

Lección 46

Datos y sesgos del muestreo

Propósito

En esta Lección, los estudiantes conocerán distintos factores que pueden influenciar el resultado de un modelo de aprendizaje automático, y cómo esta práctica puede ser corregida oportunamente cuando se conocen los errores.

Secuencia para el aprendizaje

Conocimiento inicial (10 min)

Ampliación del conocimiento (30 min)

Transferencia del conocimiento (15 min)

Objetivos

Los estudiantes serán capaces de:

- Entender los recursos y las implicaciones de un muestreo sesgado en los conjuntos de datos.

Preparación

- Sitio web Aprendizaje automático para niños listo para visualizarlo.
- Asegúrese de que cada estudiante tenga su [Bitácora de trabajo](#).

Lección en línea

Recursos

¡Atención!

Por favor, haga una copia de cada documento que planea compartir con los estudiantes.

Para los Profesores:

- Computadora(s) o tablets con conexión a internet para acceder a herramientas y recursos en línea.
 - Sitio web: Aprendizaje automático para niños

Para los estudiantes:

- Computadora(s) o tablets con conexión a internet para acceder a herramientas y recursos en línea.
 - Sitio web: Aprendizaje automático para niños

Vocabulario

- **Sesgo de muestreo:** distorsión que sufre un análisis estadístico.

Estrategia de aprendizaje

Conocimiento inicial (10 min)

En esta Lección, los alumnos considerarán las formas en que el muestreo de datos para el entrenamiento y la evaluación de datos puede afectar el resultado de un modelo de aprendizaje automático. Luego, volverán a los conjuntos de datos para su modelo "Hazme feliz" y realizarán mejoras para refinar los conjuntos de datos.

Ampliación del conocimiento (30 min)

Pregunte a los alumnos: ¿Les gustó la forma en que funcionó el modelo "Hazme feliz"? ¿Fue siempre correcta? ¿Podría haber funcionado de mejor manera?

Comente a los estudiantes que una IA no tiene opiniones ni pensamientos propios, sino que solo puede tomar decisiones con base en los datos de los que aprende. Debido a que los modelos de aprendizaje automático aprenden de los datos de entrenamiento, la calidad de la muestra de los datos de entrenamiento determina de forma directa la calidad del modelo. El sesgo en el muestreo, que daría lugar a un modelo inexacto, es causado por tener un conjunto de datos que no representa con precisión las etiquetas. Un conjunto de datos de calidad tiene las siguientes características:

- Suficientes datos. La IA necesita suficientes ejemplos para poder identificar patrones en las características de los datos. La cantidad de datos necesaria depende del objetivo de desempeño específico para la precisión en la tarea. A mayor necesidad de precisión, será mayor cantidad de datos de entrenamiento necesarios.
- Datos correctos: La IA necesita recibir suficientes tipos correctos de ejemplos para comprender todas las características correctas de los elementos a los que se les debe dar una etiqueta particular con precisión. Esto significa que, si faltan ejemplos de esa etiqueta o son confusos, la IA probablemente no pueda identificarlos correctamente en la evaluación de los datos. Algunos ejemplos de sesgo de muestreo incluyen lo siguiente:
 - Si uno entrena a la IA para identificar insectos, pero solo es entrenada con ejemplos de hormigas y escarabajos, es probable que no pueda identificar una mantis religiosa como insecto porque algunas características de los insectos faltarían en los datos. Del mismo modo, si el modelo es utilizado por muchas personas en muchos países, pero solo los insectos de un país están representados, es posible que la IA no pueda reconocer insectos de otros lugares.
 - Si se entrenó a la IA usando un conjunto de datos en el que todas las imágenes de insectos se tomaron entre la hierba, pero las imágenes de los que no eran insectos se tomaron en una variedad de lugares, la IA podría identificar la hierba como una característica de la etiqueta del insecto. Si el modelo se probó con una imagen de un perro en el césped, podría etiquetarlo como un insecto.
 - Si hay una característica que puede aplicarse a ambas etiquetas, pero solo se muestra en una, esto podría confundir a la IA. Por ejemplo, si todas las imágenes de entrenamiento de no-insectos tampoco fue de no-animales, la IA podría confundir a cualquier animal con patas, cabeza, etc. con un insecto cuando se pruebe el modelo.
 - Si los datos de entrenamiento incluyen significativamente más ejemplos de una etiqueta

que de otra, entonces la IA podría aprender que la primera etiqueta es más común y, por lo tanto, seleccionará incorrectamente esa etiqueta con más frecuencia.

Transferencia del conocimiento (15 min)

Invite a los alumnos a reflexionar sobre la experiencia colaborativa de crear el modelo “Hazme feliz”. Pídales que identifiquen los elementos que funcionaron bien de los conjuntos de datos de entrenamiento y prueba que fueron creados. Luego, pregúnteles qué se podría hacer para mejorar los resultados del modelo. Registre sus respuestas y exhíbalas al grupo. Haga que los alumnos regresen a su proyecto “Hazme feliz”, realicen las mejoras del conjunto de datos que identificaron y prueben el modelo para ver si los resultados mejoran. Haga que los alumnos continúen refinando el conjunto de datos hasta que el modelo sea capaz de identificar correcta y constantemente oraciones con un alto grado de confianza.

Opcional: Amplíe esta Lección de aprendizaje examinando los cuatro ejemplos de herramientas de aprendizaje automático en la Lección 44 y analizando las posibles fuentes e impactos del sesgo de muestreo. Quizás sea recomendable que los alumnos vayan más allá e investiguen acontecimientos de la actualidad que describan incidentes en los que el sesgo de muestreo condujo a resultados negativos en el uso de modelos de aprendizaje automático en aplicaciones. Si bien esta extensión del proyecto revelará una desventaja de usar el aprendizaje automático en las aplicaciones, los alumnos deben centrarse en la importancia de evitar el sesgo de muestreo y el papel que juegan las personas en el resultado del modelo de aprendizaje automático.

Escribir en el diario y charla rápida

El acto de escribir en sus diarios sobre lo aprendido, respecto de si les pareció útil y de lo que sintieron, ayuda a sus estudiantes a fortalecer cualquier conocimiento que hayan obtenido hoy y servir como un resumen al que puedan recurrir en el futuro.

Sugerencias para el diario:

- ¿Sobre qué se trataba la Lección de hoy?
- ¿Cómo te sentiste durante la Lección?

Sugerencias para evaluación

Se sugiere el siguiente indicador para evaluar formativamente los aprendizajes:

- Debaten sobre la importancia de evitar el sesgo de muestreo y el papel que juegan las personas en el resultado del modelo de aprendizaje automático.