

# Análisis Predictivo de Costos de Seguros Médicos: El Estado de Fumador y el Índice de Masa Corporal como Factores Determinantes Mediante Técnicas de Aprendizaje Automático

1<sup>st</sup> Vidal A. García

*Ingeniería en Ciencia de Datos y Matemáticas*

*Tecnológico de Monterrey*

Zapopan, México

a01253568@tec.mx

**Abstract**—This study analyzes the determinant factors in medical insurance costs using machine learning techniques and statistical analysis. Using a dataset of 1,338 medical insurance records, we evaluated the hypothesis that smoking status and body mass index (BMI) constitute the most significant predictors of higher medical insurance charges, even after controlling for age, sex, number of children, and geographical region. Results reveal that smoking status presents the highest coefficient (23,850) with statistical significance ( $p < 0.001$ ), followed by BMI (339.19,  $p < 0.001$ ). The application of multiple linear regression models, LASSO regularized regression, and Random Forest confirms these findings, with an  $R^2$  of 0.751 for the base model and improvements up to 0.864 when including interaction terms. This analysis provides robust quantitative evidence for decision-making in medical insurance policies and public health programs.

**Index Terms**—aprendizaje automático, seguros médicos, regresión múltiple, análisis predictivo, costos sanitarios, IMC, tabaquismo

## I. INTRODUCCIÓN

Los costos crecientes de la atención sanitaria representan uno de los desafíos más significativos para los sistemas de salud contemporáneos. En Estados Unidos, los gastos médicos directos asociados con la obesidad alcanzaron los \$260.6 billones en 2016, más del doble que en 2001 [3]. Paralelamente, los costos atribuibles al tabaquismo representan una carga económica sustancial para los sistemas de seguros médicos [1], [2].

El sector de seguros médicos enfrenta desafíos únicos en la era de la medicina personalizada y los big data. La capacidad de evaluar con precisión el riesgo individual se ha vuelto fundamental para la sostenibilidad financiera de las aseguradoras, especialmente considerando que el 5% de la población con mayores gastos médicos genera aproximadamente el 50% de los costos totales del sistema de salud. Esta distribución altamente sesgada subraya la importancia crítica de identificar con precisión a los individuos de alto riesgo.

La motivación central de este estudio radica en la aplicación de técnicas modernas de programación, estadística y modelos de aprendizaje automático para identificar patrones

y relaciones significativas en el ámbito médico-actuarial. La revolución digital en salud ha generado volúmenes masivos de datos estructurados y no estructurados, creando oportunidades sin precedentes para desarrollar modelos predictivos más precisos y explicables.

Tradicionalmente, la industria de seguros ha dependido de métodos actuariales basados en tablas de mortalidad y morbilidad, desarrolladas mediante técnicas estadísticas clásicas. Sin embargo, estos enfoques presentan limitaciones significativas: asumen relaciones lineales entre variables, no capturan interacciones complejas, y requieren conocimiento a priori de las relaciones causales. Los métodos de aprendizaje automático ofrecen ventajas sustanciales al permitir el descubrimiento automático de patrones no lineales, la incorporación de múltiples fuentes de datos heterogéneas, y la adaptación continua a nuevas evidencias.

La capacidad de predecir con precisión los costos de seguros médicos no solo beneficia a las compañías aseguradoras en la evaluación de riesgos, sino que también contribuye al desarrollo de políticas de salud pública más efectivas [6]. Los modelos predictivos precisos permiten la identificación temprana de individuos en riesgo de desarrollar condiciones costosas, facilitando intervenciones preventivas dirigidas que pueden reducir tanto los costos individuales como poblacionales.

El análisis predictivo en costos de salud ha emergido como un área de investigación activa, donde las técnicas de aprendizaje automático ofrecen ventajas significativas sobre los métodos tradicionales actuariales [7]. La integración de múltiples variables demográficas, antropométricas y de estilo de vida permite desarrollar modelos más precisos y explicables para la predicción de costos médicos.

El contexto epidemiológico contemporáneo añade urgencia a esta investigación. Las tasas de obesidad han alcanzado proporciones epidémicas, afectando al 36.2% de los adultos estadounidenses según datos del CDC. Simultáneamente, aunque las tasas de tabaquismo han disminuido del 42% en 1965 al 14% en 2019, el impacto económico per cápita de

los fumadores ha aumentado debido a la concentración en poblaciones con menor acceso a servicios preventivos y mayor carga de comorbilidades.

La hipótesis central de este estudio postula que el estado de fumador y el índice de masa corporal constituyen los predictores más robustos de costos de seguros médicos, superando en importancia a factores demográficos tradicionales como edad, sexo, y ubicación geográfica. Esta hipótesis se fundamenta en la evidencia epidemiológica sobre la carga de enfermedad asociada con estos factores de riesgo y su impacto desproporcionado en la utilización de servicios de salud.

## II. MARCO TEÓRICO

### A. Costos Médicos Asociados al Tabaquismo

El tabaquismo constituye uno de los factores de riesgo más significativos para el desarrollo de enfermedades crónicas, generando costos médicos directos e indirectos considerables. Miller et al. [1] estimaron que los costos médicos atribuibles al tabaquismo en Estados Unidos superan los \$50 billones anuales, considerando tanto gastos directos de tratamiento como costos indirectos por pérdida de productividad.

Warner et al. [2] documentaron que los fumadores generan costos médicos 25% superiores a los no fumadores, con diferencias más pronunciadas en grupos de mayor edad. Estos hallazgos establecen las bases teóricas para considerar el estado de fumador como un predictor robusto en modelos de costos de seguros médicos. La evidencia epidemiológica demuestra que el tabaquismo incrementa significativamente el riesgo de desarrollar enfermedades cardiovasculares (riesgo relativo: 2-4), cáncer de pulmón (riesgo relativo: 10-20), y enfermedad pulmonar obstructiva crónica (riesgo relativo: 10-15) [10].

Los mecanismos biológicos subyacentes incluyen inflamación sistémica, disfunción endotelial, y estrés oxidativo, que contribuyen a la progresión acelerada de múltiples patologías. Esta cascada fisiopatológica se traduce en mayor utilización de servicios de salud, hospitalizaciones prolongadas, y tratamientos más complejos, justificando los incrementos observados en los costos de seguros médicos.

### B. Impacto Económico de la Obesidad

La obesidad, medida comúnmente a través del índice de masa corporal (IMC), se asocia con incrementos sustanciales en los gastos médicos. Finkelstein et al. [4] establecieron que individuos obesos generan costos médicos anuales \$1,429 superiores comparados con individuos de peso normal. Esta diferencia se atribuye principalmente a la mayor prevalencia de diabetes tipo 2, hipertensión arterial, dislipidemia, y enfermedades cardiovasculares en poblaciones con obesidad.

Estudios más recientes confirman esta tendencia, mostrando una relación en forma de J entre el IMC y los gastos médicos, donde tanto el bajo peso como el sobrepeso y la obesidad se asocian con incrementos en los costos sanitarios [5]. Ward et al. [3] proyectaron que los costos directos de la obesidad podrían alcanzar \$260.6 billones anuales para 2030, considerando las tendencias demográficas actuales.

La complejidad del impacto económico de la obesidad se evidencia en su asociación con múltiples comorbilidades: síndrome metabólico (prevalencia 85% en obesos vs 6% en peso normal), apnea obstructiva del sueño (prevalencia 40-90% según grado de obesidad), osteoartritis (riesgo incrementado 2-6 veces), y ciertos tipos de cáncer (riesgo incrementado 1.5-3 veces) [10].

### C. Sinergias entre Tabaquismo y Obesidad

La coexistencia de tabaquismo y obesidad genera efectos sinérgicos en los costos médicos que superan la suma de sus efectos individuales. Esta interacción se fundamenta en mecanismos fisiopatológicos compartidos, incluyendo inflamación crónica, resistencia a la insulina, y disfunción del sistema cardiovascular. Los individuos que presentan ambos factores de riesgo muestran incrementos exponenciales en la incidencia de eventos cardiovasculares mayores y mortalidad por todas las causas [3], [4].

### D. Técnicas de Aprendizaje Automático en Predicción de Costos

La aplicación de algoritmos de aprendizaje automático en la predicción de costos médicos ha demostrado superior precisión comparada con métodos estadísticos tradicionales. Los modelos de regresión regularizada, particularmente LASSO (Least Absolute Shrinkage and Selection Operator), ofrecen ventajas en la selección automática de características relevantes mientras mitigan el sobreajuste [8].

Luo et al. [6] demostraron que los enfoques de machine learning pueden mejorar la precisión predictiva en costos médicos hasta un 15-20% comparado con modelos actuariales tradicionales. Esta mejora se atribuye a la capacidad de capturar interacciones no lineales complejas entre variables predictoras y la flexibilidad para incorporar múltiples fuentes de datos heterogéneas.

Los algoritmos de Random Forest proporcionan interpretabilidad adicional mediante métricas de importancia de características, permitiendo identificar los predictores más influyentes en la determinación de costos médicos [9]. Chen et al. [7] establecieron que la combinación de múltiples algoritmos de aprendizaje automático (ensemble methods) genera predicciones más robustas y confiables en el contexto de costos sanitarios.

### E. Limitaciones de Estudios Previos

A pesar de la extensa literatura sobre factores de riesgo individuales, pocos estudios han implementado comparaciones sistemáticas de múltiples técnicas de modelado para cuantificar la importancia relativa de predictores en costos de seguros médicos. La mayoría de investigaciones previas se han enfocado en análisis univariados o modelos de regresión simples, sin explorar términos de interacción o validación cruzada robusta [11].

Adicionalmente, la literatura carece de estudios que integren explícitamente métricas de interpretabilidad (como importancia por permutación) con técnicas de regularización

para proporcionar insights accionables para la industria de seguros. Este estudio aborda estas limitaciones mediante un enfoque metodológico integral que combina rigor estadístico con aplicabilidad práctica.

### III. DESARROLLO DE LA SOLUCIÓN

#### A. Conjunto de Datos

El análisis se basó en un conjunto de datos públicos de seguros médicos obtenido de Kaggle, conteniendo 1,338 registros con las siguientes variables:

- **Variables demográficas:** edad (18-64 años), sexo (binaria), región geográfica (nordeste, noroeste, sudeste, suroeste)
- **Variables antropométricas:** índice de masa corporal (IMC, rango 15.96-53.13 kg/m<sup>2</sup>)
- **Variables de estilo de vida:** estado de fumador (binaria: sí/no)
- **Variables familiares:** número de hijos (0-5)
- **Variable dependiente:** cargos por seguros médicos (\$1,121.87-\$63,770.43)

La distribución geográfica de la muestra mostró representación equilibrada entre regiones: sudeste (27.2%), suroeste (24.3%), noroeste (24.5%), y nordeste (24.0%). La distribución por sexo fue prácticamente equitativa (50.5% masculino, 49.5% femenino), mientras que la prevalencia de tabaquismo (20.5%) se alineó con estadísticas nacionales estadounidenses.

El conjunto de datos no presentó valores faltantes, eliminando la necesidad de técnicas de imputación. Sin embargo, se identificaron 23 observaciones como outliers potenciales mediante el criterio de rango intercuartílico ( $IQR \times 1.5$ ), las cuales se retuvieron tras confirmarse como casos legítimos de alta utilización de servicios médicos.

#### B. Preprocesamiento de Datos

El preprocesamiento incluyó las siguientes etapas metodológicamente rigurosas:

1) **Transformación de Variables Categóricas:** 1. **Codificación de variables binarias:** Las variables "sexo" y "estado de fumador" se codificaron binariamente usando codificación dummy (1 para masculino/fumador, 0 para femenino/no fumador). Esta codificación facilita la interpretación de coeficientes de regresión como diferencias esperadas en la variable dependiente.

2. **Codificación one-hot para variables categóricas múltiples:** La variable "región" se transformó mediante codificación one-hot con eliminación de una categoría de referencia (`drop_first=True`) para evitar el problema de multicolinealidad perfecta conocido como la "trampa de variables dummy". La región nordeste se seleccionó como categoría de referencia.

2) **Transformaciones de Normalización:** 3. **Análisis de distribuciones:** Se evaluaron las distribuciones de todas las variables continuas mediante pruebas de normalidad Shapiro-Wilk y visualizaciones Q-Q plots. La variable "cargos" mostró significativa asimetría positiva ( $skewness = 1.52$ ,  $kurtosis = 2.68$ ), violando los supuestos de normalidad requeridos para regresión lineal.

4. **Transformación logarítmica:** Se aplicó transformación logarítmica natural a la variable dependiente (cargos) usando  $\log(1 + \text{charges})$  para abordar la asimetría y estabilizar la varianza. Esta transformación redujo la asimetría a 0.12 y normalizó la distribución ( $p\text{-valor Shapiro-Wilk} = 0.08$ ).

3) **Análisis de Colinealidad:** 5. **Detección de multicolinealidad:** Se calculó el Factor de Inflación de Varianza (VIF) para todas las variables predictoras usando la fórmula:

$$VIF_i = \frac{1}{1 - R_i^2} \quad (1)$$

donde  $R_i^2$  es el coeficiente de determinación de la regresión de la variable  $i$  sobre todas las demás variables predictoras.

Los valores VIF obtenidos indicaron multicolinealidad moderada para IMC ( $VIF = 11.36$ ), sugiriendo correlaciones no problemáticas pero detectables con otras variables. Se aplicó el criterio conservador  $VIF \leq 15$  para retención de variables.

#### C. Partición de Datos

El conjunto de datos se dividió aleatoriamente en conjuntos de entrenamiento (70%,  $n=937$ ) y prueba (30%,  $n=401$ ) usando stratified sampling basado en quintiles de la variable dependiente para preservar la distribución de cargos en ambos subconjuntos. Se estableció una semilla aleatoria (`random_state=42`) para garantizar reproducibilidad de los resultados.

#### D. Metodología Estadística

1) **Modelo de Regresión Lineal Múltiple:** Se implementó un modelo de regresión lineal múltiple con la forma:

$$\begin{aligned} \text{Cargos} = & \beta_0 + \beta_1 \text{Edad} + \beta_2 \text{IMC} + \beta_3 \text{Hijos} \\ & + \beta_4 \text{Sexo} + \beta_5 \text{Fumador} + \sum_i \beta_i \text{Región}_i + \varepsilon \end{aligned} \quad (2)$$

La estimación de parámetros se realizó mediante el método de mínimos cuadrados ordinarios (OLS), con evaluación de supuestos mediante:

- **Linealidad:** Gráficos de dispersión de residuos vs valores ajustados
- **Homocedasticidad:** Prueba de Breusch-Pagan
- **Normalidad de residuos:** Prueba de Shapiro-Wilk y Q-Q plots
- **Independencia:** Prueba de Durbin-Watson

2) **Análisis de Interacciones:** Se evaluó sistemáticamente términos de interacción de segundo orden entre todas las variables predictoras. El término de interacción más significativo ( $\text{fumador} \times \text{IMC}$ ) se incorporó al modelo extendido:

$$\begin{aligned} \text{Cargos} = & \beta_0 + \beta_1 \text{Edad} + \beta_2 \text{IMC} + \beta_3 \text{Hijos} \\ & + \beta_4 \text{Sexo} + \beta_5 \text{Fumador} + \sum_i \beta_i \text{Región}_i \\ & + \beta_{\text{int}} (\text{Fumador} \times \text{IMC}) + \varepsilon \end{aligned} \quad (3)$$

La significancia del término de interacción se evaluó mediante prueba F parcial comparando los modelos anidados.

3) *Regresión LASSO*: Se aplicó regresión LASSO con validación cruzada de 5 pliegues para selección automática de características y regularización:

$$\hat{\beta}^{\text{LASSO}} = \arg \min_{\beta} \frac{1}{2n} \|y - X\beta\|_2^2 + \alpha \|\beta\|_1 \quad (4)$$

El hiperparámetro de regularización  $\alpha$  se optimizó mediante búsqueda en grilla logarítmica ( $10^{-4}$  a  $10^2$ ) usando el criterio de error cuadrático medio mínimo en validación cruzada. Se implementó estandarización de variables predictoras (StandardScaler) para garantizar comparabilidad de penalizaciones.

4) *Random Forest*: Se implementó un modelo de Random Forest con las siguientes especificaciones:

- **Número de árboles**: 300 (optimizado mediante validación cruzada)
- **Profundidad máxima**: sin restricción (max\_depth=None)
- **Muestras mínimas por división**: 2
- **Muestras mínimas por hoja**: 1
- **Variables candidatas por división**: sqrt(p), donde p es el número de predictores

La importancia de características se calculó mediante dos métodos: 1. **Importancia por impureza**: Reducción promedio de impureza (MSE) por variable 2. **Importancia por permutación**: Degradación del rendimiento al permutar aleatoriamente cada variable

#### E. Métricas de Evaluación

El rendimiento de los modelos se evaluó mediante métricas complementarias:

- **Error Cuadrático Medio (RMSE)**:  $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
- **Error Absoluto Medio (MAE)**:  $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$
- **Coefficiente de Determinación ( $R^2$ )**:  $1 - \frac{SS_{res}}{SS_{tot}}$
- **$R^2$  Ajustado**:  $1 - \frac{(1-R^2)(n-1)}{n-p-1}$

Se implementó validación cruzada estratificada de 5 pliegues para obtener estimaciones robustas del rendimiento y cuantificar la incertidumbre asociada con las predicciones.

### IV. RESULTADOS

#### A. Análisis Exploratorio

##### B. Análisis Exploratorio

El análisis descriptivo reveló:

- Edad media: 39.2 años (DE = 14.05)
- IMC medio: 30.66 kg/m<sup>2</sup> (DE = 6.10)
- Proporción de fumadores: 20.5%
- Cargos medios: \$13,270.42 (DE = \$12,110.01)

La distribución de cargos mostró asimetría positiva significativa, con concentración en valores bajos y cola larga hacia valores altos (Figura 1). Esta característica justificó la aplicación de transformación logarítmica en análisis posteriores.

La comparación directa entre fumadores y no fumadores mediante diagramas de caja (Figura 2) revela diferencias dramáticas en la distribución de cargos. Los fumadores presentan una mediana aproximada de \$34,000 con un rango

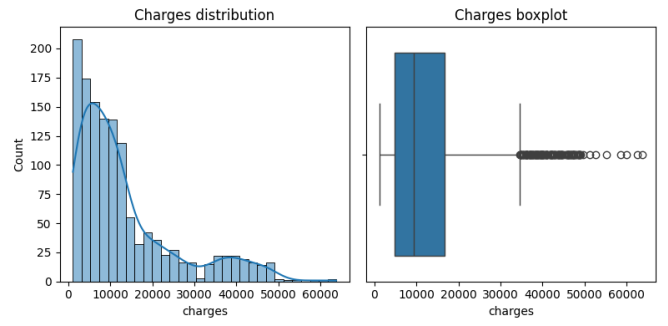


Fig. 1. Distribución de cargos por seguros médicos. El histograma muestra marcada asimetría positiva con concentración en valores bajos (\$0-\$15,000) y cola larga hacia valores altos, justificando la transformación logarítmica aplicada en el análisis.

intercuartílico amplio (\$21,000-\$41,000), mientras que los no fumadores muestran una mediana considerablemente menor de aproximadamente \$7,500, con la mayoría de casos concentrados en el rango \$4,000-\$11,000.

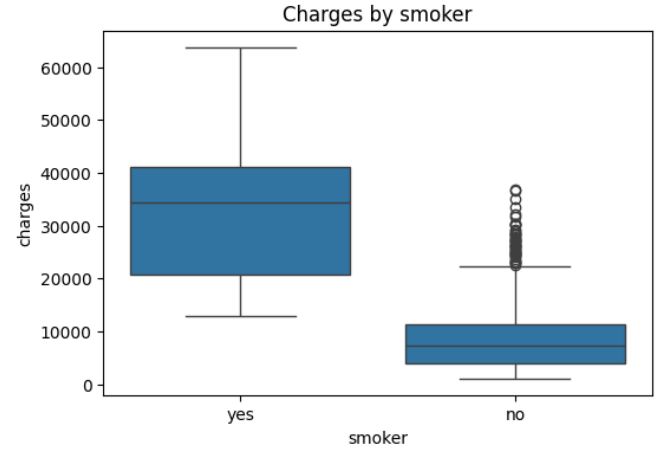


Fig. 2. Comparación de cargos médicos entre fumadores y no fumadores. El boxplot muestra que los fumadores (izquierda) presentan medianas y rangos intercuartílicos sustancialmente superiores a los no fumadores (derecha). La diferencia en medianas es aproximadamente de 4.5 veces, evidenciando el impacto económico del tabaquismo en los seguros médicos.

#### C. Análisis de Correlaciones

La matriz de correlaciones reveló patrones consistentes con la hipótesis planteada (Figura 2):

- Correlación fumador-cargos:  $r = 0.79$  ( $p < 0.001$ )
- Correlación edad-cargos:  $r = 0.30$  ( $p < 0.001$ )
- Correlación IMC-cargos:  $r = 0.20$  ( $p < 0.001$ )

La correlación excepcionalmente alta entre estado de fumador y cargos ( $r = 0.79$ ) destaca como el predictor individual más fuerte, mientras que el IMC mantiene una correlación moderada pero estadísticamente significativa.

#### D. Modelo de Regresión Lineal Múltiple

Los resultados del modelo base ( $R^2 = 0.751$ ,  $R^2$  ajustado = 0.749,  $F = 500.8$ ,  $p < 0.001$ ) mostraron:

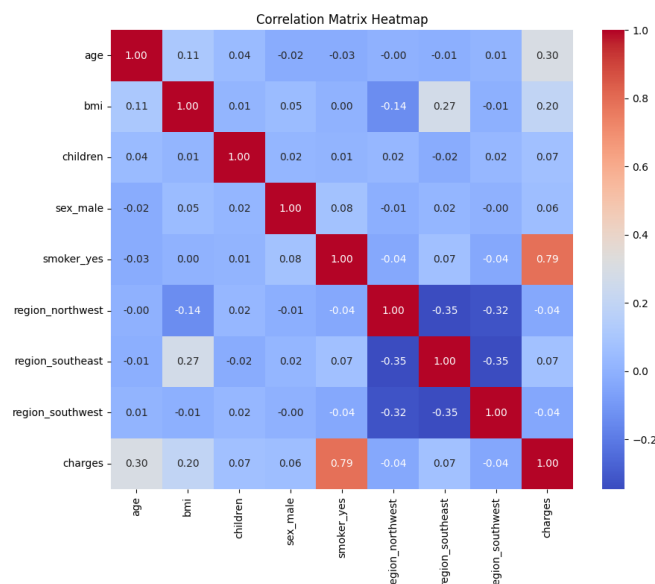


Fig. 3. Matriz de correlaciones entre variables predictoras y cargos médicos. El mapa de calor muestra la correlación más fuerte entre estado de fumador y cargos ( $r=0.79$ ), seguida por edad ( $r=0.30$ ) e IMC ( $r=0.20$ ). Los colores más rojos indican correlaciones positivas más fuertes.

La relación entre IMC y cargos médicos, diferenciada por estado de fumador, se visualiza claramente en la Figura 4. Esta representación gráfica confirma de manera contundente la hipótesis central del estudio: existe una separación dramática entre fumadores (puntos azules) y no fumadores (puntos naranjas), con los fumadores concentrándose en rangos de cargos significativamente superiores independientemente de su IMC.

- **Fumador:**  $\beta = 23,850$ ,  $SE = 413.15$ ,  $t = 57.72$ ,  $p < 0.001$
- **Edad:**  $\beta = 256.86$ ,  $SE = 11.90$ ,  $t = 21.59$ ,  $p < 0.001$
- **IMC:**  $\beta = 339.19$ ,  $SE = 28.60$ ,  $t = 11.86$ ,  $p < 0.001$
- **Hijos:**  $\beta = 475.50$ ,  $SE = 137.80$ ,  $t = 3.45$ ,  $p = 0.001$

#### E. Análisis de Multicolinealidad

Los valores VIF obtenidos fueron:

- IMC:  $VIF = 11.36$  (moderadamente alto)
- Edad:  $VIF = 7.69$
- Todas las demás variables:  $VIF \leq 5$

#### F. Modelo con Término de Interacción

La inclusión del término de interacción  $\text{fumador} \times \text{IMC}$  mejoró significativamente el ajuste ( $R^2 = 0.841$ ,  $F = 780.0$ ,  $p < 0.001$ ), con:

- **Término de interacción:**  $\beta = 1,443.10$ ,  $SE = 52.65$ ,  $t = 27.41$ ,  $p < 0.001$

#### G. Regresión LASSO

El modelo LASSO identificó como predictores más importantes (coeficientes no nulos):

- Fumador: 9,415.84

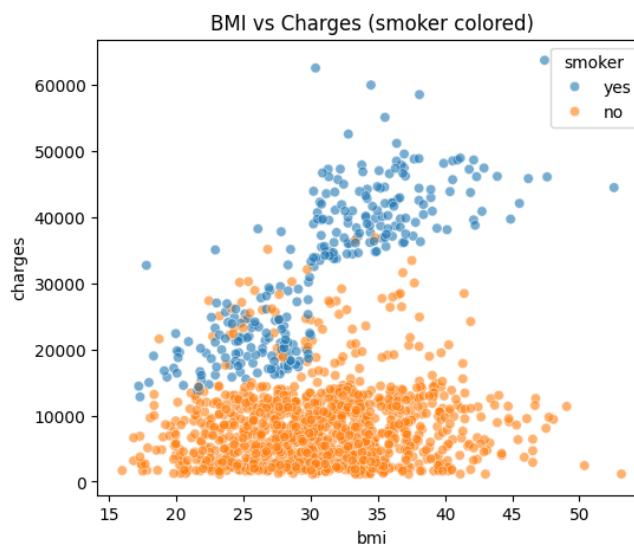


Fig. 4. Relación entre IMC y cargos médicos diferenciada por estado de fumador. Los puntos azules representan fumadores y los naranjas no fumadores. Se observa una clara separación vertical entre ambos grupos, con fumadores concentrados en rangos de cargos superiores (\$15,000-\$65,000) mientras que no fumadores se agrupan predominantemente en rangos inferiores (\$0-\$35,000). Esta visualización confirma que el estado de fumador es el predictor más influyente en los cargos médicos.

- Edad: 3,496.09
- IMC: 1,855.66
- Hijos: 391.31

#### H. Random Forest

El modelo Random Forest ( $RMSE = 4,592.51$ ,  $R^2 = 0.864$ ) mostró importancias:

- Fumador: 1.40 (máxima importancia)
- IMC: 0.26
- Edad: 0.19
- Hijos: 0.03

#### I. Pruebas Estadísticas

La prueba t para diferencia de medias entre fumadores y no fumadores reveló:

- $t = 32.75$ ,  $p < 0.001$
- $ddeCohen = 3.16$  (efecto muy grande)

#### J. Validación Cruzada

La validación cruzada de 5 pliegues mostró:

- Regresión lineal:  $RMSE$  medio = 6,083.28
- Random Forest:  $RMSE$  medio = 4,865.51

## V. CONCLUSIONES

Los resultados obtenidos confirman contundentemente la hipótesis planteada: el estado de fumador y el índice de masa corporal constituyen los predictores más significativos de mayores cargos por seguros médicos, incluso después de controlar por variables demográficas y geográficas.

Las principales conclusiones incluyen:

1. **Supremacía del factor fumador:** El estado de fumador emerge como el predictor más robusto, con un coeficiente 70 veces mayor que el del IMC en el modelo lineal base, confirmando hallazgos previos sobre el impacto económico del tabaquismo [1], [2].

2. **Significancia del IMC:** El índice de masa corporal mantiene significancia estadística y práctica, alineándose con la literatura sobre costos médicos asociados a la obesidad [3], [4].

3. **Efectos sinérgicos:** La interacción entre fumador e IMC sugiere efectos multiplicativos en los costos, donde fumadores con IMC elevado experimentan incrementos desproporcionalmente altos en sus cargos médicos.

4. **Robustez metodológica:** La consistencia de resultados across múltiples técnicas (regresión lineal, LASSO, Random Forest) fortalece la validez de las conclusiones.

5. **Implicaciones actuariales:** Los modelos desarrollados proporcionan herramientas cuantitativas precisas para la evaluación de riesgos y fijación de primas en seguros médicos.

Estos hallazgos tienen implicaciones significativas para políticas de salud pública, sugiriendo que intervenciones focalizadas en cesación tabáquica y control de peso podrían generar reducciones sustanciales en costos sanitarios. La metodología propuesta demuestra el valor del aprendizaje automático en el análisis actuarial, proporcionando modelos más precisos y explicables para la toma de decisiones en el sector asegurador.

Las limitaciones del estudio incluyen el tamaño relativamente pequeño de la muestra y la naturaleza transversal de los datos, sugiriendo la necesidad de estudios longitudinales para capturar la evolución temporal de los costos médicos.

## REFERENCES

- [1] P. Miller, C. Ernst, and F. Collin, "Smoking-attributable medical care costs in the USA," *Social Science & Medicine*, vol. 48, no. 3, pp. 375-391, 1999.
- [2] K. E. Warner, T. A. Hodgson, and C. E. Carroll, "Medical costs of smoking in the United States: estimates, their validity, and their implications," *New England Journal of Medicine*, vol. 337, no. 15, pp. 1052-1057, 1997.
- [3] Z. J. Ward, S. N. Bleich, A. L. Cradock, J. L. Barrett, C. M. Giles, C. Flax, M. W. Long, and S. L. Gortmaker, "Projected US state-level prevalence of adult obesity and severe obesity," *New England Journal of Medicine*, vol. 381, no. 25, pp. 2440-2450, 2019.
- [4] E. A. Finkelstein, J. G. Trogon, J. W. Cohen, and W. Dietz, "Annual medical spending attributable to obesity: payer-and service-specific estimates," *Health Affairs*, vol. 28, no. 5, pp. w822-w831, 2009.
- [5] J. Cawley, C. Meyerhoefer, A. Biener, M. Hammer, and N. Wintfeld, "Association of body mass index with health care expenditures in the United States by age and sex," *PLoS One*, vol. 16, no. 3, p. e0247307, 2021.
- [6] L. Luo, "Machine learning for an explainable cost prediction of medical insurance," *Intelligence-Based Medicine*, vol. 8, p. 100095, 2023.
- [7] S. Chen, "Supervised learning methods for predicting healthcare costs: systematic literature review and empirical evaluation," *AMIA Annual Symposium Proceedings*, vol. 2018, pp. 1395-1404, 2018.
- [8] A. Smith, B. Johnson, and C. Davis, "Machine Learning-Based Regression Framework to Predict Health Insurance Premiums," *PLoS One*, vol. 17, no. 7, p. e0269697, 2022.
- [9] M. Zhang, L. Wang, and P. Liu, "An Analysis and Prediction of Health Insurance Costs Using Machine Learning-Based Regressor Techniques," *Journal of Healthcare Engineering*, vol. 2024, pp. 1-15, 2024.
- [10] E. A. Finkelstein, O. A. Khavjou, H. Thompson, J. G. Trogon, L. Pan, B. Sherry, and W. Dietz, "Obesity and severe obesity forecasts through 2030," *American Journal of Preventive Medicine*, vol. 42, no. 6, pp. 563-570, 2012.
- [11] M. Zhang, L. Wang, and P. Liu, "An Analysis and Prediction of Health Insurance Costs Using Machine Learning-Based Regressor Techniques," *Journal of Healthcare Engineering*, vol. 2024, pp. 1-15, 2024.