

Este documento proporciona una guía paso a paso para llevar a cabo un proyecto de aprendizaje automático de extremo a extremo, utilizando como ejemplo la predicción de precios de viviendas en California. El documento se enfoca en el proceso, no en el negocio de bienes raíces. Los pasos principales son:

- **Observar el panorama general:** Definir el objetivo del negocio y cómo se usará el modelo. En este caso, el modelo predecirá el precio medio de la vivienda de un distrito y se utilizará en un sistema de aprendizaje automático posterior para determinar si vale la pena invertir en un área. El documento también explica la importancia de comprender el problema a resolver y el tipo de aprendizaje a aplicar: en este caso, es un problema de aprendizaje supervisado de regresión.
- **Obtener los datos:** Se proporcionan enlaces a repositorios de datos abiertos y se explica cómo descargar un conjunto de datos de precios de viviendas de California desde un repositorio de GitHub. Se enfatiza la importancia de automatizar la descarga de datos y se muestra cómo cargar los datos en un DataFrame de Pandas.
- **Explorar y visualizar los datos:** Se describe cómo usar varios métodos de Pandas para obtener información sobre los datos, como `.head()`, `.info()`, `.value_counts()`, `.describe()` e `.hist()`. También se muestra cómo visualizar los datos geográficos utilizando gráficos de dispersión y cómo usar el color y el tamaño de los puntos para representar diferentes atributos. Se aborda la correlación entre atributos, con el uso de `.corr()` y `scatter_matrix()`. Se explica la importancia de observar la correlación no lineal. Se muestra cómo crear nuevas combinaciones de atributos y evaluar su correlación con el valor medio de la vivienda.
- **Preparar los datos para los algoritmos de aprendizaje automático:** Se explica cómo manejar los valores faltantes, se utiliza `SimpleImputer` para completar los valores faltantes utilizando la mediana. Se muestra cómo codificar los atributos categóricos, utilizando `OrdinalEncoder` y `OneHotEncoder`. Se explica la importancia de la escalada de atributos y cómo utilizar `MinMaxScaler` y `StandardScaler`. Se muestra cómo transformar los atributos utilizando `FunctionTransformer`, incluyendo el cálculo de logaritmos y el uso de la función de base radial gaussiana (RBF). También se introduce cómo crear transformadores personalizados con las clases `BaseEstimator` y `TransformerMixin`, además de cómo usar la clase `ClusterSimilarity`. Se explica cómo utilizar la clase `Pipeline` para encadenar múltiples transformaciones y cómo usar la clase `ColumnTransformer` para aplicar diferentes transformaciones a diferentes columnas. Se define un pipeline que realiza imputación, escalado, creación de ratio features, cluster similarity features, transformaciones logarítmicas y escalado estándar.
- **Seleccionar un modelo y entrenarlo:** Se explica cómo utilizar el pipeline de preprocesamiento y cómo entrenar modelos de regresión lineal, árboles de decisión y bosques aleatorios. También se menciona la importancia de la validación cruzada para obtener una estimación más fiable del rendimiento del modelo.
- **Afinar el modelo:** Se presenta cómo utilizar la clase `GridSearchCV` para buscar la mejor combinación de hiperparámetros utilizando validación cruzada. También se explica cómo utilizar `RandomizedSearchCV` para buscar hiperparámetros de forma aleatoria. Se menciona cómo usar métodos de ensamble y cómo analizar los mejores modelos y sus errores. Se explica la importancia de evaluar el rendimiento del modelo en diferentes categorías de datos.
- **Presentar la solución:** Se recomienda crear documentación y presentaciones claras para resaltar lo aprendido. Se discute el balance entre las expectativas del modelo y los expertos.

- **Lanzar, monitorear y mantener el sistema:** Se explica cómo guardar y cargar el modelo utilizando la librería `joblib`. Se menciona cómo desplegar el modelo como un servicio web y cómo utilizar plataformas en la nube como Google Vertex AI. Se enfatiza la importancia de monitorear el rendimiento del modelo en vivo y de actualizar los datos y los modelos periódicamente. Se discute la necesidad de tener copias de seguridad de los modelos y de los conjuntos de datos. Se enfatiza que construir un sistema de aprendizaje automático implica bastante trabajo en infraestructura.

En resumen, el documento proporciona una guía práctica y completa para desarrollar proyectos de aprendizaje automático, desde la obtención de los datos hasta el lanzamiento y el mantenimiento del modelo. Destaca la importancia del preprocesamiento de los datos y la necesidad de comprender el problema y los datos con los que se está trabajando, también la importancia de usar `Pipeline` y `ColumnTransformer` para encadenar transformaciones de datos. El documento proporciona ejemplos de código que el usuario puede ejecutar.