

Fundamentos Teóricos del Aprendizaje Automático: Scikit-learn y TensorFlow

Dr. Darío Ezequiel Díaz

24 de abril de 2025

Índice

1. Introducción al Aprendizaje Automático	3
2. Preprocesamiento de Datos	3
2.1. Normalización y Estandarización	4
2.1.1. StandardScaler	4
2.1.2. MinMaxScaler	4
2.2. Codificación de Variables Categóricas	4
2.2.1. OneHotEncoder	4
3. Modelos Fundamentales de Scikit-learn	5
3.1. Regresión Lineal	5
3.2. Regresión Logística	5
3.3. Árboles de Decisión	6
4. Validación y Evaluación de Modelos	6
4.1. División de Conjuntos de Datos	6
4.2. Validación Cruzada	7
4.3. Matrices de Confusión y Métricas de Evaluación	7
5. Introducción a TensorFlow y Keras	7
5.1. Fundamentos de Redes Neuronales	8
5.2. Arquitecturas de Redes Neuronales en Keras	8
5.2.1. Modelos Secuenciales	8
5.2.2. API Funcional	9
5.3. Entrenamiento de Redes Neuronales	9
6. Redes Neuronales Convolucionales	9
6.1. Operación de Convolución	9
6.2. Arquitectura CNN Básica	10
6.3. Transfer Learning	10

7. Aprendizaje por Refuerzo	10
7.1. Fundamentos Teóricos	10
7.2. Q-Learning	11
8. Conclusiones y Perspectivas	11

1. Introducción al Aprendizaje Automático

El aprendizaje automático constituye actualmente uno de los pilares fundamentales para la transformación digital de nuestra sociedad, abarcando un amplio espectro de aplicaciones que abarcan desde sistemas de recomendación hasta algoritmos de conducción autónoma. Esta disciplina, fundamentada en sólidos principios matemáticos y estadísticos, ha experimentado una notable evolución conceptual desde sus orígenes, cuando Arthur Samuel la definió en 1959 como “el campo de estudio que proporciona a las computadoras la capacidad de aprender sin ser explícitamente programadas” [15].

La esencia del aprendizaje automático radica en su capacidad para identificar patrones en conjuntos de datos voluminosos y heterogéneos, permitiendo así construir modelos matemáticos cuya complejidad trasciende la programación explícita. Estos modelos adquieren la facultad de generalizar conocimiento a partir de experiencias previas, realizando predicciones o tomando decisiones sobre nuevos datos no contemplados durante su fase de entrenamiento [6].

Desde una perspectiva formal, podemos clasificar los enfoques principales del aprendizaje automático en tres paradigmas esenciales:

- **Aprendizaje Supervisado:** Orientado hacia la predicción, utiliza datos etiquetados para establecer una correspondencia entre entradas y salidas deseadas. Las implementaciones algorítmicas más representativas incluyen regresiones, árboles de decisión y redes neuronales.
- **Aprendizaje No Supervisado:** Centrado en el descubrimiento de estructuras latentes en datos no etiquetados, facilitando la identificación de patrones subyacentes. Técnicas como el agrupamiento (clustering), reducción de dimensionalidad y detección de anomalías se inscriben dentro de este paradigma.
- **Aprendizaje por Refuerzo:** Basado en el principio de recompensa-penalización, donde un agente aprende a tomar decisiones secuenciales óptimas maximizando una función de recompensa acumulativa. Este enfoque resulta particularmente valioso para problemas de toma de decisiones secuenciales en entornos dinámicos.

El desarrollo contemporáneo del aprendizaje automático ha sido posible gracias a la confluencia de tres factores determinantes: la disponibilidad de datos masivos, el avance en capacidad computacional y la sofisticación de algoritmos matemáticos. Este progreso ha propiciado la consolidación de herramientas técnicas como Scikit-learn y TensorFlow, que constituyen el eje central de este módulo académico.

2. Preprocesamiento de Datos

El preprocesamiento de datos representa una fase crítica en el desarrollo de modelos de aprendizaje automático eficientes, incidiendo significativamente en su capacidad predictiva y rendimiento general. Esta etapa preliminar implica la transformación sistemática de datos crudos en formatos optimizados que maximizan la efectividad de los algoritmos subsecuentes [5].

2.1. Normalización y Estandarización

La normalización y estandarización constituyen técnicas fundamentales que abordan la heterogeneidad en las escalas de medición entre variables, un requisito imperativo para múltiples algoritmos de aprendizaje automático, especialmente aquellos basados en gradiente [8].

2.1.1. StandardScaler

La estandarización mediante **StandardScaler** transforma cada variable para que presente una media aritmética $\mu = 0$ y desviación estándar $\sigma = 1$. Matemáticamente, para cada valor x_i en un conjunto de datos, la transformación se define como:

$$z_i = \frac{x_i - \mu}{\sigma} \quad (1)$$

Esta transformación resulta particularmente valiosa para algoritmos sensibles a la escala de los datos, como máquinas de vectores de soporte (SVM), regresión logística y redes neuronales, contribuyendo a la estabilidad y eficiencia en la convergencia del gradiente durante la optimización [10].

2.1.2. MinMaxScaler

Por su parte, **MinMaxScaler** implementa una normalización que reescala linealmente cada característica al intervalo $[0, 1]$ o cualquier otro rango especificado. La transformación se expresa como:

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (2)$$

Esta técnica preserva las relaciones entre los valores originales mientras contiene todos los datos dentro de límites predefinidos, resultando especialmente útil para algoritmos que requieren entradas no negativas o acotadas, como las redes neuronales con funciones de activación sigmoideas [14].

2.2. Codificación de Variables Categóricas

Las variables categóricas, fundamentales en numerosos contextos analíticos, requieren transformaciones específicas para su integración en modelos matemáticos de aprendizaje automático [7].

2.2.1. OneHotEncoder

La codificación one-hot representa cada categoría mediante un vector binario donde únicamente la posición correspondiente a la categoría presente contiene el valor 1, mientras las restantes posiciones se establecen en 0. Formalmente, para una variable categórica con k categorías posibles, la transformación genera k nuevas variables binarias, evitando así imponer una ordinalidad artificial entre categorías nominales.

Esta técnica resulta crucial para algoritmos que presuponen relaciones numéricas entre valores, como regresiones y redes neuronales, aunque puede incrementar significativamente la dimensionalidad del espacio de características, especialmente con variables de alta cardinalidad [18].

3. Modelos Fundamentales de Scikit-learn

Scikit-learn constituye una biblioteca de código abierto para Python que implementa un conjunto exhaustivo de algoritmos de aprendizaje automático, caracterizándose por su interfaz consistente, documentación rigurosa y eficiencia computacional. Esta herramienta facilita tanto la implementación como la evaluación de modelos predictivos en entornos de investigación y producción [13].

3.1. Regresión Lineal

La regresión lineal representa uno de los modelos predictivos más fundamentales, orientado a establecer relaciones lineales entre variables independientes (predictores) y una variable dependiente (objetivo) continua [4].

Matemáticamente, para un conjunto de m observaciones con n características, el modelo se expresa como:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_n x_{in} + \varepsilon_i \quad (3)$$

donde:

- y_i representa el valor objetivo para la observación i .
- x_{ij} representa el valor de la característica j para la observación i .
- $\beta_0, \beta_1, \dots, \beta_n$ constituyen los parámetros del modelo.
- ε_i representa el término de error aleatorio.

La estimación de los parámetros β se realiza comúnmente mediante el método de mínimos cuadrados ordinarios (OLS), que minimiza la suma de los cuadrados de los residuos:

$$\min_{\beta} \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \min_{\beta} \sum_{i=1}^m (y_i - \beta_0 - \sum_{j=1}^n \beta_j x_{ij})^2 \quad (4)$$

3.2. Regresión Logística

A diferencia de la regresión lineal, la regresión logística aborda problemas de clasificación binaria modelando la probabilidad de pertenencia a una clase específica [9]. El modelo transforma la combinación lineal de predictores mediante la función logística (sigmoide):

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}} = \frac{1}{1 + e^{-\beta^T x}} \quad (5)$$

La estimación de parámetros en la regresión logística se realiza mediante máxima verosimilitud, maximizando:

$$L(\beta) = \prod_{i=1}^m P(y_i|x_i, \beta) = \prod_{i=1}^m P(y_i = 1|x_i, \beta)^{y_i} \cdot (1 - P(y_i = 1|x_i, \beta))^{1-y_i} \quad (6)$$

O equivalentemente, minimizando la log-verosimilitud negativa:

$$-\log L(\beta) = -\sum_{i=1}^m y_i \log(p_i) + (1 - y_i) \log(1 - p_i) \quad (7)$$

donde $p_i = P(y_i = 1|x_i, \beta)$.

3.3. Árboles de Decisión

Los árboles de decisión implementan un enfoque no paramétrico que particiona recursivamente el espacio de características, creando una estructura jerárquica de reglas de decisión [2]. Cada nodo interno representa una condición sobre una característica específica, cada rama corresponde a un resultado de dicha condición, y cada hoja proporciona una predicción para la variable objetivo.

La construcción del árbol se basa en la optimización de métricas como:

- **Impureza de Gini** para problemas de clasificación:

$$Gini(t) = 1 - \sum_{i=1}^k p(i|t)^2 \quad (8)$$

- **Entropía** como medida alternativa de impureza:

$$Entropía(t) = -\sum_{i=1}^k p(i|t) \log_2 p(i|t) \quad (9)$$

- **Error Cuadrático Medio** para problemas de regresión:

$$MSE(t) = \frac{1}{N_t} \sum_{i \in t} (y_i - \bar{y}_t)^2 \quad (10)$$

donde $p(i|t)$ representa la proporción de instancias pertenecientes a la clase i en el nodo t , y \bar{y}_t es el valor medio de las instancias en dicho nodo.

4. Validación y Evaluación de Modelos

La evaluación rigurosa constituye un aspecto crucial del desarrollo de modelos de aprendizaje automático, proporcionando estimaciones realistas sobre su capacidad de generalización a datos no observados durante el entrenamiento [8].

4.1. División de Conjuntos de Datos

La función `train_test_split` implementa una estrategia fundamental que segmenta los datos disponibles en conjuntos de entrenamiento y prueba, permitiendo simular el rendimiento del modelo en un contexto de producción. Esta metodología facilita la detección de problemas como el sobreajuste (overfitting), donde el modelo memoriza patrones específicos del conjunto de entrenamiento en detrimento de su capacidad generalizadora.

4.2. Validación Cruzada

La validación cruzada (`cross_val_score`) extiende el concepto de división de datos mediante particiones múltiples, proporcionando una evaluación más robusta del rendimiento esperado. La técnica de k-fold divide los datos en k subconjuntos (folds), utilizando iterativamente $k - 1$ subconjuntos para entrenamiento y el restante para validación, promediando finalmente las métricas obtenidas.

Matemáticamente, para un modelo M y una función de pérdida L , la puntuación de validación cruzada se calcula como:

$$CV_{score} = \frac{1}{k} \sum_{i=1}^k L(M_{-i}, D_i) \quad (11)$$

donde M_{-i} representa el modelo entrenado en todos los subconjuntos excepto D_i , y $L(M_{-i}, D_i)$ es la pérdida evaluada en D_i .

4.3. Matrices de Confusión y Métricas de Evaluación

Las matrices de confusión proporcionan una visualización tabulada del rendimiento de un clasificador, contraponiendo predicciones y valores reales. Para un problema binario, la matriz define cuatro categorías fundamentales:

	Predicción Positiva	Predicción Negativa
Valor Real Positivo	Verdadero Positivo (VP)	Falso Negativo (FN)
Valor Real Negativo	Falso Positivo (FP)	Verdadero Negativo (VN)

Cuadro 1: Estructura general de una matriz de confusión para clasificación binaria

A partir de estos componentes, se derivan métricas evaluativas como:

- **Exactitud (Accuracy):** $\frac{VP+VN}{VP+FP+VN+FN}$
- **Precisión:** $\frac{VP}{VP+FP}$
- **Sensibilidad (Recall):** $\frac{VP}{VP+FN}$
- **F1-Score:** $2 \cdot \frac{Precisión \cdot Recall}{Precisión + Recall}$

5. Introducción a TensorFlow y Keras

TensorFlow representa actualmente uno de los marcos de referencia más sofisticados para el desarrollo de modelos de aprendizaje automático, especialmente en el ámbito del aprendizaje profundo [1]. Este ecosistema, inicialmente desarrollado por Google, facilita tanto la implementación como el despliegue de arquitecturas complejas mediante operaciones tensoriales en estructuras de grafo computacional.

5.1. Fundamentos de Redes Neuronales

Las redes neuronales artificiales constituyen modelos computacionales inspirados en la estructura biológica del sistema nervioso, compuestos por unidades interconectadas (neuronas) organizadas en capas secuenciales [6]. Cada neurona artificial implementa una transformación matemática que combina linealmente sus entradas, aplicando posteriormente una función de activación no lineal:

$$y = f\left(\sum_{i=1}^n w_i x_i + b\right) = f(W^T X + b) \quad (12)$$

donde:

- $X = (x_1, x_2, \dots, x_n)$ representa el vector de entrada
- $W = (w_1, w_2, \dots, w_n)$ contiene los pesos sinápticos
- b es el término de sesgo (bias)
- f corresponde a la función de activación no lineal

Las funciones de activación más comunes incluyen:

- **ReLU (Rectified Linear Unit):** $f(x) = \max(0, x)$
- **Sigmoide:** $f(x) = \frac{1}{1+e^{-x}}$
- **Tangente hiperbólica:** $f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- **Softmax:** $f(x_i) = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}}$ (para salidas multiclase)

5.2. Arquitecturas de Redes Neuronales en Keras

Keras, integrada actualmente en el ecosistema TensorFlow, proporciona una API de alto nivel para la especificación y entrenamiento de redes neuronales, simplificando significativamente su implementación [3]. La biblioteca soporta dos paradigmas arquitectónicos principales:

5.2.1. Modelos Secuenciales

El modelo secuencial (**Sequential**) facilita la definición de redes donde cada capa recibe como entrada la salida de la capa precedente, resultando particularmente adecuado para arquitecturas lineales como perceptrones multicapa (MLP):

```
1 model = Sequential([
2     Dense(64, activation='relu', input_shape=(10,)),
3     Dense(32, activation='relu'),
4     Dense(16, activation='relu'),
5     Dense(1, activation='sigmoid')
6 ])
```


5.2.2. API Funcional

Para arquitecturas más complejas con conexiones no lineales, bifurcaciones o múltiples entradas/salidas, Keras proporciona la API Funcional:

```
1 input_layer = Input(shape=(10,))
2 x = Dense(64, activation='relu')(input_layer)
3 x = Dense(32, activation='relu')(x)
4 output = Dense(1, activation='sigmoid')(x)
5 model = Model(inputs=input_layer, outputs=output)
```

5.3. Entrenamiento de Redes Neuronales

El entrenamiento de redes neuronales requiere la especificación de tres componentes esenciales:

- **Función de pérdida** que cuantifica la discrepancia entre predicciones y valores reales
- **Optimizador** que implementa un algoritmo para minimizar la función de pérdida
- **Métricas** que evalúan el rendimiento durante el entrenamiento

```
1 model.compile(
2     optimizer='adam',           # Algoritmo de optimización
3     loss='binary_crossentropy', # Función de pérdida
4     metrics=['accuracy']        # Métrica a monitorear
5 )
```

El proceso de optimización generalmente implementa variantes del descenso de gradiente estocástico (SGD), donde los parámetros θ se actualizan iterativamente mediante:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} J(\theta_t) \quad (13)$$

donde η representa la tasa de aprendizaje y $\nabla_{\theta} J(\theta_t)$ es el gradiente de la función de pérdida respecto a los parámetros.

6. Redes Neuronales Convolucionales

Las Redes Neuronales Convolucionales (CNN) constituyen una clase especializada de arquitecturas para el procesamiento eficiente de datos con estructura reticular, particularmente imágenes [11]. Su diseño incorpora principios de conectividad local y compartición de parámetros, reduciendo significativamente el número de parámetros respecto a redes totalmente conectadas equivalentes.

6.1. Operación de Convolución

La operación fundamental en estas redes implica la convolución discreta entre un tensor de entrada (por ejemplo, una imagen) y un conjunto de filtros o kernels aprendibles. Para una imagen bidimensional y un filtro 2D, la operación se define como:

$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) \cdot K(m, n) \quad (14)$$

donde I representa la imagen, K el kernel, y (i, j) las coordenadas de salida.

6.2. Arquitectura CNN Básica

Una arquitectura CNN típica incluye los siguientes componentes secuenciales:

- **Capas Convolucionales** que aplican filtros para extraer características locales
- **Funciones de Activación** no lineales, predominantemente ReLU
- **Capas de Pooling** (Max o Average) que reducen la dimensionalidad espacial
- **Capas Totalmente Conectadas** para la clasificación final

6.3. Transfer Learning

El Transfer Learning representa un paradigma metodológico donde modelos preentrenados en grandes conjuntos de datos (como ImageNet) se adaptan a tareas específicas con conjuntos de datos más limitados [12]. Esta aproximación aprovecha representaciones genéricas aprendidas previamente, requiriendo significativamente menos datos y recursos computacionales para alcanzar altos niveles de rendimiento.

Comúnmente, el proceso implica:

- Cargar una arquitectura preentrenada (como MobileNetV2)
- Congelar las primeras capas que contienen características genéricas
- Reemplazar y reentrenar las capas finales para la tarea específica
- Opcionalmente, realizar un "fine-tuning" de capas intermedias

7. Aprendizaje por Refuerzo

El aprendizaje por refuerzo (RL) establece un paradigma fundamentalmente distinto dentro del aprendizaje automático, orientado hacia la optimización de comportamientos secuenciales mediante interacciones con un entorno dinámico [16]. A diferencia de enfoques supervisados que requieren ejemplos etiquetados, los agentes de RL aprenden mediante experimentación directa, recibiendo señales de recompensa que guían su comportamiento hacia objetivos específicos.

7.1. Fundamentos Teóricos

Formalmente, el aprendizaje por refuerzo se modela como un Proceso de Decisión de Markov (MDP) definido por la tupla $\langle S, A, P, R, \gamma \rangle$, donde:

- S representa el espacio de estados posibles
- A corresponde al conjunto de acciones disponibles
- $P : S \times A \times S \rightarrow [0, 1]$ define la función de transición de probabilidad
- $R : S \times A \times S \rightarrow \mathbb{R}$ especifica la función de recompensa
- $\gamma \in [0, 1]$ constituye el factor de descuento temporal

El objetivo central consiste en aprender una política $\pi : S \rightarrow A$ que maximice la recompensa acumulada descontada esperada:

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1}) \right] \quad (15)$$

7.2. Q-Learning

Q-Learning representa uno de los algoritmos fundamentales en el aprendizaje por refuerzo, implementando una estrategia de aprendizaje libre de modelo (model-free) para estimar la función de valor acción-estado $Q(s, a)$, que representa la recompensa esperada al tomar la acción a en el estado s y seguir posteriormente la política óptima [17].

El algoritmo actualiza iterativamente sus estimaciones mediante la ecuación:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (16)$$

donde:

- $\alpha \in [0, 1]$ representa la tasa de aprendizaje
- r_t corresponde a la recompensa inmediata
- $\gamma \max_a Q(s_{t+1}, a)$ constituye la estimación de recompensa futura óptima

Esta regla de actualización, derivada de la ecuación de Bellman para optimalidad, garantiza la convergencia asintótica hacia la función Q óptima bajo ciertas condiciones técnicas, permitiendo extraer directamente la política óptima como:

$$\pi^*(s) = \arg \max_a Q(s, a) \quad (17)$$

8. Conclusiones y Perspectivas

Este módulo ha proporcionado una fundamentación teórica completa de los principios subyacentes al aprendizaje automático, con especial énfasis en las herramientas Scikit-learn y TensorFlow. Las técnicas y metodologías exploradas constituyen pilares fundamentales para abordar problemas complejos de análisis de datos y predicción en diversos dominios científicos y aplicados.

La comprensión profunda de estos conceptos resulta esencial para desarrollar modelos robustos y eficientes, permitiendo no solo su implementación técnica sino también su interpretación crítica y optimización metodológica. Los conocimientos adquiridos establecen una base sólida para explorar técnicas avanzadas como el aprendizaje por refuerzo, tema que se desarrollará extensivamente en módulos subsecuentes.

La integración de perspectivas matemáticas, estadísticas y computacionales proporciona un enfoque transdisciplinario que caracteriza el paradigma contemporáneo de la ciencia de datos, potenciando el desarrollo de soluciones innovadoras a problemáticas complejas de relevancia socioeconómica.

Referencias

- [1] Abadi, M., et al. (2016). TensorFlow: A system for large-scale machine learning. 12th USENIX Symposium on Operating Systems Design and Implementation, 265-283.
- [2] Breiman, L. (2017). Classification and regression trees. Routledge.
- [3] Chollet, F. (2018). Deep learning with Python. Manning Publications.
- [4] Fox, J. (2015). Applied regression analysis and generalized linear models. Sage Publications.
- [5] García, S., Luengo, J., & Herrera, F. (2015). Data preprocessing in data mining. Springer.
- [6] Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep learning. MIT Press.
- [7] Hancock, J. T., & Khoshgoftaar, T. M. (2020). Survey on categorical data for neural networks. Journal of Big Data, 7(1), 1-41.
- [8] Hastie, T., Tibshirani, R., & Friedman, J. (2009). The elements of statistical learning: data mining, inference, and prediction. Springer.
- [9] Hosmer Jr, D. W., Lemeshow, S., & Sturdivant, R. X. (2013). Applied logistic regression. John Wiley & Sons.
- [10] Kuhn, M., & Johnson, K. (2019). Feature engineering and selection: A practical approach for predictive models. CRC Press.
- [11] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.
- [12] Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. IEEE Transactions on Knowledge and Data Engineering, 22(10), 1345-1359.
- [13] Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825-2830.
- [14] Raschka, S., & Mirjalili, V. (2019). Python machine learning. Packt Publishing.
- [15] Samuel, A. L. (1959). Some studies in machine learning using the game of checkers. IBM Journal of Research and Development, 3(3), 210-229.
- [16] Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT Press.
- [17] Watkins, C. J., & Dayan, P. (1992). Q-learning. Machine Learning, 8(3-4), 279-292.
- [18] Zheng, A., & Casari, A. (2018). Feature engineering for machine learning: principles and techniques for data scientists. O'Reilly Media.