

Apuntes Semana 5 - 18/3/25

John Sánchez C.

Tecnológico de Costa Rica

Abstract—Este documento presenta un resumen estructurado de los conceptos fundamentales en aprendizaje automático cubiertos en clase. Se abordan temas como regresión lineal, descenso del gradiente, sesgo y varianza, validación de modelos, y estrategias para evitar el sobreajuste y subajuste.

I. RESPUESTAS QUIZ 2

1) ¿Qué es un mínimo global y mínimo local en una función?

Un mínimo local es un punto donde la función alcanza un valor menor en comparación con sus alrededores, pero puede haber puntos más bajos en otras regiones. Un mínimo global es el punto más bajo en toda la función.

2) Describa la diferencia entre un problema de clasificación versus un problema de regresión lineal.

Un problema de clasificación tiene un conjunto limitado de clases y busca asignar cada entrada a una de ellas. Un problema de regresión lineal predice un valor continuo basado en una función matemática.

3) Describa paso a paso el algoritmo KNN.

- Dado un nuevo punto de datos, calcular la distancia con todos los elementos en el conjunto de entrenamiento.
- Seleccionar los k puntos más cercanos.
- Realizar una votación entre los k vecinos para determinar la clase del nuevo punto.

4) Calcule la derivada respecto a x de la función $x^3 + 2xy + y^2 + 5x + 3y + 7$.

Aplicando derivadas parciales:

$$\frac{d}{dx}(x^3 + 2xy + y^2 + 5x + 3y + 7) = 3x^2 + 2y + 5. \quad (1)$$

II. INTRODUCCIÓN

El aprendizaje automático se basa en la capacidad de los modelos de identificar patrones en datos y generalizar a nuevas muestras. En esta clase, se presentan conceptos esenciales para entender el proceso de entrenamiento y validación de modelos.

III. REPASO - REGRESIÓN LINEAL Y DESCENSO DEL GRADIENTE

La regresión lineal es un modelo utilizado para predecir valores continuos. Se optimiza mediante el descenso del gradiente, que ajusta los parámetros minimizando la función de pérdida L .

Resultados de derivadas parciales de L

$$\bullet \frac{\partial L}{\partial w} = \frac{1}{N} \sum 2((wx_i + b) - y_i) \cdot x_i$$

$$\bullet \frac{\partial L}{\partial w} = \frac{1}{N} \sum 2(f_{w,b}(x_i) - y_i) \cdot x_i$$

$$\bullet \frac{\partial L}{\partial b} = \frac{1}{N} \sum 2((wx_i + b) - y_i)$$

$$\bullet \frac{\partial L}{\partial b} = \frac{1}{N} \sum 2(f_{w,b}(x_i) - y_i)$$

Fig. 1.

A. Derivadas parciales b y w

- Derivadas parciales de b y w

B. Actualización de Parámetros

- $w = w - \alpha \nabla L$
- α es el learning rate
- Valores pequeños de α requieren muchas iteraciones
- Valores grandes pueden hacer que el modelo no converja
- Se recomienda iniciar con un α de 0.001 e ir ajustando

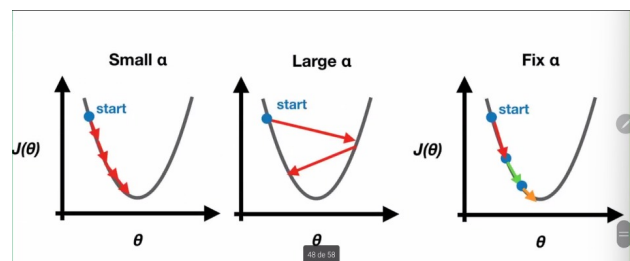


Fig. 2. Comportamiento de α

- Epoch:** Recorrer todo el conjunto de datos de entrenamiento, es un hiperparámetro que se debe definir.
- Batch:** Procesamiento de datos por bloques, subconjunto del total de datos.

C. Algoritmos para aplicar el Descenso del Gradiente

- Batch Gradient Descent:** Usa todo el dataset para cada actualización.
- Stochastic Gradient Descent:** Modifica los parámetros por cada muestra, rápido pero ruidoso.
- Mini-Batch Gradient Descent:** Combina las dos anteriores, mejora estabilidad.

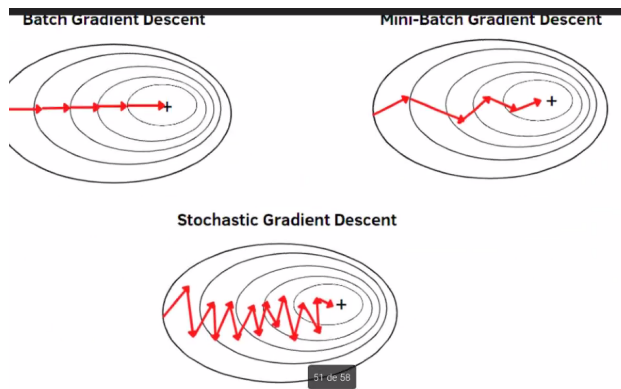


Fig. 3. Algoritmos

IV. SESGO Y VARIANZA

El balance entre sesgo y varianza es clave en el rendimiento del modelo:

- **Alto Sesgo:** Modelo simple, no aprende correctamente (underfitting).
- **Alta Varianza:** Modelo demasiado ajustado a los datos de entrenamiento (overfitting).

V. VALIDACIÓN DE MODELOS

Para evaluar el modelo, se dividen los datos en:

- **Training set:** Para ajustar los parámetros. El modelo identifica patrones basados en estos datos de forma generalizada.
- **Testing set:** Para medir el rendimiento final. Se utiliza este set para probar con datos que nunca se han visto antes. Como un "examen" para el modelo.
- **Validation set:** Este set es un subconjunto del training set. Cada época se hace una prueba de validación para ver si se hace overfitting antes de terminar el entrenamiento. Debería comportarse como el training pero un poco más equivocado.

Técnicas para subdividir el set de validación incluyen:

- **Random Sampling:** División aleatoria. Ideal para datos balanceados, no agrega sesgo pero con datos imbalanceados no se puede tomar representaciones de todas las clases.
- **Stratified Sampling:** Mantiene distribuciones representativas. Ideal para datos imbalanceados, asegura representación de todas las clases.
- **K-Fold Cross-Validation:** Utiliza particiones alternadas del dataset. Se divide el subconjunto en k partes, se reserva k-1 partes para validación.

VI. POSIBLES ESCENARIOS

- **Bajo error en training, bajo error en testing:** El escenario ideal para el modelo, evita el ruido existente en los datos y puede generalizar (aprendió la tarea).
- **Bajo error en training, alto error en testing:** Hay overfitting. El modelo no generaliza, hay alta varianza y está demasiado ajustado a los datos de entrenamiento.

- **Alto error en training, alto error en testing:** Hay underfitting. El modelo no está aprendiendo nada de los datos, modelo muy simple, alto sesgo.
- **Bias Variance tradeoff:** Necesitamos un modelo de baja varianza y bajo sesgo.

VII. BIAS VARIANCE TRADEOFF

A. Alto Bias

El modelo comete muchos errores en training (underfitting). Modelo asume mucho en el training set, no usa todos los features del modelo, modelo es simple. Como evitarlo:

- Usar un modelo más complejo
- Features del training set no son confiables

B. Alta Varianza

El modelo se ajusta mucho a los datos de entrenamiento (overfitting), no generaliza, sucede con datos de alta dimensionalidad y pocos ejemplos. Como evitarlo:

- Usar un modelo más simple.
- Reducir la dimensionalidad
- Obtener más ejemplos en el training set
- Aplicar técnicas de regularización (hacen el modelo más simple).

VIII. CONCLUSIONES

Comprender la relación entre sesgo y varianza, así como seleccionar el método adecuado de optimización y validación, es esencial para el desarrollo de modelos de aprendizaje automático que generalicen correctamente.