones entre solsticios para definir el año. En las últimas décadas se ha visto un aumento de interés del uso de la estadística en astronomía, impulsado por la presencia de grandes conjuntos de datos en todos los campos de la astronomía. Por tal motivo, se ha llegado a que estas disciplinas se complementen para formar una rama llamada la astroestadística. La astronomía moderna produce datos que requieren de herramientas estadísticas para ser explorados. La investigación en astronomía ha visto un cambio de paradigma en los últimos años, tratando habitualmente la minería de datos con procesos complejos que exigen un conjunto muy diverso de técnicas estadísticas. Entre algunas técnicas estadísticas aplicadas en astronomía se pueden destacar la estimación de parámetros, muy útil para estimación de parámetros cosmológicos y parámetros orbitales de cuerpos celestes. El análisis multivariado, para estudio cúmulos globulares, estudio de rayos cósmicos y gamma ray burns. Las series de tiempo, de alta relevancia en el estudio de manchas solares y variabilidad de rayos X. Los modelos de mixtura para fotometría galáctica y pertenencia de estrellas, entre otros. Una de las investigaciones en astronomía es la pertenencia de estrellas en cúmulos abiertos. Este estudio es de gran importancia en astronomía para comprender rasgos de la evolución estelar y edad de cúmulos. En este artículo se desarrolla un estudio de pertenencia de estrellas analizando los movimientos propios, centrándonos al cúmulo de las Hyades ubicado en la constelación de Tauro. Usando una muestra de 2678 estrellas tomada del catálogo de Hipparcos, se utiliza el método de mezclas de densidades gaussianas multivariadas para encontrar cuales de estas estrellas pertenecen al cúmulo de las Hyades y de esta forma generar el diagrama Hertzsprung-Russell (H-R) para revelar propiedades muy importantes del mismo. Este artículo se organiza de la siguiente forma: En la sección 2 se comenta acerca del estudio de la pertenencia de estrellas en cúmulos abiertos a partir de movimientos propios y se habla de la importancia del diagrama H-R en el estudio de la astronomía estelar. En la sección 3 se discute el método de clasificación estadística basada en mezcla de gaussianas. En la sección 4 se implementa una aplicación utilizando el conjunto de estrellas mencionadas presentan los resultados, detección de outliers, cómo a través del método de mezcla se analizan las variables de estudio para determinar las posibles estrellas que pertenecen al cúmulo, algunas características de la clasificación, el diagrama H-R,

la construcción de filtros y comparación de resultados. Finalmente en la sección 5 se describen las conclusiones y futuros trabajos alrededor del tema.

2. Pertenencia de estrellas y diagrama Hertzsprung-Russell (H-R)

Los cúmulos abiertos son regiones que contienen de diez hasta centenares de estrellas. La distancias de estos cúmulos pueden ser obtenidos por métodos fotométricos o espectroscópicos. Para cúmulos cercanos como las Hyades se utiliza el método de paralaje cinético, donde se supone que las estrellas que pertenecen al cúmulo tienen la misma velocidad espacial en promedio respecto al sol. Sin embargo, el estudio de la pertenencia de estrellas en cúmulos abiertos ha sido muy compleja (Karttunen H. et al. 2007). A través del estudio de la pertenencia de estrellas en un cúmulo se puede obtener las características de la distribución estelar y la evolución de la galaxia donde se encuentra el cúmulo. Para determinar si una estrella pertenece al cúmulo se utiliza los siguientes métodos: método fotométrico cuya limitación es debida a la absorción interestelar, método de velocidades radiales que tiene dificultad en la medición por efecto Doppler y método de movimientos propios. Este último es muy preciso cuando el cúmulo no se encuentra lejos de nosotros. El movimiento propio de una estrella se define como el cambio angular en la posición de una estrella, respecto a la línea de visión del observador, medida en arco-segundos por año. Es una medida indirecta de la velocidad transversal de la estrella con respecto a la Tierra. Después de saber la pertenencia de las estrellas en el cúmulo, se procede a elaborar el diagrama de Hertzsprung-Russell (diagrama H-R) con estas estrellas y de este diagrama se infiere las propiedades del cúmulo, dinámica y edad. El diagrama H-R¹, es un diagrama estadístico en el que las estrellas están clasificadas en base a su temperatura y luminosidad. El diagrama está hecho sobre un sistema en el que se dispone la temperatura superficial de la estrella sobre el eje horizontal, en sentido decreciente de izquierda a derecha y la luminosidad sobre el eje vertical, en sentido creciente de abajo hacia arriba (ver figura 1). Aquí se observa que la mayor parte de las estrellas están ubicadas sobre una diagonal que cruza el diagrama conocida como secuencia principal. En esta región, se ubica las estrellas más jóvenes (las cuales están quemando Hidrógeno en su núcleo) y en la cual pasan el mayor tiempo de su vida. Las estrellas azules de gran masa y luminosidad se encuentran en la parte superior izquierda. Las estrellas amarillas medianas como el Sol, se encuentran en el centro y las rojas pequeñas están ubicadas en la parte inferior derecha. Además de la secuencia principal, existe una rama de las gigantes rojas ubicadas a la derecha de la secuencia principal que se caracterizan por tener gran tamaño, brillo y baja temperatura superficial. Finalmente las enanas blancas, en la parte inferior del diagrama son de estrellas de baja luminosidad pero de altas temperaturas, y pequeñas comparadas con el Sol.

¹Ideado por E. Hertzsprung y H. N Russell entre 1905 y 1913.

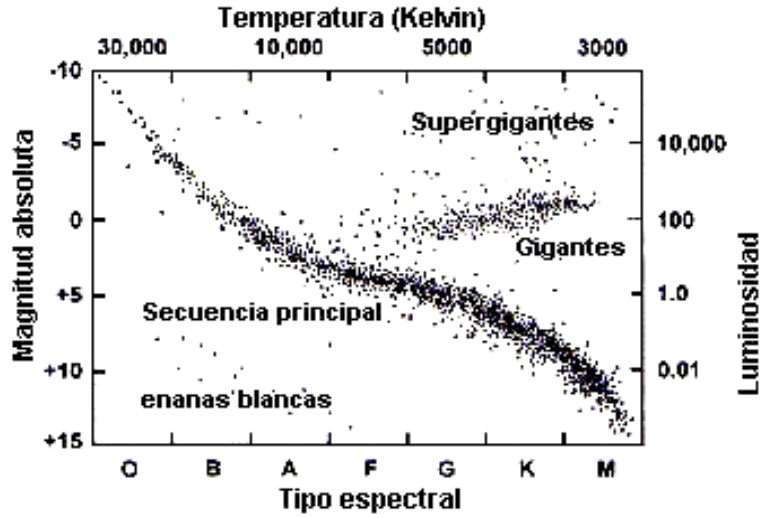


Figura 1: Diagrama H-R

3. Clasificación usando modelos gaussianos

La clasificación de grupos o análisis de conglomerados (cluster analysis) es una metodología que consiste en ubicar objetos, ítems, individuos, etc, dentro de ciertos grupos denominados conglomerados, de tal forma que en cada grupo, los objetos sean semejantes entre sí y, entre grupos, sean diferentes. Existen muchas técnicas de este tipo, en particular las clasificaciones apoyadas en modelos (Fraley et al. 2012). Esta última, considera la agrupación usando modelos normales multivariados y se describe a continuación.

Sea X una variable p -dimensional observada en el conjunto de datos, $\phi(x)$ su función de densidad de la mezcla de normales multivariadas y sea $\{x_i; i = 1, \dots, n\}$ las observaciones de X correspondientes a una muestra aleatoria simple de la población objeto en estudio.

Una clasificación usando modelos, asume que los datos provienen de una función de densidad mixta dada por

$$\phi(x) = \sum_{k=1}^G \tau_k \phi_k(x), \quad (1)$$

donde $\phi_k(x)$ es la función de densidad de las observaciones en el grupo k , τ_k es la probabilidad de que una observación haga parte de la componente k -ésima ($\tau_k \in (0, 1)$ y $\sum_{k=1}^G \tau_k = 1$). Cada componente es usualmente modelada a partir de una función de densidad normal multivariada. Cada componente se caracteriza por un vector de medias μ_k y una matriz de covarianzas Σ_k , cuya función de densidad

viene dada por

$$\phi_k(x_i; \mu_k, \Sigma_k) = (2\pi)^{-\frac{p}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x_i - \mu_k)^t \Sigma_k^{-1} (x_i - \mu_k) \right\}. \quad (2)$$

La matriz de covarianza Σ_k determina las características geométricas tales como forma, volumen, orientación de cada uno de los grupos, a partir de la descomposición matricial de la siguiente manera

$$\Sigma_k = \lambda_k D_k A_k D_k^t,$$

donde D_k , es la matriz ortogonal de vectores propios, A_k es la matriz diagonal cuyos elementos son los valores propios de Σ_k , y λ_k es un valor escalar. La orientación de las componentes principales de Σ_k es determinada por D_k , mientras A_k determina la forma de los contornos de densidad; λ_k especifica el volumen correspondiente al elipsoide, proporcional a $\lambda_k^d \|A\|$, con d la dimensión de los datos. Las características de las distribuciones son usualmente estimadas a partir de los

Tabla 1: Parametrizaciones de la matriz de covarianzas Σ_k

| identifier | Model | Distribution | Volumen | Shape | Orientation |
|------------|---------------------------|--------------|----------|----------|-----------------|
| E | | (univariate) | equal | | |
| V | | (univariate) | variable | | |
| EII | λI | Spherical | equal | equal | NA |
| VII | $\lambda_k I$ | Spherical | variable | equal | NA |
| EEI | λA | Diagonal | equal | equal | coordinate axes |
| VEI | $\lambda_k A$ | Diagonal | variable | equal | coordinate axes |
| EVI | λA_k | Diagonal | equal | variable | coordinate axes |
| VVI | $\lambda_k A_k$ | Diagonal | variable | variable | coordinate axes |
| EEE | $\lambda D A D^t$ | Ellipsoidal | equal | equal | equal |
| EEV | $\lambda D_k A D_k^t$ | Ellipsoidal | equal | equal | variable |
| VEV | $\lambda_k D_k A D_k^t$ | Ellipsoidal | variable | equal | variable |
| VVV | $\lambda_k D_k A_k D_k^t$ | Ellipsoidal | variable | variable | variable |

datos, y pueden variar entre cluster. Todas las parametrizaciones son consideradas en la Tabla 1. Por ejemplo, un modelo EVI denota un modelo en el cual el volumen de todos los conglomerados es igual (E “equal”), la forma de los conglomerados puede variar (V “varying”) y la orientación es idéntica (I “identity”) (Fraley et al. 2012).

La verosimilitud para los datos consiste en asumir que las n observaciones provienen de un modelo de mezclas finitas de G normal multivariadas, es decir

$$\prod_{i=1}^n \sum_{k=1}^G \tau_k \phi_k(x_i; \mu_k, \Sigma_k).$$

Para un numero fijo de componentes G , los parámetros del modelo τ_k , μ_k , y Σ_k pueden ser estimados usando el algoritmo EM (Esperanza y Maximización) (Dasgupta & Raftery 1998, Fraley & Raftery 1998).

Luego de estimar los parámetros de las G componentes, se puede determinar con el teorema de Bayes, la probabilidad posterior de que una observación x pertenezca al k -ésima grupo mediante

$$P(x \in \text{clase } k) = \frac{\tau_k \phi_k(x)}{\sum_{l=1}^G \tau_l \phi_l(x)}, \quad (3)$$

donde G es el número de clases con distribución de probabilidad $\phi_k(\cdot)$ y la proporción de miembros de la población τ_k que están en la k -ésima clase.

4. Aplicación

Inicialmente se realizó una breve descripción de las variables y conjunto de datos a utilizar. Posteriormente se realiza una identificación de estrellas atípicas. Después se utiliza la librería `mclust` creada por Fraley et al. (2012) del paquete estadístico R Core Team (2013) para clasificar las estrellas en diferentes grupos, para luego identificar la secuencia de estrellas que pertenecen al cúmulo de las Hyades. Por último, se caracterizan los resultados estadísticamente y se elabora el diagrama H-R descrito en la sección 2.

4.1. Descripción de los datos

Utilizando 2678 estrellas del catalogo de Hipparcos², bajo el criterio de que el ángulo paraláctico este entre 20° y 25° y el grupo de estrellas esté a una distancia entre 40 y 50 pc. Además, no se tienen en cuenta estrellas que carezcan de información en las variables utilizadas. En la tabla 2 se describen las variables para cada estrella obtenidas a través de la base de datos de Hipparcos.

Tabla 2: Variables a utilizar

| Variable | Descripción |
|----------|---|
| Vmag | Magnitud de banda Visual. |
| RA | Ascension Recta (grados). |
| DE | Declinación (grados). |
| Plx | Ángulo Paraláctico (mas = milliarcseconds). |
| pmRA | Movimiento propio en RA (mas/yr). |
| pmDE | Movimiento propio en DE (mas/yr). |
| e.Plx | Error de medición en Plx (mas). |
| B-V | Color de la estrella (mag). |

De las variables anteriormente mencionadas, solamente se tendrán en cuenta las que están relacionadas con los movimientos propios de las estrellas (pmRA, pmDE). Para el diagrama H-R se tienen en cuenta el color (B-V), magnitud (Vmag)

²<http://heasarc.gsfc.nasa.gov/W3Browse/all/hipparcos.html>

y ángulo paraláctico (Plx). Por último, para procesos de filtros se usaran las coordenadas espaciales de las estrellas (RA, DE).

4.2. Detección de estrellas atípicas

Con los datos descritos anteriormente, se depura la base eliminando aquellas estrellas cuyos movimientos propios no se comportan igual al resto del conjunto de estrellas.

Brieva & Uribe (1985) realiza un proceso de depuración utilizando filtros para una aplicación similar al cúmulo de estrellas NGC654, para detectar estrellas atípicas. También, Fraley & Raftery (2002) sugiere un método alternativo para detectar outliers. Por simpleza se utilizó el procedimiento propuesto por Johnson & Wichern (1998), el cual consiste en calcular la distancia de mahalanobis

$$d_i^2 = (x_i - \bar{x})^t S^{-1} (x_i - \bar{x}) \quad i = 1, 2, \dots, n,$$

luego se compara estos valores con un valor crítico de la tabla de la distribución $F_{(1-\alpha, p, n-p-1)}$, donde p es el número de variables, n el número de observaciones y $\alpha = 1 - (1 - 0.0027)^p$. Para nuestro caso se encontraron 58 estrellas, las cuales se omitieron para este trabajo.

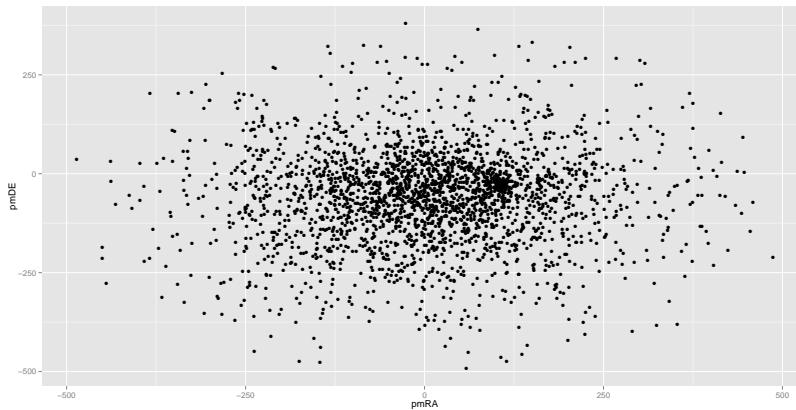


Figura 2: Diagrama de dispersión de los movimientos propios de 2620 estrellas del catalogo de Hipparcos sin observaciones atípicas

En la Figura 2, se observa el diagrama de dispersión de los movimientos propios del catalogo de estrellas sin observaciones atípicas. Nótese que los movimientos propios están muy agrupados en la parte central, razón por la cual no se observa claramente cuantos grupos de estrellas se lograrían obtener. En la tabla 3 se describen los resultados estadísticos de los movimientos propios de este conjunto de estrellas.

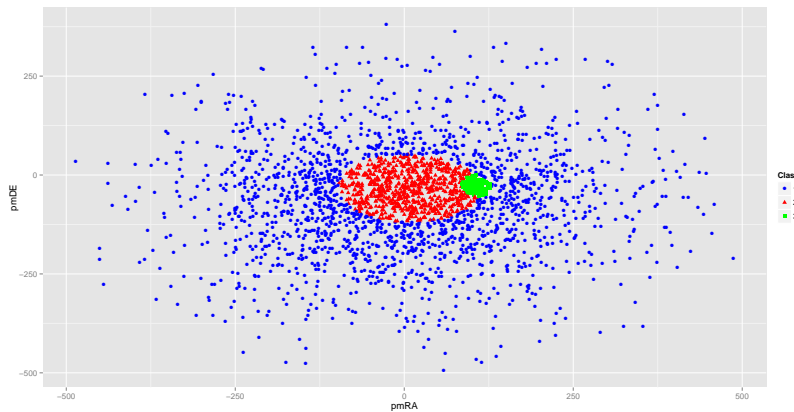


Figura 4: Diagrama de dispersión de los movimientos propios según los grupos de clasificación obtenidos

Con el resultado anterior se puede observar en la Figura 4 como se agrupan las estrellas en los tres grupos según sus movimientos propios.

Los tres grupos tienen distribuciones normales bivariadas totalmente diferentes en volumen y orientación. Por otro lado se observa que las estrellas en el grupo de color azul (clase 1, ●) son las estrellas más dispersas, mientras que las estrellas que se ubican en el grupo de color rojo (clase 2, ▲) presentan menos dispersión. Sin embargo, las estrellas en el grupo de color verde (clase 3, ■) presentan muy poca dispersión con respecto a los dos grupos de estrellas anteriores. Entonces se tiene un grupo de estrellas (clase 3) mucho más compacto en sus movimientos propios.

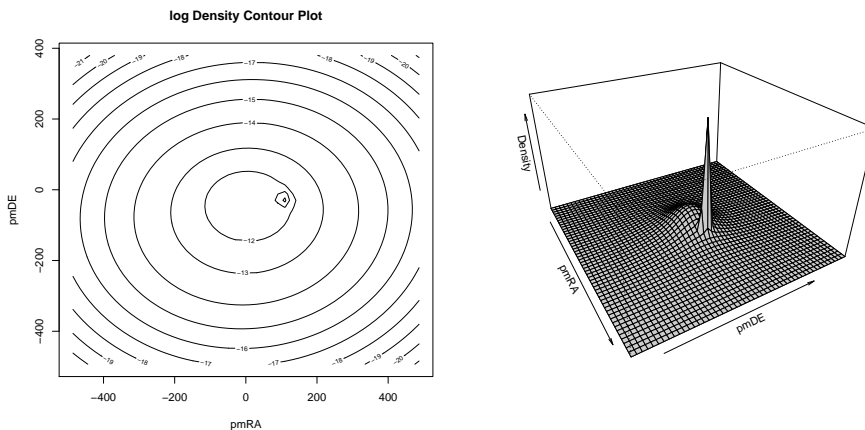


Figura 5: Diagrama de los contornos de la función de densidad y gráfico en 3D de la función de densidad obtenida.

En la Figura 5 se observa la función de densidad de la mezcla de distribuciones normales bivariadas obtenidas. Se observa que la clase 3 es un grupo muy compacto en sus movimientos propios, mientras que los otros grupos tienen una dispersión más alta. A continuación se caracterizan estadísticamente cada uno de los grupos obtenidos.

4.4. Caracterización de los grupos de estrellas obtenidos

Al utilizar este método se clasifican 1770 estrellas en la clase 1, 717 estrellas en la clase 2 y 133 estrellas en la clase 3. Cada clase tienen las siguientes probabilidades $\tau_1 = 0.678$, $\tau_2 = 0.280$ y $\tau_3 = 0.041$. Las distribuciones de ϕ_1 , ϕ_2 y ϕ_3 tienen vectores de medias y matrices de covarianzas dadas por:

| | pmRA | pmDE |
|------------|----------|----------|
| μ_1^t | 1.17 | -68.67 |
| μ_2^t | 6.72 | -40.71 |
| μ_3^t | 105.80 | -26.71 |
| Σ_1 | 29581.58 | 822.32 |
| | 822.32 | 19627.19 |
| Σ_2 | 6157.33 | -86.53 |
| | -86.53 | 4067.98 |
| Σ_3 | 93.95 | -10.36 |
| | -10.36 | 136.39 |

Al observar los resultados anteriores notamos que las covarianzas de la clase 1, son las únicas positivas, mientras que las restantes son negativas. Por otro lado, las covarianzas de la clase 3 son mucho más pequeñas que los otros grupos de estrellas. Al calcular las correlaciones entre los movimientos propios de los grupos se observan que los valores son muy pequeños (0.034, -0.02, -0.09), lo cual corrobora que estos son independientes como se esperaba físicamente.

El diagrama de box-plot de la Figura 6, muestra que el grupo de estrellas de la clase 3 tiene muy poca dispersión. Por otro lado, también observamos que los tres grupos tienen comportamientos muy simétricos.

En la tabla 4 se describen los estadísticos descriptivos de los movimientos propios de cada uno de los grupos obtenidos.

Obsérvese que los coeficientes de asimetría y kurtosis son cercanos a cero, esto nos da entender que los movimientos propios en cada grupo tienden a ser simétricos. El coeficiente de variabilidad resulta ser más alto en el grupo uno, esto indica que los movimientos propios tienen mucha más variación en este grupo. Mientras, que el grupo tres, el coeficiente de variación es mucho más pequeño, indicando una dispersión mínima en este grupo de estrellas.

Se ha encontrado además que los movimientos propios, tienen una menor dispersión en la clase 3. De esta forma se entiende que todas las estrellas en esta clase tienen

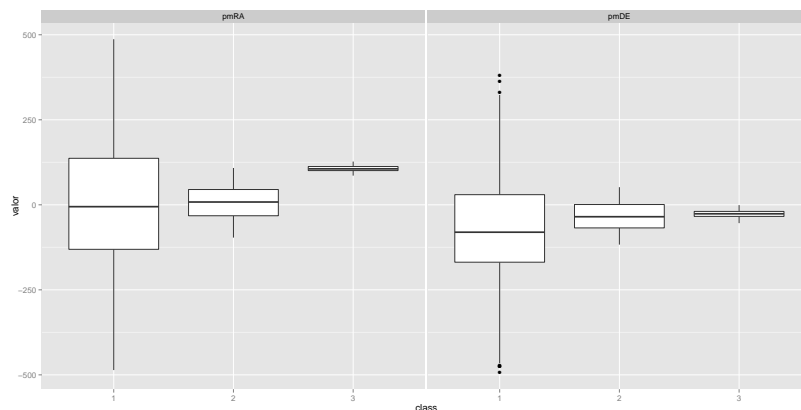


Figura 6: Diagrama de box-plot de los movimientos propios según grupos de clasificación

Tabla 4: Estadísticos de los movimientos propios en Declinación (pmDE) y Ascensión Recta (pmRA)

| Variable: pmDE | | | | | | | | |
|----------------|---------|---------|---------|---------|----------|----------|------|--|
| Grupos | mean | sd | IQR | cv | skewness | kurtosis | n | |
| 1 | -71.452 | 144.708 | 198.578 | 2.025 | 0.102 | -0.071 | 1770 | |
| 2 | -34.477 | 41.400 | 68.770 | 1.201 | 0.020 | -0.987 | 717 | |
| 3 | -27.298 | 11.147 | 14.640 | 0.408 | -0.112 | -0.354 | 133 | |
| Variable: pmRA | | | | | | | | |
| Grupos | mean | sd | IQR | cv | skewness | kurtosis | n | |
| 1 | -0.498 | 174.736 | 267.618 | 351.162 | 0.088 | -0.567 | 1770 | |
| 2 | 7.347 | 48.643 | 77.130 | 6.621 | -0.030 | -0.946 | 717 | |
| 3 | 106.174 | 9.197 | 11.940 | 0.087 | -0.051 | -0.277 | 133 | |

poca variabilidad. Desde el punto de vista estelar, indica que las estrellas de este grupo, pertenecen al cúmulo abierto de las Hyades. Por otra parte, en la clase 1 se encuentra una alta variabilidad en los movimientos propios. Esto indica que cada una de estas estrellas pertenece al background o foreground del cúmulo. Por último, la clase 2 se observa una gran dispersión respecto a la clase 3 pero menor a la clase 1. De esta forma se llega a un resultado importante ya que a través de este grupo se obtiene una especie de datos atípicos que indican un sesgo de estas estrellas a pertenecer o no al cúmulo. Analizando este grupo se encuentra que algunas estrellas pueden pertenecer al cúmulo, pero que debido a sus características que difieren del resto de estrellas, no pudieron ser categorizadas como clase 3, es decir estrellas tales como gigantes, sistemas binarios, entre otros.

4.5. Diagrama H-R

Después de encontrar las estrellas que pertenecen al cúmulo de Hyades usando el método estadístico mencionado anteriormente, se procede a ubicar estas estrellas en el diagrama H-R. El resultado obtenido se muestra en la figura 7. La luminosidad fue calculada usando la expresión dada por

$$\log(L) = (15 - V_{\text{mag}} - 5 \cdot \log_{10}(\text{Plx}))/2.5. \quad (4)$$

En este diagrama se observa que el cúmulo de las Hyades contiene cuatro estrellas del grupo de las gigantes rojas las cuales se encuentran localizadas en la parte superior del diagrama.

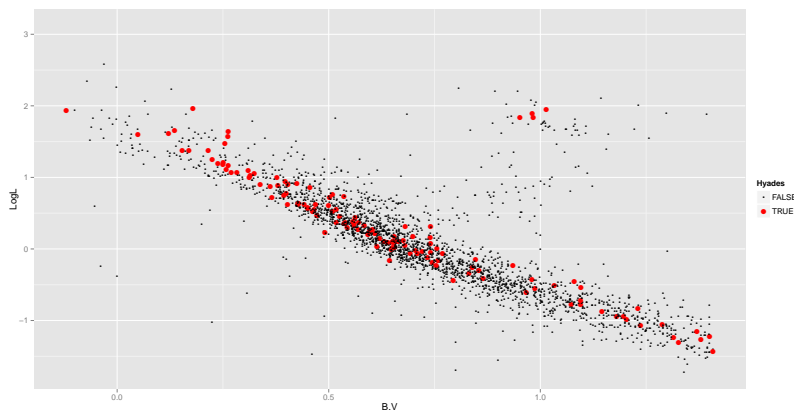


Figura 7: Diagrama H-R obtenido para estrellas pertenecientes al cúmulo de Hyades.

Por otra parte, el cúmulo contiene en su mayoría sus estrellas en la secuencia principal, indicando que este es un cúmulo joven (635 millones de años). En el diagrama se muestra en color rojo las estrellas del grupo tres obtenidas durante

la clasificación y de las cuales se concluyen altamente pertenecientes al cúmulo. Las estrellas mostradas en este grupo concuerdan con los resultados encontrados por Perryman et al. (1998). Para el grupo dos, se realizara un filtro o un análisis estadístico adicional para determinar si algunas estrellas de este grupo, pertenecen al cúmulo de las Hyades. Algunas estrellas de este grupo tienen movimientos propios estadísticamente diferentes respecto al conjunto, debido a su masa o también a que forman sistemas binarios. El grupo restante simplemente experimenta una dispersión grande en sus movimientos propios indicando una gran variabilidad y por tanto no pertenecen al cúmulo.

4.6. Construcción de filtros y comparación

En la Figura 8 se consideran las variables (RA, DE) de las 717 estrellas del grupo 2 y 133 del grupo 3 durante el proceso de clasificación.

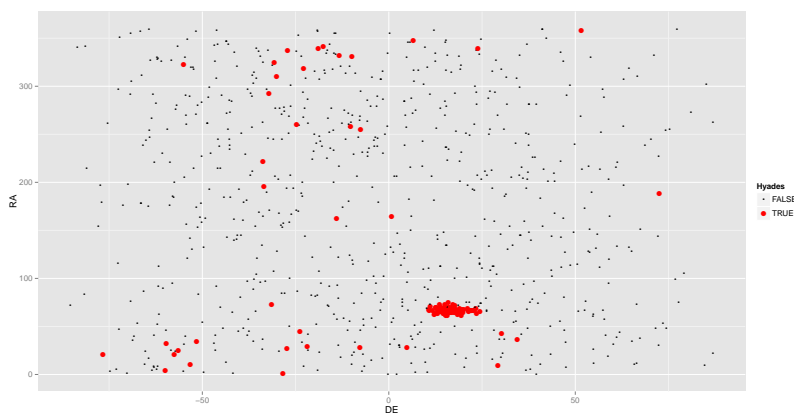


Figura 8: Diagrama de dispersión de las variables (RA, DE) según pertinencia al cúmulo de las Hyades

Se observa la posición donde se encuentra el cúmulo de las Hyades, de esta forma se puede pensar en un filtro para determinar las estrellas en el cúmulo de las Hyades. Para ello se implementa un árbol de clasificación con la función `rpart` de la librería `mvpart` creada por De'ath (2013) del paquete estadístico R Core Team (2013)³. Las variables implementadas en el árbol de clasificación son (RA, DE), donde se determina si la estrella pertenece o no al cúmulo de las Hyades encontradas en el proceso de clasificación.

En la Figura 9 se observa que la gran mayoría de las estrellas del cúmulo de las Hyades se ubican en el nodo 9. Siguiendo el recorrido del árbol se encuentra que $60.54 \leq RA < 72.97$ y $10.46 \leq DE < 22.93$.

³Para la visualización se utiliza la librería `partykit` creada por Hothorn & Zeileis (2013)

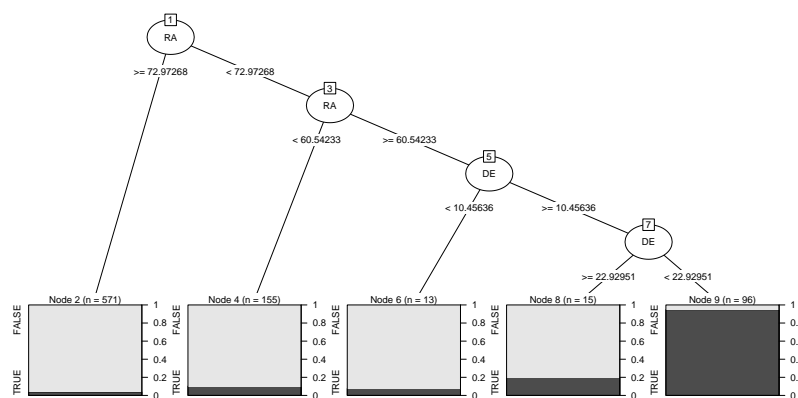


Figura 9: Árbol de clasificación de las variables (RA, DE) según pertinencia al cúmulo de las Hyades

Tabla 5: Matriz de confusión de la clasificación según filtros implementados

| Predicción/Hyades | Falso | Verdadero |
|-------------------|-------|-----------|
| Falso | 712 | 42 |
| Verdadero | 5 | 91 |

En la Tabla 5, se muestra que sólo 5 estrellas que pertenecían al grupo dos pueden ser catalogadas como estrellas del cúmulo de las Hyades. Por otro lado, de las 133 estrellas del cúmulo de Hyades, solo 91 estrellas se encuentran con los filtros implementados. La tasa de mala clasificación es de 5.5 %.

Perryman et al. (1998), realiza un estudio observacional del cúmulo de las Hyades basado en distancias, estructuras, dinámicas y edad de las estrellas pertenecientes a este cúmulo. Para ello implementa la lectura de una muestra de 282 estrellas del catalogo de Hipparcos. Teniendo en cuenta la ecuación (1) de la función de densidad mixta y los parámetros estimados para la clasificación obtenida (ver sección 4.4), se clasifica estas estrellas utilizando la ecuación (3) y los filtros descritos en la sección 4.6, para comparar los resultados. Para ello, se implementa la lectura de las variables anteriormente mencionadas, utilizando el número de la estrella en el catalogo de Hipparcos (HIP)⁴.

En el diagrama H-R mostrado en la Figura 10 se observa 5 grupos. Los cuales se describen a continuación:

- El grupo denominado **FALSE**, son aquellas 54 estrellas que tanto en la propuesta como en el trabajo de Perryman et al. (1998) no se consideran pertenecientes al cúmulo de las Hyades.

⁴Si el lector desea ver los resultados intermedios se recomienda ver el blog Bitácoras en Estadística

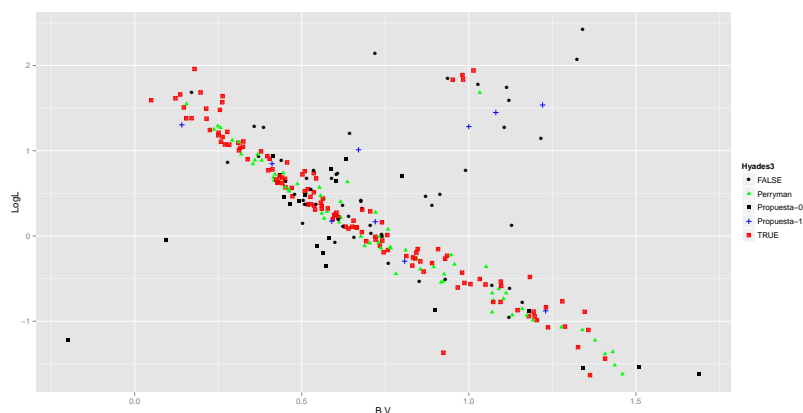


Figura 10: Diagrama H-R obtenido para estrellas pertenecientes al cúmulo de Hyades comparando los resultados obtenidos para el conjunto de Perryman et al. (1998).

- El grupo denominado **Perryman**, son 71 estrellas detectadas por Perryman et al. (1998) las cuales se consideran del cúmulo de Hyades; en nuestro trabajo no se consideran del cúmulo de las Hyades.
- El grupo denominado como **Propuesta-0**, son 21 estrellas las cuales se proponen como falsas; en el trabajo de Perryman et al. (1998) no se catalogaron.
- El grupo denominado como **Propuesta-1**, son 10 estrellas las cuales se proponen pertenecientes al cúmulo de las Hyades; en el trabajo de Perryman et al. (1998) eran falsas.
- El grupo denominado como **TRUE**, son 126 estrellas las cuales se consideran del cúmulo de las Hyades tanto en la propuesta de este trabajo como en el trabajo de Perryman et al. (1998). Este último grupo es el más numeroso, indicando una alta concordancia entre las dos técnicas.

5. Conclusiones

En este artículo se estudió una de las aplicaciones de la estadística en el área de la astronomía, utilizando un método de clasificación usando modelos gaussianos. El objetivo principal del trabajo era encontrar la pertenencia de estrellas al cúmulo de las Hyades analizando el movimiento propio de las estrellas. Los datos fueron tomados de la base de datos de Hipparcos. Usando el método de clasificación se encontró tres grupos en los cuales de acuerdo a la dispersión en los movimientos propios, se catalogó como perteneciente y no perteneciente al cúmulo. El primer grupo contiene 133 estrellas cuya correlación en sus velocidades es muy alta, indicando una alta probabilidad de pertenencia al cúmulo. El segundo grupo contiene

717 estrellas donde la dispersión es más alta, sin embargo, algunas de estas estrellas tienen un movimiento propio similar al primer grupo. Esto indica que los miembros de dicho grupo pueden ser catalogados como outliers, por lo tanto el uso de algunos filtros tales como en la ascensión recta (RA), declinación (DE) y variable $e\text{-Plx}$ deben ser impuestos a este grupo para poder catalogar las estrellas que pueden pertenecer al cúmulo. Para ello, se usó las variables (RA, DE) para la realización de un filtro con el árbol de clasificación con la función `rpart`. Con este filtro se encontró que sólo 5 estrellas que pertenecían al grupo dos pueden ser catalogadas como estrellas del cúmulo de las Hyades. Por otro lado, de las 133 estrellas, solo 91 estrellas pertenecen al cúmulo de las Hyades. Por último, el tercer grupo contiene una gran dispersión en los datos de movimientos propios indicando que los miembros de este grupo no pertenecen al cúmulo. Después de determinar cuales estrellas pertenecen al cúmulo se elaboró el diagrama H-R para estas estrellas encontrando la figura 7. En este gráfico se observa que la mayoría de estas estrellas siguen la secuencia principal (lugar donde se encuentran la mayor parte de su vida), concluyendo que este cúmulo es joven. Se observa algunas estrellas atípicas (outliers) que se ubican fuera de la secuencia principal y que corresponde gigantes rojas. Por otra parte, al comparar los resultados obtenidos junto con los encontrados en la literatura se puede decir que el método de clasificación basada en modelos gaussianos es bastante útil para determinar la pertenencia de estrellas en cúmulos abiertos y se pueden clasificar de forma adecuada datos que sean compactos en sus variables de estudio. Como trabajos futuros se pretende utilizar otro tipo de técnicas de clasificación paramétricas y no paramétricas y comparar los resultados con los obtenidos en este trabajo. También se pretenderá aislar la secuencia principal de Hyades en el diagrama H-R y determinar su ajuste mediante técnicas de regresión no paramétrica.

Agradecimientos

Los autores agradecen al profesor Antonio Uribe de la Universidad Nacional de Colombia y a la profesora Luz Ángela García de la Fundación Universitaria los Libertadores, por sus importantes aportes y comentarios a este trabajo. El trabajo fue elaborado en el semillero de investigación en Astronomía, de la Fundación Universitaria los Libertadores.

Referencias

- Brieva, E. & Uribe, A. (1985), ‘Una aplicación del método de máxima verosimilitud en astronomía galáctica’, *Revista Colombiana de Estadística*.
- Dasgupta, A. & Raftery, A. E. (1998), ‘Detecting features in spatial point processes with clutter via model-based clustering’, *Journal of the American Statistical Association*.

- De'ath, G. (2013), *mvpart: Multivariate partitioning*. rpart by Terry M Therneau and Beth Atkinson. R port of rpart by Brian Ripley ripley@stats.ox.ac.uk. Some routines from vegan – Jari Oksanen jari.oksanen@oulu.fi Extensions and adaptations of rpart to mvpart by Glenn Deáth.
*<http://CRAN.R-project.org/package=mvpart>
- Fraley, C. & Raftery, A. E. (1998), ‘How many clusters? which clustering method? answers via model-based cluster analysis’, *Computer Journal* .
- Fraley, C. & Raftery, A. E. (2002), ‘Model-based Clustering, Discriminant Analysis and Density Estimation’, *Journal of the American Statistical Association* **97**, 611–631.
- Fraley, C., Raftery, A. E., Murphy, T. B. & Scrucca, L. (2012), mclust version 4 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation, Technical report no. 597, Department of Statistics, University of Washington.
- Hothorn, T. & Zeileis, A. (2013), *partykit: A Toolkit for Recursive Partytioning*.
*<http://CRAN.R-project.org/package=partykit>
- Johnson, R. & Wichern, D. (1998), *Applied Multivariate Statistical Analysis*, 4 edn, Prentice Hall.
- Karttunen H., K. P., Oja H., P. M. & Donner, K. J. (2007), *Fundamental Astronomy*, 5 edn, Springer-Verlag GmbH.
- Perryman, M. A. C., Brown, A. G. A., Lebreton, Y., Gómez, A., Turon, C., Cayrel de Strobel, G., Mermilliod, J. C., Robichon, N., Kovalevsky, J. & Crifo, F. (1998), ‘The Hyades: distance, structure, dynamics, and age’, *Astronomy and Astrophysics* **331**, 81–120.
- R Core Team (2013), *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria.
*<http://www.R-project.org/>