



Predicción de la quema de calorías durante el ejercicio

19 de febrero de 2023



Uso de técnicas de Machine Learning: Método Jerárquico + XGB Regressor

AUTOR: Robledo, Nallely.^{a,*}

^a Universidad Autónoma de Nuevo León, Facultad de Fisicomatemáticas
Pedro de Alba S/N, Niños Héroes, Ciudad Universitaria, San Nicolás de los Garza, N.L.

To cite this article: Robledo, Nallely. 2023. Calorie burn prediction during exercise. Universidad Autónoma de Nuevo León, Facultad de Ciencias Fisicomatemáticas 00, 1-5. <https://doi.org/10.4995/riai.2020.7133>

Resumen

Si la gente respondiera honestamente a la pregunta '¿Cuáles son las razones por las que haces ejercicio?', una respuesta frecuente sería quemar calorías. De hecho, según el Departamento de Salud y Servicios Humanos de EE. UU. (1992), el 26 por ciento de los adultos estadounidenses entre 20 y 74 años tienen sobrepeso, lo que demuestra claramente el impacto de esta preocupación nacional.

Se sabe que la reducción de la grasa corporal puede revertir varios procesos de enfermedades (p. ej., diabetes tipo II, enfermedades cardíacas, etc.), el ejercicio aumenta el gasto calórico total y también maximiza la pérdida de grasa corporal y el mantenimiento o aumento de masa muscular, la participación en el ejercicio es una estrategia muy consecuente y gratificante para perder grasa corporal y mejorar su salud.

El ejercicio como medio para quemar calorías ha sido reconocido por la industria del fitness. Hay muchos tipos de modalidades de ejercicio que se comercializan con el reclamo de "quemar más calorías", y el consumidor se pregunta qué es lo que determina la cantidad de calorías quemadas durante el ejercicio. Esta situación es la razón fundamental para escribir este artículo.

Calorie burn prediction during exercise -

Abstract

If people were to honestly answer the question 'What are the reasons you exercise?' a frequent answer would be to burn calories. In fact, according to the US Department of Health and Human Services (1992), 26 percent of American adults ages 20-74 are overweight, clearly demonstrating the impact of this national concern.

It is known that reducing body fat can reverse various disease processes (eg, type II diabetes, heart disease, etc.), exercise increases total caloric expenditure and also maximizes body fat loss and maintenance or gaining muscle mass, engaging in exercise is a very consistent and rewarding strategy for losing body fat and improving your health.

Exercise as a means of burning calories has been recognized by the fitness industry. There are many types of exercise modalities that are marketed as "burning more calories." and the consumer wonders what determines the number of calories burned during exercise. This situation is the fundamental reason for writing this article.

1. Introducción

La principal directriz para la creación de este artículo es representar la solución a la interrogante: *¿Cuál es el factor de mayor influencia en la quema de calorías durante el ejercicio?*

Se evalúan diferentes variables que involucran las características físicas de la persona que realiza el ejercicio, contra otras variables, que describen la manera de hacer el ejercicio. Las

variables de estudio serán: Edad, Sexo, Peso, Estatura, Pulso cardíaco, Temperatura y Duración.

En reposo, el cuerpo gasta energía para mantener las funciones de las células que son esenciales para la vida. El bombeo continuo de sangre por parte del corazón exige energía, al igual que la ventilación continua (movimiento de aire hacia adentro y hacia afuera) de los pulmones. Además, mantener un entorno de soporte vital dentro y alrededor de las células

* Autor para correspondencia: autor1@ceautomatica.es
Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0)

requiere una descomposición constante de ciertas moléculas liberadoras de energía. Esta energía también se utiliza para formar las moléculas necesarias para reparar las células, almacenar energía (glucógeno y triglicéridos), combatir infecciones y procesar los nutrientes obtenidos de la digestión. Estas funciones exigentes de energía se combinan para formar la tasa metabólica basal del cuerpo, que puede variar de aproximadamente 800 a 1500 Kcal dependiendo del tamaño del cuerpo y la ingesta calórica total (cantidad ingerida de alimentos)..

2. Marco teórico

El trifosfato de adenosina (ATP) es la molécula principal que el cuerpo utiliza como medio para utilizar la energía química para realizar el trabajo celular. El ejercicio aumenta el gasto calórico del cuerpo, ya que la contracción muscular implica la necesidad de formar y descomponer ATP repetidamente. La energía liberada por la descomposición del ATP alimenta la contracción del músculo esquelético, lo que aumenta las demandas de energía del cuerpo y aumenta el gasto calórico. Las investigaciones han demostrado que durante el ejercicio el aumento del gasto calórico se debe casi en su totalidad a la contracción del músculo esquelético; el equilibrio se debe a un aumento en las demandas de energía del corazón y los músculos utilizados durante la ventilación.

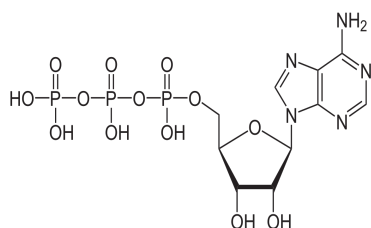


Figura 1: Estructura química del trifosfato.

La investigación evaluará las variables en consideración para al quema de calorías y se definirá a través de modelos matemáticos los de mayor influencia en la quema de calorías. Al final de esta investigación se concluirá de manera lógica si estas variables están influidas por el efecto de el ATP.

Los modelos matemáticos de Machine Learning seleccionados para el análisis: Método Jerárquico y XGBoost, están sustentados bajo la misma teoría de clasificación; que es la de los árboles de decisión. Esto, apesar de que siguen procedimientos totalmente adversos.

En general, los árboles de decisión clasifican datos a partir de su separación en regiones y obtienen una clasificación a partir de las cotas que limitan las regiones. Una vez obtenidas dichas regiones, la función de predicción.

3. Machine Learning

Esta documento hace una recopilación de conjuntos apropiados para enseñar a nuestros modelos de aprendizaje automático para que logre saber cuál es la cantidad de calorías que el individuo gasta para quemar.

Usaremos el método jerárquico y Linear Regression como modelos de aprendizaje automático para comparar y luego evaluar estos modelos. La herramienta es Google Colab, el cual es un servicio basado en la nube.

3.1. Selección de características

La primera fase del análisis de características, se presentará a través de un ANOVA, esta herramienta es de gran utilidad, pues es una fórmula estadística que sirve para comparar las varianzas entre las medias (o el promedio) de diferentes grupos. También se utiliza para determinar si existe alguna diferencia entre las medias de los diferentes grupos.

A continuación se presentan los resultados del valor de F, ordenando de las variables que tienen un valor más alto, al menos representativo.

Tabla 1: Resultados de ANOVA, F-Value por variable

Variable	F Value
Duration	157053.43
Heart Rate	62387.94
Body Temp	31855.44
Age	18.904356
Weight	366.25
Gender	7.5
Height	4.61

En base a esta evaluación se interpreta que un valor F alto, indica alta relación lineal; valores menores, lo contrario. Por lo tanto asumimos que existen 3 variables con mayor relación con la variable de respuesta *Burned Calories* y estas son Duration, Heart Rate y Body Temp.

Se busca una segunda opción para la selección de características, esto con la intención de revisar si existen resultados coincidentes, el método seleccionado es el de información mutua.

La información mutua mide la cantidad de información transferida cuando x^i = (variable de interés) es transmitido y y^i = (variable de respuesta) es recibido.

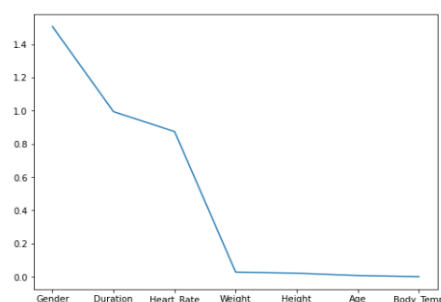


Figura 2: Gráfico de Información mutua entre variables

Con el gráfico anterior podemos observar que las variables Heart Rate, Weight y Body Temp comparten información mutua relevante, en comparación con el resto de las variables.

En base a lo anterior se cree que la eliminación de las variables Age y Gender pudiera ser conveniente para el agrupamiento de características.

3.2. Agrupamiento de características y análisis de grupos

En esta sección se realizará el análisis de grupos o también conocido como clustering, es la tarea de agrupar objetos por similitud, en conjuntos de manera que los miembros del mismo grupo tengan características similares.

Es la tarea principal de la minería de datos exploratoria y es una técnica común en el análisis de datos estadísticos. Se puede realizar a través del aprendizaje automático No Supervisado y el Supervisado.

En primer instancia es importante mencionar que para esta sección del análisis, las variables ya fueron filtradas por medio de la selección de características precedente, por lo que solo se tomarán en cuenta aquellas que aportan mayor información al modelo.

A. Método jerárquico

Los datos son escalados, debido a la diferencia de unidades entre las variables. Para esto se hace uso de la librería `sklearn.preprocessing` importando `StandardScaler`.

■ Definir cantidad de clústers

Se realiza el gráfico de dendrograma, el cual de acuerdo a su estructura muestra los datos en subcategorías que se van dividiendo en otros hasta llegar al nivel de detalle deseado.

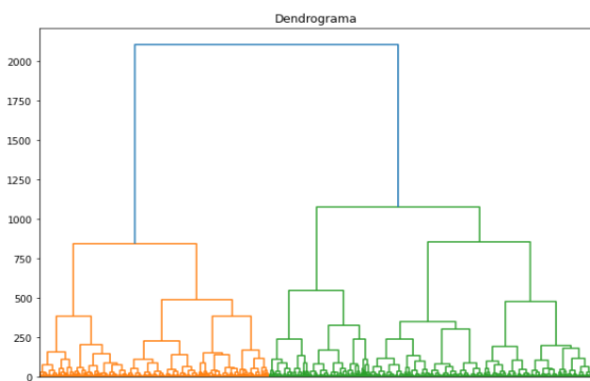


Figura 3: Gráfico Dendrograma

■ Visualización de clústers

Para esta sección se hace uso de la librería `sklearn.cluster` importando la `AgglomerativeClustering`.

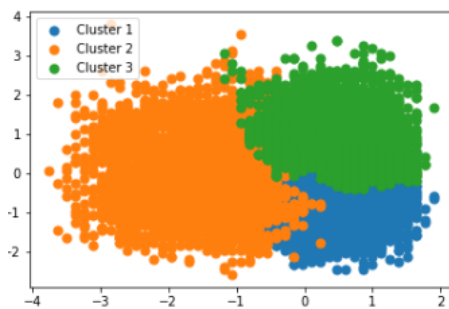


Figura 4: Visualización de clústers en método jerárquico

■ Evaluación de características de los clústers

En el punto anterior, se visualizó la unión de los clusters

en el plano 2D. Posterior a esto, se evalúan las características relevantes de los grupos. Esto nos servirá para entender las comunales de estos modelos. El primer paso es asignar a cada observación su respectivo cluster. Posteriormente se crea un nuevo conjunto de datos que contenga solamente la variable de interés y la etiqueta de cluster. Y por último se selecciona un diagrama de boxplot, es una herramienta muy utilizada para la evaluación de la estadística descriptiva.

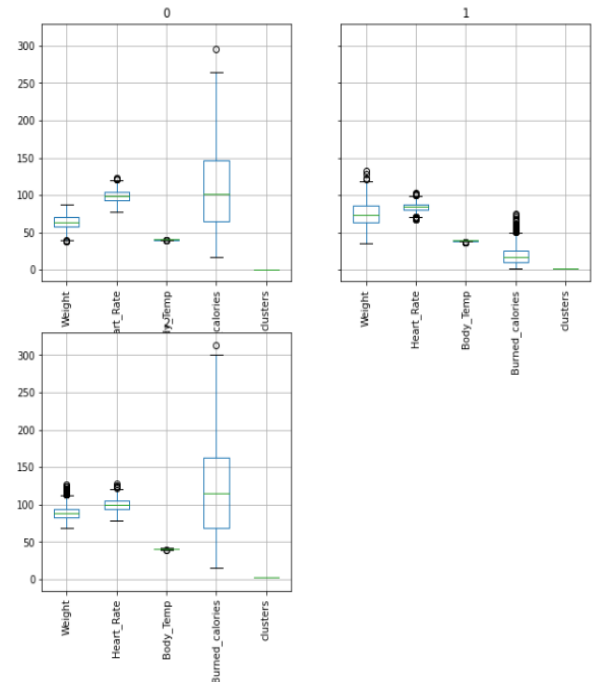


Figura 5: Boxplot de grupos (Método Jerárquico)

■ Desempeño de los modelos o clústers

Para esto se realiza un diseño de experimentos. La herramienta a utilizar es un ANOVA, esta herramienta se explica a detalle previamente ya que también fue utilizada para la Selección de características. Posteriormente una prueba de Tukey's se usa para hacer las comparaciones por pares entre las medias de cada grupo mientras controlamos la tasa de error por familia.

B. Método RGBBoost Regression

Este es un algoritmo de aprendizaje automático supervisado que utiliza una técnica de ensamblado de árboles de decisión para mejorar la precisión de la predicción.

Los datos son escalados nuevamente por su diferencia de unidades, utilizando las mismas herramientas que en el método jerárquico.

Se crean las variables para X y Y de prueba y entrenamiento. En este método se utilizan los errores residuales del modelo inicial para entrenar un segundo modelo. El segundo modelo se enfoca en corregir los errores del primer modelo. Se repite el proceso para cada modelo subsiguiente, utilizando los errores residuales del modelo anterior para entrenar el siguiente modelo.

■ Desempeño del modelo.

Para conocer el desempeño o conocer la exactitud de predicción brindado por el modelo:

Se evalúa la media de error absoluto. MAE es una métrica que mide el promedio de la diferencia absoluta entre las predicciones y los valores reales. Se hace uso de la librería sklearn en la sección de metrics. Se puede expresar matemáticamente como:

$$\text{MAE}(y, \hat{y}) = \frac{\sum_{i=0}^{N-1} |y_i - \hat{y}_i|}{N}$$

El resultado de MAE obtenido da 0.04224.

Después se evalúa la eficacia de un modelo, a través de otra métrica. En este caso utilizamos Kfold, para uso se importo de la librería sklearn.modelselection. En esta técnica, se divide el conjunto de datos en k subconjuntos o "folds" de aproximadamente el mismo tamaño. Posteriormente, el modelo se entrena k veces, cada vez utilizando k-1 subconjuntos para entrenamiento y el subconjunto restante para validación. El resultado de KFold obtenido da un score de 0.91.

Por último se grafica el comportamiento de la variable dependiente de prueba y la de predicción. La visualización se hace en un plot, la función se importo de la librería matplotlib.pyplot.

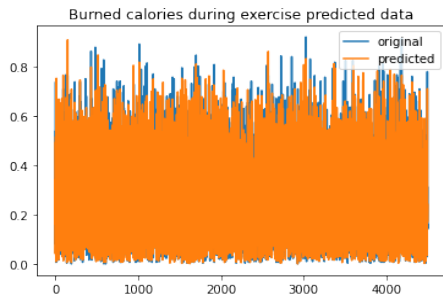


Figura 6: Visualización gráfica del desempeño de la predicción)

4. Conclusiones

El objetivo de este artículo fue explicar la relación entre el ejercicio y el gasto calórico durante mediante la demostración de un estudio de caso. Si el compromiso con la actividad física relacionada con la salud es el objetivo para usted, estas son las recomendaciones:

PENDIENTE

Esperamos haberle proporcionado la información suficiente para que pueda tomar decisiones óptimas en la realización de ejercicio físico. Cabe mencionar que los beneficios fisiológicos a largo plazo del ejercicio regular de la parte superior e inferior del cuerpo no se han dilucidado completamente en los resultados de la investigación.

Agradecimientos

Este trabajo ha sido realizado con el apoyo de los maestros de la Facultad de Ciencias Fisicomatéticas de la UANL.

5. Apéndice

PENDIENTE

6. Referencias

[1] <https://www.kaggle.com/datasets/aadhavvignesh/calories-burned-during-exercise-and-activities>