

# Práctica 2

## Ejercicio 1

**Descripción del dataset. ¿Por que es importante y que pregunta/problema pretende responder?**

Este dataset está formado 27 variables y 1260 observaciones. Estas variables son:

1. Timestamp: momento de presentación de respuestas
2. Age: edad
3. Gender: género
4. Country: país
5. state: estado. ¿Si vives en los Estados Unidos, cual es el estado o el territorio dónde vives?
6. self\_employed: auto-empleado. ¿Es autónomo (auto-empleado)?
7. family\_history: historia familiar. ¿Tiene antecedentes de enfermedad mental en la familia?
8. treatment: tratamiento. ¿Ha sido tratado por una enfermedad mental?
9. work\_interfere: ¿Si tiene una enfermedad mental, siente que interfiere con su trabajo?
10. no\_employees: ¿número de empleados. ¿Cuántos empleados tiene su compañía u organización?
11. remote\_work: ¿realiza teletrabajo (fuera de la oficina) al menos el 50% del tiempo?
12. tech\_company: ¿su empleador primario es una organización o empresa de tecnología?
13. benefits: ¿su empleador provee beneficios de salud mental?
14. care\_options: ¿conoce las opciones de cuidado mental de la compañía médica que el empleador provee?
15. wellness\_program: ¿Su empleador ha mencionado alguna vez que tiene un programa de bienestar mental para sus empleados?
16. seek\_help: ¿Su empleador proporciona recursos para saber más sobre aspectos de salud mental y cómo encontrar ayuda?
17. anonymity: ¿Está protegida su privacidad si elige acogerse a ventajas de salud mental o recursos de tratamiento de abusos de sustancias?
18. leave: ¿Le sería fácil, acogerse a una baja por situación de salud mental?
19. mental\_health\_consequence: ¿Cree que hablar de un aspecto de salud mental con su empleador, tendría consecuencias negativas?
20. phys\_health\_consequence: ¿Cree que hablar de un aspecto de salud física con su empleador, tendría consecuencias negativas?
21. coworkers: ¿Estaría dispuesto a hablar con sus compañeros de un aspecto de salud mental?
22. supervisor: ¿Estaría dispuesto a hablar con sus supervisores de un aspecto de salud mental?
23. mental\_health\_interview: ¿Mencionaría un aspecto de salud mental con un potencial empleador en una entrevista?

24. phys\_health\_interview: ¿Mencionaría un aspecto de salud física con un potencial empleador en una entrevista?
25. mental\_vs\_physical: ¿Siente que su empleador se toma la salud mental como un aspecto importante de la salud?
26. obs\_consequence: ¿Ha oído u observado consecuencias negativas para sus compañeros de trabajo que se encuentren en situación de enfermedad mental en su puesto de trabajo?
27. comments: comentarios adicionales

El dataset, de 2014 (facilitado por Open Sourcing Mental Illness), procede de una encuesta que mide las actitudes sobre salud mental y la frecuencia de desórdenes mentales en puestos de trabajo extraído en un contexto de Tecnologías de Información. Es de especial interés dado que aspectos como el uso intensivo de las tecnologías de la información está dando lugar también a nuevas enfermedades, también de tipo mental, como así se pone de manifiesto en la literatura (ver, por ejemplo: Gentile, D., Coyne, S., & Bricolo, F.(2012-12-31). Pathological Technology Addictions: What Is Scientifically Known and What Remains to Be Learned. In The Oxford Handbook of Media Psychology: Oxford University Press).

Para realizar un trabajo de forma correcta, el trabajador debe estar en situación de condiciones mentales normales.

Teniendo en cuenta que la Organización Mundial de la Salud informa que la salud mental no más que una actitud de bienestar para que la persona sea capaz de desarrollar sus capacidades, de afrontar el estrés del día a día, que en su trabajo se observe una productividad y que sea capaz de aportar a la comunidad. Luego mirándolo de forma positiva, la salud mental es el pilar de un funcionamiento correcto tanto a nivel individual como a nivel comunidad.

No hay que olvidar que durante nuestro día a día nos encontramos con diferentes situaciones tanto a nivel persona como laboral que nos provocan estrés, esto está dentro de unos baremos de la normalidad y en ningún caso debe considerarse como un problema a tratar.

El hecho de sentir estrés no es malo, siempre y cuando sea en unas cantidades que nos permitan en todo momento tener un nivel de sensatez mental adecuado y un positivo rendimiento a nivel de conducta como cognitivo. Se afirma que el estrés agudo, de poca duración, pone en predisposición el cerebro para un mejor rendimiento.

Si lugar a dudas el estrés lleva a las personas a tener problemas de salud, relaciones insuficientes y una baja productividad laboral. Con lo que conlleva aspectos negativos tanto personalmente como profesionalmente. Visiblemente esto se observa con facilidad ya que el individuo se enfada constantemente con los que están más cerca.

Solamente, en la Unión Europea, las enfermedades relacionadas con los músculos del esqueleto superan al estrés laboral.

Una persona con estrés tiene los siguientes síntomas fatiga, tensión muscular, variación en el apetito, bruxismo, cambios en el estímulo sexual, mareos y dolores de

cabeza. Psicológicamente estos factores pueden ser la irritabilidad, nerviosismo, falta de energía y ganas de llorar.

La cuestión que podemos llegar a responder es si el trato es el mismo laboralmente en la enfermedad física que en la enfermedad mental.

Pretendemos por tanto con ello responder a la siguiente pregunta/problema: ¿se trata de igual modo en el contexto laboral a las enfermedades físicas y mentales?

Por las variables existentes en el conjunto de datos y a partir de estas preguntas previas, deducimos que hay dos aspectos que se podrían tratar: la existencia de enfermedad mental, y las actitudes hacia ésta por las personas en el puesto de trabajo. Nosotros nos centraremos en el segundo aspecto, buscando respuestas en cuanto al trato (o consideración) de igualdad (o no) entre enfermedades físicas y enfermedades mentales.

Los datos corresponden a la encuesta realizada durante el 2014. La licencia que tiene toda esta información es Creative Commons Attribution-ShareAlike 3.0 Unported License.

No debemos olvidar que todo proyecto analítico en ciencia de datos tiene las siguientes fases:

1. Se trata de encontrar la cuestión que deseamos resolver.
2. Consiste en la recogida y almacenamiento de los datos. Conocer de dónde se han extraído los datos y el formato de almacenamiento.
3. Limpieza de datos. Los datos son preparados para el análisis. Para ello es muy posible que se produzca eliminaciones, transformaciones, etc.
4. En esta etapa se produce el estudio de los datos y un aprendizaje de forma automática.
5. Aquí nos encontramos con el estudio de establecer la forma visual más eficiente para la representación de los datos.
6. Resolvemos la cuestión que se planeó en la primera fase del proyecto.

Sin olvidar la peculiaridad y necesidades de cada proyecto, no todos tienen que llevar a cabo las 6 fases anteriormente nombradas de manera estricta y única. A veces es necesario que alguna fase se repita de manera iterativa.

## Ejercicio 2

### Limpieza de los datos.

#### 2.1 Selección de los datos de interés a analizar. ¿Cuáles son los campos más relevantes para responder al problema?

Dado que el dataset, contiene las siguientes variables y preguntas asociadas a su explicación:

1. Timestamp: momento de presentación de respuestas
2. Age: edad

3. Gender: género
4. Country: país
5. state: estado. ¿Si vives en los Estados Unidos, cual es el estado o el territorio donde vives?
6. self\_employed: auto-empleado. ¿Es autónomo (auto-empleado)?
7. family\_history: historia familiar. ¿Tiene antecedentes de enfermedad mental en la familia?
8. treatment: tratamiento. ¿Ha sido tratado por una enfermedad mental?
9. work\_interfere: ¿Si tiene una enfermedad mental, siente que interfiere con su trabajo?
10. no\_employees: número de empleados. ¿Cuántos empleados tiene su compañía o organización?
11. remote\_work: ¿realiza teletrabajo (fuera de la oficina) al menos el 50% del tiempo?
12. tech\_company: ¿su empleador primario es una organización o empresa de tecnología?
13. benefits: ¿su empleador provee beneficios de salud mental?
14. care\_options: ¿conoce las opciones de cuidado mental de la compañía médica que el empleador provee?
15. wellness\_program: ¿Su empleador ha mencionado alguna vez que tiene un programa de bienestar mental para sus empleados?
16. seek\_help: ¿Su empleador proporciona recursos para saber más sobre aspectos de salud mental y cómo encontrar ayuda?
17. anonymity: ¿Está protegida su privacidad si elige acogerse a ventajas de salud mental o recursos de tratamiento de abusos de sustancias?
18. leave: ¿Le sería fácil, acogerse a una baja por situación de salud mental?
19. mental\_health\_consequence: ¿Cree que hablar de un aspecto de salud mental con su empleador, tendría consecuencias negativas?
20. phys\_health\_consequence: ¿Cree que hablar de un aspecto de salud física con su empleador, tendría consecuencias negativas?
21. coworkers: ¿Estaría dispuesto a hablar con sus compañeros de una aspecto de salud mental?

22. supervisor: ¿Estaría dispuesto a hablar con sus supervisores de un aspecto de salud mental?
23. mental\_health\_interview: ¿Mencionaría un aspecto de salud mental con un potencial empleador en una entrevista?
24. phys\_health\_interview: ¿Mencionaría un aspecto de salud física con un potencial empleador en una entrevista?
25. mental\_vs\_physical: ¿Siente que su empleador se toma la salud mental como un aspecto importante de la salud?
26. obs\_consequence: ¿Ha oído u observado consecuencias negativas para sus compañeros de trabajo que se encuentren en situación de enfermedad mental en su puesto de trabajo?
27. comments: comentarios adicionales

De estas variables, dado que algunas de ellas no son directamente asociadas al objetivo de nuestro trabajo, debido a las razones previamente expuestas, prescindimos de las siguientes 3 variables:

1. Timestamp
5. state
6. comments

Así pues, tenemos 1259 observaciones y 27 variables de las que nos quedamos con 24 variables que, a priori, podrían ser útiles para nosotros.

Los campos más importantes para resolver el problema serían

1. Age
2. treatment
3. mental\_vs\_physical

En este dataset nos encontramos con un conjunto de variables que son cuantitativas y cualitativas.

Las cualitativas son las que tienen su origen en características o categorías. Mientras que la variable cuantitativa hace referencia a un valor de naturaleza numérica, estas pueden ser discretas (corresponden a un valor numérico entero) y continuos (toman cualquier valor existente en un intervalo).

La forma de analizar estos datos es diferente, la primera de ella es la ordenación, un dato cualitativo no puede ordenarse de manera numérica.

Para obtener información de datos cualitativos partimos de distribuciones de frecuencias, en la cual podemos observar el número de veces que sucede una categoría o nivel de la variable cualitativa.

En variables cuantitativas la distribución de frecuencia nos proporciona una zona visible más espesa donde se establecen el mayor número de observaciones y una zona más liviana donde nos encontramos con muy pocas observaciones.

En el dataset que nos ocupa la única variable cuantitativa discreta es Age el resto son variables cualitativas.

## **2.2 ¿Los datos contienen ceros o elementos vacíos? ¿Y valores extremos? ¿Cómo gestionarás cada uno de estos casos?**

Tenemos 1259 observaciones. Vemos que, en primer lugar, hay niveles y valores inadecuados en algunas de las variables. Es necesario estandarizarlos.

Cuando hablamos de un dato cero tenemos siempre en mente una asociación a un valor numérico. No hay que olvidar que si el dato es de carácter numérico el valor cero es el que mejor se adapta.

Un dato vacío existe cuando se carece de observación. Este es de utilidad cuando nos encontramos con cadena de caracteres, si añadimos un espacio en blanco el dato pierde el carácter de vacío.

En el momento de la lectura del fichero hemos especificado `na.strings = "NA"` con lo cual cualquier elemento vacío ha sido rellenado con "NA".

Comprobamos que variables tienen datos perdidos. Las únicas variables que contienen valores vacíos son `self_employed` con un total de 18 (Valor TRUE) y `work_Interfere` con 254 (Valor TRUE). Como estas variables no son de interés para nuestro estudio hemos decidido no eliminar las observaciones con valor "NA". En caso de que hubiéramos deseado tener estas variables como parte importante del estudio hubiéramos procedido a la eliminación de los registros u observaciones donde los valores de estas variables fueran "NA", esta decisión principalmente está fundamentada en el hecho de que son variables cualitativas.

Las variables `no_employees`, `family_history`, `remote_work`, `tech_company`, `benefits`, `care_options`, `wellness_program`, y `treatment` no tiene NA.

Se entiende por dato atípico como una observación fuera de la normalidad de la variable, una observación con una desviación tan grande de las otras observaciones que incluso podemos poner en duda si ha sido producido por los mismos mecanismos que las anteriores. El punto en común es lo alejado que esta del resto de las observaciones de la variable.

Los motivos por los cuales aparecen los datos atípicos pueden ser:

- 1.Outliers o datos atípicos cuyo origen está en la equivocación de los datos.
- 2.Valores atípicos u outlets con un propósito.
- 3.Valores atípicos u outlets cuyo origen son errores del muestreo.
- 4.Valores atípicos u outlets de errores en la estandarización.
- 5.Valores atípicos u outlets por asumir distribuciones erróneas.
- 6.Valor atípico u outlets cuyo origen es el muestreo correcto de la población.
- 7.Outliers o datos atípicos que proporcionan orígenes de nuevas investigaciones.



Los datos atípicos pueden tener efectos peligrosos en los diferentes análisis estadísticos que realicemos, con ellos presentes se puede llegar a aumentar el error de la varianza y hacer disminuir los resultados de las pruebas estadísticas.

Las únicas variables que poseen datos atípicos son Age y Gender.

En general, hemos acordado eliminar las observaciones cuyas respuestas no sean adecuadas dentro de los límites de lo aceptable para las preguntas formuladas, dado que si para Age o Gender el encuestado está optando una actitud que no refleja un comportamiento serio para estas variables, entendemos que tampoco es fiable la respuesta que da en el resto de la encuesta.

Así pues, eliminaremos las observaciones que contienen respuestas inadecuadas o fuera de rango aceptable.

Para la variable Gender debería haber 2 niveles (Male, Female), cuando hay 49. Así pues, para el variable género, cambiamos sus valores correspondientes por M y F respectivamente.

En aquellas observaciones que no es posible determinar o que la respuesta no es adecuada procedemos a eliminar la observación en lugar de asignar NA ( "A little about you", "Agender", "All", "Enby", "fluid", "Genderqueer", "Nah", "Neuter", "non-binary", "p", "queer", "queer/she/they", "Trnas woman", "Trans-female"). Nos quedamos por tanto con 1245 observaciones.

Exploramos a continuación los valores de la variable Age, variable cuantitativa. Dado que hemos visto que hay respuestas inadecuadas, como edades con signos negativos o valores que no pueden corresponderse con edad del encuestado, entendemos que son inadecuados y es preciso corregirlos (-29000, 329000, 1.000e+11, -1.726e+03, 5.000e+00, 8.000e+00, 1.100e+01, -1.000e+00)

Por tanto, del total de observaciones nos quedamos con 1240.

## Ejercicio 3

### Análisis de los datos.

#### 3.1 Selección de los grupos de datos que se quieren analizar/comparar.

Vamos a investigar:

- El hecho de recibir tratamiento tiene algo que ver con la edad, es decir, si existen diferencias en la variable Age según la variable treatment (Tratamiento)
- Dependiendo de la edad del individuo como percibe este el hecho de que la organización de igual importancia a la salud mental vs salud física, es decir, si existen diferencias en la variable Age según la variable mental\_vs\_physical

#### 3.2 Comprobación de la normalidad y homogeneidad de la varianza. Si es necesario (y posible), aplicar transformaciones que normalicen los datos.

El análisis de la normalidad o contrastes de normalidad, investigan cuanto de lejos está la distribución de los valores observados con respecto a una distribución normal con la misma media y desviación típica.

Para este análisis inicialmente podemos realizar unos estudios de manera gráfica. Vamos a comenzar observando si existe diferencias significativas según la edad del individuo para pensar que la organización da igual importancia a la salud mental o a la salud física.

Los estudios gráficos que utilizamos son Box Plot, histogramas, función de densidad, Normal Q-Q Plot, grafico para comprobar la homogeneidad de la varianza.

Realizamos el estudio de la normalidad mediante los contrastes de hipótesis.

Tenemos diferentes test de hipótesis:

- Test de Shapiro-Wilk: Para muestras de tamaño menor de 50

- Test de Kolmogorov-Smirnov

- Lilliefors: Da por hecho que la media y la varianza son desconocidas. Se considera que cuando tenemos muestras con tamaño superior a 50 es la alternativa de Shapiro-Wilk

- Test Jarque-Bera: Esta da valor a la alejancia que existe entre los coeficientes de asimetría y curtosis de los esperados por una distribución normal.

Todos estos test tenemos como hipótesis nula que los datos proceden de una distribución normal y la hipótesis alternativa que no lo hacen. El pvalor nos da la probabilidad de tener una distribución como la observada siempre y cuando los datos proceden de una población con distribución normal. Al estar hablando de pvalor, hay que tener en cuenta que a mayor tamaño de la muestra más finos son los test y es más sencillo encontrar evidencias en contra de  $H_0$ . De igual manera, a mayor tamaño de la muestra menos sensibles son los test paramétricos en falta de normalidad.

No realizamos el test de Shapiro-Wilk ya que nuestra muestra tiene un tamaño mayor a 50.

Vamos a utilizar el test de Kolmogorov-Smirnov, para estudiar si una muestra proviene de una población con una distribución de media y desviación típica específica.

Si no podemos asumir normalidad este hecho nos influya en los test de hipótesis paramétricos y en los modelos de regresión luego los estimadores calculados por mínimos cuadrados no serán eficientes y tanto los intervalos de confianza de los parámetros del modelo como contrastes significativos serán únicamente aproximados y no exactos.

Si tenemos en cuenta el teorema del límite central el cual necesita que las poblaciones de las que procede la muestra sea una normal, no las muestras. Si la muestra se distribuye según una normal está claro que la población también lo hará. Puede ocurrir que la muestra no se distribuye según una norma, pero si conocemos que la población se distribuye según una normal, entonces los contrastes paramétricos si son válidos. El Teorema del Límite Central permite simplificar los requisitos de normalidad cuando las muestras son grandes.

A continuación, estudiamos la homogeneidad de la varianza u homocedasticidad, se está considerando que la varianza es constante en los diferentes niveles.



Tenemos diferentes test para evaluar la distribución de la varianza. En todos ellos estamos considerando como hipótesis nula que la varianza es la misma en todos los grupos y como hipótesis alternativa que no lo es.

-F-Test. Razón de varianzas: Es recomendado siempre y cuando se tenga la certeza de que las poblaciones se distribuyen con normalidad. Luego es muy sensible en caso de no cumplir normalidad

-Test de Levene: Se puede utilizar en el caso de tener más de dos poblaciones. Permite elegir entre diferentes estadísticos de centralidad. Lo cual tiene relevancia a la hora de realizar el contraste de homocedasticidad según se tenga distribuciones normales o no.

-Test de Bartlett: Es muy sensible si no existe normalidad. Permite realizar el contraste para muestras de diferente tamaño.

-Test de Brown-Forsyth: Se basa en el test de Levene pero únicamente se utiliza la mediana como medida de centralidad.

-Test de Fligner-Killeen: Es el idóneo cuando no se cumple la condición de normalidad en las poblaciones. Es un test no paramétrico donde la comparativa de las varianzas se realizan basándonos en la mediana.

Al tener muestras de diferentes tamaños utilizaremos el test de Bartlett, aunque teniendo en cuenta los resultados anteriormente no sería el más idóneo ya que este es muy sensible si no existe normalidad.

Por lo tanto, tenemos un conjunto de la variable edad donde consideran que la organización da mayor importancia a la salud mental vs salud física y otro conjunto de la variable edad donde consideran que la organización no da mayor importancia a la salud mental vs salud física.

Comencemos con el análisis de la normalidad para el conjunto de valores de la variable Age que consideran que la organización da mayor importancia a la salud mental vs salud física.

Observamos primero los datos de una manera gráfica en este análisis comprobamos que existe cierta asimetría a la derecha, ya que las colas no son de igual longitud. Este gráfico fracciona los datos en 4 partes de igual frecuencia, es decir, cada grupo contiene mas o menos el mismo número de observaciones. Pero la ocupación de estos es diferente. El primer grupo (desde el valor mas pequeño hasta Q1) los valores de la variable Age donde los individuos consideran que la organización da mayor importancia a la salud mental que a la salud física va desde 18 hasta 27. El último grupo (desde Q3 hasta el máximo valor) desde 35 hasta 65. El 50% de los individuos observados tienen Age entre Q1 y Q3.

Al comprobar para este grupo tanto la función de densidad como su histograma confirmamos que se produce cierta asimetría a la derecha.

Vamos a utilizar también el gráfico de los cuantiles teóricos, Graficos Q-Q. Estos consisten en la comparación de los cuantiles de la distribución observada con los cuantiles teóricos de la distribución normal. Cuanto más se asemejen a una normal, mas alineados están los puntos a una recta. Al observar este gráfico vemos como los puntos situados en la parte mas alta (derecha) se alejan de la recta por lo tanto ya vamos con disposición de que no existe normalidad.

Si realizamos los test de normalidad de Kolmogorov-Smirnov, Lilliefors y Jarque-Bera procedemos a rechazar la hipótesis nula de normalidad ya que en todos ellos obtenemos un  $p\text{valor} < 0.05$ .

Realizamos la transformación  $\sqrt[3]{(1/\text{Age\_Mental})}$  para eliminar la asimetría de la derecha y suavizar la gráfica e intentar que esta transformación de la variable acepte un contraste de normalidad. Realizando las mismas pruebas gráficas y de contrastes podemos afirmar que procedemos a rechazar la hipótesis nula de normalidad ya que en todos ellos obtenemos un  $p\text{valor} < 0.05$ .

Continuamos con el análisis de la normalidad para el conjunto de valores de la variable Age que consideran que la organización da menor importancia a la salud mental vs salud física.

El primer grupo (desde el valor mas pequeño hasta Q1) los valores de la variable Age va desde 18 hasta 27. El último grupo va (desde Q3 hasta el máximo valor) desde 37 hasta 61. Podemos observar que la longitud desde el mínimo hasta Q1 es diferente a la de Q3 al máximo, por lo que podemos decir que no existe simetría con respecto a la mediana, por tanto podemos hablar de asimetría. El 50% de los individuos observados tienen Age entre Q1 y Q3.

Al comprobar para este grupo tanto la función de densidad como su histograma confirmamos que se produce cierta asimetría a la derecha.

Vamos a utilizar también el gráfico de los cuantiles teóricos, Gráficos Q-Q. Estos consisten en la comparación de los cuantiles de la distribución observada con los cuantiles teóricos de la distribución normal. Cuanto más se asemejen a una normal, mas alineados están los puntos a una recta. Al observar este gráfico vemos como los puntos situados en la parte mas alta (derecha) se alejan de la recta por lo tanto ya vamos con disposición de que no existe normalidad.

Si realizamos los test de normalidad anteriormente nombramos procedemos a rechazar la hipótesis nula de normalidad ya que en todos ellos obtenemos un  $p\text{valor} < 0.05$ .

Realizamos la transformación  $\sqrt[3]{(1/\text{Age\_Física})}$  para eliminar la asimetría de la derecha y suavizar la gráfica e intentar que esta transformación de la variable acepte un contraste de normalidad. Realizando las mismas pruebas gráficas y de contrastes podemos afirmar que procedemos a rechazar la hipótesis nula de normalidad ya que en todos ellos obtenemos un  $p\text{valor} < 0.05$ .

Para el análisis de la homogeneidad de la varianza para el conjunto de la variable edad donde consideran que la organización da mayor importancia a la salud mental vs salud física y el conjunto de la variable edad donde consideran que la organización da menor importancia a la salud mental vs salud física. Realizamos el test de Bartlett ya que consideramos que es el más idóneo ya que las muestras tienen diferentes tamaño aunque sea muy sensible si no existe normalidad. Procedemos a concluir que el test no haya diferencias significativas entre las varianzas de los dos grupos. Ya que el  $p\text{valor} > 0.05$ . Este resultado también lo confirmamos gráficamente.

Tenemos también un conjunto de la variable edad donde han tenido tratamiento relacionado con la salud mental y otro conjunto de la variable edad que no ha tenido tratamiento alguno relacionado con la salud mental.

Para estos grupos vamos a estudiar su normalidad y su homogeneidad de la varianza.

Comenzamos con la homogeneidad de la varianza, para su estudio realizamos el test de Bartlett con su resultado podemos concluir que si hay diferencias entre las varianzas de los grupos debido a que el pvalor  $< 0.05$ . Este hecho también lo podemos ver gráficamente.

Para el estudio de la normalidad utilizamos las gráficas y contrastes anteriormente nombrados en el caso anterior. En todos ellos podemos proceder a rechazar la hipótesis nula de normalidad (en contraste por pvalor  $< 0.05$  y gráficamente por cierta simetría)

En estos grupos realizamos la transformación  $\sqrt{1/\text{Age\_Tratamiento}}$  y  $\sqrt{1/\text{Age\_NTratamiento}}$  y aun así no conseguimos poder afirmar la normalidad.

### 3.3 Aplicación de pruebas estadísticas (tantas como sea posible) para comparar los grupos de datos.

Vamos a realizar el estudio ANOVA de un factor (one-way ANOVA o independent samples ANOVA) para investigar si existen diferencias en la edad entre los individuos que han tenido tratamiento o no de salud mental.

Nos hubiera gustado realizar un modelo de regresión logística pero finalmente este no ha sido posible.

## Ejercicio 4

### Representación de los resultados a partir de tablas y gráficas.

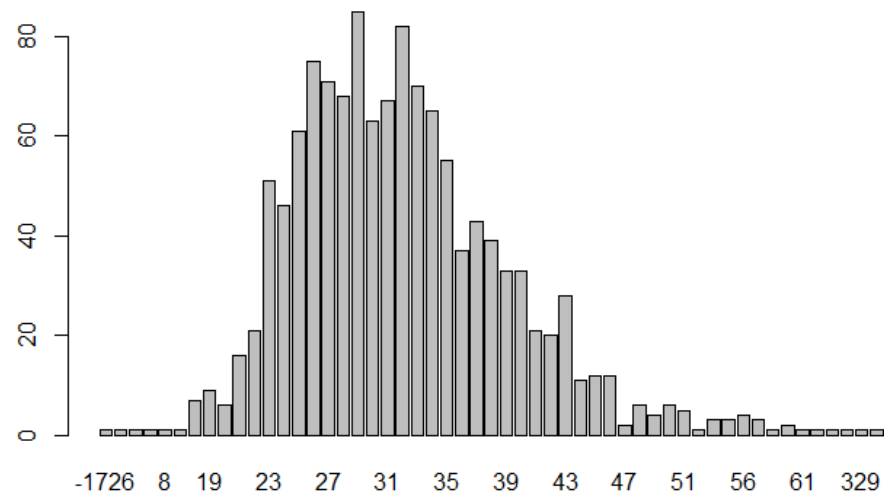
```
sapply(surveyMentalHealth,class)
## Age Gender
## "numeric" "factor"
## Country self_employed
## "factor" "factor"
## family_history treatment
## "factor" "factor"
## work_interfere no_employees
## "factor" "factor"
## remote_work tech_company
## "factor" "factor"
## benefits care_options
## "factor" "factor"
## wellness_program seek_help
## "factor" "factor"
## anonymity leave
## "factor" "factor"
## mental_health_consequence phys_health_consequence
## "factor" "factor"
## coworkers supervisor
## "factor" "factor"
## mental_health_interview phys_health_interview
## "factor" "factor"
## mental_vs_physical obs_consequence
## "factor" "factor"
```

No debemos olvidar que la transformación entre los diferentes tipos de datos es una labor ineludible en la limpieza de datos. Hay que tener siempre en mente que estas transformaciones conllevan un riesgo principal, que no es otro que la pérdida de datos al transformar un tipo de dato en otro. Recordemos que los principales factores que dan lugar a esta situación son:

- Mismo tipo de dato con transformación en diferente tamaño.
- Transformación con cota de exactitud diferente.

En el caso que nos ocupa todas las variables están definidas de forma correcta.

En la observación de los datos atípicos de la variable Age comenzamos visualizándolos mediante la gráfica

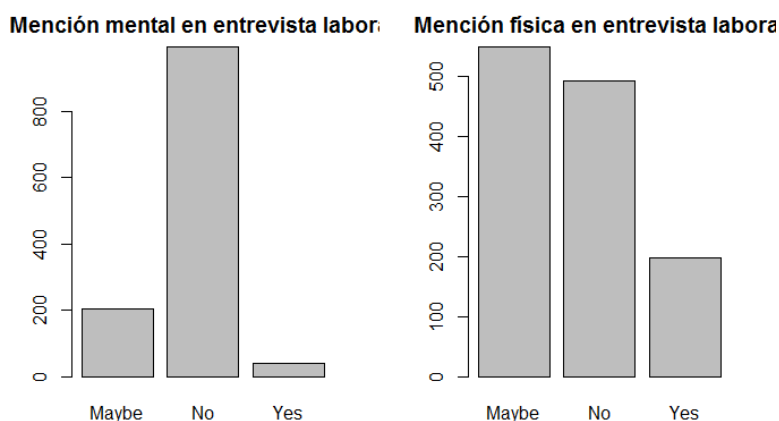


En la variable Gender obtenemos los siguientes niveles que nos incitan a un estudio mas en detalle de los datos atípicos para esta variable

```
## [1] "A little about you"
## [2] "Agender"
## [3] "All"
## [4] "Androgyne"
## [5] "cis-female/femme"
## [6] "Cis Female"
## [7] "cis male"
## [8] "Cis Male"
## [9] "Cis Man"
## [10] "Enby"
## [11] "f"
## [12] "F"
## [13] "femail"
## [14] "Femake"
## [15] "female"
## [16] "Female"
```

```
## [17] "Female "
## [18] "Female (cis)"
## [19] "Female (trans)"
## [20] "fluid"
## [21] "Genderqueer"
## [22] "Guy (-ish) ^_^"
## [23] "m"
## [24] "M"
## [25] "Mail"
## [26] "maile"
## [27] "Make"
## [28] "Mal"
## [29] "male"
## [30] "Male"
## [31] "Male-ish"
## [32] "Male "
## [33] "Male (CIS)"
## [34] "male leaning androgynous"
## [35] "Malr"
## [36] "Man"
## [37] "msle"
## [38] "Nah"
## [39] "Neuter"
## [40] "non-binary"
## [41] "ostensibly male, unsure what that really means"
## [42] "p"
## [43] "queer"
## [44] "queer/she/they"
## [45] "something kinda male?"
## [46] "Trans-female"
## [47] "Trans woman"
## [48] "woman"
## [49] "Woman"
```

Realizamos una comparativa entre la Mención mental en entrevista laboral y la Mención física en entrevista laboral.



Para la variable cuantitativa Age estudiamos los valores del mínimo, Q1, Mediana, Media, Q3 y máximo.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	18.00	27.00	31.00	32.08	36.00	72.00

Y de los siguientes indicadores

Media\_Age32.076739

Mediana\_Age31.000000

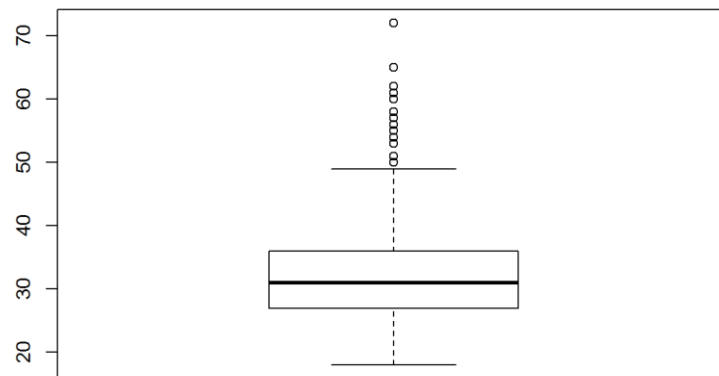
Media\_Recortada\_Age31.655723

Desviacion\_estandar\_Age7.288272

RIC\_Age9.000000

Desviacion\_Absoluta\_Mediana\_Age5.930400

De igual forma realizamos la representación del Box-Plot.

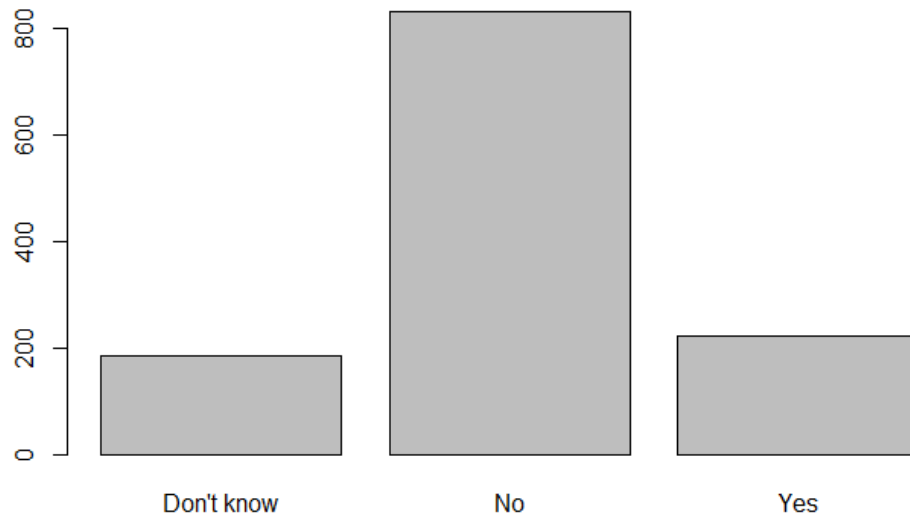


Para el resto que son variables cualitativas se realiza un análisis gráfico mediante gráficos de barras, en los cuales podemos observar los niveles de cada variable y el número de observaciones que pertenece a cada uno de ellos.

Wellness-program

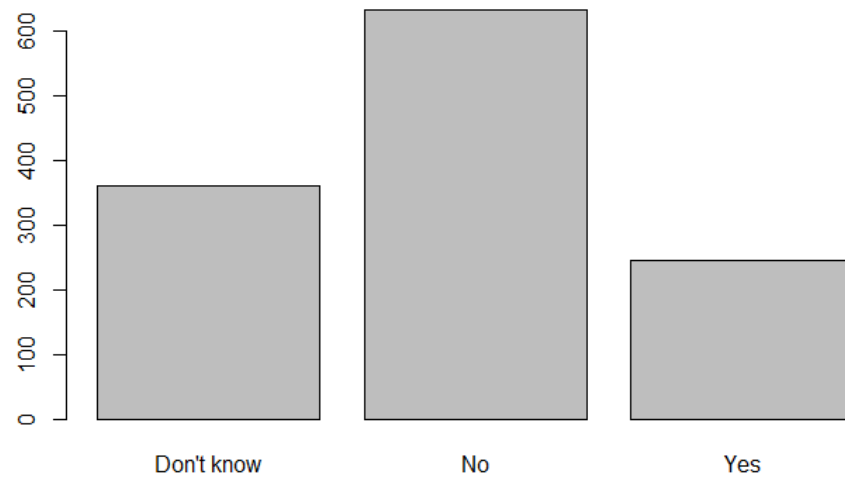


### Información del conocimiento de programas específicos



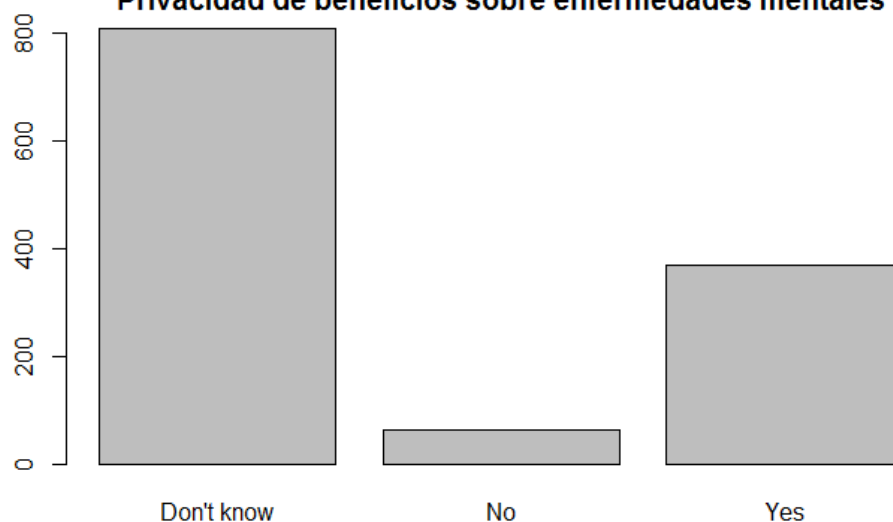
### Seek\_help

#### Información de recursos y ayuda desde la organización



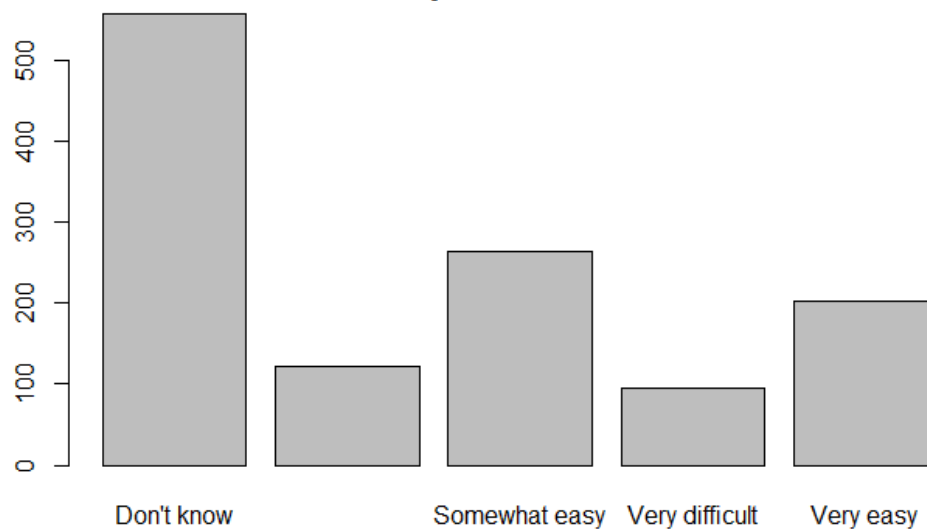
### Anonymity

### Privacidad de beneficios sobre enfermedades mentales



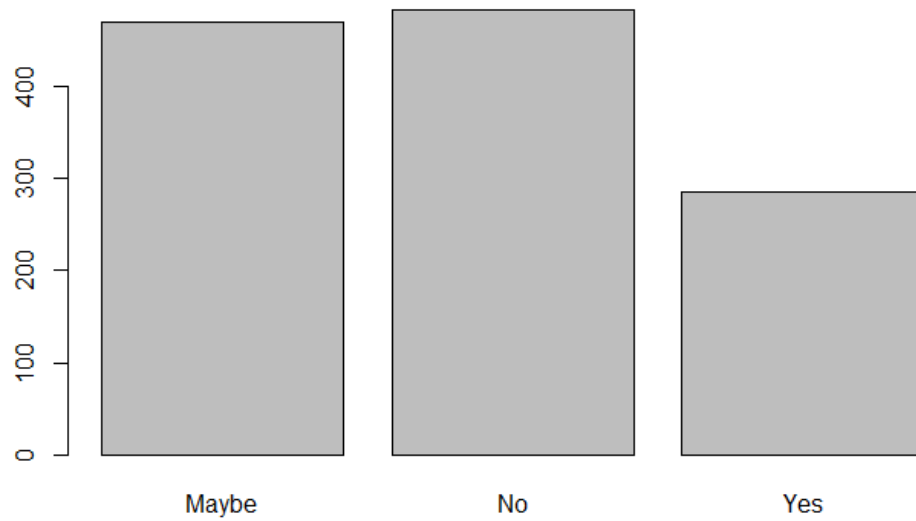
Leave

### Posibilidad de baja en enfermedades mentales



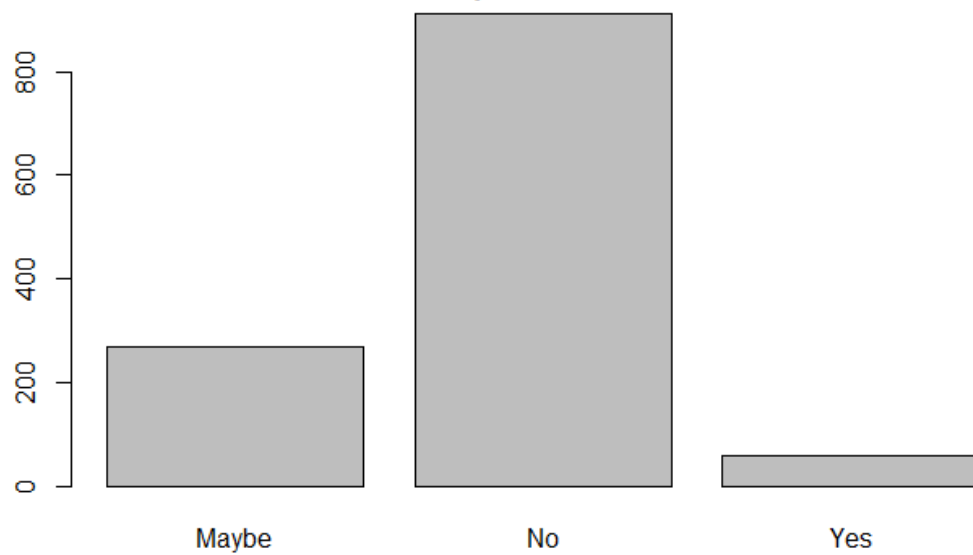
Mental\_health\_consequence

### Consecuencias por hablar de salud mental



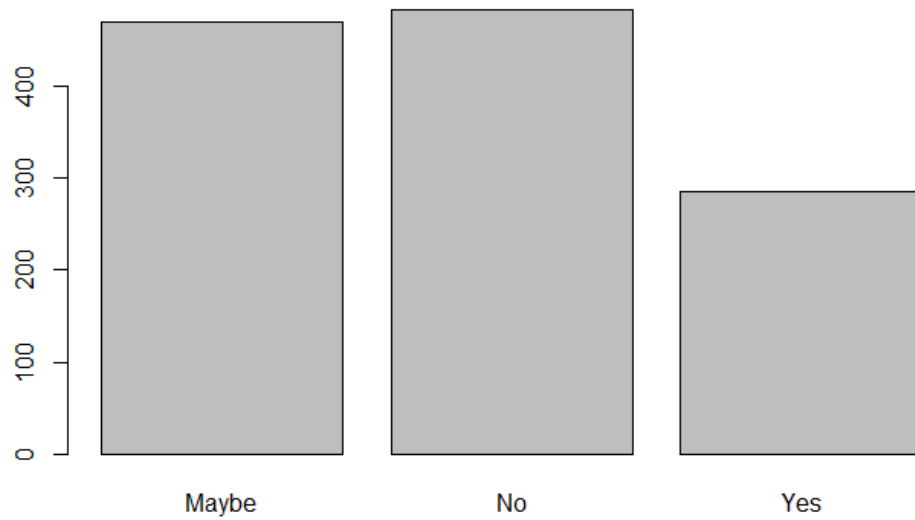
Phys\_health\_consequence

### Consecuencias por hablar de salud física



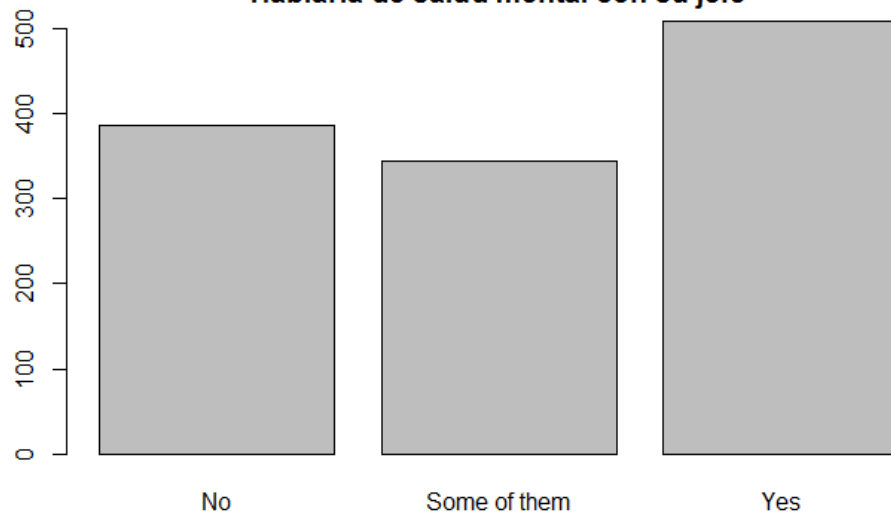
coworkers

### Hablaria de salud mental con compa eros



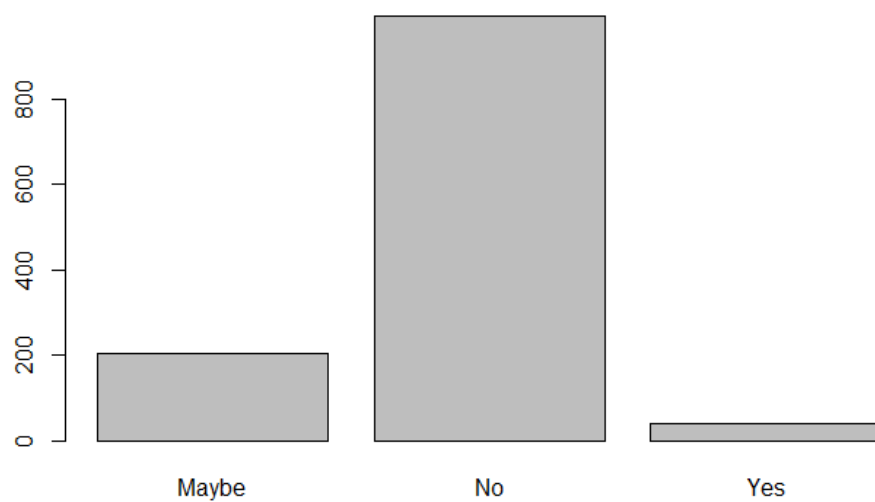
supervisor

### Hablaria de salud mental con su jefe



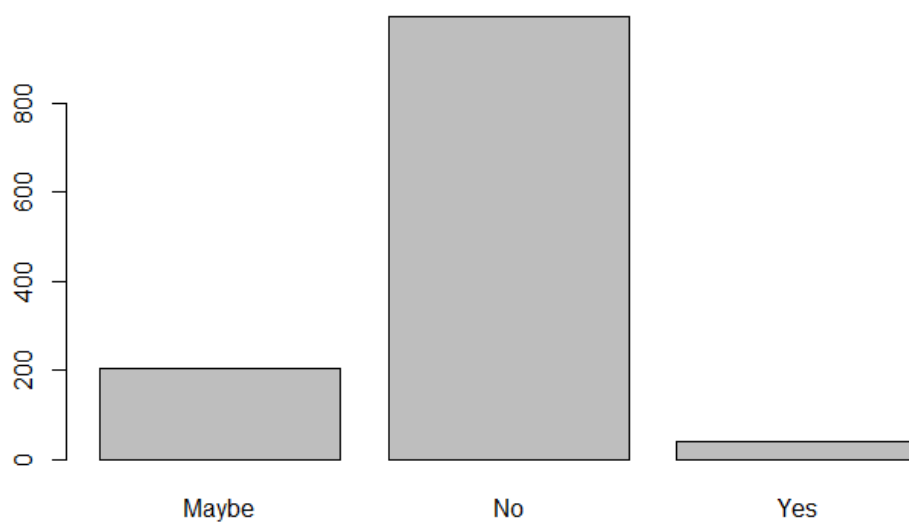
mental\_health\_interview

### Hablaria de salud mental en una entrevista laboral

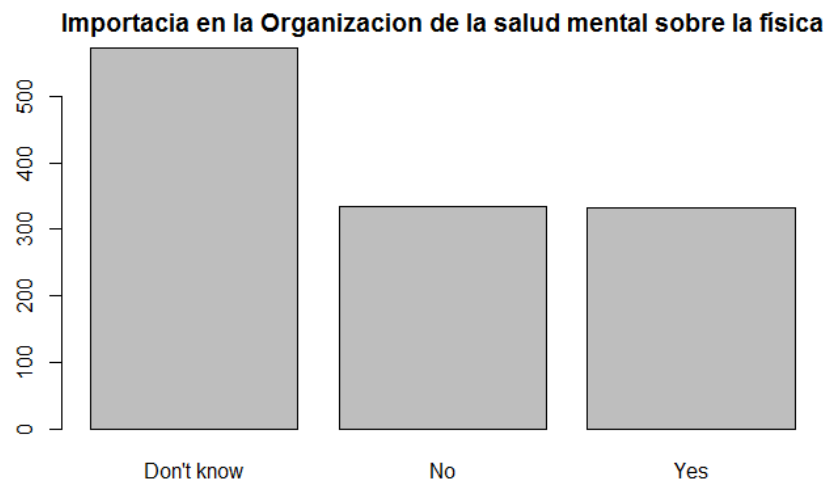


Phys\_health\_interview

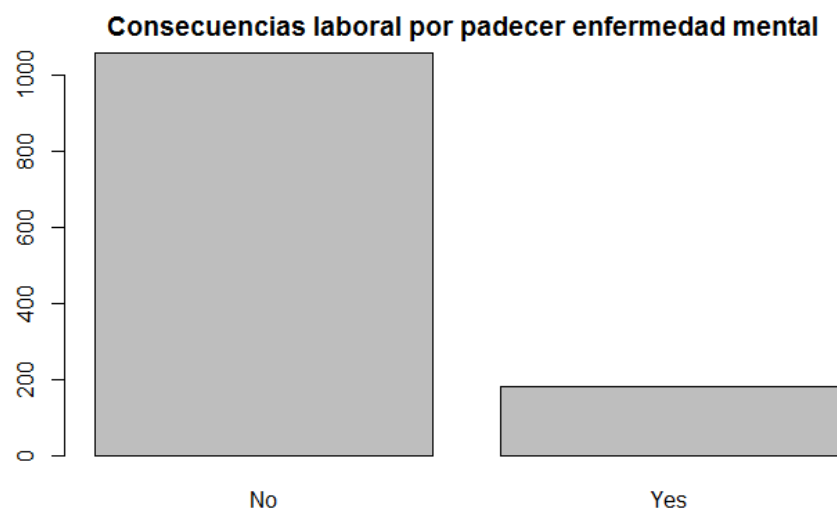
### Hablaria de salud física en una entrevista laboral



mental\_vs\_physical



obs\_consequence



En varias de las variables hay un porcentaje relevante de observaciones que indican incertidumbre de los propios encuestados, por lo que el análisis debería tomarse con cautela. Ello nos ayudaría a sacar conclusiones a priori sobre los observados, dado que muestra un grado de incertidumbre del sujeto, pero no permite extraer otras conclusiones sobre el entorno que pretendemos analizar.



Teníamos previsto realizar análisis de regresión logística para buscar la probabilidad e influencia entre las diferentes variables. Sin embargo, a la vista del análisis de los datos previo visto gráficamente, entendemos que en los mismo hay importantes carencias que, en aspectos que van desde la adecuada representatividad de la muestra hasta el control de respuestas que se ha observado, podrían implicar la invalidez de cualquier análisis.

Entendemos que es especialmente positivo para esta práctica, en este caso, la comprobación de la importancia del reprocesado y preparación de los datos y comprobación y validación de las respuestas en la muestra en relación con aspectos como la validez de las observaciones, niveles, o representatividad, así como sus valores. Ello nos ha permitido que -antes de entrar a hacer análisis estadísticos más avanzados para extraer conclusiones-, (ya que pueden por si mismos asegurar los posteriores pasos a tratar en el análisis, (como podría ser otras pruebas más avanzadas, como regresión múltiple, por ejemplo) podamos comprender la importancia del análisis de los datos en la muestra para garantizar su validez desde la recogida del dato. Lejos de ser un aspecto accesorio, ha resultado ser crítico. Desafortunadamente esto por otro lado, impide entrar a realizar análisis más profundos para responder a la pregunta que habíamos planteado.

Teniendo en cuenta que los grupos de datos que vamos estudiar corresponde a:

1. Dependiendo de la edad del individuo como percibe el hecho de que la organización de igual importancia a la salud mental vs salud física.
2. El hecho de recibir tratamiento tiene algo que ver con la edad

Comenzamos con el grupo 1 para ello realizamos los siguientes estudios

Calculo de la media

Recibe el nombre de Age\_Mental el grupo de observaciones correspondientes a las edades que perciben que la organización da mayor importancia a la enfermedad mental vs enfermedad física.

Age\_Fisica corresponde al grupo de observaciones correspondientes a las edades que perciben que la organización da mayor importancia a la enfermedad física vs enfermedad mental.

Cálculo de la media

`mean(Age_Mental)`

31.84985

`mean(Age_Fisica)`

32.4491

Cálculo de la mediana

`median(Age_Mental)`

31

`median(Age_Fisica)`

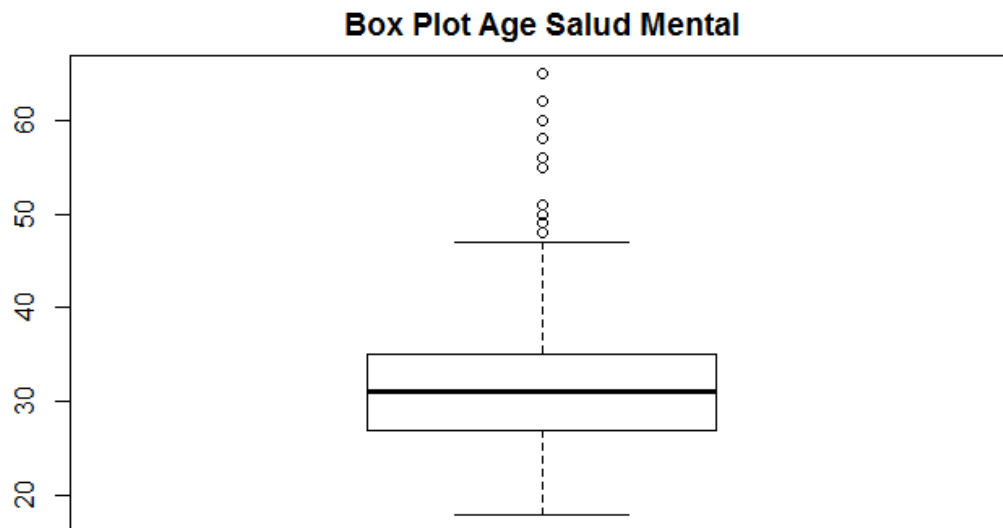
32

Cálculo de los cinco números de Tukey para Age\_Mental

Mínimo     18

Q1	27
Mediana	31
Q3	35
Máximo	65

Gráfico boxplot para Age\_Mental



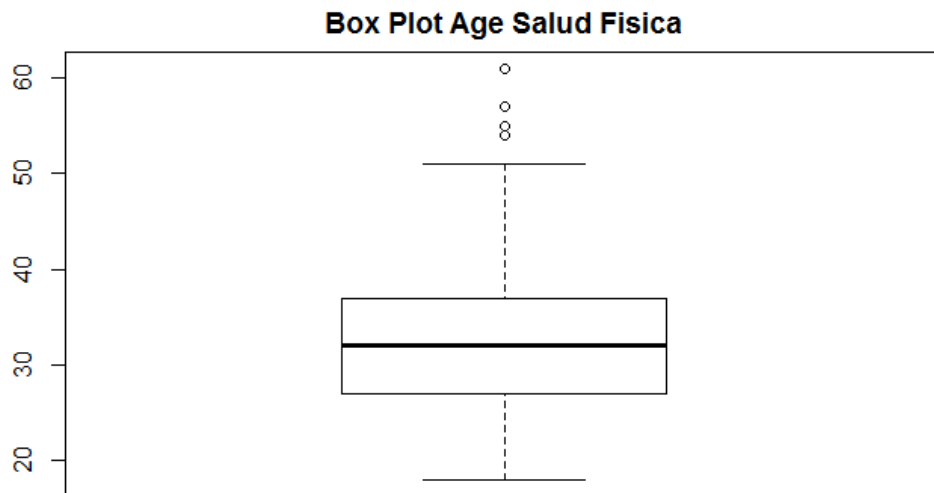
En un gráfico de Boxplot podemos estudiar la simetría, detectar outlier e incluso contrastar algunas hipótesis de la distribución. El gráfico fracciona los datos en 4 partes de igual frecuencia, es decir, cada grupo contiene más o menos el mismo número de observaciones. Pero la ocupación de estos es diferente. El primer grupo (desde el valor más pequeño hasta Q1) el valor de la variable Age va desde 18 hasta 27. El último grupo va (desde Q3 hasta el máximo valor) desde 35 hasta 65. Podemos observar que la longitud desde el mínimo hasta Q1 es diferente a la de Q3 al máximo, por lo que podemos decir que no existe simetría con respecto a la mediana, por tanto, podemos hablar de asimetría. El 50% de los individuos observados tienen Age entre Q1 y Q3.

Cálculo de los cinco números de Tukey para Age\_Física

Mínimo	18
Q1	27
Mediana	32
Q3	37

Máximo 61

Gráfico boxplot para Age\_Fisica



En un gráfico de Boxplot podemos estudiar la simetría, detectar outlier e incluso contrastar algunas hipótesis de la distribución. El gráfico fracciona los datos en 4 partes de igual frecuencia, es decir, cada grupo contiene más o menos el mismo número de observaciones. Pero la ocupación de estos es diferente. El primer grupo (desde el valor más pequeño hasta Q1) el valor de la variable Age va desde 18 hasta 27. El último grupo va (desde Q3 hasta el máximo valor) desde 37 hasta 61. Podemos observar que la longitud desde el mínimo hasta Q1 es diferente a la de Q3 al máximo, por lo que podemos decir que no existe simetría con respecto a la mediana, por tanto, podemos hablar de asimetría. El 50% de los individuos observados tienen Age entre Q1 y Q3.

Realizamos una representación de un histograma y superponemos una curva normal o función de densidad estimada para que se pueda ver la forma de la gráfica.

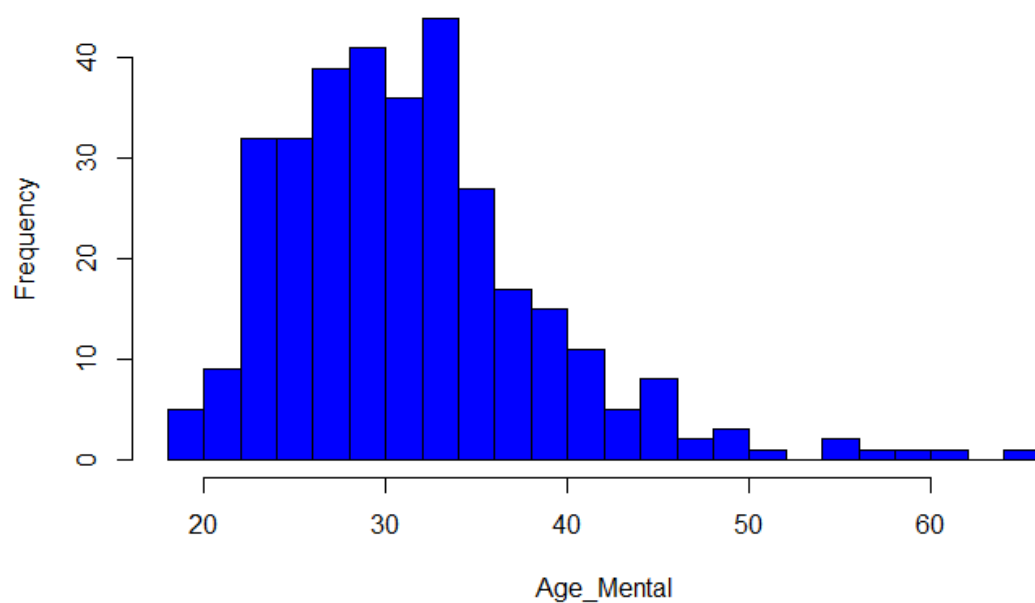
Representamos el histograma de la variable Age de la muestra. Para calcular el número de clases

que necesitamos realizamos el siguiente cálculo  $k=1+3,3*\log(n)$  ó  $k=\sqrt[n]{n}$ .

El número de intervalos correspondientes tanto a Age\_Mental como Age\_Fisica es 18.

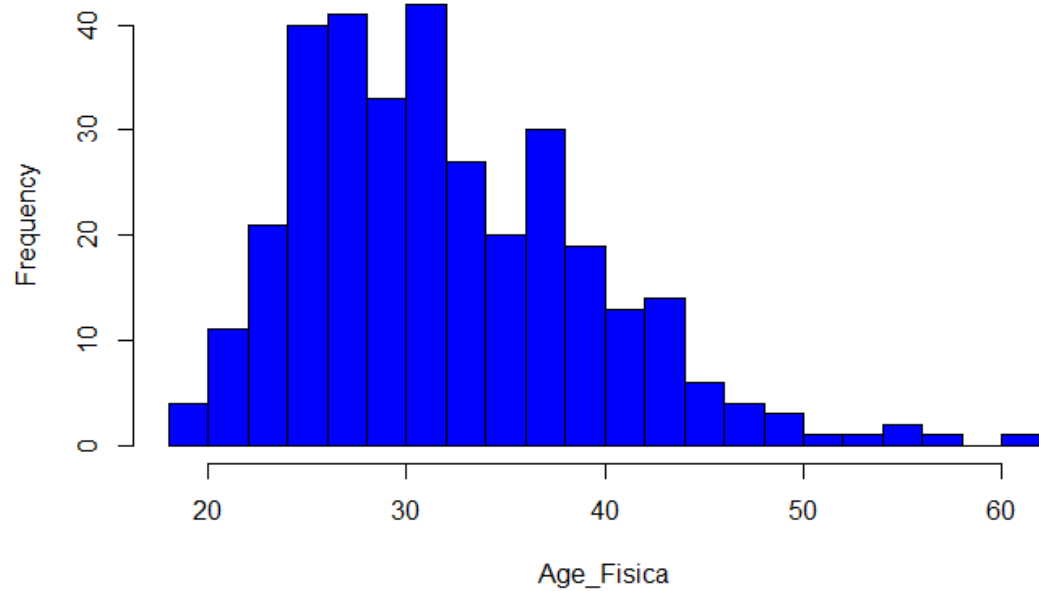
El histograma para Age\_Mental

### Edad individuos con Organizacion mas valor a la Salud Mental



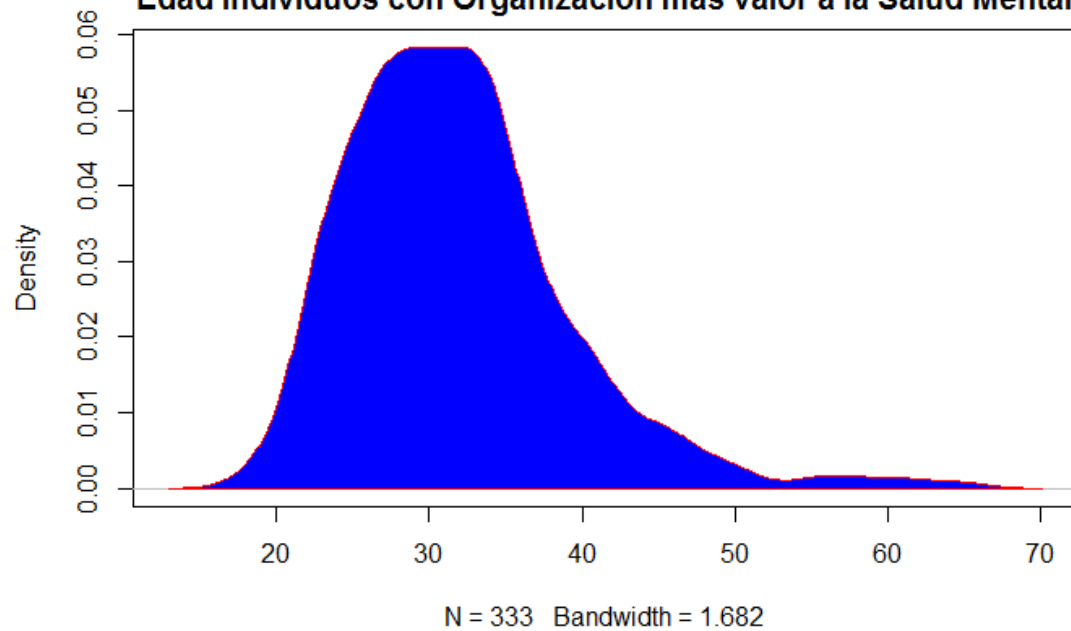
El histograma para Age\_Fisica

### Edad individuos con Organizacion mas valor a la Salud Fisica

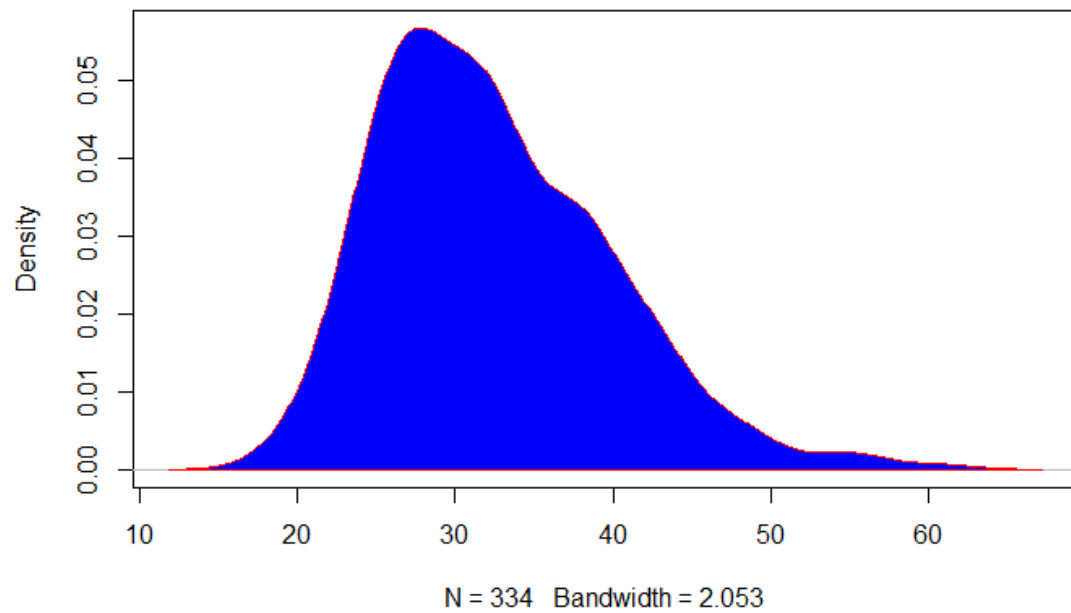


Calculamos la función de densidad de cada uno para finalmente superponerla sobre el histograma anteriormente representado

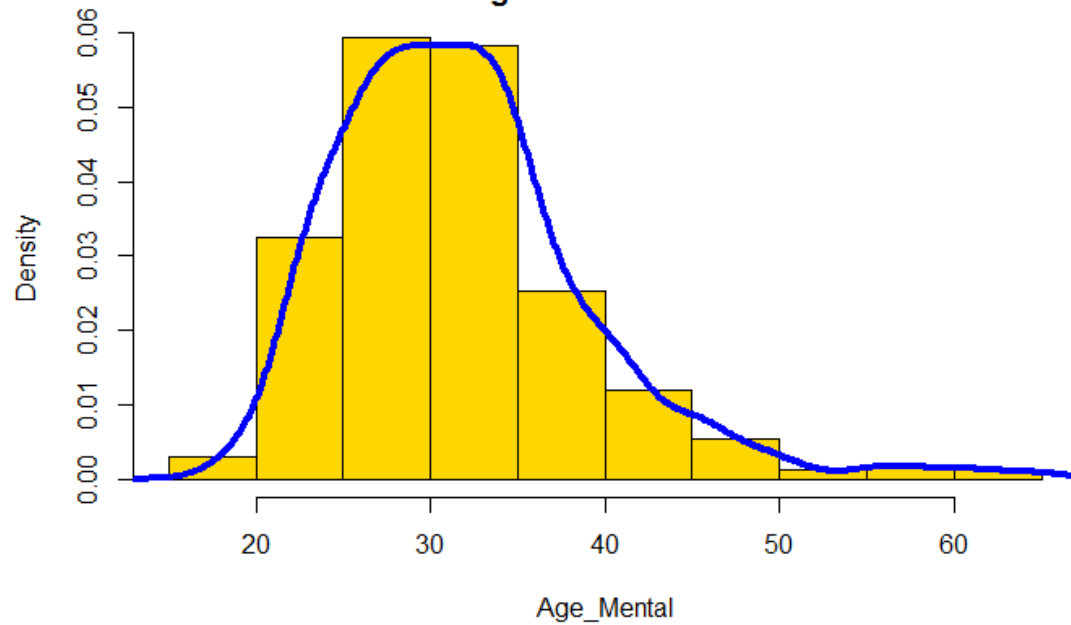
### Edad Individuos con Organización mas valor a la Salud Mental



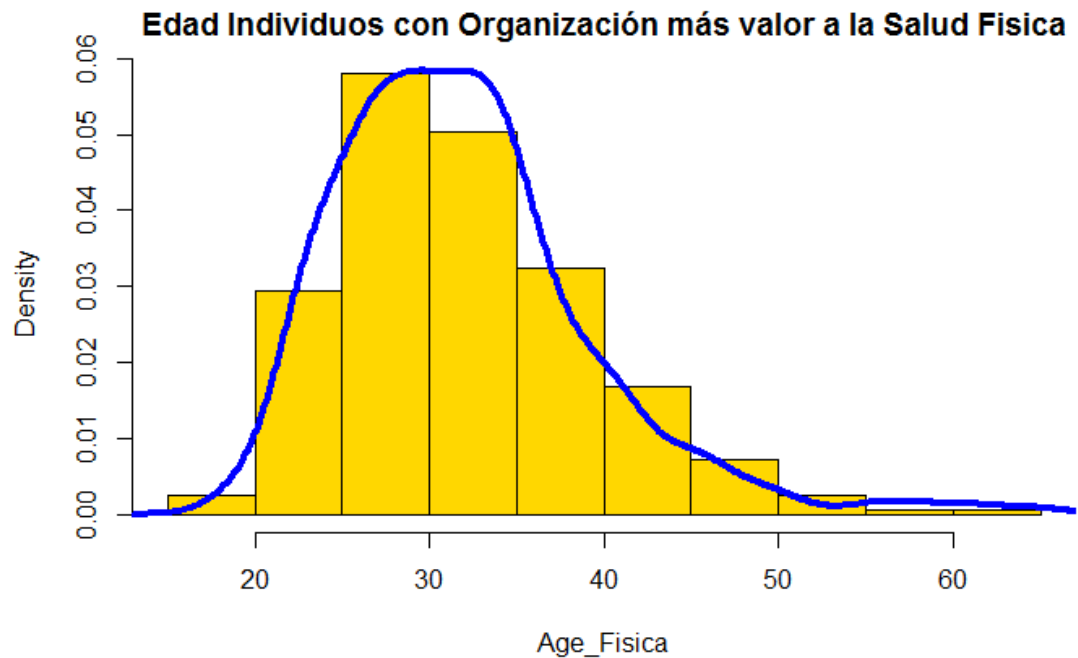
### Edad Individuos con Organización mas valor a la Salud Fisica



### Edad individuos con Organización mas valor a la Salud Mental

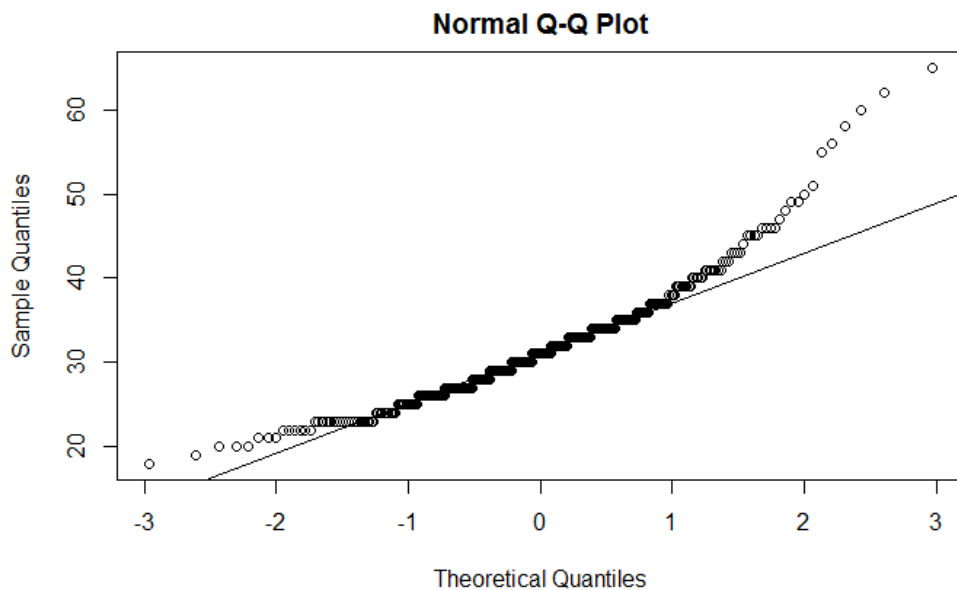




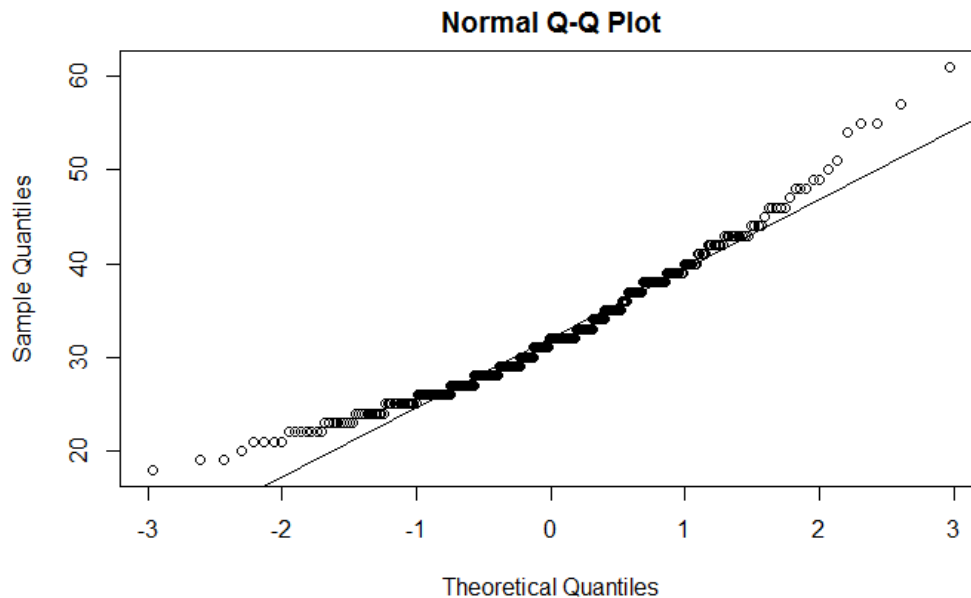


Vamos a utilizar también el gráfico de los cuantiles teóricos (Gráficos Q-Q). Estos consisten en la comparación de los cuantiles de la distribución observada con los cuantiles teóricos de la distribución normal. Cuanto más se asemejen a una normal, más alineados están los puntos a una recta.

Age\_Mental



Age\_Fisica



Para el estudio de la normalidad utilizamos el test de Kolmogorov-Smirnov  
Age\_Mental obtenemos el siguiente resultado

ties should not be present for the kolmogorov-Smirnov test  
One-sample kolmogorov-Smirnov test

data: Age\_Mental  
D = 0.098747, p-value = 0.003024  
alternative hypothesis: two-sided

Age\_Fisica obtenemos el siguiente resultado

ties should not be present for the kolmogorov-Smirnov test  
One-sample kolmogorov-Smirnov test

data: Age\_Fisica  
D = 0.099406, p-value = 0.002718  
alternative hypothesis: two-sided

Como ya hemos dicho anteriormente el test de Kolmogorov-Smirnov acepta que conoce la media y varianza poblacional, lo que hace que dicho test sea conservador y poco potente. Así tenemos el test de Lilliefors, en este caso se acepta que la media y la varianza son desconocidas.

Age\_Mental

Lilliefors (kolmogorov-Smirnov) normality test

data: (x = Age\_Mental)  
D = 0.098747, p-value = 2.646e-08

## Age\_Fisica

Lilliefors (Kolmogorov-Smirnov) normality test

```
data: (x = Age_Fisica)
D = 0.099406, p-value = 1.903e-08
```

Podemos tener en cuenta también el test de normalidad de Jarque-Bera, este no pide estimación de los parámetros con los que podemos caracterizar una normal. Este lo que hace es saber lo que se alejan los coeficientes de asimetría y curtosis de una distribución normal.

## Age\_Mental

Jarque-Bera test for normality

```
data: Age_Mental
JB = 180.38, p-value < 2.2e-16
```

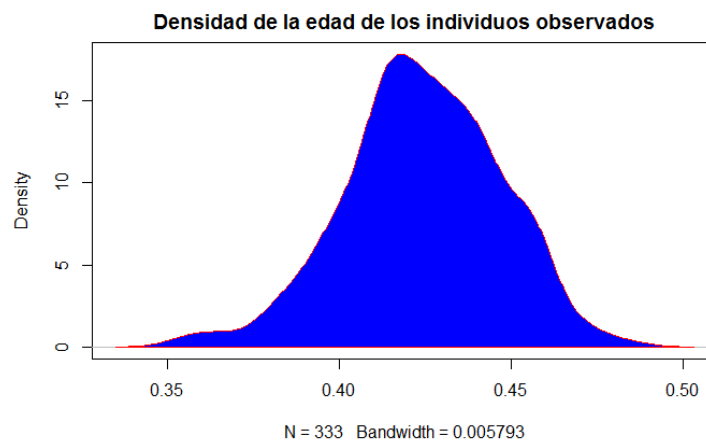
## Age\_Fisica

Jarque-Bera test for normality

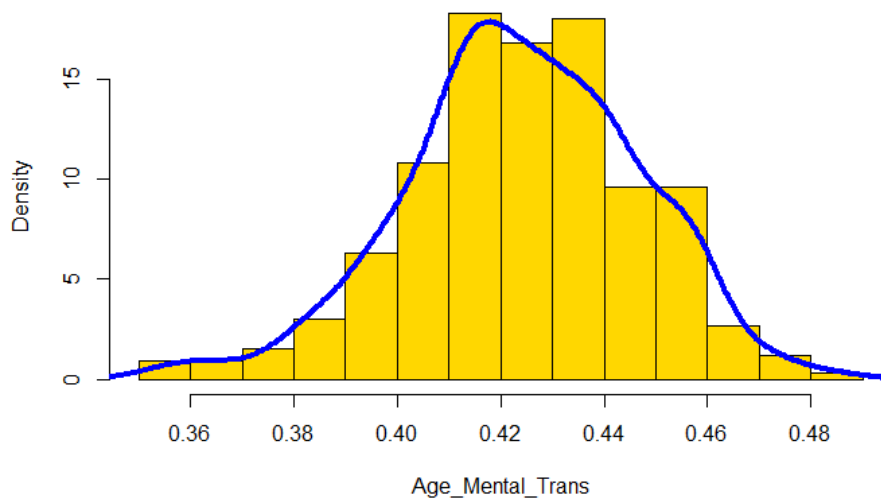
```
data: Age_Fisica
JB = 40.943, p-value < 2.2e-16
```

Procedemos a rechazar la hipótesis nula de normalidad ya que en todos los test obtenemos un p-valor < 0.05

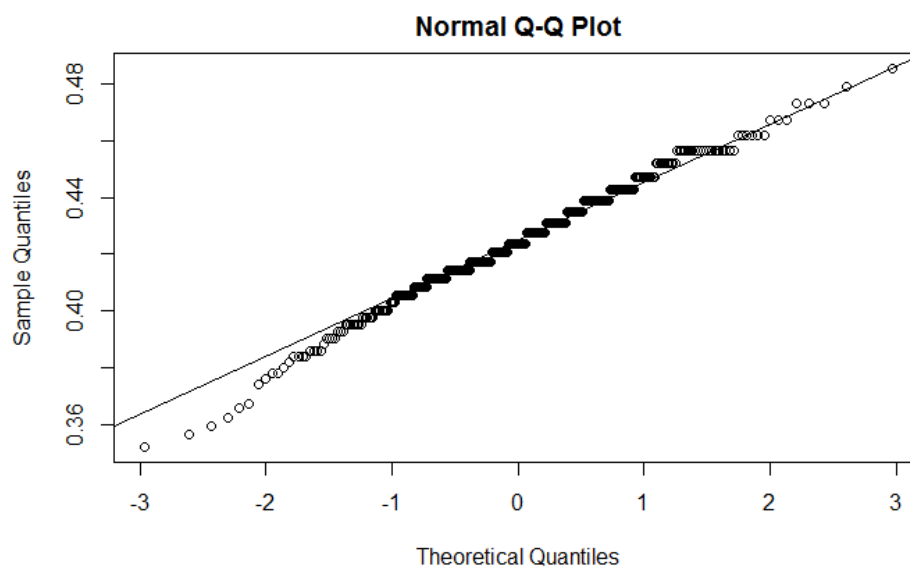
Realizamos la transformación  $\sqrt{1/\text{Age\_Mental}}$  y con ella calculamos su correspondiente función de densidad y la superposición de esta sobre su histograma



Histograma y densidad del valor de la edad de los individuos observado:



Observamos el gráfico Normal Q-Q Plot para la transformación de Age\_Mental

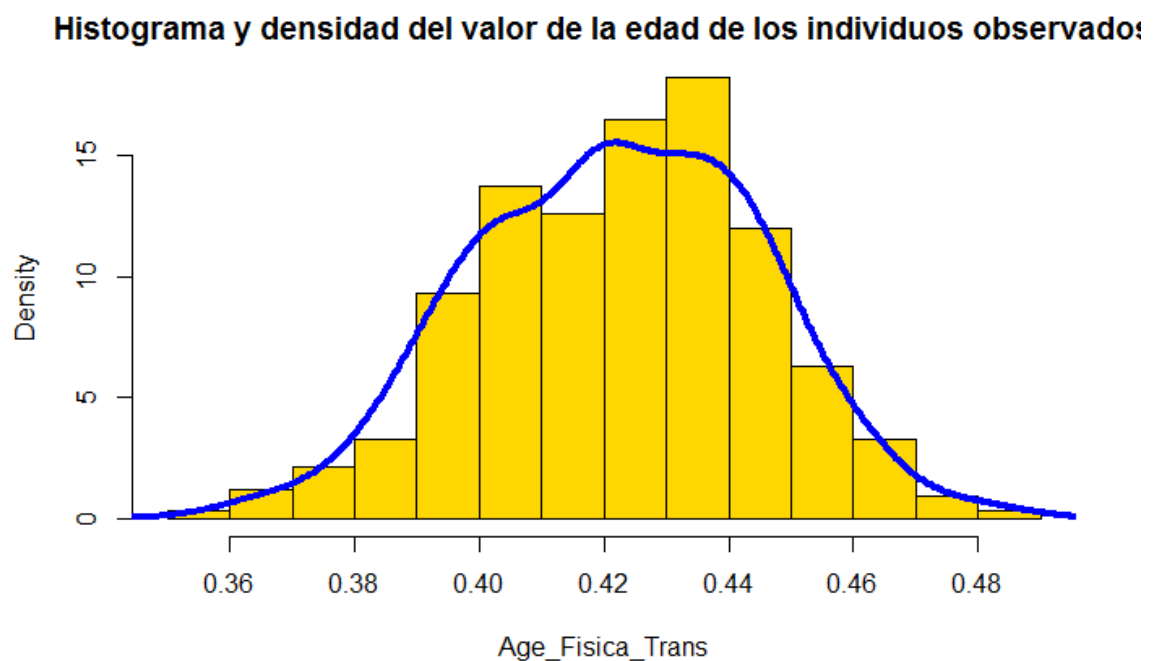
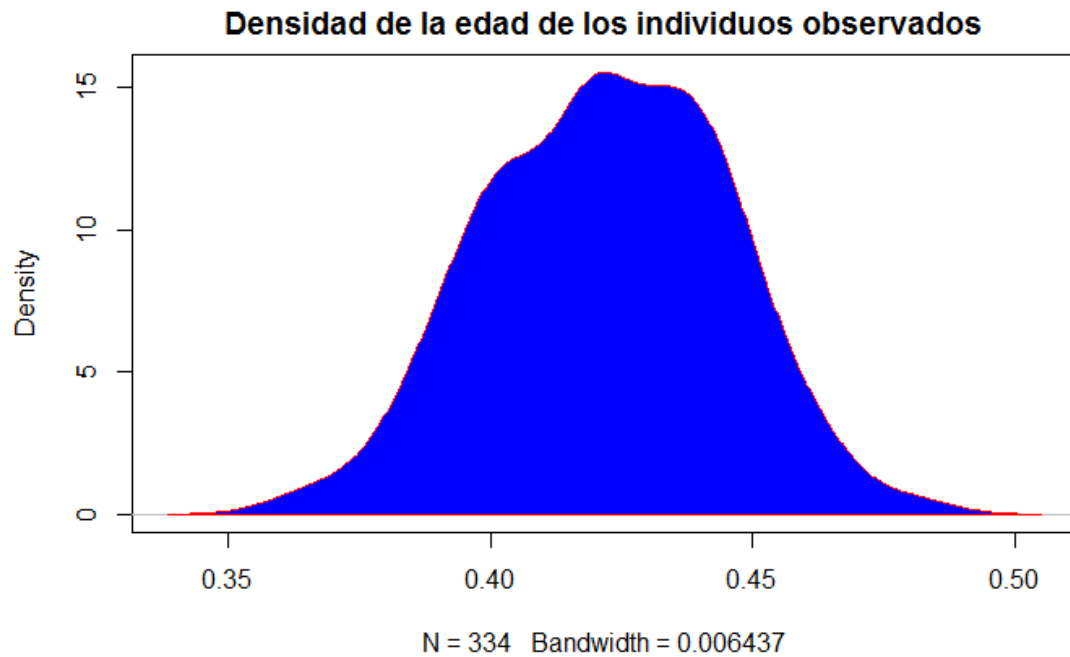


Lilliefors (Kolmogorov-Smirnov) normality test

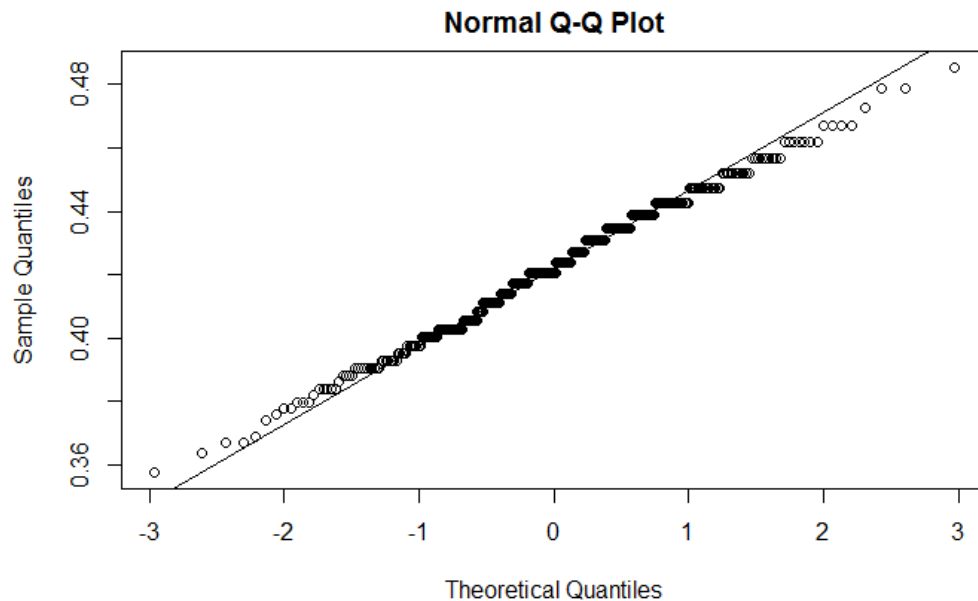
data: (x = Age\_Mental\_Trans)  
D = 0.04929, p-value = 0.04973

Donde continuamos rechazando la hipótesis nula de normalidad ya que  $p\text{-value} < 0.05$   
Pero el valor del pvalor se ha mejorado bastante.

Realizamos la transformación  $\sqrt{1/\text{Age\_Física}}$  y con ella calculamos su correspondiente función de densidad y la superposición de esta sobre su histograma



Observamos el gráfico Normal Q-Q Plot para la transformación de Age\_Fisica



Lilliefors (Kolmogorov-Smirnov) normality test

```
data: (x = Age_Fisica_Trans)
D = 0.059944, p-value = 0.005716
```

Procedemos a rechazar la hipótesis nula de normalidad ya que el pvalor<0.05

A continuación, estudiamos la homogeneidad de la varianza u homocedasticidad, se está considerando que la varianza es constante en los diferentes niveles.

Tenemos diferentes test para evaluar la distribución de la varianza. En todos ellos estamos considerando como hipótesis nula que la varianza es la misma en todos los grupos y como hipótesis alternativa que no lo es.

Al tener muestras de diferentes tamaños utilizaremos el test de Bartlett, aunque teniendo en cuenta los resultados anteriormente no sería el más idóneo ya que este es muy sensible si no existe normalidad.

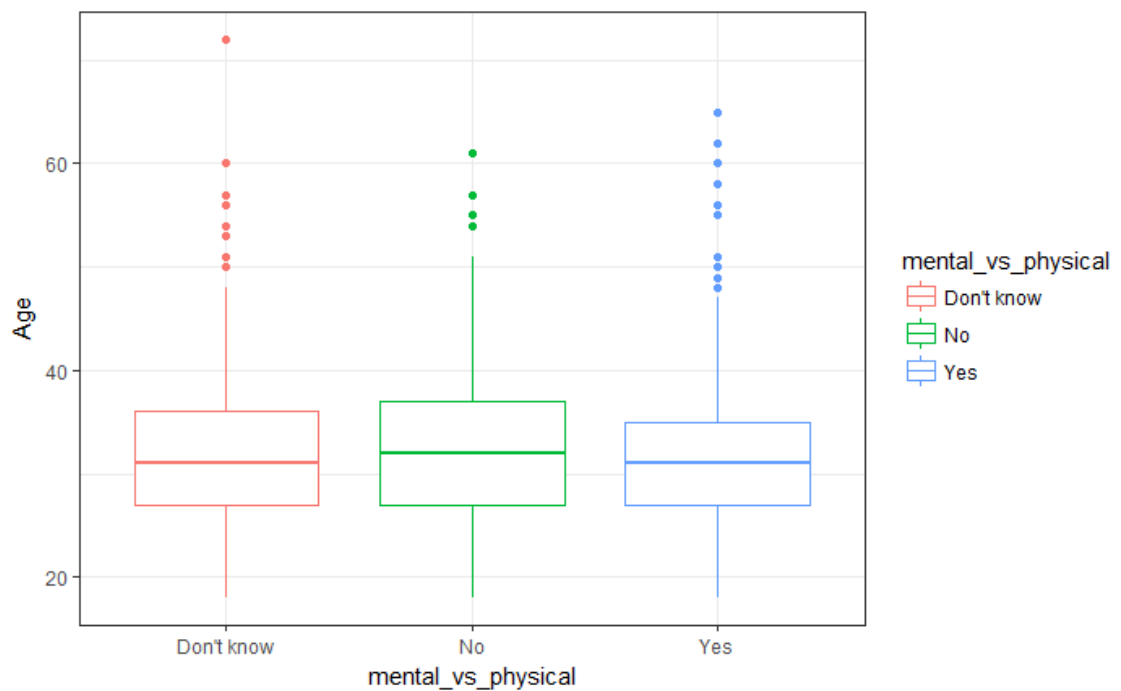
Bartlett test of homogeneity of variances

```
data: list(Age_Mental, Age_Fisica)
Bartlett's K-squared = 2.3584e-06, df = 1, p-value = 0.9988
```

Podemos concluir que el test no hay diferencias significativas entre las varianzas de los dos grupos, pvalor>0.05

Gráficamente también lo podemos visualizar





Vamos a continuar observando si existe diferencias significativas según la edad del individuo de haber recibido o no tratamiento.

Definimos por Age\_Tratamiento aquellas observaciones que corresponden a las edades de los individuos que han recibido tratamiento mientras que Age\_NTratamiento corresponden a las edades de los individuos que no han recibido tratamiento.

1.Cálculo de la media

```
mean(Age_Tratamiento)
```

```
32.6886
```

```
mean(Age_NTratamiento)
```

```
31.53485
```

2.Cálculo de la mediana

```
median(Age_Tratamiento)
```

```
32
```

```
median(Age_NTratamiento)
```

```
31
```

Cálculo de los cinco números de Tukey para

Age\_Tratamiento

Mínimo

```
18
```

Q1

```
27
```

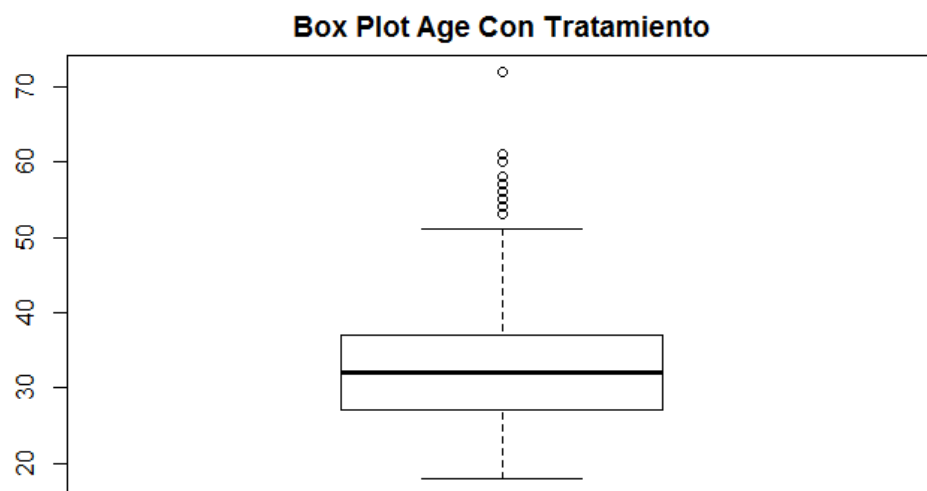
Mediana

```
32
```

Q3

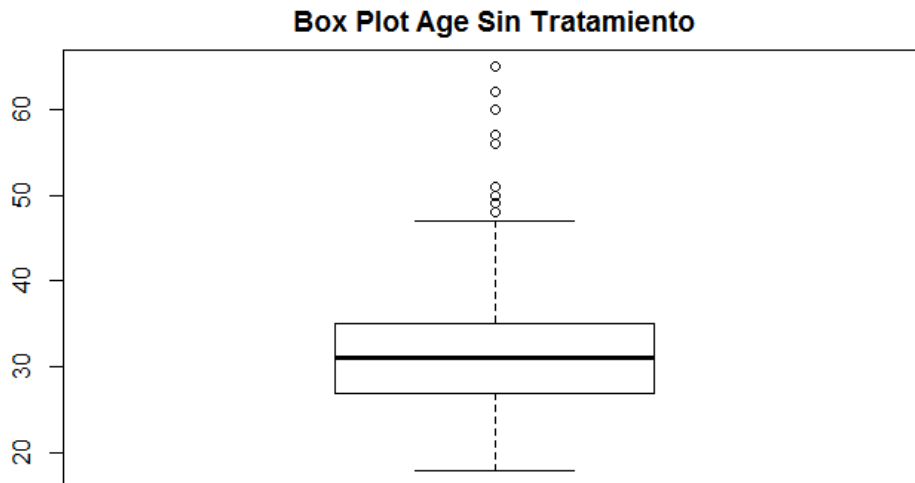
37  
Máximo  
72  
Para Age\_NTratamiento  
Mínimo  
18  
Q1  
27  
Mediana  
31  
Q3  
35  
Máximo  
65

Gráfico Box-Plot para la variable Age\_Tratamiento



En un gráfico de Boxplot podemos estudiar la simetría, detectar outlier e incluso contrastar algunas hipótesis de la distribución. El gráfico fracciona los datos en 4 partes de igual frecuencia, es decir, cada grupo contiene más o menos el mismo número de observaciones. Pero la ocupación de estos es diferente. El primer grupo (desde el valor más pequeño hasta Q1) el valor de la variable Age va desde 18 hasta 27. El último grupo va (desde Q3 hasta el máximo valor) desde 37 hasta 72. Podemos observar que la longitud desde el mínimo hasta Q1 es diferente a la de Q3 al máximo, por lo que podemos decir que no existe simetría con respecto a la mediana, por tanto, podemos hablar de asimetría. El 50% de los individuos observados tienen Age entre Q1 y Q3.

Gráfico Box-Plot para la variable Age\_NTratamiento



En un gráfico de Boxplot podemos estudiar la simetría, detectar outlier e incluso contrastar algunas hipótesis de la distribución. El gráfico fracciona los datos en 4 partes de igual frecuencia, es decir, cada grupo contiene más o menos el mismo número de observaciones. Pero la ocupación de estos es diferente. El primer grupo (desde el valor más pequeño hasta Q1) el valor de la variable Age va desde 18 hasta 27. El último grupo va (desde Q3 hasta el máximo valor) desde 35 hasta 75. Podemos observar que la longitud desde el mínimo hasta Q1 es diferente a la de Q3 al máximo, por lo que podemos decir que no existe simetría con respecto a la mediana. Por tanto, podemos hablar de asimetría. El 50% de los individuos observados tienen Age entre Q1 y Q3.

Realizamos una representación de un histograma y superponemos una curva normal o función de densidad estimada para que se pueda ver la forma de la gráfica.

Calculamos el número de intervalos

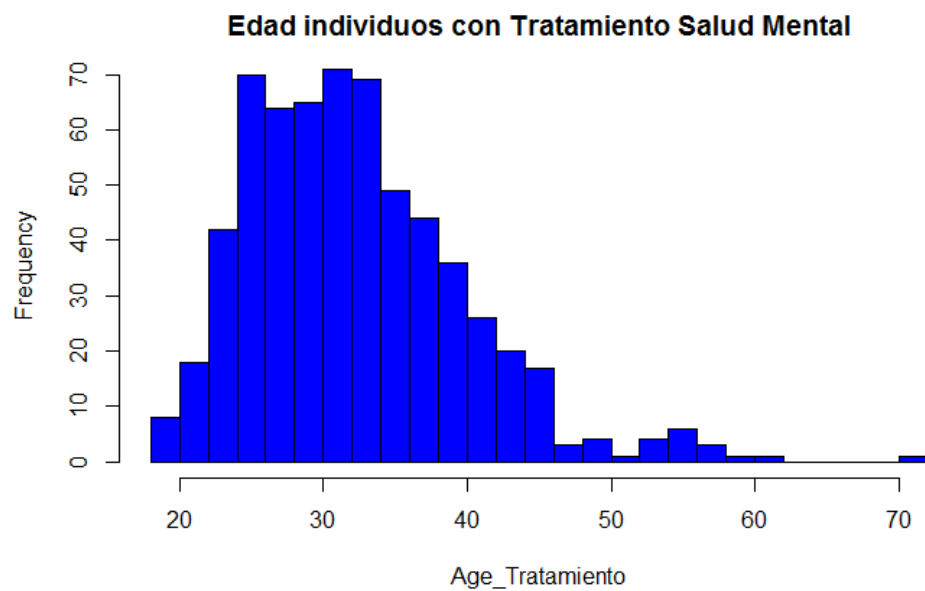
Age\_Tratamiento

25

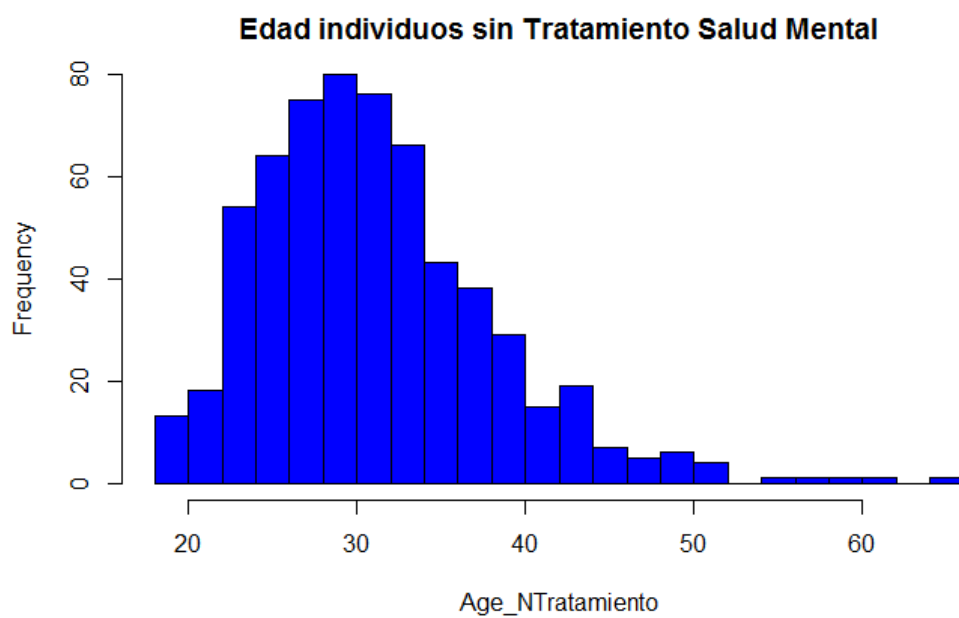
Age\_NTratamiento

25

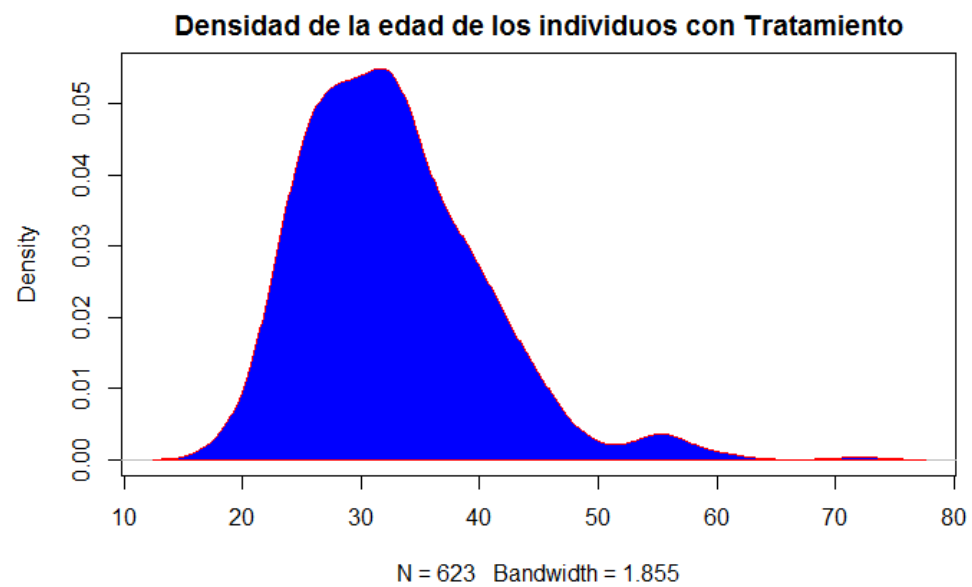
Visualizamos el histograma de Age\_Tratamiento



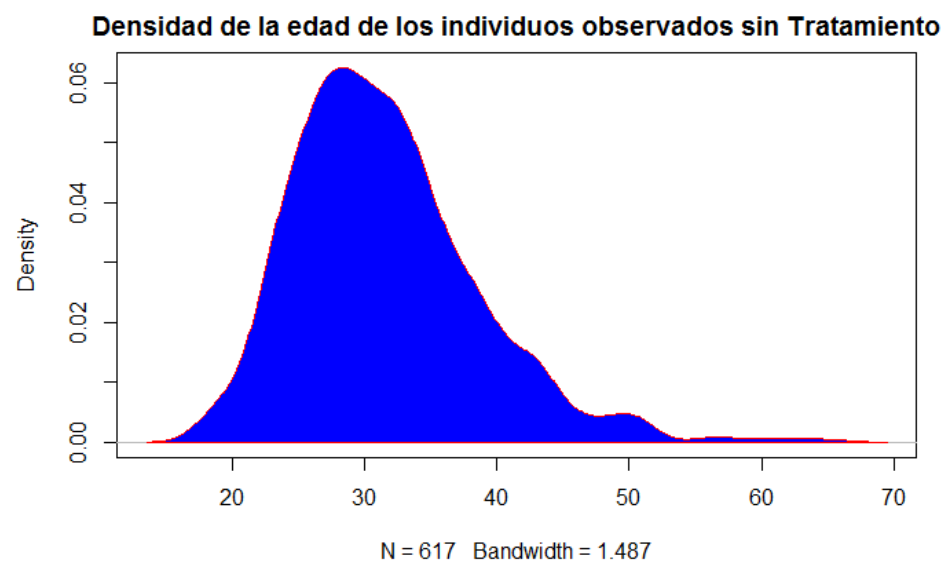
Histograma de Age\_NTratamiento



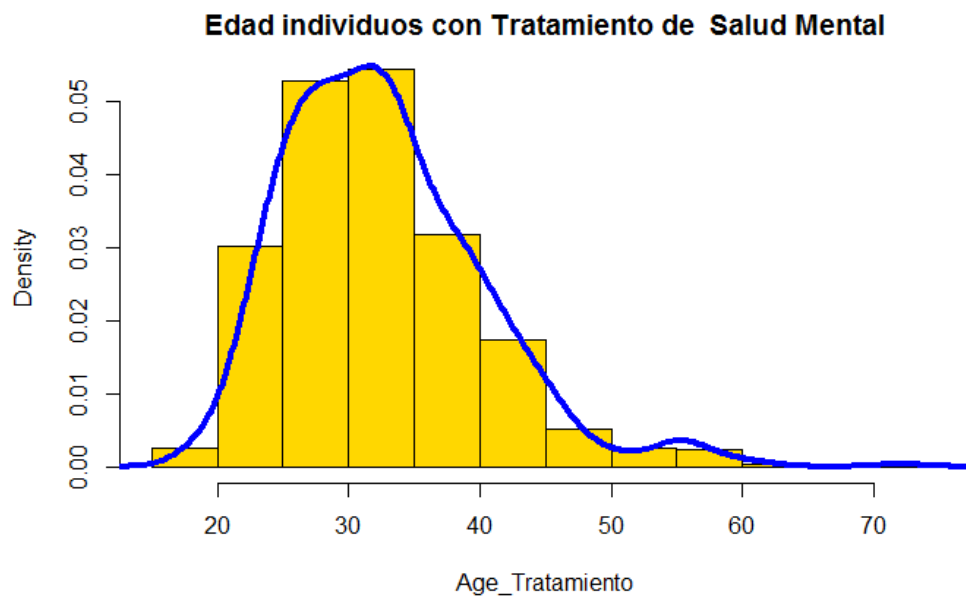
Calculamos las funciones de densidad  
Age\_Tratamiento



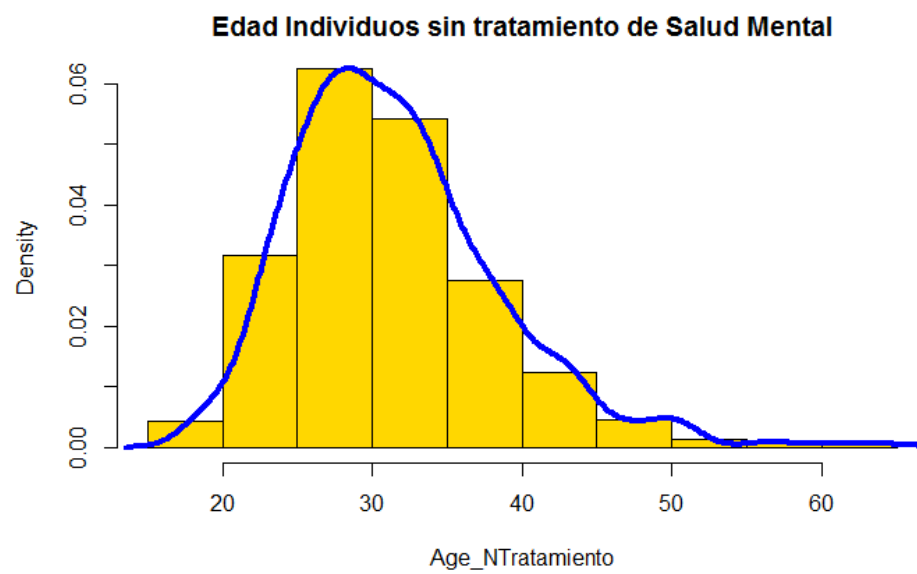
Función de densidad para Age\_NTratamiento



Realizamos la superposición del histograma y la función de densidad  
Age\_Tratamiento

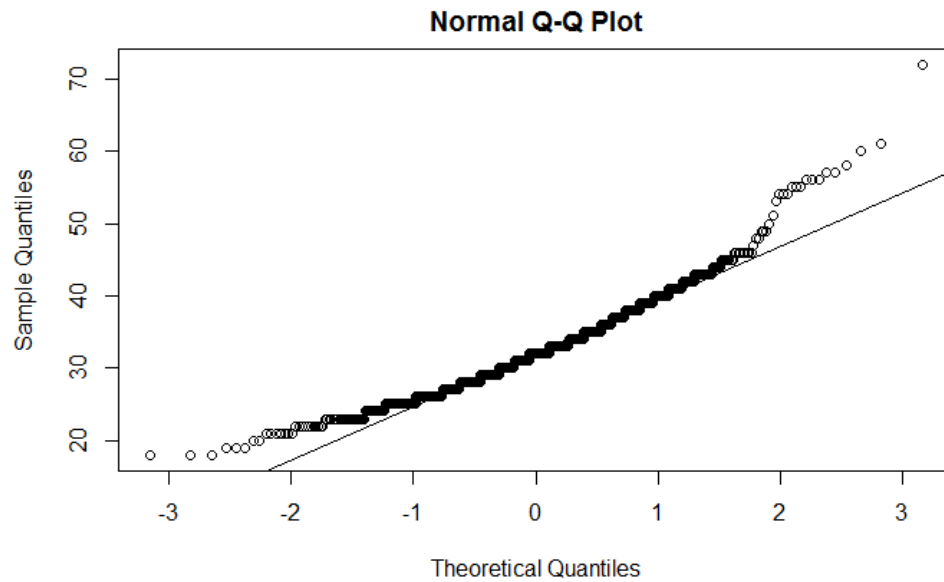


Age\_NTratamiento

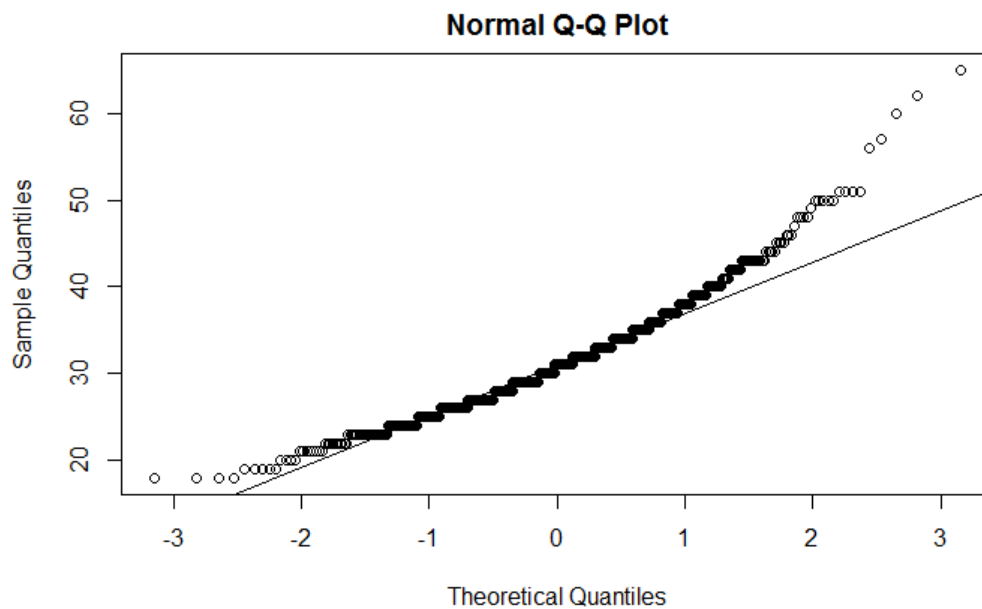


Vamos a utilizar también el gráfico de los cuantiles teóricos (Gráficos Q-Q).

### Age\_Tratamiento



### Age\_NTratamiento



Realizamos los contrastes para Age\_Tratamiento

ties should not be present for the kolmogorov-Smirnov test  
One-sample kolmogorov-Smirnov test

```
data: Age_Tratamiento
D = 0.088718, p-value = 0.0001101
alternative hypothesis: two-sided
```

Lilliefors (Kolmogorov-Smirnov) normality test

```
data: (x = Age_Tratamiento)
D = 0.088718, p-value = 1.166e-12
```

Jarque-Bera test for normality

```
data: Age_Tratamiento
JB = 190.72, p-value < 2.2e-16
```

Procedemos a rechazar la hipótesis nula de normalidad ya que en todos los test obtenemos un  $p\text{-value} < 0.05$

Age\_NTratamiento

ties should not be present for the Kolmogorov-Smirnov test  
One-sample Kolmogorov-Smirnov test

```
data: Age_NTratamiento
D = 0.089244, p-value = 0.0001078
alternative hypothesis: two-sided
```

```
data: (x = Age_NTratamiento)
D = 0.089244, p-value = 1.085e-12
```

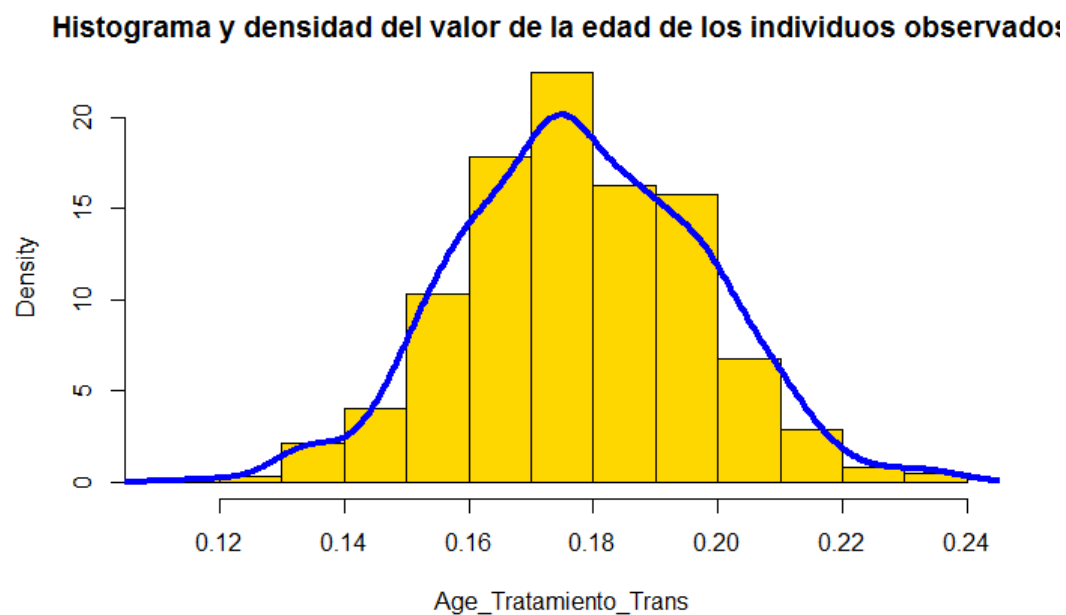
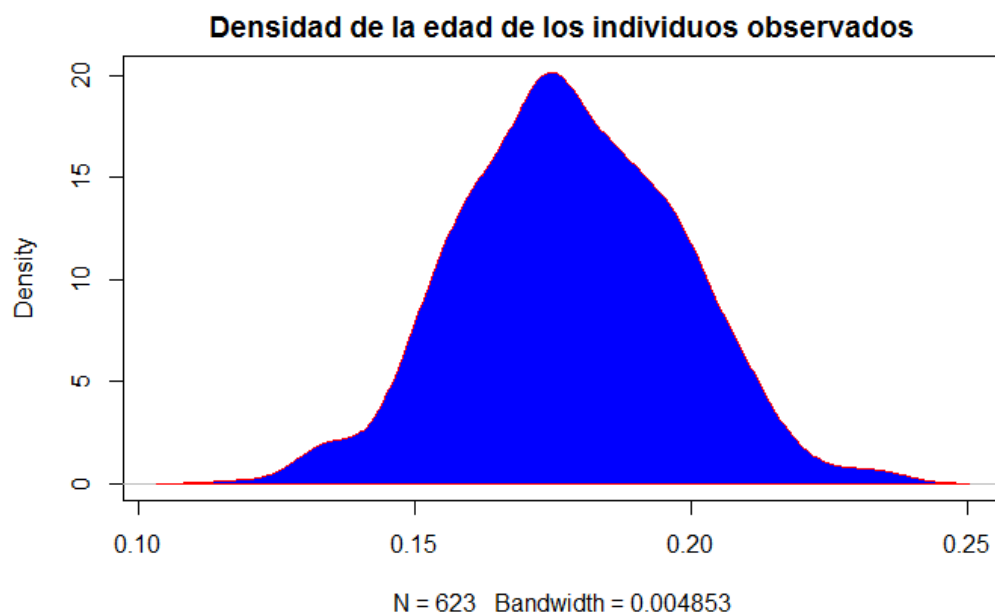
Jarque-Bera test for normality

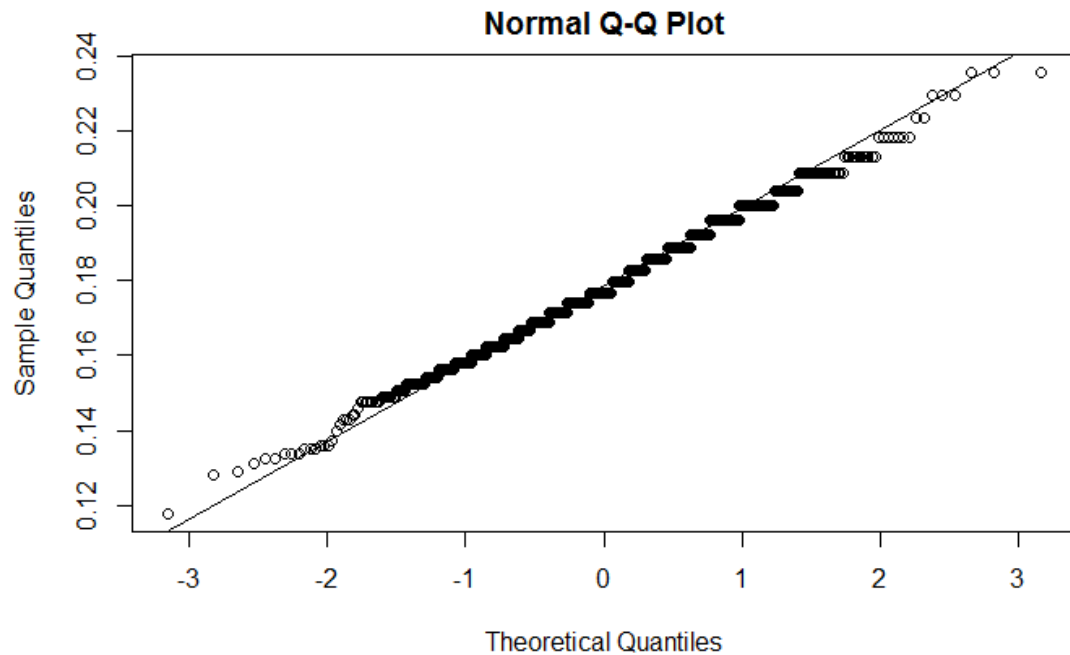
```
data: Age_NTratamiento
JB = 186.91, p-value < 2.2e-16
```

Procedemos a rechazar la hipótesis nula de normalidad ya que en todos los test obtenemos un  $p\text{-value} < 0.05$

Si realizamos la transformación  $\sqrt{1/\text{Age\_Tratamiento}}$  tenemos las siguientes gráficas y contrastes





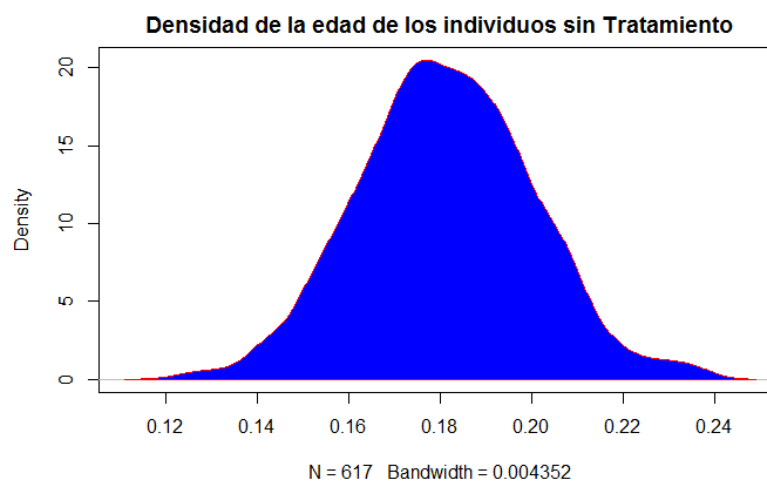


Lilliefors (Kolmogorov-Smirnov) normality test

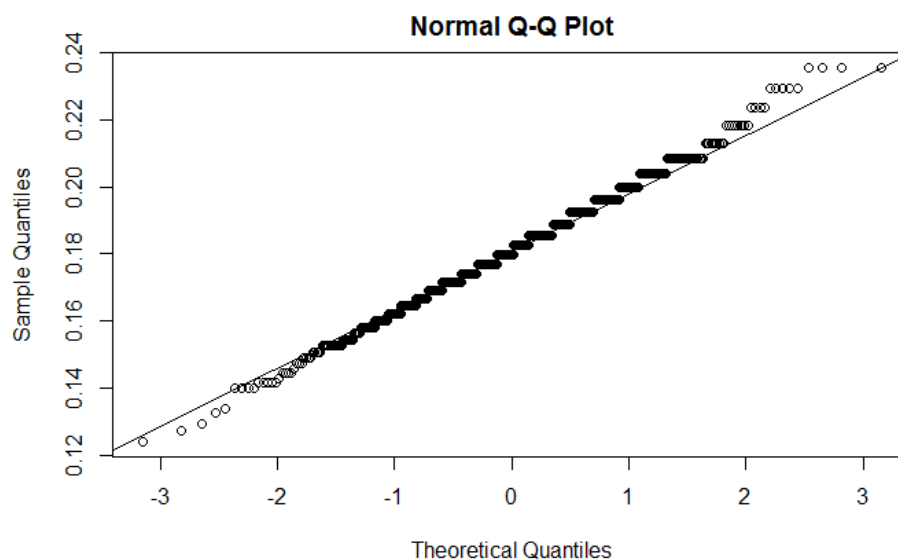
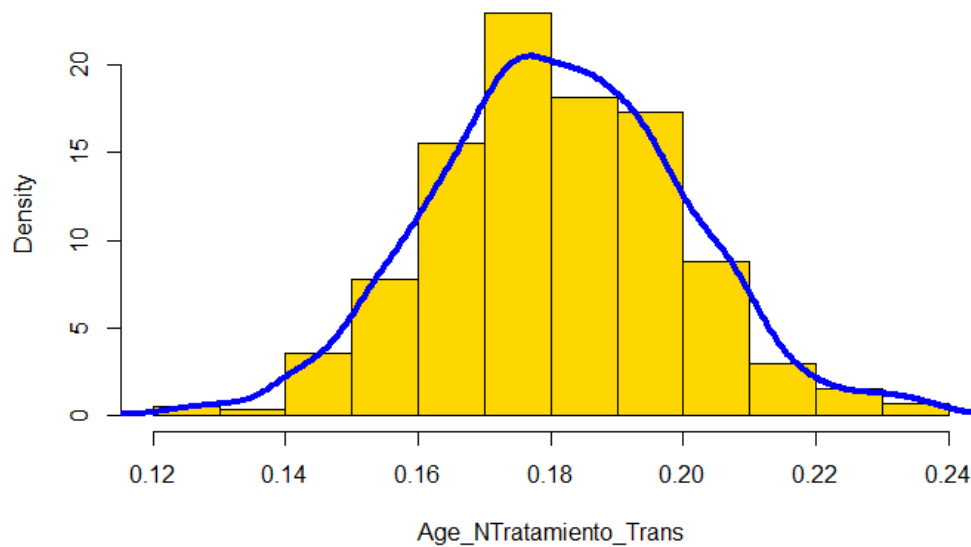
data: (x = Age\_Tratamiento\_Trans)  
D = 0.050462, p-value = 0.0006927

Rechazamos hipótesis nula de normalidad ya que el pvalor<0.05

Realizamos la misma transformación para la variable Age\_NTratamiento  
 $\sqrt{1/\text{Age\_NTratamiento}}$  obteniendo las siguientes gráficas y contrastes



### Histograma y densidad del valor de la edad de los individuos sin Tratamier



Lilliefors (Kolmogorov-Smirnov) normality test

data: (x = Age\_NTratamiento\_Trans)  
D = 0.041727, p-value = 0.01249

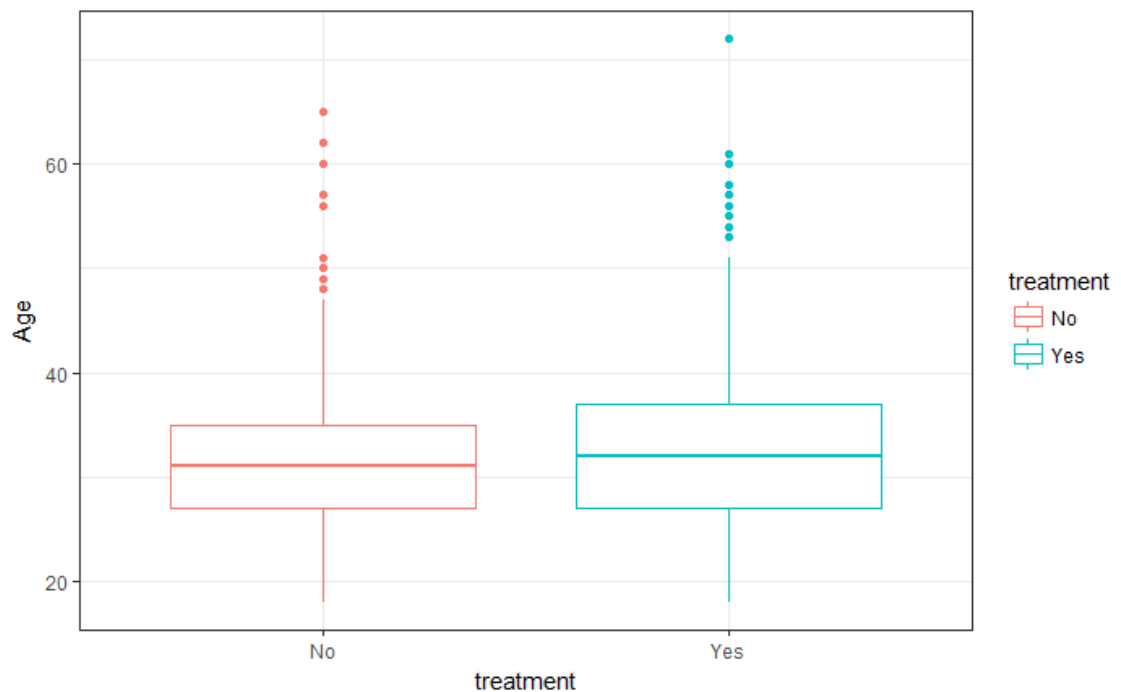
Procedemos a rechazar la hipótesis nula ya que el pvalor<0.05

A continuación, procedemos a estudiar la homogeneidad de la varianza u homocedasticidad.

Bartlett test of homogeneity of variances

```
data: list(Age_Tratamiento, Age_NTratamiento)
Bartlett's K-squared = 4.2709, df = 1, p-value = 0.03877
```

Debido a que el  $p\text{-value} < 0.05$  podemos concluir que si hay diferencias significativas entre las varianzas de los dos grupos. Grafica también lo podemos observar



Vamos a realizar el estudio ANOVA de un factor (one-way ANOVA o independent samples ANOVA) para investigar si existen diferencias en la edad entre los individuos que han tenido tratamiento o no de salud mental.

```
Call:
aov(formula = fit)
```

```
Terms:
          treatment Residuals
Sum of Squares    412.65  65439.09
Deg. of Freedom         1    1238
```

```
Residual standard error: 7.2704
Estimated effects may be unbalanced
```

```
Df Sum Sq Mean Sq F value Pr(>F)
treatment    1    413   412.6   7.807 0.00529 **
Residuals 1238  65439    52.9
---
```

signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

En este caso hemos encontrado cambio significativo de la variable treatment (si han recibido o no tratamiento relacionado con la salud mental) ya que el pvalor ha sido menor que 0.05

Continuamos con el estudio ANOVA de un factor (one-way ANOVA o independent samples ANOVA) para investigar si existen diferencias en la edad entre los individuos que consideran que las empresas le dan más importancia a la salud mental vs salud física.

```
Call:
aov(formula = fit2)

Terms:
              mental_vs_physical Residuals
Sum of Squares              61.69 65790.05
Deg. of Freedom                2    1237

Residual standard error: 7.292816
Estimated effects may be unbalanced
```

Df	Sum Sq	Mean Sq	F value	Pr(>F)
mental_vs_physical	2	62	30.85	0.58
Residuals	1237	65790	53.19	0.56

En este caso no hemos encontrado ningún cambio significativa de la variable mental\_vs\_physical ya que el pvalor ha sido mayor que 0.05

## Ejercicio 5

**Resolución del problema. A partir de los resultados obtenidos. ¿cuáles son las conclusiones? ¿Los resultados permiten responder al problema?**

Una primera aproximación a la respuesta de nuestra pregunta, podría ser la que viene derivada de la comparación en un contexto de entrevista laboral, de trabajo entre la enfermedad física y mental, así como la comparativa de consecuencias negativas por hablar de saludo mental con el empleador.

Si tenemos en cuenta que los encuestados consideran que las respuestas sobre consecuencias de hablar de salud mental vs física con el empleador son (quizá 477, No 490, Sí 292) versus (quizá 273, no 925 y si 61), para el ámbito mental y físico, respectivamente y que asimismo los encuestados en una entrevista de trabajo y en relación con la salud mental vs física son (quizá 207, No 1008, Sí 44) versus (Quizá

557, No 500, Si 202), para el ámbito mental y físico, respectivamente, tal y como se aprecia en los gráficos previos, podríamos inducir cierta condición de estigmatización de la enfermedad mental, versus enfermedad física. Por otra parte, el grado de incertidumbre que los datos presentan en muchas de las respuestas de los encuestados apuntarían a la falta de conocimiento del sujeto; es decir, no solo a aspectos externos derivados del entorno de trabajo o la organización, sino también a aspectos de voluntad o interés asociados al propio sujeto encuestado, lo cual podría ser un indicador más a tener en cuenta.

Sin embargo, es necesario ser sumamente cautelosos al extraer esta conclusión, ya que otros factores que no aparecen en la encuesta podrían ser determinantes, como el tipo de trabajo objetivo, el perfil del encuestado y el tipo de trabajo que desempeña, etc. ya que esto podría condicionar la respuesta. Si además, tenemos en cuenta las observaciones realizadas a lo largo de la práctica sobre los datos, como la relativa al número de observaciones por país, o la relativa a la variable `work_interfere`, debemos concluir que los resultados no permiten responder al problema o pregunta planteada, por varios motivos, a los que hemos llegado a partir del análisis de los datos. Entendemos que existen errores importantes en el diseño del proceso para la obtención de los datos, debido a la falta de completitud del mismo. Derivado de ello, y a partir de los datos extraídos, el conjunto de datos presenta un problema de validez de la muestra y no permite responder a la pregunta.

## Ejercicio 6

**Código:** Hay que adjuntar el código, preferiblemente en R, con el que se ha realizado la limpieza, análisis y representación de los datos. Si lo preferís, también podéis trabajar en Python.

```
```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = TRUE)
```
```

```
```{r load_libraries, include=FALSE}
library(knitr)
```
```

```
```{r}
library(rockchalk)
library(nortest)
library(normtest)
library(moments)
library(car)
```

```
library(ggplot2)
library(dplyr)
...

```{r read}
surveyMentalHealth<-read.csv("survey.csv", sep=",",na.strings = "NA")
#Mostramos las primeras filas
head(surveyMentalHealth)
...

```{r}
#Eliminación de las variables
surveyMentalHealth$Timestamp<-NULL
surveyMentalHealth$state<-NULL
surveyMentalHealth$comments<-NULL
#summary(surveyMentalHealth)
...

```{r}
colnames(surveyMentalHealth)
...

```{r str}
str(surveyMentalHealth)
...

```{r }
#Número de fila del fichero.
nrow(surveyMentalHealth)
...

```{r}
sapply(surveyMentalHealth,class)
...

```

```

```{r}
#En el momento de la lectura del fichero establecemos
#que si se encuentra un valor perdido los asigne por NA (na.strings = "NA")
table(is.na(surveyMentalHealth$Age))
```

```{r}
table(is.na(surveyMentalHealth$Gender))
```

```{r}
table(is.na(surveyMentalHealth$Country ))
```

```{r}
table(is.na(surveyMentalHealth$self_employed))
```

```{r}
table(is.na(surveyMentalHealth$family_history))
```

```{r}
table(is.na(surveyMentalHealth$treatment))
```

```{r}
table(is.na(surveyMentalHealth$work_interfere))
```

```{r}
table(is.na(surveyMentalHealth$no_employees))
```

```{r}
table(is.na(surveyMentalHealth$remote_work))
```

```{r}

```



```

table(is.na(surveyMentalHealth$tech_company))
```

```{r}
table(is.na(surveyMentalHealth$benefits))
```

```{r}
table(is.na(surveyMentalHealth$care_options))
```

```{r}
table(is.na(surveyMentalHealth$wellness_program))
```

```{r}
table(is.na(surveyMentalHealth$seek_help))
```

```{r}
table(is.na(surveyMentalHealth$anonymity))
```

```{r}
table(is.na(surveyMentalHealth$leave))
```

```{r}
table(is.na(surveyMentalHealth$mental_health_consequence))
```

```{r}
table(is.na(surveyMentalHealth$phys_health_consequence))
```

```{r}
table(is.na(surveyMentalHealth$coworkers))
```

```{r}

```

```

table(is.na(surveyMentalHealth$supervisor))
```

```{r}
table(is.na(surveyMentalHealth$mental_health_interview))
```

```{r}
table(is.na(surveyMentalHealth$phys_health_interview))
```

```{r}
table(is.na(surveyMentalHealth$mental_vs_physical))
```

```{r}
table(is.na(surveyMentalHealth$obs_consequence))
```

```{r}
surveyMentalHealth$Age[which(surveyMentalHealth$Age==0) ]

```

**Age**

```{r}
barplot(table(surveyMentalHealth$Age))
```

```{r}
datos_atipicos<-subset(surveyMentalHealth[1:1], surveyMentalHealth$Age<16 |
surveyMentalHealth$Age>75)
datos_atipicos
```

```

```
```{r}
surveyMentalHealth_clean<-subset(surveyMentalHealth, surveyMentalHealth$Age>16
& surveyMentalHealth$Age<75)
nrow(surveyMentalHealth_clean)
```
```

```
```{r}
summary(surveyMentalHealth_clean$Age)
```
```

```
```{r}
boxplot(surveyMentalHealth_clean$Age)
```
```

```
```{r}
#Media Aritmetica
Media_Age<-mean(surveyMentalHealth_clean$Age)
#Mediana
Mediana_Age<-median(surveyMentalHealth_clean$Age)
#Media Recortada
Media_Recortada_Age<-mean(surveyMentalHealth_clean$Age, trim=0.05)
#Desviación estándar
Desviacion_estandar_Age<-sd(surveyMentalHealth_clean$Age)
#Rango Intercuartilico (RIC)
RIC_Age<-IQR(surveyMentalHealth_clean$Age)
#Desviación Absoluta Respecto de la Mediana
Desviacion_Absoluta_Mediana_Age<-mad(surveyMentalHealth_clean$Age)
#Tabla
kable(rbind(Media_Age,Mediana_Age,Media_Recortada_Age,
            Desviacion_estandar_Age,RIC_Age,
            Desviacion_Absoluta_Mediana_Age))
```

**Gender**
```

```
```{r}
levels(surveyMentalHealth_clean$Gender)
```
```

```
``{r}
surveyMentalHealth_clean$Gender<-as.character(surveyMentalHealth_clean$Gender)
surveyMentalHealth_clean$Gender<-
replace(surveyMentalHealth_clean$Gender,surveyMentalHealth_clean$Gender=="Cis
Female"|surveyMentalHealth_clean$Gender=="cis-
female/femme"|surveyMentalHealth_clean$Gender=="f"|surveyMentalHealth_clean$G
ender=="femail"|surveyMentalHealth_clean$Gender=="Femake"|surveyMentalHealth_
clean$Gender=="female"|surveyMentalHealth_clean$Gender=="Female"|surveyMental
Health_clean$Gender=="Female "|surveyMentalHealth_clean$Gender=="Female
(cis)"|surveyMentalHealth_clean$Gender=="Female
(trans)"|surveyMentalHealth_clean$Gender=="Woman"|surveyMentalHealth_clean$Ge
nder=="woman","F")
```

```
surveyMentalHealth_clean$Gender<-
replace(surveyMentalHealth_clean$Gender,surveyMentalHealth_clean$Gender=="Mal
e"|surveyMentalHealth_clean$Gender=="male"|surveyMentalHealth_clean$Gender=="
Cis
Male"|surveyMentalHealth_clean$Gender=="m"|surveyMentalHealth_clean$Gender=="
Mail"|surveyMentalHealth_clean$Gender=="Make"|surveyMentalHealth_clean$Gender
=="male leaning
androgynous"|surveyMentalHealth_clean$Gender=="Malr"|surveyMentalHealth_clean$
Gender=="msle"|surveyMentalHealth_clean$Gender=="ostensibly male, unsure what
that really means"|surveyMentalHealth_clean$Gender=="something kinda
male?"|surveyMentalHealth_clean$Gender=="Androgyne"|surveyMentalHealth_clean$
Gender=="cis male"|surveyMentalHealth_clean$Gender=="Cis
Man"|surveyMentalHealth_clean$Gender=="Guy (-ish)
^_^"|surveyMentalHealth_clean$Gender=="maile"|surveyMentalHealth_clean$Gender=
=="Mal"|surveyMentalHealth_clean$Gender=="Male
(CIS)"|surveyMentalHealth_clean$Gender=="Male-
ish"|surveyMentalHealth_clean$Gender=="Man"|surveyMentalHealth_clean$Gender==
"Male ","M")
```

```
surveyMentalHealth_clean$Gender<-
replace(surveyMentalHealth_clean$Gender,surveyMentalHealth_clean$Gender=="A
little about
you"|surveyMentalHealth_clean$Gender=="Agender"|surveyMentalHealth_clean$Gend
er=="All"|surveyMentalHealth_clean$Gender=="Enby"|surveyMentalHealth_clean$Gen
der=="fluid"|surveyMentalHealth_clean$Gender=="Genderqueer"|surveyMentalHealth_
clean$Gender=="Nah"|surveyMentalHealth_clean$Gender=="Neuter"|surveyMentalHe
alth_clean$Gender=="non-
binary"|surveyMentalHealth_clean$Gender=="p"|surveyMentalHealth_clean$Gender==
"queer"|surveyMentalHealth_clean$Gender=="queer/she/they"|surveyMentalHealth_cle
```

```
an$Gender=="Trans woman"|surveyMentalHealth_clean$Gender=="Trans-
female",NA)
```

```
surveyMentalHealth_clean$Gender<-as.factor(surveyMentalHealth_clean$Gender)
```
```

```
```{r}
levels(surveyMentalHealth_clean$Gender)
```
```

```
```{r}
surveyMentalHealth_clean<-subset(surveyMentalHealth_clean,
surveyMentalHealth_clean$Gender!="NA")
nrow(surveyMentalHealth_clean)
```
```

**\*\*Country\*\***

```
```{r}
levels(surveyMentalHealth_clean$Country)
summary(surveyMentalHealth_clean$Country)
```
```

```
```{r}
barplot(table(surveyMentalHealth_clean$Country),
main="Países")
```
```

**\*\*family\_history\*\***

```
```{r}
levels(surveyMentalHealth_clean$family_history)
summary(surveyMentalHealth_clean$family_history)

```
```

```

```{r}
barplot(table(surveyMentalHealth_clean$family_history),
        main="Antecedentes familiares")
```

**treatment**

```{r}
levels(surveyMentalHealth_clean$treatment)
summary(surveyMentalHealth_clean$treatment)

```

```{r}
barplot(table(surveyMentalHealth_clean$treatment),
        main="Ha sido tratado de alguna enfermedad mental")
```

**work_interfere**

```{r}
levels(surveyMentalHealth_clean$work_interfere)
summary(surveyMentalHealth_clean$work_interfere)

```

```{r}
barplot(table(surveyMentalHealth_clean$work_interfere),
        main="La enfermedad mental interfiere en su trabajo")
```

**no_employees**

```{r}
levels(surveyMentalHealth_clean$no_employees)
summary(surveyMentalHealth_clean$no_employees)

```

```{r}
barplot(table(surveyMentalHealth_clean$no_employees),
        main="Número de empleados de la compañía u organizacion")

```

...

**\*\*remote\_work\*\***

```{r}

levels(surveyMentalHealth\_clean\$remote\_work)

summary(surveyMentalHealth\_clean\$remote\_work)

...

```{r}

barplot(table(surveyMentalHealth\_clean\$remote\_work),  
main="Teletrabajo al menos el 50% del tiempo")

...

**\*\*tech\_company\*\***

```{r}

levels(surveyMentalHealth\_clean\$tech\_company)

summary(surveyMentalHealth\_clean\$tech\_company)

...

```{r}

barplot(table(surveyMentalHealth\_clean\$tech\_company),  
main="La Organización es Tecnológica")

...

**\*\*benefits\*\***

```{r}

levels(surveyMentalHealth\_clean\$benefits)

summary(surveyMentalHealth\_clean\$benefits)

...

```{r}

barplot(table(surveyMentalHealth\_clean\$benefits),  
main="La Organización provee de beneficios de salud Mental")

...

**\*\*care\_options\*\***

```

```{r}
levels(surveyMentalHealth_clean$care_options)
summary(surveyMentalHealth_clean$care_options)

...

```{r}
barplot(table(surveyMentalHealth_clean$care_options),
          main="Conoce Opciones de cuidado mental de su compañía médica")
...

**wellness_program**

```{r}
levels(surveyMentalHealth_clean$wellness_program)

...

```{r}
summary (surveyMentalHealth_clean$wellness_program)

...

```{r}
barplot(table(surveyMentalHealth_clean$wellness_program),
          main="Información del conocimiento de programas especificos")
...

**seek_help**

```{r}
levels(surveyMentalHealth_clean$seek_help)
...

```{r}
summary (surveyMentalHealth_clean$seek_help)

```



```

```{r}

```
barplot(table(surveyMentalHealth_clean$seek_help),
  main="Información de recursos y ayuda desde la organización")
```

```

**\*\*anonymity\*\***

```{r}

```
levels(surveyMentalHealth_clean$anonymity)
```

```

```{r}

```
summary (surveyMentalHealth_clean$anonymity)
```

```

```{r}

```
barplot(table(surveyMentalHealth_clean$anonymity),
  main="Privacidad de beneficios sobre enfermedades mentales")
```

```

**\*\*leave\*\***

```{r}

```
levels(surveyMentalHealth_clean$leave)
```

```

```{r}

```
summary (surveyMentalHealth_clean$leave)
```

```

```{r }

```
barplot(table(surveyMentalHealth_clean$leave),
  main="Posibilidad de baja en enfermedades mentales")
```

```
'''
```

```
**mental_health_consequence**
```

```
'''{r}
```

```
levels(surveyMentalHealth_clean$mental_health_consequence)
```

```
'''
```

```
'''{r}
```

```
summary (surveyMentalHealth_clean$mental_health_consequence)
```

```
'''
```

```
'''{r}
```

```
barplot(table(surveyMentalHealth_clean$mental_health_consequence),  
          main="Consecuencias por hablar de salud mental")
```

```
'''
```

```
**phys_health_consequence**
```

```
'''{r}
```

```
levels(surveyMentalHealth_clean$phys_health_consequence)
```

```
'''
```

```
'''{r}
```

```
summary (surveyMentalHealth_clean$phys_health_consequence)
```

```
'''
```

```
'''{r}
```

```
barplot(table(surveyMentalHealth_clean$phys_health_consequence),  
          main="Consecuencias por hablar de salud física")
```

```
'''
```

```
**coworkers**
```

```

```{r}
levels(surveyMentalHealth_clean$coworkers)
```

```{r}
summary (surveyMentalHealth_clean$coworkers)

```

```{r}
barplot(table(surveyMentalHealth_clean$mental_health_consequence),
          main="Hablaria de salud mental con compañeros")
```

**supervisor**

```{r}
levels(surveyMentalHealth_clean$supervisor)
```

```{r}
summary (surveyMentalHealth_clean$supervisor)

```

```{r}
barplot(table(surveyMentalHealth_clean$supervisor),
          main="Hablaria de salud mental con su jefe")
```

**mental_health_interview**

```{r}
levels(surveyMentalHealth_clean$mental_health_interview)
```

```{r}
summary (surveyMentalHealth_clean$mental_health_interview)

```

```

```

```{r}
barplot(table(surveyMentalHealth_clean$mental_health_interview),
  main="Hablaria de salud mental en una entrevista laboral")
```

**phys_health_interview**

```{r}
levels(surveyMentalHealth_clean$phys_health_interview)
```

```{r}
summary (surveyMentalHealth_clean$phys_health_interview)

```

```{r}
barplot(table(surveyMentalHealth_clean$mental_health_interview),
  main="Hablaria de salud física en una entrevista laboral")
```

**mental_vs_physical**

```{r}
levels(surveyMentalHealth_clean$mental_vs_physical)
```

```{r}
summary (surveyMentalHealth_clean$mental_vs_physical)

```

```{r }
barplot(table(surveyMentalHealth_clean$mental_vs_physical),
  main="Importacia en la Organizacion de la salud mental sobre la física")
```

**obs_consequence**

```{r}

```

```

levels(surveyMentalHealth_clean$obs_consequence)
```

```{r}
summary (surveyMentalHealth_clean$obs_consequence)
```

```{r}
barplot(table(surveyMentalHealth_clean$obs_consequence),
  main="Consecuencias laboral por padecer enfermedad mental ")
```

```{r}
#Gráfica comparativa 1
par(mfrow=c(1,2))
barplot(table(surveyMentalHealth_clean$mental_health_interview),
  main="Mención mental en entrevista laboral")
barplot(table(surveyMentalHealth_clean$phys_health_interview),
  main="Mención física en entrevista laboral")
```

```{r}
#Gráfica comparativa 2
par(mfrow=c(1,2))
barplot(table(surveyMentalHealth_clean$mental_health_consequence),
  main="Mental: ¿consecuencias negativas?")
barplot(table(surveyMentalHealth_clean$phys_health_consequence),
  main="Física: ¿consecuencias negativas?")
```

```{r}
#Conjunto de valores de la variable Edad que consideran que la organización da
mayor importancia a la salud mental vs salud física.
Age_Mental<-subset(surveyMentalHealth_clean$Age,
surveyMentalHealth_clean$mental_vs_physical=="Yes")
```

```

```

```{r}
Age_Fisica<-subset(surveyMentalHealth_clean$Age,
surveyMentalHealth_clean$mental_vs_physical=="No")
```

```{r}
#Calculo Media
mean(Age_Mental)
mean(Age_Fisica)
```

```{r}
#Calculo Mediana
median(Age_Mental)
median(Age_Fisica)
```

```{r}
#Sumario de los cinco números (Mínimo, Q1, Mediana, Q3, Maximo)
fivenum(Age_Mental)
```

```{r}
#Sumario de los cinco números (Mínimo, Q1, Mediana, Q3, Maximo)
fivenum(Age_Fisica)
```

```{r}
#Diagrama de caja (Boxplot)
boxplot(Age_Mental, main="Box Plot Age Salud Mental")
```

```{r}
#Diagrama de caja (Boxplot)
boxplot(Age_Fisica, main="Box Plot Age Salud Fisica")
```

```{r}

```

```
#Calculamos el numero de intervalor
k_Age_Mental<- round(sqrt(length(Age_Mental)))
k_Age_Mental
```

```{r}
#Calculamos el numero de intervalor
k_Age_Fisica<- round(sqrt(length(Age_Fisica)))
k_Age_Fisica
```

```{r}
hist(Age_Mental ,main="Edad individuos Organizacion mas valor a la Salud Mental",
      breaks=k_Age_Mental, col="blue")
```

```{r}
hist(Age_Fisica ,main="Edad individuos con Organizacion mas valor a la Salud Fisica",
      breaks=k_Age_Fisica, col="blue")
```

```{r}
hh_Age_Mental<-hist(Age_Mental ,main="Edad individuos con Organizacion mas valor
a la Salud Mental",
      breaks=k_Age_Mental, col="blue")
hh_Age_Mental
```

```{r}
hh_Age_Fisica<-hist(Age_Fisica ,main="Edad individuos con Organizacion mas valor a
la Salud Fisica",
      breaks=k_Age_Fisica, col="blue")
hh_Age_Fisica
```

```{r}
#Calculo de la función de densidad
den_Age_Mental<- density(Age_Mental)
```

```

plot(den_Age_Mental ,main="Edad Individuos con Organización mas valor a la Salud
Mental")
polygon(den_Age_Mental , col="blue", border="red")

...

```{r}
#Calculo de la función de densidad
den_Age_Fisica<- density(Age_Fisica)
plot(den_Age_Fisica ,main="Edad Individuos con Organización mas valor a la Salud
Fisica")
polygon(den_Age_Fisica , col="blue", border="red")

...

```{r}
#Superposición de las gráficas
hist(Age_Mental ,main="Edad individuos con Organización mas valor a la Salud
Mental",
      col="gold",freq=FALSE)
lines(den_Age_Mental,col="blue",lwd=4)
...

```{r}
#Superposición de las gráficas
hist(Age_Fisica ,main="Edad Individuos con Organización más valor a la Salud Fisica",
      col="gold",freq=FALSE)
lines(den_Age_Mental ,col="blue",lwd=4)
...

```{r}
qqnorm(Age_Mental)
qqline(Age_Mental)
...

```{r}
qqnorm(Age_Fisica)
qqline(Age_Fisica)
...

```



```

```{r}
ks.test(x=Age_Mental,"pnorm", mean(Age_Mental), sd(Age_Mental))
```

```{r}
ks.test(Age_Fisica,"pnorm", mean(Age_Fisica), sd(Age_Fisica))
```

```{r}

lillie.test((x=Age_Mental))
```

```{r}

lillie.test((x=Age_Fisica))
```

```{r}
jb.norm.test(x=Age_Mental)
```

```{r}
jb.norm.test(x=Age_Fisica)
```

```{r}
Age_Mental_Trans<-(sqrt(sqrt(1/Age_Mental)))
```

```{r}
#Calculo de la función de densidad
den_Age_Mental_Trans<- density(Age_Mental_Trans)
plot(den_Age_Mental_Trans ,main="Densidad de la edad de los individuos
observados")
polygon(den_Age_Mental_Trans , col="blue", border="red")

```

```

```

```{r}
#Superposición de las gráficas
hist(Age_Mental_Trans ,main="Histograma y densidad del valor de la edad de los
individuos observados",
      col="gold",freq=FALSE)
lines(den_Age_Mental_Trans ,col="blue",lwd=4)
```

```{r}
qqnorm(Age_Mental_Trans)
qqline(Age_Mental_Trans)
```

```{r}

lillie.test((x=Age_Mental_Trans))
```

```{r}
Age_Fisica_Trans<-(sqrt(sqrt(1/Age_Fisica)))
```

```{r}
#Calculo de la función de densidad
den_Age_Fisica_Trans<- density(Age_Fisica_Trans)
plot(den_Age_Fisica_Trans ,main="Densidad de la edad de los individuos
observados")
polygon(den_Age_Fisica_Trans , col="blue", border="red")
```

```{r}
#Superposición de las gráficas
hist(Age_Fisica_Trans ,main="Histograma y densidad del valor de la edad de los
individuos observados",
      col="gold",freq=FALSE)
lines(den_Age_Fisica_Trans ,col="blue",lwd=4)
```

```{r}

```

```
qqnorm(Age_Fisica_Trans)
qqline(Age_Fisica_Trans)
'''

'''{r}

lillie.test((x=Age_Fisica_Trans))
'''

'''{r}
bartlett.test(list(Age_Mental, Age_Fisica))
'''

'''{r}
ggplot(surveyMentalHealth_clean, aes(x = mental_vs_physical, y = Age, colour
=mental_vs_physical)) + geom_boxplot() + theme_bw()
'''

'''{r}
#Conjunto de valores de la variable Edad que consideran que la organización da
mayor importancia a la salud mental vs salud física.
Age_Tratamiento<-subset(surveyMentalHealth_clean$Age,
surveyMentalHealth_clean$treatment=="Yes")
'''

'''{r}
Age_NTratamiento<-subset(surveyMentalHealth_clean$Age,
surveyMentalHealth_clean$treatment=="No")
'''

'''{r}
#Calculo Media
mean(Age_Tratamiento)
mean(Age_NTratamiento)
'''
```

```

```{r}
#Calculo Mediana
median(Age_Tratamiento)
median(Age_NTratamiento)
```

```{r}
#Sumario de los cinco números (Mínimo, Q1, Mediana, Q3, Maximo)
fivenum(Age_Tratamiento)
```

```{r}
#Sumario de los cinco números (Mínimo, Q1, Mediana, Q3, Maximo)
fivenum(Age_NTratamiento)
```

Gráfico de Boxplot
```{r}
#Diagrama de caja (Boxplot)
boxplot(Age_Tratamiento, main="Box Plot Age Con Tratamiento")
```

```{r}
#Diagrama de caja (Boxplot)
boxplot(Age_NTratamiento, main="Box Plot Age Sin Tratamiento")
```

```{r}
#Calculamos el numero de intervalo
k_Age_Tratamiento<- round(sqrt(length(Age_Tratamiento)))
k_Age_Tratamiento
```

```{r}
#Calculamos el numero de intervalo
k_Age_NTratamiento<- round(sqrt(length(Age_NTratamiento)))
k_Age_NTratamiento

```

```
```
```

```
```{r}
hist(Age_Tratamiento ,main="Edad individuos con Tratamiento Salud Mental",
      breaks=k_Age_Tratamiento, col="blue")
```
```

```
```{r}
hist(Age_NTratamiento ,main="Edad individuos sin Tratamiento Salud Mental",
      breaks=k_Age_NTratamiento, col="blue")
```
```

```
```{r}
hh_Age_Tratamiento<-hist(Age_Tratamiento ,main="Edad individuos con Tratamiento
Salud Mental",
      breaks=k_Age_Tratamiento, col="blue")
hh_Age_Tratamiento
```
```

```
```{r}
hh_Age_NTratamiento<-hist(Age_NTratamiento ,main="Edad individuos sin
Tratamiento Salud Mental",
      breaks=k_Age_NTratamiento, col="blue")
hh_Age_NTratamiento
```
```

```
```{r}
#Calculo de la función de densidad
den_Age_Tratamiento<- density(Age_Tratamiento)
plot(den_Age_Tratamiento ,main="Densidad de la edad de los individuos con
Tratamiento")
polygon(den_Age_Tratamiento , col="blue", border="red")
```
```

```
```{r}
#Calculo de la función de densidad
den_Age_NTratamiento<- density(Age_NTratamiento)
```

```

plot(den_Age_NTratamiento ,main="Densidad de la edad de los individuos observados
sin Tratamiento")
polygon(den_Age_NTratamiento , col="blue", border="red")

...

```{r}
#Superposición de las gráficas
hist(Age_Tratamiento ,main="Edad individuos con Tratamiento de Salud Mental",
      col="gold",freq=FALSE)
lines(den_Age_Tratamiento,col="blue",lwd=4)
...

```{r}
#Superposición de las gráficas
hist(Age_NTratamiento ,main="Edad Individuos sin tratamiento de Salud Mental",
      col="gold",freq=FALSE)
lines(den_Age_NTratamiento ,col="blue",lwd=4)
...

```{r}
qqnorm(Age_Tratamiento)
qqline(Age_Tratamiento)
...

```{r}
qqnorm(Age_NTratamiento)
qqline(Age_NTratamiento)
...

```{r}
ks.test(x=Age_Tratamiento,"pnorm", mean(Age_Tratamiento), sd(Age_Tratamiento))
...

```{r}
ks.test(Age_NTratamiento,"pnorm", mean(Age_NTratamiento), sd(Age_NTratamiento))
...

```{r}

```

```

lillie.test((x=Age_Tratamiento))
```

```{r}

lillie.test((x=Age_NTratamiento))
```

```{r}
jb.norm.test(x=Age_Tratamiento)
```

```{r}
jb.norm.test(x=Age_NTratamiento)
```

```{r}
Age_Tratamiento_Trans<-((sqrt(1/Age_Tratamiento)))
```

```{r}
#Calculo de la función de densidad
den_Age_Tratamiento_Trans<- density(Age_Tratamiento_Trans)
plot(den_Age_Tratamiento_Trans ,main="Densidad de la edad de los individuos
observados")
polygon(den_Age_Tratamiento_Trans , col="blue", border="red")
```

```{r}
#Superposición de las gráficas
hist(Age_Tratamiento_Trans ,main="Histograma y densidad del valor de la edad de los
individuos observados",
      col="gold",freq=FALSE)
lines(den_Age_Tratamiento_Trans ,col="blue",lwd=4)
```

```{r}
qqnorm(Age_Tratamiento_Trans)

```

```

qqline(Age_Tratamiento_Trans)
'''

'''{r}

lillie.test((x=Age_Tratamiento_Trans))
'''

'''{r}
Age_NTratamiento_Trans<-((sqrt(1/Age_NTratamiento)))
'''

'''{r}
#Calculo de la función de densidad
den_Age_NTratamiento_Trans<- density(Age_NTratamiento_Trans)
plot(den_Age_NTratamiento_Trans ,main="Densidad de la edad de los individuos sin
Tratamiento")
polygon(den_Age_NTratamiento_Trans , col="blue", border="red")

'''

'''{r}
#Superposición de las gráficas
hist(Age_NTratamiento_Trans ,main="Histograma y densidad del valor de la edad de
los individuos sin Tratamiento",
      col="gold",freq=FALSE)
lines(den_Age_NTratamiento_Trans ,col="blue",lwd=4)
'''

'''{r}
qqnorm(Age_NTratamiento_Trans)
qqline(Age_NTratamiento_Trans)
'''

'''{r}

lillie.test((x=Age_NTratamiento_Trans))
'''

```



```

```{r}
bartlett.test(list(Age_Tratamiento, Age_NTratamiento))
```

```{r}
ggplot(surveyMentalHealth_clean, aes(x = treatment, y = Age, colour = treatment)) +
  geom_boxplot() + theme_bw()
```

```{r}
fit=lm(Age~ treatment, surveyMentalHealth_clean)
aov(fit)
```

```{r}
summary(aov(fit))
```

```{r}
fit2=lm(Age~ mental_vs_physical, surveyMentalHealth_clean)
aov(fit2)
```

```{r}
summary(aov(fit2))
```

```{r}
#Se guardan los cambios realizados en el fichero rodriguez_inmuebles_clean.csv
write.csv(surveyMentalHealth_clean, file="survey_clean.csv")
```

```

**Realizada por:** Carlos E. Jimenez Gomez  
M<sup>a</sup> Sonia Rodríguez Cepedano