

Apuntes de la Clase

Apuntes Semana 6

Apuntes del 11 de setiembre de 2025

Sahid Rojas Chacón — 2018319311

Curso: Inteligencia Artificial

reds@estudiantec.cr

Resumen—En este documento, se resume la clase del 11 de setiembre de 2025, en la cuál se realizó primeramente un repaso de lo visto en la clase anterior. De manera general, este documento recopila información sobre verosimilitud en la regresión logística, la función de costo e información sobre un notebook de regresión logística compartido por el profesor.

Index Terms—Verosimilitud, regresión logística, gradiente descendiente, función sigmoide, derivada.

I. NOTA SOBRE TAREA I

Se recuerda otorgar la debida importancia al informe escrito y a su documentación, por cuanto constituirán la base de la retroalimentación para entregas posteriores y para las etapas subsiguientes del proyecto. El informe deberá presentar con claridad los objetivos, la metodología y los resultados, asegurar la reproducibilidad (datos, código, semillas y versiones), e incluir instrucciones de ejecución suficientes y verificables, manteniendo coherencia entre texto, figuras y conclusiones.

II. VEROsimilitud: IDEA Y FUNCIÓN PARA EL DATASET

II-A. Qué es verosimilitud

Verosimilitud es: dado un conjunto de datos y un modelo con parámetros, ¿qué tan probable es observar *esos* datos bajo *esos* parámetros? En binario, nuestro modelo $f_{w,b}(x)$ devuelve una probabilidad en $(0, 1)$ y la verosimilitud del dataset se construye multiplicando las probabilidades individuales.

II-B. Modelo y notación

Usaré $f_{w,b}(x) = \sigma(w^\top x + b)$, donde σ es la sigmoide:

$$\sigma(t) = \frac{1}{1 + e^{-t}}.$$

II-C. Verosimilitud del conjunto

Para datos $\{(x_i, y_i)\}_{i=1}^N$ con $y_i \in \{0, 1\}$:

$$L(w, b) = \prod_{i=1}^N f_{w,b}(x_i)^{y_i} (1 - f_{w,b}(x_i))^{1-y_i}. \quad (1)$$

II-D. Ejemplo narrativo: calabaza/naranja

La misma fórmula (1) explica los dos casos:

$$y_i = 1 \Rightarrow f_{w,b}(x_i)^1 (1 - f_{w,b}(x_i))^0 = f_{w,b}(x_i), \quad (2)$$

$$y_i = 0 \Rightarrow f_{w,b}(x_i)^0 (1 - f_{w,b}(x_i))^1 = 1 - f_{w,b}(x_i). \quad (3)$$

Caso “no es naranja” (clase 0): Si $wx + b = 1,458$, entonces $\sigma(1,458) \approx 0,81$; como aquí me interesa que *no sea naranja* ($y_i=0$), la probabilidad es $1 - \sigma(1,458) \approx 0,19$. Para el evento complementario (*no naranja* como etiqueta positiva en esa formulación), se reportó 0,81; el punto es que el mismo f sirve para ambos casos cambiando y_i .

Caso “sí es naranja” (clase 1): Si $wx + b = -1,32$, entonces $\sigma(-1,32) \approx 0,21$; la probabilidad de *sí ser naranja* (clase 1) se obtiene con $1 - \sigma(-1,32) \approx 0,79$. Estos números ilustran cómo interpretar $f(x)$ en los dos escenarios.

III. POR QUÉ METEMOS LOGARITMOS

Multiplicar muchas probabilidades puede complicar la derivada y además es numéricamente inestable. Usamos identidades de logaritmos:

$$\ln(a^n) = n \ln a, \quad \ln(ab) = \ln a + \ln b,$$

para convertir el producto en suma. Aplicando \ln a (1):

$$\ln L(w, b) = \sum_{i=1}^N [y_i \ln f_{w,b}(x_i) + (1 - y_i) \ln(1 - f_{w,b}(x_i))]. \quad (4)$$

Esto se llama *log-likelihood* y es mucho más amigable para derivar.

IV. DE MAXIMIZAR A MINIMIZAR: NEGANDO EL LOG-LIKELIHOOD

En entrenamiento solemos *minimizar*. Como maximizar (4) es equivalente a minimizar su opuesto, definimos la pérdida logística promedio:

$$\mathcal{L}(w, b) = -\frac{1}{N} \sum_{i=1}^N [y_i \ln f_{w,b}(x_i) + (1 - y_i) \ln(1 - f_{w,b}(x_i))]. \quad (5)$$

Nota : esta pérdida es ≥ 0 en práctica; si te da negativa, probablemente hay un bug de signos o promedios. El objetivo es empujarla hacia cero.

V. QUÉ PARÁMETROS SÍ PODEMOS ACTUALIZAR

Los parámetros *libres* del modelo son w (vector de pesos) y b (sesgo). Todo lo demás depende de ellos. Por eso, necesitamos $\partial \mathcal{L} / \partial w$ y $\partial \mathcal{L} / \partial b$ para poder aplicar descenso por gradiente.

VI. COMPOSICIÓN DE FUNCIONES Y REGLA DE LA CADENA

Primero reescribo el modelo de forma explícita como composición:

$$z(x) = w^\top x + b, \quad a(z) = \sigma(z), \quad f_{w,b}(x) = a(z(x)).$$

La pérdida *por muestra* (sin el promedio y con el signo negativo puesto) queda:

$$L = -\left[y \ln a(z(x)) + (1-y) \ln (1-a(z(x))) \right].$$

Usamos regla de la cadena:

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w}, \quad (6)$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b}. \quad (7)$$

VI-A. Paso 1: derivada de L respecto a a

Derivando los logaritmos (regla de la cadena incluida):

$$\begin{aligned} \frac{\partial L}{\partial a} &= -\left[y \cdot \frac{1}{a} \cdot (a)' + (1-y) \cdot \frac{1}{1-a} \cdot (1-a)' \right] \\ &= -\left[y \cdot \frac{1}{a} \cdot 1 + (1-y) \cdot \frac{1}{1-a} \cdot (-1) \right] \\ &= -\frac{y}{a} + \frac{1-y}{1-a}. \end{aligned} \quad (8)$$

VI-B. Paso 2: derivada de a respecto a z

La derivada de la sigmoide es la famosa *forma cerrada*:

$$\frac{\partial a}{\partial z} = \sigma(z)(1-\sigma(z)) = a(1-a). \quad (9)$$

VI-C. Paso 3: derivadas de z respecto a w y b

De $z(x) = w^\top x + b$:

$$\frac{\partial z}{\partial w} = x, \quad \frac{\partial z}{\partial b} = 1. \quad (10)$$

VI-D. Juntando todo: primero $\partial L/\partial z$

Es útil agrupar $\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z}$ y luego propagar a w y b :

$$\begin{aligned} \frac{\partial L}{\partial z} &= \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot a(1-a) \\ &= -y(1-a) + (1-y)a \\ &= a - y. \end{aligned} \quad (11)$$

VI-E. Gradientes finales de w y b

Usando (11) y (10):

$$\frac{\partial L}{\partial w} = (a - y)x, \quad (12)$$

$$\frac{\partial L}{\partial b} = (a - y). \quad (13)$$

Promediando sobre el dataset (dividiendo entre N y sumando en i) recuperamos las expresiones habituales de la pérdida promedio (5).

VII. ACTUALIZACIÓN DE PARÁMETROS (DESCENSO POR GRADIENTE)

Con *learning rate* $\alpha > 0$:

$$w \leftarrow w - \alpha \frac{1}{N} \sum_{i=1}^N (a_i - y_i) x_i, \quad (14)$$

$$b \leftarrow b - \alpha \frac{1}{N} \sum_{i=1}^N (a_i - y_i). \quad (15)$$

(Nota : si la curva de pérdida sube o se vuelve errática, probá reducir α .)

VIII. CÓDIGO

En esta sección explico únicamente el *código fuente* y cómo implementa la teoría vista en clase: verosimilitud \rightarrow log-likelihood \rightarrow log-loss, derivadas por regla de la cadena y actualización de los parámetros w y b . Las Figuras 1–3 corresponden a tres capturas del archivo `main.py`.

VIII-A. Definiciones del modelo: clase, sigmoide y pérdida (`cap01`)

La Figura 1 muestra la clase `LogisticRegressionAI`. En el constructor (`__init__`) se fijan los **hiperparámetros** `lr` (tasa de aprendizaje) y `epochs` (épocas de entrenamiento), y se inicializan los **parámetros entrenables** $w \in \mathbb{R}^d$ y $b \in \mathbb{R}$, que son los únicos valores que el algoritmo ajusta.

El método `sigmoid` implementa la activación logística

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad z = w^\top x + b,$$

que mapea la combinación lineal z a una probabilidad $a = \sigma(z) \in (0, 1)$.

La función `binary_crossentropy_loss` codifica la **log-loss** (negativa del log-likelihood promedio):

$$\begin{aligned} \mathcal{L}(w, b) &= -\frac{1}{N} \sum_{i=1}^N [y_i \ln a_i + (1-y_i) \ln (1-a_i)], \\ a_i &= \sigma(w^\top x_i + b). \end{aligned} \quad (16)$$

Antes de aplicar los logaritmos, el código realiza *clipping* con $\varepsilon = 10^{-15}$ para evitar $\log(0)$ y mejorar la estabilidad numérica. Minimizar (16) es equivalente a maximizar la verosimilitud $\prod_i a_i^{y_i} (1-a_i)^{1-y_i}$.

VIII-B. Entrenamiento vectorizado y predicción (`cap02`)

La Figura 2 concentra el corazón del aprendizaje: el método `fit` (bucle de entrenamiento) y `predict` (inferencia).

Forward (vectorizado): Con $X \in \mathbb{R}^{N \times d}$ y $w \in \mathbb{R}^d$, se calcula

$$z = Xw + b\mathbf{1}, \quad a = \sigma(z),$$

usando multiplicación matricial de numpy (`np.dot`). Este paso implementa la composición $x \rightarrow z \rightarrow a$ vista en clase.

```

class LogisticRegressionAI:

    def __init__(self, lr=0.0001, epochs=1000):
        self.epochs = epochs
        self.lr = lr
        self.w = None
        self.b = None

    # activación sigmoid
    def sigmoid(self, x):
        return 1 / (1 + np.exp(-x))

    # función de costo (log-loss)
    def binary_cross_entropy_loss(self, y_true, y_pred):
        epsilon = 1e-15 # evita log(0)
        y_pred = np.clip(y_pred, epsilon, 1 - epsilon)
        loss = (-1 / len(y_true)) * np.sum(
            y_true * np.log(y_pred) + (1 - y_true) * np.log(1 - y_pred)
        )
        return loss

```

Figura 1. **cap01:** Clase LogisticRegressionAI. Hiperparámetros (lr, epochs), parámetros w, b , sigmoid y pérdida (16).

```

def fit(self, X, y):
    n_samples, n_features = X.shape
    # z = Xw + b
    self.w = np.random.rand(n_features)
    self.b = 0.0

    for epoch in range(self.epochs):
        # forward
        z = np.dot(X, self.w) + self.b
        a = self.sigmoid(z)

        # gradientes promedio
        dw = (1 / n_samples) * np.dot(X.T, (a - y))
        db = (1 / n_samples) * np.sum(a - y)

        # actualización
        self.w = self.w - self.lr * dw
        self.b = self.b - self.lr * db

        # logging (espaciado para que no sature la terminal)
        if epoch % 100 == 0:
            loss = self.binary_cross_entropy_loss(y, a)
            print(f"[fit] epoch={epoch:4d} loss={loss:.4f}")

    def predict(self, X):
        z = np.dot(X, self.w) + self.b
        a = self.sigmoid(z)
        return (a >= 0.5).astype(int)

```

Figura 2. **cap02:** fit (forward $z \rightarrow a$, gradientes dw/db , actualización) y predict (umbral 0,5).

Gradientes (regla de la cadena): De la derivación teórica se obtiene

$$\frac{\partial \mathcal{L}}{\partial z} = a - y, \quad \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{N} X^T (a - y), \quad \frac{\partial \mathcal{L}}{\partial b} = \frac{1}{N} \mathbf{1}^T (a - y).$$

En el código aparecen como $dw = (1/n_samples) * X.T @ (a - y)$ y $db = (1/n_samples) * np.sum(a - y)$.

Actualización (descenso por gradiente): Con tasa $\alpha = lr$:

$$w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}, \quad b \leftarrow b - \alpha \frac{\partial \mathcal{L}}{\partial b},$$

que en el código se implementa como $w = w - lr * dw$ y $b = b - lr * db$. Al disminuir (16) se incrementa la verosimilitud del modelo.

Predictión: El método predict repite $z = Xw + b$ y $a = \sigma(z)$ y umbraliza con 0,5 para devolver etiquetas binarias $\hat{y} \in \{0, 1\}$, consistente con decidir por la clase más probable.

```

acc_manual = accuracy_score(y_test, y_hat)
print("\n--- RESULTADOS MODELO DESDE CERO ---")
print("Accuracy: {:.2f}")
print(classification_report(y_test, y_hat, zero_division=0))

acc_sk = accuracy_score(y_test, y_pred)
print("\n--- RESULTADOS SKLEARN ---")
print("Accuracy: {:.2f}")
print(classification_report(y_test, y_pred, zero_division=0))

```

Figura 3. **cap03:** Flujo completo: imports, generación del dataset con make_classification, división entrenamiento/prueba, entrenamiento del modelo propio y bloque análogo con sklearn.

VIII-C. Bloque principal: imports, dataset y flujo del entrenamiento (cap03)

Como se muestra en la Figura 3, se realiza la importación de librerías necesarias (numpy y módulos de sklearn). Luego, se crea un **dataset de clasificación** con make_classification, que genera datos sintéticos controlados para problemas binarios. En este caso, se indica:

`n_samples = 1000, n_features = 2, n_informative = 2,`

Esto produce dos *features* informativas sin redundancia y un solo clúster por clase, coherente con los ejemplos del curso.

A continuación, se realiza la **división entrenamiento/prueba** con train_test_split, manteniendo la proporción de clases mediante stratify=y y fijando random_state para reproducibilidad. Seguidamente, se **instancia** y **entrena** el modelo implementado desde cero:

`modelo_manual = LogisticRegressionAI(lr=0.001, epochs=1000)`

ejecutando el bucle explicado en la subsección anterior (forward, gradientes y actualización). Finalmente, el bloque incluye LogisticRegression de sklearn para **replicar el mismo enfoque** con la librería estándar, lo que sirve como referencia y valida que la implementación manual respeta la teoría.

VIII-D. Resumen código ↔ teoría

Verosimilitud → log-loss. El código minimiza (16), que es $-\log L$; así, “minimizar la pérdida” equivale a “maximizar la verosimilitud”.

Regla de la cadena. La ruta $x \rightarrow z \rightarrow a \rightarrow \mathcal{L}$ da $\frac{\partial \mathcal{L}}{\partial z} = a - y$, base de los gradientes vectorizados dw y db.

Parámetros actualizables. Solo w y b cambian; el resto (sigmoid, datos) son transformaciones/entradas fijas conforme a la formulación del modelo.

IX. CONCLUSIONES

La clave para no perderse es mirar la composición $x \rightarrow z \rightarrow a \rightarrow \mathcal{L}$ y empujar las derivadas con la regla de la cadena. El uso de logaritmos cambia productos por sumas y, al negar el log-likelihood, pasamos a minimizar una función estable y derivable. Con los gradientes compactos $(a - y)x$ y $(a - y)$, actualizar w y b se vuelve mecánico con descenso por gradiente.

APÉNDICE: FÓRMULAS RELEVANTES

Sigmoide y derivada:

$$\sigma(t) = \frac{1}{1 + e^{-t}}, \quad \sigma'(t) = \sigma(t)(1 - \sigma(t)).$$

Verosimilitud y log-likelihood:

$$L = \prod_i f(x_i)^{y_i} (1-f(x_i))^{1-y_i}, \quad \ln L = \sum_i [y_i \ln f(x_i) + (1-y_i) \ln(1-f(x_i))].$$

Pérdida logística promedio (lo que minimizo):

$$\mathcal{L}(w, b) = -\frac{1}{N} \sum_i [y_i \ln f(x_i) + (1 - y_i) \ln(1 - f(x_i))].$$

Gradientes (por muestra):

$$\begin{aligned} \frac{\partial L}{\partial a} &= -\frac{y}{a} + \frac{1-y}{1-a}, & \frac{\partial a}{\partial z} &= a(1-a), & \frac{\partial z}{\partial w} &= x, & \frac{\partial z}{\partial b} &= 1, \\ \frac{\partial L}{\partial z} &= a - y, & \frac{\partial L}{\partial w} &= (a - y)x, & \frac{\partial L}{\partial b} &= (a - y). \end{aligned}$$