

UNIVERSIDADE DE SANTIAGO DE COMPOSTELA
ESCOLA TÉCNICA SUPERIOR DE ENXEÑARÍA



**Plataforma Web para la Validación de Experimentación en
Aprendizaje Automático y Minería de Datos**

TRABAJO DE FIN DE GRADO

Realizado por:
Adrián Canosa Mouzo

Dirigido por:
Ismael Rodríguez Fernández
Alberto J. Bugarín Díz
Manuel Mucientes Molina

Índice general

1. Introducción	1
1.1. Objetivos del proyecto	2
1.2. Organización del documento	2
2. Contraste de hipótesis	3
2.1. Hipótesis nula y alternativa	3
2.2. Estadístico de contraste	5
2.3. Decisiones y tipos de error	6
2.4. Intervalos de confianza	7
2.5. Decisión final y Concepto de p-valor	8
2.6. Etapas en la resolución de un contraste de hipótesis	9
2.7. Tests paramétricos y no paramétricos	10

Índice de figuras

2.1. Funciones de distribución de probabilidad.	5
2.2. Decisiones y tipos de error.	7
2.3. Regiones de aceptación y rechazo.	8
2.4. Punto crítico.	9

Capítulo 1

Introducción

En el contexto tecnológico actual, en donde el Big Data es un recurso cada vez más utilizado, el rol del analista de datos (data scientist) se está convirtiendo en una profesión emergente y de elevada demanda. Un analista de datos es aquel profesional que reúne, analiza e interpreta los datos obtenidos con el objetivo de sacar ciertas conclusiones de ellos y así tomar diferentes decisiones, con las que aumentar la productividad en una organización. Relacionado con el analista de datos, un nuevo rol emergente en muchas empresas / organizaciones es el de CDO (“*Chief Data Officer*”). Este rol es el responsable de la gestión y la utilización de la información como un activo para toda la empresa. El analista de datos combina diferentes habilidades, especialmente las técnicas de la minería de datos y del aprendizaje automático (DM&ML).

Según Mitchell [1], una definición de aprendizaje automático sería la siguiente: Un programa de ordenador aprende a partir de una experiencia E a realizar una tarea T (de acuerdo con una medida de rendimiento P), si su rendimiento al realizar T, medido con P, mejora gracias a la experiencia E. La minería de datos, por otra parte es un campo de las ciencias de la computación referido al proceso que trata de descubrir patrones en grandes volúmenes de conjuntos de datos [2]. Para ello utiliza, entre otros métodos, análisis matemático mediante la estadística para deducir estos patrones y tendencias que existen en los datos. Por lo general, estos patrones no pueden ser detectados mediante exploración tradicional debido a la complejidad o la gran cantidad de datos.

Una de las tareas más importantes que se deben llevar a cabo en el aprendizaje automático es la validación de resultados obtenidos por los algoritmos de aprendizaje. El método estándar más aceptado en la actualidad es el de la aplicación de test estadísticos sobre los experimentos, que, entre otras utilidades, apoyan la toma de decisiones (por ejemplo la elección del algoritmo más adecuado).

Este proyecto se centra en crear y desarrollar una plataforma para asistir al analista de datos en el proceso de validación de resultados. En este proyecto se extenderá una librería de tests estadísticos, se crearan servicios web para facilitar su consulta y se desarrollara una interfaz web que hará uso de estos servicios. El objetivo es que el analista pueda introducir en la web los datos obtenidos mediante experimentación y seleccionar el test estadístico que desee utilizar para que, de forma automática, el sistema muestre los resultados de la aplicación del test. Así, el sistema permitirá de un modo fácil y centralizado la validación de resultados mediante el uso de tests estadísticos.

La herramienta se incorporará en la lista de aplicaciones disponibles a través de la web del CiTIUS para su acceso. El impacto y difusión del resultado del proyecto tiene el potencial de ser amplio, ya que en la actualidad no existe ninguna herramienta que centralice la aplicación de los test estadísticos de mayor utilidad para la validación de aprendizaje automático y que resulte fácil de usar.

1.1. Objetivos del proyecto

El proyecto se centra en crear y desarrollar una plataforma web para asistir al analista de datos en el proceso de validación de resultados obtenidos de diferentes algoritmos de aprendizaje. Para ello, habrá que realizar las siguientes tareas:

1. Completar y extender una librería de test estadísticos, actualmente implementada en el lenguaje Python.
2. Estudiar y conocer las definiciones de los tests estadísticos de uso habitual en aprendizaje automático.
3. Crear los servicios web en Python, basados en REST, que hagan disponible el acceso a los tests estadísticos vía web.
4. Desarrollar una interfaz web (HTML + JavaScript) para facilitar el uso de los tests sobre los datos introducidos por el analista de datos.

1.2. Organización del documento

Sección que incluirá la descripción de los distintos apartados del documento.

Capítulo 2

Contraste de hipótesis

El contraste de hipótesis, o lo que se conoce como tests estadísticos, se engloba en el ámbito de la Inferencia Estadística, que es la parte de la estadística que estudia cómo sacar conclusiones generales (sujetas a un determinado grado de fiabilidad o significancia) para toda la población a partir del estudio de una muestra. En nuestro caso, se tratará de sacar conclusiones de los resultados obtenidos por diferentes algoritmos (muestra) para determinar, por ejemplo, si los algoritmos tienen un rendimiento significativamente diferente y por lo tanto no se pueden considerar iguales (población).

El **contraste de hipótesis** es uno de los problemas más comunes dentro de la inferencia estadística. En él se contrasta una hipótesis estadística. Por ejemplo:

Un ingeniero de software afirma que la media de los resultados obtenidos por un algoritmo de aprendizaje automático es 10. ¿Se podría desmentir la afirmación del ingeniero?

El planteamiento del contraste sería el siguiente (μ indica media poblacional):

$$\begin{aligned}\mu &= 10 \\ \mu &\neq 10\end{aligned}$$

Para tomar una decisión (desmentir o no la afirmación), hay que basarse en los datos de una muestra, para comprobar si en efecto la media de los resultados es 10 (media muestral). Para ello, podría establecer una regla de decisión sobre la cual se basaría nuestra decisión final. Por ejemplo: si la media obtenida está próxima a la afirmada por el ingeniero (10), entonces se podría afirmar que dice la verdad. Si por el contrario la muestra nos proporciona una media muy distinta a 10, entonces se puede afirmar que la evidencia desmiente la afirmación del ingeniero sobre el algoritmo en cuestión. Esto supone un problema y es el hecho de cuándo considerar que la media es lo suficientemente distinta como para determinar que la afirmación del ingeniero es errónea. Por ejemplo si la media de la muestra es 8.5, ¿se podría desmentir la afirmación inicial? El contraste de hipótesis nos proporciona una forma de establecer este criterio y poder rechazar o aceptar la afirmación inicial.

2.1. Hipótesis nula y alternativa

En todo contraste de hipótesis siempre se dan dos posibilidades o hipótesis, las cuales se representan con los siguientes símbolos:

$$\begin{aligned}H_0 &: \text{Hipótesis nula} \\ H_1 &: \text{Hipótesis alternativa}\end{aligned}$$

- H_0 : Es la hipótesis que se supone cierta de partida, es decir, es la hipótesis que establece que lo que indica la muestra es solamente debido a la variación aleatoria entre la muestra y la población.
- H_1 : Es la hipótesis alternativa y es la que reemplazará a la hipótesis nula si ésta es rechazada. H_1 establece que lo que indica la muestra es verdadero, ya que representa a toda la población.

A modo de ejemplo, supongamos que unos programadores están trabajando en la optimización de un algoritmo de aprendizaje. El objetivo es mejorar el algoritmo de forma que los resultados que proporcione sean menores de 100. Se toma una muestra de los resultados obtenidos por el nuevo algoritmo optimizado y se observa que la media de la muestra es de 92. Si no hubiera incertidumbre en la media muestral, entonces se podría concluir que la modificación reduciría los resultados a 92. Sin embargo, siempre existe incertidumbre en la media muestral. La media poblacional en realidad será poco mayor o menor a 92.

Los programadores están preocupados de que el nuevo algoritmo en realidad no mejore al anterior, es decir, que la media poblacional pudiera ser mayor o igual a 100. Quieren saber si esta preocupación está justificada. Se ha observado una muestra con media de 92 y existen dos posibles interpretaciones o, como se ha mencionado más arriba dos tipos de hipótesis que serán contrastadas más adelante mediante un determinado test estadístico:

1. La media poblacional es mayor o igual a 100 (la media muestral es, por tanto, menor debido sólo a la variación aleatoria de la media poblacional). El nuevo algoritmo no mejorará al anterior.
2. La media poblacional es menor que 100, y la media muestral lo refleja. El nuevo algoritmo sí mejorará al anterior.

La primera interpretación sería la hipótesis nula o H_0 . La segunda, la hipótesis alternativa o H_1 , como bien se comentó más arriba.

En este caso, los programadores están preocupados de que la hipótesis nula sea cierta. Un test estadístico o prueba de hipótesis hallará una medida cuantitativa de la factibilidad de la hipótesis nula (denominado estadístico de contraste, que para este ejemplo viene dado por la media obtenida en la muestra) y se podrá decir a los programadores (después de que el test tome la decisión) si su preocupación está o no justificada. Por tanto, a modo de resumen este ejemplo nos proporciona dos hipótesis:

$$H_0 : \mu \geq 100 \text{ vs. } H_1 : \mu < 100$$

La realización de un contraste de hipótesis no consiste en decidir cuál de las dos hipótesis (H_0 , H_1) es más creíble, sino en decidir si la muestra proporciona o no suficiente evidencia para descartar H_0 . Para realizar la prueba de hipótesis o test estadístico se pone la hipótesis nula en juicio, es decir se empieza suponiendo que H_0 es verdadera. Se podría poner como analogía el supuesto de “*En un juicio, el acusado siempre es inocente hasta que se demuestre lo contrario.*” Esto es:

H_0 : El acusado es inocente

H_1 : El acusado es culpable

y, mientras no se tenga suficiente evidencia para aceptar H_1 , hay que creer que lo que dice H_0 es cierto. La muestra aleatoria proporcionará la evidencia. Si el juicio (test o prueba de hipótesis) determina que el acusado es inocente, sólo se puede decir que no se tiene suficiente evidencia para asegurar que el acusado es culpable, mientras que si aceptamos la hipótesis alternativa, se estará bastante seguro de que el acusado sí es culpable.

2.2. Estadístico de contraste

Los tests estadísticos o pruebas de hipótesis, calculan internamente una medida cuantitativa que proporciona la factibilidad de la hipótesis nula. Esta medición se extrae de la muestra proporcionada. Por ejemplo, si queremos contrastar la hipótesis de que la media poblacional es 5, un estadístico a calcular puede ser la media de una muestra. En este caso, la muestra viene determinada por los resultados obtenidos por los algoritmos y cada uno de los tests tiene una forma particular de hallar este estadístico mediante una fórmula que lo caracteriza. Estos estadísticos siguen una determinada distribución de probabilidad. Por ejemplo en esta plataforma web, los tests implementados harán uso de estadísticos que siguen distribuciones como:

- Distribución normal (p. ej. test de Wilcoxon).
- Distribución chi cuadrado χ^2 (p. ej. test de Friedman).
- Distribución f de Fisher-Snedecor (p. ej. test de Iman-Davenport).
- Distribución t de Student (p. ej. t-test).

La figura 2.1, nos muestra el aspecto que presentan las distintas distribuciones de probabilidad. Las distribuciones dependen de ciertos parámetros para determinar su forma ($\mu, \sigma^2 \dots$):

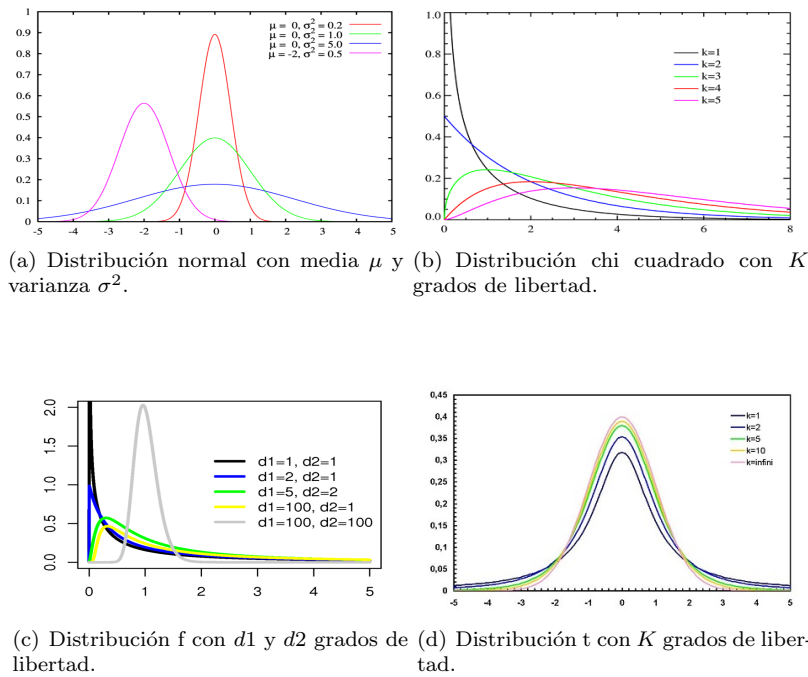


Figura 2.1: Funciones de distribución de probabilidad.

Como podemos ver en la figura 2.1, la distribución normal presenta μ y σ^2 como parámetros. Éstos indican media y varianza respectivamente. La varianza, es una medida de dispersión que indica cómo se distribuye la población. Por ejemplo: en una distribución normal de media 0 y varianza 1, que es la línea roja en la figura 2.1(a), aproximadamente el 68% de la población se encuentra en el intervalo $[-1, 1]$, ya que el área bajo la curva es de 0.68. Por tanto, la probabilidad de que un individuo de la población

se encuentre en ese intervalo es del 68 %. Si un estadístico sigue una distribución normal con media μ y varianza σ^2 , se expresa como:

$$\text{Estadístico} \sim N(\mu, \sigma^2)$$

En las distribuciones chi cuadrado y t de Student se habla del parámetro K o grados de libertad ($d1$ y $d2$ en la distribución f de Fisher-Snedecor). La media y la varianza de estas tres distribuciones vendrán determinadas por el parámetro K . Cuando se habla de grados de libertad se está refiriendo al número de valores que se pueden elegir libremente en una muestra. Por ejemplo: una muestra con dos datos y media 5 si el primer dato toma el valor 4 entonces necesariamente el segundo dato debe de ser 6 (para lograr la media de 5). En este caso, se tienen:

$$N - 1 \text{ grados de libertad, donde } N \text{ es el tamaño de la muestra.}$$

Se hallan con la fórmula $N - R$, donde N es el número de individuos en la muestra cuyo valor puede ser elegido de forma libre y R es el número de sujetos cuyo valor dependerá del valor que tengan los individuos de la muestra que son libres. También se puede representar por $K - R$, donde K es el número de grupos (cuando intervienen grupos y no sujetos individuales).

En nuestro caso, N viene determinado por el número de resultados obtenidos por los algoritmos (número de filas de la matriz de la muestra de datos) y K por el número de algoritmos o variables relacionadas que tiene la muestra de datos con la que se están aplicando los tests (número de columnas de la matriz). Cada test que use el parámetro de grados de libertad lo calcula de acuerdo a su fórmula característica para el estadístico.

Todas las distribuciones de la figura 2.1 son continuas, pues se puede tomar cualquier valor dentro de un intervalo, a diferencia de las distribuciones discretas. Por otra parte, en la distribución t de Student a medida que aumentan los grados de libertad se tiende más a una distribución normal estandarizada (de $\mu = 0$ y $\sigma^2 = 1$).

Las distribuciones de probabilidad que puedan seguir un estadístico nos dan un valor diferente de probabilidad para cada valor diferente del estadístico. Este valor de probabilidad indica cómo de probable que es obtener ese valor del estadístico siendo la hipótesis nula cierta. Por ejemplo, si es cierta la hipótesis nula de que la media de una población es 5, es más probable que obtengamos una media de una muestra igual a 4.5 que a 3.

2.3. Decisiones y tipos de error

Cuando se lleva a cabo un contraste de hipótesis sólo se pueden tomar dos decisiones. Los datos de la muestra, que en este proyecto vendrá dada por los resultados obtenidos por los algoritmos, evidenciarán qué decisión se debe tomar:

1. Aceptar la hipótesis nula (H_0) (Rechazar la hipótesis alternativa H_1)
2. Rechazar H_0 (Aceptar la hipótesis alternativa)

Sin embargo, cuando se toma la decisión se pueden cometer dos tipos de error. La figura 2.2, nos muestra las decisiones y los dos tipos de errores que se pueden cometer:

		Decisión	
		No se rechaza H_0	Se rechaza H_0
Realidad	H_0 es verdadera	Decisión correcta	Error de tipo I
	H_0 es falsa	Error de tipo II	Decisión correcta

Figura 2.2: Decisiones y tipos de error.

La probabilidad de “Error tipo I” se denota por α y se denomina nivel de significación:

$$P(\text{“Error tipo I”}) = P(\text{Rechazar } H_0 | H_0 \text{ es cierta}) = \alpha$$

El nivel de significación consiste en la probabilidad de rechazar la hipótesis nula H_0 cuando verdaderamente es cierta. Este valor α es un parámetro que debe seleccionar la persona que quiere realizar un test estadístico en base a cómo de importante considere rechazar H_0 cuando es cierta. Normalmente es del 5 %, lo que implicará que 5 de cada 100 veces se acepta la hipótesis alternativa cuando la cierta es la hipótesis nula. Cuanto menor sea el nivel de significación, cada vez es más difícil rechazar la hipótesis nula. Es decir, si queremos equivocarnos menos veces, necesitamos mucha más evidencia para justificar el rechazo. Si es grande es más fácil aceptar la hipótesis alternativa cuando en realidad es falsa.

Por otra parte, la probabilidad de “Error tipo II” se denota por β :

$$P(\text{“Error tipo II”}) = P(\text{Aceptar } H_0 | H_0 \text{ es falsa}) = \beta$$

Este error β consiste en la probabilidad de aceptar la hipótesis nula H_0 cuando verdaderamente es falsa. Por último, cabe destacar el concepto de “Potencia”.

$$P(\text{“Potencia”}) = P(\text{Rechazar } H_0 | H_0 \text{ es falsa}) = 1 - \beta.$$

La potencia es la probabilidad de detectar que una hipótesis es falsa. Los tests estadísticos o pruebas de hipótesis implementados en el presente proyecto se caracterizan por su potencia, siendo esta fija, y dejando como parámetro libre el nivel de significación. Así, cuanto mayor es el nivel de potencia, mejor será el test, ya que se rechazarán más hipótesis nulas cuando se deben rechazar (mayor habilidad en aceptar correctamente hipótesis alternativas).

En este proyecto se pondrá el énfasis en el nivel de significación, ya que es la hipótesis alternativa la que se quiere probar y no se quiere aceptar si en realidad no es cierta, es decir, si aceptamos la hipótesis alternativa queremos equivocarnos con un margen de error muy pequeño. Obviamente, lo ideal sería que tanto α como β fuesen nulos y que no se cometiese ningún error, o que ambos valores fuesen muy pequeños. Como no se pueden disminuir ambos errores a la vez, se controla el “Error tipo I”.

2.4. Intervalos de confianza

El nivel de significación fijado divide en dos regiones el conjunto de posibles valores del estadístico de contraste: la región de aceptación y la región de rechazo o región crítica. Se denomina región de aceptación a la región que conduce a la aceptación de H_0 y región de rechazo a la región que conduce al rechazo de H_0 en favor de H_1 . Aquí surge el concepto de “Cola”, que indica la porción o porciones de una distribución de probabilidad en la cual se rechaza la hipótesis nula.

La determinación de las regiones de aceptación o de rechazo depende de cómo se establezca la hipótesis alternativa H_1 . Por ejemplo, si hablamos de un contraste en el que se esté contrastando una determinada media (μ_0) se podría establecer como H_1 que la media en realidad sea menor, mayor o distinta a (μ_0):

- Media menor (test unilateral con cola a la izquierda):

$$\begin{aligned} H_0 : \mu &= \mu_0 \\ H_1 : \mu &< \mu_0 \end{aligned}$$

- Media mayor (test unilateral con cola a la derecha):

$$\begin{aligned} H_0 : \mu &= \mu_0 \\ H_1 : \mu &> \mu_0 \end{aligned}$$

- Media distinta (test bilateral o de dos colas):

$$\begin{aligned} H_0 : \mu &= \mu_0 \\ H_1 : \mu &\neq \mu_0 \end{aligned}$$

En la figura 2.3, podemos ver cómo quedarían establecidos los intervalos para el ejemplo:

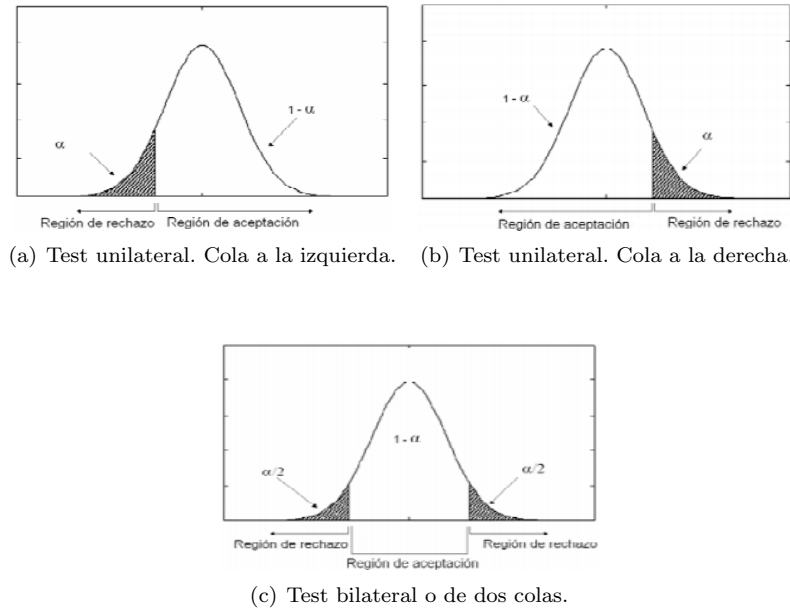


Figura 2.3: Regiones de aceptación y rechazo.

2.5. Decisión final y Concepto de p-valor

Si el valor del estadístico cae en la región de aceptación, se acepta la hipótesis nula, ya que no existen razones suficientes para rechazar H_0 con el nivel de significación dado. Por tanto, en este caso se diría que el contraste es estadísticamente no significativo, es decir, no existe evidencia estadísticamente significativa en favor de H_1 . La figura 2.4 nos muestra el punto crítico: si el estadístico obtenido por el test o prueba de hipótesis es 5 y el punto crítico es 4.5, se rechaza H_0 ya que el estadístico pertenece a la región de rechazo.

La decisión de rechazar o aceptar la hipótesis nula, se puede determinar también mediante el p-valor, que es el parámetro utilizado para realizar los tests estadísticos en este proyecto. El p-valor proporciona



Figura 2.4: Punto crítico.

un forma más eficiente de determinar si el contraste es o no estadísticamente significativo, ya que no sería necesario recalcular regiones de aceptación y rechazo cada vez que el usuario de los tests cambia de nivel de significación.

2.6. Etapas en la resolución de un contraste de hipótesis

En un contraste de hipótesis siempre se siguen una serie de pasos definidos. Como se ha ido viendo a lo largo del capítulo, los pasos para la realización de una prueba de hipótesis o test estadístico son las siguientes:

1. Especificación de la hipótesis nula H_0 y de la hipótesis alternativa H_1 .
2. Suponer que H_0 es verdadera (el test sirve para que a partir de la muestra de datos podamos rechazar H_0 en beneficio de H_1).
3. Calcular un estadístico de prueba o estadístico de contraste. Este estadístico se usa para evaluar la fuerza de la evidencia en contra de H_0 (medir la discrepancia entre la hipótesis y la muestra).
4. Establecer un nivel de significación α en base a cómo de importante se considere rechazar H_0 cuando realmente es verdadera.
5. El nivel de significación fijado divide en dos regiones el conjunto de posibles valores del estadístico de contraste: la región de aceptación y la región de rechazo o región crítica.
6. Si el valor del estadístico cae en la región de rechazo, se rechaza la hipótesis nula, ya que esto evidencia que los datos obtenidos de la muestra no son compatibles con H_0 . Por tanto, en este caso se diría que el contraste es estadísticamente significativo, es decir, existe evidencia estadísticamente significativa en favor de H_1 .
7. Si el valor del estadístico cae en la región de aceptación, se acepta la hipótesis nula, ya que no existen razones suficientes para rechazar H_0 con el nivel de significación dado. Por tanto, en este caso se diría que el contraste es estadísticamente no significativo, es decir, no existe evidencia estadísticamente significativa en favor de H_1 .
8. La decisión de rechazar o aceptar la hipótesis nula, se puede determinar también mediante el p-valor, que es el parámetro utilizado para realizar los tests estadísticos en este proyecto. El p-valor proporciona un forma más eficiente de determinar si el contraste es o no estadísticamente significativo, ya que no sería necesario recalcular regiones de aceptación y rechazo cada vez que el usuario de los tests cambia de nivel de significación.

2.7. Tests paramétricos y no paramétricos

Bibliografía

- [1] Tom M. Mitchell, *Machine Learning*, McGraw Hill, 1997.
- [2] Oded Maimon and Lior Rokach, *Data Mining and Knowledge Discovery Handbook*, Springer, New York, 2010.