

TEORÍA

1. *En la empresa GA, en el área de compras necesitan CLASIFICAR y organizar los correos que llegan a la bandeja de entrada entre 4 tipos de correos (Compras cementos, Compras energía, Compras concretos y correos generales o de otra índole). Esta tarea se le encomienda a usted, gracias a su rol puede solicitar al área interesada los recursos humanos que necesite para llevar a cabo este proyecto, también puede solicitar en tecnología todo lo que necesite, además tiene las bandejas de entrada de correos históricos de los analistas que reciben estas solicitudes con aproximadamente: 5500 correos de compras cementos, 2700 correos de compras de energía, 1100 correos de compras concretos y 12876 correos generales o de otra índole. Explique cómo resolvería este problema, metodología, algoritmos, modelos, arquitectura del proyecto etc.*

1.1. Identificación del problema

Este caso representa la necesidad de analizar un texto e identificar patrones que puedan ayudar a clasificarlo en una de las 4 posibles opciones. Como se describe en el problema, se goza con una base de datos con 22,176 correos clasificados cada uno de ellos en la clase correspondiente, teniendo en cuenta que hay un desbalance en la clasificación de la data para entrenamiento.

1.2. Metodología

Recolección de datos : Se toma los correos con los que se dispone para el entrenamiento

Exploración : Se identifican comportamientos y necesidades de limpieza de datos, así como también la distribución de las clases y requerimientos adicionales como procesamiento en diferentes idiomas

Limpieza de datos : Este es uno de los pasos más importantes, pues se realiza toda la limpieza de texto que permita el acceso a los datos en su forma más pura posible (ejemplo la eliminación de caracteres especiales, signos de puntuación, lematización, eliminación de palabras que no aporten valor, convertir a minúsculas, etc)

Definición de posibles modelos : Actualmente hay varios modelos que han sido pre-entrenados con una gran cantidad de texto y permite definir de allí características para la búsqueda de un contexto para personalizar la solución de clasificación en nuestro caso. Uno de los modelos que ha tenido mucho auge últimamente es BERT

Definición del conjunto de entrenamiento y test : Normalmente se toma entre un 20% de los datos para test y el resto para entrenamiento, sin embargo, este muestreo debe estar basado en el desbalance de las clases

Definición de métricas para evaluar el desempeño del modelo : Dichas métricas deben tener un castigo por el desbalance de clases, puede ser f1 macro o f1 para cada clase

Optimización : Ajuste de los diferentes hiperparámetros de los modelos evaluados para obtener los óptimos de cada uno de ellos

Elección del modelo : Se toma el modelo que mejor se ajuste al problema, tomando como criterio de decisión las métricas de desempeño para casos desbalanceados y matriz de confusión para evaluar la distribución de las clasificaciones

Despliegue : Puesta en servicio del mejor modelo

Monitoreo : Es muy factible que se presente alteraciones del comportamiento de los correos recibidos dado cambios de situación país, criterios de negocio, crecimiento/decrecimiento de líneas de negocio, cambio de proveedores, etc ; por lo que es importante realizar monitoreo del modelo desplegado para hacer ajustes si es del caso

1.3. Definición de recursos

De acuerdo a la metodología planteada, se debe evaluar los recursos humanos y de infraestructura para llevar a cabo cada una de las etapas del proceso y luego de ello, asegurarse de su mantenimiento en el tiempo

1.4. Puesta en marcha del proyecto

Se definen alcances, tareas, tiempos y seguimiento de cada uno de los recursos asignados con el propósito de ejecutar a cabalidad el proyecto.

2. Seis meses después de haber desplegado un modelo de regresión en producción, los usuarios se dan cuenta que las predicciones que este está dando no son tan acertadas, se le encarga a usted que revise que puede estar sucediendo.

¿Cree que el modelo esté sufriendo Drift?

Efectivamente puede estar sufriendo Drift ya que podría ser muy sensible a cualquiera o varias de las siguientes situaciones :

- **Condiciones de mercado** : Ya sean internas o externas que generen crecimiento/decrecimiento de líneas de negocio y por ende se mueva la distribución de los correos entrantes
- **Dinámica de los correos entrantes** : Por cambios de proveedores, cambios en la clasificación de los correos ya sea por políticas internas o errores humanos.
- **Cambios en vocabulario o aparición de nuevos patrones** : Referenciado a la naturaleza del texto procesado

¿Cómo puede validarlo?

Inicialmente debe hacerse un sondeo sobre alguna de las hipótesis planteadas en el ítem anterior, indagar con los equipos afectados para identificar posibles cambios en las reglas de negocio o definiciones del modelo y con ello revisar del modelo lo siguiente :

- Evaluar cambios en las distribuciones de las clases actuales vs las de entrenamiento.
- Revisar cambios en las métricas y si hay un comportamiento focalizado en alguna de las clases.

¿De ser así, que haría usted para corregir esto? Explique sus respuestas.

- Reentrenar el modelo con la data verificada y correctamente etiquetada si se identifica cambios en las distribuciones de las clases, reglas de negocio, vocabulario, proveedores y demás factores que puedan estar alterando la clasificación
- Considerar arquitecturas dinámicas que permitan ajustarse continuamente
- Monitoreo continuo para alertar oportunamente posibles desviaciones
- Entrenamiento periódico para garantizar la funcionalidad de la herramienta

3. Su equipo de trabajo está trabajando en un chatbot con generación de texto utilizando el modelo GPT-3.5, según cómo funciona este modelo, ¿cómo haría usted para hacer que las respuestas del chatbot estén siempre relacionadas a conseguir cierta información particular del usuario y no empiece a generar texto aleatorio sobre cualquier tema? Explique su respuesta.

En la configuración del chatbot puede especificarse el contexto a necesidad del equipo de trabajo. En dicho contexto se definen las reglas de las respuestas y su filtro de acuerdo a las preguntas ingresadas en caso de que correspondan o no a la finalidad del chatbot

También pueden modificarse otros parámetros como temperature que controlan las respuestas a ser más precisas, se puede limitar la extensión de la respuesta para obligar una respuesta más contundente y análisis de historial de conversaciones.