

Creacion de características por agrupamiento y separación

Agrupamiento de datos

Agrupamiento

Método groupby()

Operaciones Comunes con groupby()

Método get_group()

Casos de Uso Comunes

Buenas Prácticas

Metodos y Funciones

1. Lectura y Escritura de Datos

2. Exploración de Datos

3. Limpieza de Datos

4. Transformación de Datos

5. Filtrado y Selección

6. Agregación y Resumen

Escalamiento

Escalamiento del máximo absoluto

Escalamiento min-max

Estandarización

Cuando la relación entre x y y no es lineal.

Explicación sobre potencias y polinomios:

Ejemplo de regresión polinómica:

Conclusión:

Logaritmos y su aplicación en ML

Logaritmos: Explicación y Aplicaciones

Propiedades de los Logaritmos

Tipos Comunes de Logaritmos

Transformación de Datos:

Agrupamiento de datos

Agrupamiento

El agrupamiento en pandas mediante `groupby()` y `get_group()` son herramientas fundamentales para el análisis y manipulación de datos estructurados.

Método `groupby()`

Definición: Es una operación que implica uno o más de los siguientes pasos:

- Dividir los datos en grupos según uno o más criterios
- Aplicar una función a cada grupo de manera independiente
- Combinar los resultados en una estructura de datos

Sintaxis básica:

```
df.groupby('columna')  
df.groupby(['columna1', 'columna2'])
```

Operaciones Comunes con `groupby()`

- Agregación: `.agg()`, `.sum()`, `.mean()`, `.count()`
- Transformación: `.transform()`
- Filtrado: `.filter()`
- Iteración: `.groups`, `.get_group()`

Método `get_group()`

Definición: Permite extraer un grupo específico de un objeto `GroupBy`

Ejemplo de uso:

```
# Crear un groupby object  
grouped = df.groupby('categoria')  
  
# Obtener un grupo específico  
grupo_especifico = grouped.get_group('valor_categoria')
```

Casos de Uso Comunes

- Análisis de ventas por región o período
- Estadísticas por categorías
- Segmentación de datos para análisis específicos
- Cálculos agregados por grupo

Buenas Prácticas

- Verificar la existencia de grupos antes de usar `get_group()`
 - Considerar el impacto en memoria al trabajar con grandes datasets
 - Utilizar métodos de agregación eficientes
 - Documentar las operaciones de agrupamiento para mantenibilidad
-

Metodos y Funciones

Métodos Esenciales para Manipulación de Datos:

1. Lectura y Escritura de Datos

- `pd.read_csv()`, `pd.read_excel()`: Importar datos

```
# Leer un archivo CSV
df = pd.read_csv('datos.csv')
# Leer un archivo Excel
df = pd.read_excel('datos.xlsx', sheet_name='Hoja1')
```

- `to_csv()`, `to_excel()`: Exportar datos

```
# Guardar en CSV
df.to_csv('output.csv', index=False)
# Guardar en Excel
df.to_excel('output.xlsx', sheet_name='Resultados')
```

2. Exploración de Datos

- `head()`, `tail()`: Visualizar primeras/últimas filas

```
# Ver primeras 5 filas
print(df.head())
# Ver últimas 3 filas
print(df.tail(3))
```

- `info()`, `describe()`: Información y estadísticas

```
# Información del DataFrame
df.info()
# Estadísticas descriptivas
print(df.describe())
```

3. Limpieza de Datos

- `dropna()`, `fillna()`: Manejo de valores nulos

```
# Eliminar filas con valores nulos
df_limpio = df.dropna()
# Rellenar valores nulos con un valor específico
df['columna'].fillna(0, inplace=True)
```

- `drop_duplicates()`: Eliminar duplicados

```
# Eliminar filas duplicadas
df_sin_duplicados = df.drop_duplicates()
# Eliminar duplicados basados en columnas específicas
df_sin_duplicados = df.drop_duplicates(subset=['columna1', 'columna2'])
```

4. Transformación de Datos

- `apply()`, `map()`: Aplicar funciones

```
# Aplicar función a una columna
df['nueva_columna'] = df['columna'].apply(lambda x: x * 2)
# Mapear valores
mapeo = {'A': 1, 'B': 2, 'C': 3}
df['valores_mapeados'] = df['columna'].map(mapeo)
```

5. Filtrado y Selección

- `loc[], iloc[]`: Selección de datos

```
# Selección por etiquetas
resultado = df.loc[df['edad'] > 25]
# Selección por posición
primeras_filas = df.iloc[0:5, [0, 2]]
```

6. Agregación y Resumen

- `value_counts()`, `pivot_table()`

```
# Contar valores únicos
conteo = df['categoria'].value_counts()
# Crear tabla dinámica
tabla_pivot = df.pivot_table(
    values='ventas',
    index='región',
    columns='producto',
    aggfunc='sum'
)
```

Funciones Auxiliares:

- `numpy.where()`: Condiciones lógicas

```
# Crear columna basada en condición
df['estado'] = np.where(df['valor'] > 100, 'Alto', 'Bajo')
```

- pandas.cut(), pandas.qcut(): Discretización

```
# Discretización en intervalos iguales
df['grupos'] = pd.cut(df['valor'], bins=4)
# Discretización por cuantiles
df['cuartiles'] = pd.qcut(df['valor'], q=4)
```

Escalamiento

El escalamiento de datos es una técnica fundamental en el preprocesamiento de datos que consiste en transformar las variables numéricas para que estén en una escala similar o comparable.

¿Por qué es importante el escalamiento?

- Mejora el rendimiento de muchos algoritmos de machine learning
- Evita que variables con rangos más grandes dominen sobre otras en los cálculos
- Acelera la convergencia en algoritmos basados en gradiente
- Permite comparar características en una misma escala

Esta técnica es especialmente crucial cuando trabajamos con algoritmos sensibles a la magnitud de los datos, como K-means, redes neuronales o métodos basados en distancias.

Escalamiento del máximo absoluto

Fórmula matemática:

$$X_{scaled} = \frac{X}{|X_{max}|}$$

Donde:

- X es el valor original
- $|X_{\{max\}}|$ es el valor absoluto del máximo

Implementación en Python usando scikit-learn:

```
from sklearn.preprocessing import MaxAbsScaler

# Crear y ajustar el escalador
scaler = MaxAbsScaler()
X_scaled = scaler.fit_transform(X)

# Para un solo valor
# X_scaled = X / np.abs(X).max()
```

Este método escala los datos al rango [-1, 1] preservando el cero y es útil para datos dispersos.

Escalamiento min-max

Fórmula matemática:

$$X_{scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Donde:

- X es el valor original
- $X_{\{min\}}$ es el valor mínimo
- $X_{\{max\}}$ es el valor máximo

Implementación en Python usando scikit-learn:

```
from sklearn.preprocessing import MinMaxScaler

# Crear y ajustar el escalador
scaler = MinMaxScaler()
X_scaled = scaler.fit_transform(X)
```

```
# Para un solo valor
# X_scaled = (X - X.min()) / (X.max() - X.min())
```

Este método escala los datos al rango [0, 1] y es útil cuando se necesita acotar todos los valores a un rango específico.

Estandarización

Fórmula matemática:

$$X_{scaled} = \frac{X - \mu}{\sigma}$$

Donde:

- X es el valor original
- μ es la media de la distribución
- σ es la desviación estándar

Implementación en Python usando scikit-learn:

```
from sklearn.preprocessing import StandardScaler

# Crear y ajustar el escalador
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Para un solo valor
# X_scaled = (X - X.mean()) / X.std()
```

Este método transforma los datos para que tengan media 0 y desviación estándar 1, y es especialmente útil cuando los datos siguen una distribución normal.

Cuando la relación entre x y y no es lineal.

Cuando esto ocurre, se recurre a una forma más flexible de modelar la relación entre ambas variables: los **polinomios**.

Explicación sobre potencias y polinomios:

1. Polinomios de grado k:

- En estos casos, en lugar de una relación simple $y=x$, se usa una función de **grado k** de la forma:

$$y = w_0 + w_1x + w_2x^2 + \dots + w_kx^k$$

- Los coeficientes w_0, w_1, \dots, w_k son constantes que se ajustan para que el modelo se adapte a los datos observados.

2. Transformación de los datos:

- Cuando tenemos datos donde **x** es la variable independiente y **y** es la variable dependiente, y la relación no es lineal, podemos **transformar** las variables. En lugar de usar solo x, utilizamos x^1, x^2, \dots, x^k como nuevas características.
- **Ejemplo:** Para un conjunto de datos con variables x y y, podemos crear nuevas variables como x^2, x^3, \dots y entrenar un modelo de regresión con estas variables para ajustar el modelo no lineal.

3. Uso de regresión polinómica:

- La **regresión polinómica** es una técnica efectiva para modelar relaciones no lineales. Esta técnica permite ajustar modelos más complejos, haciendo que los datos se ajusten a una curva en lugar de una línea recta.
- **Ejemplo práctico:** Cuando la relación entre x y y sigue una curva (como una parábola), puedes emplear un **polinomio cuadrático** o **polinomios de mayor grado** para capturar esta forma curva.

4. Ventajas y desventajas:

- **Ventaja:** La regresión polinómica ofrece mayor flexibilidad en el ajuste de los datos.
- **Desventaja:** Un grado demasiado alto puede causar **sobreajuste** (overfitting), donde el modelo captura el ruido en lugar de la verdadera

relación entre las variables.

Ejemplo de regresión polinómica:

Para un conjunto de datos donde x y y muestran una relación cuadrática, un modelo polinómico de grado 2 sería apropiado:

$$y = w_0 + w_1x + w_2x^2$$

Donde el modelo ajustará los parámetros w_0 , w_1 y w_2 para minimizar el error en los datos observados.

Conclusión:

El uso de **potencias** y **polinomios** es fundamental para modelar relaciones no lineales entre variables. Esta técnica permite que el modelo capture mejor la complejidad de los datos, aunque requiere una evaluación cuidadosa para evitar el sobreajuste.

Logaritmos y su aplicación en ML

Logaritmos: Explicación y Aplicaciones

Un **logaritmo** es la operación inversa de la **exponenciación**. Si tienes una ecuación de la forma:

$$y = bx$$

El logaritmo de y con base b es el valor de x :

$$\log_b(y) = x$$

Esto significa que $bx = y$, donde b es la base del logaritmo, x es el exponente, y y es el valor.

Propiedades de los Logaritmos

Algunas propiedades clave de los logaritmos incluyen:

1. **Propiedad del Producto:** $\log_b(x \cdot y) = \log_b(x) + \log_b(y)$ Los logaritmos permiten transformar productos en sumas.
2. **Propiedad del Cociente:** $\log_b(x/y) = \log_b(x) - \log_b(y)$ Los logaritmos permiten transformar divisiones en restas.
3. **Propiedad de la Potencia:** $\log_b(x^n) = n \cdot \log_b(x)$ Los logaritmos permiten transformar exponentes en multiplicaciones.
4. **Cambio de Base:** $\log_b(x) = \log_k(x) / \log_k(b)$ Esto te permite convertir logaritmos entre diferentes bases.

Tipos Comunes de Logaritmos

1. Logaritmo Natural (ln):

El

logaritmo natural usa la base e , que es un número irracional aproximadamente igual a 2.71828. Este tipo de logaritmo es común en **cálculo** y **estadística**. $\ln(x) = \log_e(x)$

2. Logaritmo en Base 10 (log):

El

logaritmo en base 10 es comúnmente utilizado en **ciencias** y **estadística**. $\log(x) = \log_{10}(x)$

3. Logaritmo en Base 2:

El

logaritmo en base 2 es utilizado en **informática** y **teoría de la información**. $\log_2(x)$

Transformación de Datos:

En

estadística y **machine learning**, los logaritmos son herramientas valiosas para transformar distribuciones sesgadas o con colas largas, como las distribuciones **log-normales**, en distribuciones más simétricas. Esta transformación mejora la efectividad de los modelos estadísticos.

- **Ejemplo:** Cuando un conjunto de datos de ingresos presenta valores muy dispersos, la aplicación de un **logaritmo** puede normalizar su distribución.