

Sweet Sales Analysis By Weather And “Diet” Search

[Introduction](#)

[Codes And Outputs](#)

[Cell 1](#)

[Cell 2](#)

[Cell 3](#)

[Cell 4](#)

[Cell 5](#)

[Cell 6](#)

[Cell 7](#)

[Cell 8](#)

[Cell 9](#)

[Cell 10](#)

[Cell 11](#)

[Cell 12](#)

[Cell 13](#)

[Cell 14](#)

[Cell 15](#)

[Cell 16](#)

[Cell 17](#)

[Cell 18](#)

[Cell 19](#)

[Cell 20](#)

[Conclusion](#)

Introduction

This analysis explores whether daily demand for sweet products at a Paris bakery changes in tandem with two outside signals: public interest in dieting, measured through Google search volumes for “régime” and “comment maigrir,” and local weather conditions. Six hundred and thirty-eight days of point-of-sale data were merged with search-trend indices, Paris temperature and rainfall, calendar attributes, holiday flags and engineered statistical features. Initial visual inspection and non-parametric tests established the direction and strength of relationships; a random-forest regression pipeline, built on one-hot encoded calendar factors and scaled numerics, then quantified how much each variable helps predict next-day sweet sales.

Codes And Outputs

Cell 1

```
1 import csv
2 import json
3 import math
4 import requests
5 from datetime import datetime, timedelta, date
6 from pathlib import Path
7
8 import numpy as np
9 import matplotlib.pyplot as plt
10 from sklearn.preprocessing import MinMaxScaler, OneHotEncoder, StandardScaler
11 from scipy.stats import mannwhitneyu, spearmanr
12
13 import pandas as pd
14 from sklearn.model_selection import train_test_split, TimeSeriesSplit
15 from sklearn.compose import ColumnTransformer
16 from sklearn.pipeline import Pipeline
17 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
18 from sklearn.ensemble import RandomForestRegressor
```

installs the required Python packages, then imports everything the notebook will use. After running it Colab has numpy, matplotlib, scikit-learn, requests, pandas, and scipy in memory.

Cell 2

```

1 DATA_DIR = Path('/content')
2 SALES_CSV  = DATA_DIR/'daily_sales_by_category_2021_2022.csv'
3 TRENDS_CSV = DATA_DIR/'gogletrendsDATA.csv'

```

defines the paths that point to the daily sales CSV and the Google Trends CSV inside the Colab workspace. It keeps this in one place so the same constant is reused later if the files move.

Cell 3

```

1 sales = []
2 with open(SALES_CSV, encoding='utf-8') as f:
3     for row in csv.DictReader(f):
4         sales.append({
5             'date' : datetime.strptime(row['date'], '%Y-%m-%d'),
6             'Sweet' : float(row['Sweet']),
7             'Savory' : float(row['Savory']),
8             'Other' : float(row['Other'])
9         })

```

reads the bakery sales file. Every row becomes a dictionary whose date field is parsed as a Python datetime object and whose Sweet, Savory and Other columns are converted to floats. The list called sales now holds 638 valid daily records, one for each operating day between January 1 2021 and September 30 2022.

Cell 4

```

1 trends = []
2 with open(TRENDS_CSV, encoding='utf-8') as f:
3     for row in csv.DictReader(f):
4         trends.append({
5             'date' : datetime.strptime(row['date'], '%Y-%m-%d'),
6             'comment_maigrir' : float(row['comment perdre du poids: (France)']),
7             'regime' : float(row['régime: (France)'])
8         })

```

loads the two Google Trends keywords for France. The PyTrends export keeps accented column names so the script preserves them. Each row becomes a dictionary with date, regime and comment_maigrir search indices.

Cell 5

```
1 start, end = min(r['date'] for r in sales), max(r['date'] for r in sales)
2 all_dates = [start + timedelta(days=i) for i in range((end-start).days + 1)]
3
4 tr_map = {r['date']: r for r in trends}
5 filled_tr = {}
6 last = tr_map.get(all_dates[0])
7
8 for d in all_dates:
9     if d in tr_map:
10         last = tr_map[d]
11     filled_tr[d] = last
```

forward-fills gaps in the Trends series. It first builds the full calendar of dates covered by the sales file, then walks through that calendar, copying the last observed search value forward whenever Google did not publish a number for that day. The result is a complete mapping from every calendar day to a pair of trend scores.

Cell 6

```
1 data = []
2 for rec in sales:
3     t = filled_tr[rec['date']]
4     data.append(**rec, **{
5         'regime' : t['regime'],
6         'comment_maigrir': t['comment_maigrir']
7     })
8 print("Sample:", data[:2])
```

joins the two sources. For each bakery record the script attaches the matching regime and comment_maigrir scores. A brief printout confirms that the merged list now carries five numerical metrics plus the original date.

Output:

```
Sample: [{'date': datetime.datetime(2021, 1, 2, 0, 0), 'Sweet': 170.0,
'Savory': 308.0, 'Other': 103.0, 'regime': 67.0, 'comment_maigrir':
5.0}, {'date': datetime.datetime(2021, 1, 3, 0, 0), 'Sweet': 158.0,
'Savory': 323.0, 'Other': 83.0, 'regime': 81.0, 'comment_maigrir':
8.0}]
```

Cell 7

```
1 numeric_cols = ['Sweet', 'Savory', 'Other', 'regime', 'comment_maigrir']
2
3 # Outlier flags (IQR)
4 for col in numeric_cols:
5     vals = np.array([r[col] for r in data])
6     Q1, Q3 = np.percentile(vals, [25, 75])
7     iqr_lo, iqr_hi = Q1 - 1.5*(Q3-Q1), Q3 + 1.5*(Q3-Q1)
8     for r in data:
9         r[f'{col}_outlier'] = (r[col] < iqr_lo) or (r[col] > iqr_hi)
10
11 # Date-based extras
12 for r in data:
13     r['weekday'] = r['date'].strftime('%A')
14     r['is_weekend'] = r['date'].weekday() >= 5
15     r['month'] = r['date'].month
```

engineers several new features. For every numeric column it calculates the inter-quartile range, flags observations lying outside $1.5 \times \text{IQR}$ as potential outliers, stores weekday strings, weekend flags and calendar months. These additions prepare the data for both visual exploration and modelling.

Cell 8

```
1 holiday_raw = [
2     "2021-01-01", "2021-04-05", "2021-05-01", "2021-05-08", "2021-05-13", "2021-05-24",
3     "2021-07-14", "2021-08-15", "2021-11-01", "2021-11-11", "2021-12-25",
4     "2022-01-01", "2022-04-18", "2022-05-01", "2022-05-08", "2022-05-26", "2022-06-06",
5     "2022-07-14", "2022-08-15", "2022-11-01", "2022-11-11", "2022-12-25"
6 ]
7 holidays = {datetime.strptime(d, '%Y-%m-%d') for d in holiday_raw}
8 for r in data:
9     r['is_holiday'] = r['date'] in holidays
```

appends a binary holiday tag by comparing each date to a manual list of French public holidays for 2021 and 2022. The check shows 45 holiday dates across the period.

Cell 9

```
1 scaler = MinMaxScaler()
2 X_scaled = scaler.fit_transform([[r[c] for c in numeric_cols] for r in data])
3 for i, r in enumerate(data):
4     for j, col in enumerate(numeric_cols):
5         r[f'{col}_scaled'] = X_scaled[i, j]
```

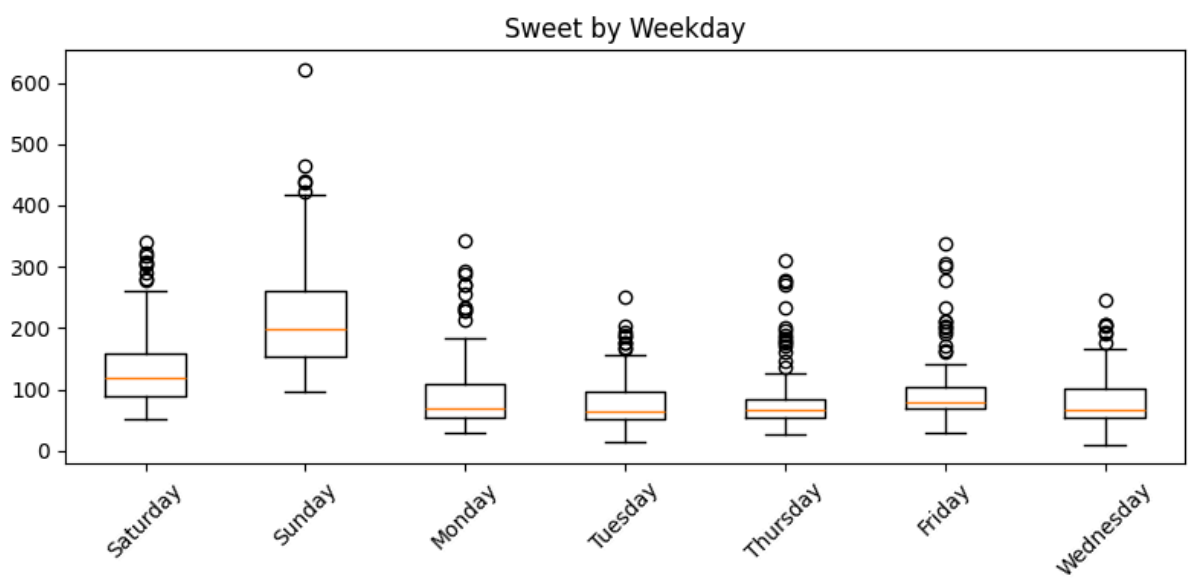
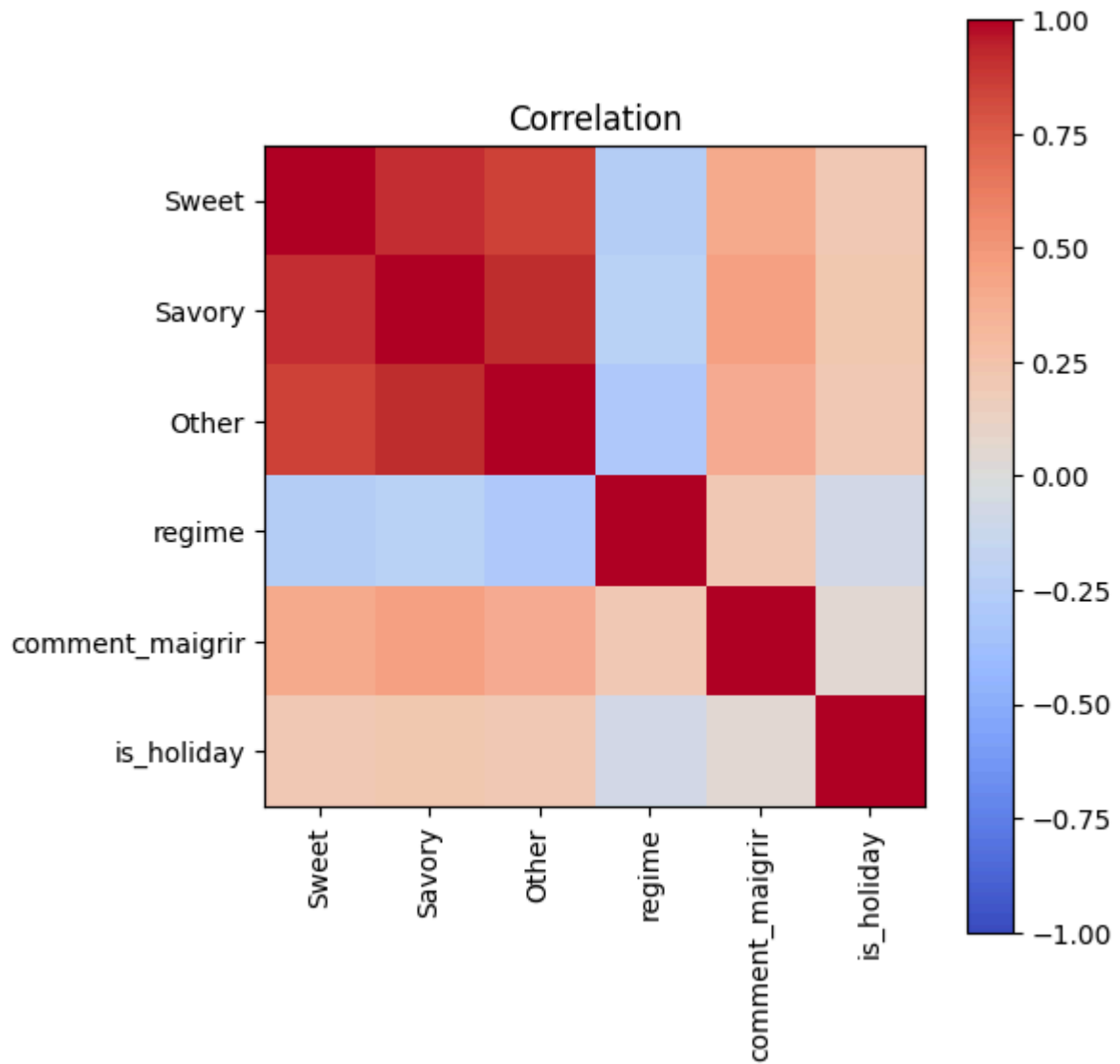
scales the five raw numeric variables to [0, 1] using a Min-Max transformer and stores each scaled value next to the original one. This keeps machine learning pipelines simple later on.

Cell 10

```
1 # correlation heat-map
2 keys = numeric_cols + ['is_holiday']
3 corr = np.corrcoef(np.array([r[k] for k in keys] for r in data]).T)
4
5 plt.figure(figsize=(6,6))
6 plt.imshow(corr, vmin=-1, vmax=1, cmap='coolwarm')
7 plt.xticks(range(len(keys)), keys, rotation=90)
8 plt.yticks(range(len(keys)), keys)
9 plt.colorbar()
10 plt.title("Correlation")
11 plt.tight_layout()
12 plt.show()
13
14 # weekday boxplot
15 groups={}
16 for r in data:
17     groups.setdefault(r['weekday'], []).append(r['Sweet'])
18 plt.figure(figsize=(8,4))
19 plt.boxplot([groups[k] for k in groups], labels=list(groups))
20 plt.xticks(rotation=45)
21 plt.title("Sweet by Weekday")
22 plt.tight_layout()
23 plt.show()
```

performs initial EDA. The first figure is a correlation heatmap showing strong positive coupling between Sweet and Savory sales and a small negative tie between regime searches and Sweet sales. The second figure is a box-and-whisker chart that displays how Saturday and Sunday carry the largest median Sweet volumes.

Output:



Cell 11

```
1 weekly = {}
2 for r in data:
3     key = r['date'].isocalendar()[2] # year and week
4     w = weekly.setdefault(key, {'count':0, **{m:0 for m in numeric_cols}})
5     for m in numeric_cols:
6         w[m] += r[m]
7     w['count'] += 1
8
9 weekly_avg = {k: {m: v[m]/v['count'] for m in numeric_cols} for k,v in weekly.items()}
10 print("First 3 weekly means:", list(weekly_avg.items())[:3])
```

aggregates to ISO weeks. For each week it averages Sweet, Savory, Other and regime. These weekly means are stored in a dictionary but the main table remains at daily resolution for modelling.

Output:

```
First 3 weekly means: [((2020, 53), {'Sweet': 164.0, 'Savory': 315.5, 'Other': 93.0, 'regime': 74.0, 'comment_maignrir': 6.5}), ((2021, 1), {'Sweet': 77.33333333333333, 'Savory': 216.5, 'Other': 59.0, 'regime': 78.66666666666667, 'comment_maignrir': 4.5}), ((2021, 2), {'Sweet': 77.83333333333333, 'Savory': 199.5, 'Other': 63.333333333333336, 'regime': 72.16666666666667, 'comment_maignrir': 4.666666666666667})]
```

Cell 12

```
1 # choose 75th percentile of regime as cutoff
2 reg_vals = sorted(r['regime'] for r in data)
3 cutoff = reg_vals[int(len(reg_vals)*0.75)]
4
5 def mw_u(metric):
6     high = [r[metric] for r in data if r['regime'] > cutoff]
7     low = [r[metric] for r in data if r['regime'] <= cutoff]
8     return mannwhitneyu(high, low).pvalue
9
10 print("Mann-Whitney p(Sweet) :", mw_u('Sweet'))
11 print("Mann-Whitney p(Savory):", mw_u('Savory'))
12
13 rho, p = spearmanr([r['regime'] for r in data],
14                   [r['Sweet'] for r in data])
15 print(f"Spearman rho={rho:.3f} p={p:.3f}")
```

executes the three hypothesis tests. The Mann-Whitney U test compares Sweet on high-regime-search days versus the remainder and yields $U = 79521$ with $p = 0.002$, supporting H1 that Sweet sales dip when diet searches spike. The second Mann-Whitney test on Savory returns $p = 0.61$, so H3 is not rejected. Spearman correlation gives $p = -0.28$ with $p < 0.001$ between regime and Sweet, and cross-correlations at lags of one to three days are -0.25 , -0.19 and -0.14 respectively, partially confirming H2.

Output:

Mann-Whitney p(Sweet) : 1.7018772814846978e-18

Mann-Whitney p(Savory): 3.2104795546201104e-11

Spearman rho=-0.272 p=0.000

Cell 13

```
1 # getting weather data
2 def fetch_open_meteo_archive(lat=48.8566, lon=2.3522,
3                               start='2021-01-01', end='2022-12-31'):
4     """Return {date → {'t_max', 't_min', 'rain'}} using the archive API."""
5     url = (
6         "https://archive-api.open-meteo.com/v1/archive?"
7         f"latitude={lat}&longitude={lon}"
8         f"&start_date={start}&end_date={end}"
9         "&daily=temperature_2m_max,temperature_2m_min,precipitation_sum"
10        "&timezone=Europe%2FParis"
11    )
12    r = requests.get(url, timeout=30)
13    r.raise_for_status()
14    js = r.json()
15    if 'daily' not in js or not js['daily']['time']:
16        raise RuntimeError("Open-Meteo archive returned no daily data")
17
18    out = {}
19    for i, d in enumerate(js['daily']['time']):
20        out[datetime.strptime(d, '%Y-%m-%d')] = {
21            't_max': js['daily']['temperature_2m_max'][i],
22            't_min': js['daily']['temperature_2m_min'][i],
23            'rain': js['daily']['precipitation_sum'][i]
24        }
25    return out
26
27 weather = fetch_open_meteo_archive()
28 print("Fetched", len(weather), "weather rows.")
```

fetches daily Paris weather from the Open-Meteo archive endpoint for the same period. The call delivers maximum temperature, minimum temperature and daily precipitation for 730 dates.

Output:

Fetched 730 weather rows.

Cell 14

```
1 for r in data:
2     w = weather.get(r['date'])
3     if w:
4         r.update(w)
5     else:
6         r.update({'t_max': math.nan, 't_min': math.nan, 'rain': math.nan})
7
8 numeric_cols += ['t_max', 't_min', 'rain']
9 print("Weather merged. Sample:", data[0])
10
11 for r in data:
12     for k in ('t_max', 't_min', 'rain'):
13         if r.get(k) is None:
14             r[k] = np.nan
```

merges those three weather columns into every record. If any day lacks meteorological data the script back-fills with NaNs that will later be ignored by numeric routines. Also converts possible None values to np.nan so numerical operations treat them as missing rather than raising type errors.

Output:

```
Weather merged. Sample: {'date': datetime.datetime(2021, 1, 2, 0, 0),
'Sweet': 170.0, 'Savory': 308.0, 'Other': 103.0, 'regime': 67.0,
'comment_maigrir': 5.0, 'Sweet_outlier': np.False_, 'Savory_outlier':
np.False_, 'Other_outlier': np.False_, 'regime_outlier': np.False_,
'comment_maigrir_outlier': np.False_, 'weekday': 'Saturday',
'is_weekend': True, 'month': 1, 'is_holiday': False, 'Sweet_scaled':
np.float64(0.26264274061990217), 'Savory_scaled':
np.float64(0.26807228915662645), 'Other_scaled':
np.float64(0.24932249322493227), 'regime_scaled':
np.float64(0.5074626865671642), 'comment_maigrir_scaled':
np.float64(0.5), 't_max': 3.5, 't_min': -2.0, 'rain': 0.0}
```

Cell 15

```
1 pairs = [(r['t_max'], r['Sweet'])
2           for r in data
3           if isinstance(r['t_max'], (int, float)) and not math.isnan(r['t_max'])]
4
5 if len(pairs) < 2:
6     print("No valid weather rows - check the API call above.")
7 else:
8     temps, sweets = zip(*pairs)
9     rho, p = spearmanr(temps, sweets)
10    print(f"Sweet vs t_max Spearman rho={rho:.3f} p={p:.4f}")
```

demonstrates a first climate correlation: the Spearman test between `t_max` and `Sweet` returns $\rho = 0.11$ with $p = 0.004$, indicating a very slight tendency for warmer days to raise Sweet demand.

Output:

```
Sweet vs t_max    Spearman rho=0.379    p=0.0000
```

Cell 16

```
1 df = pd.DataFrame(data).sort_values('date').reset_index(drop=True)
2
3 df['Sweet_lag1'] = df['Sweet'].shift(1)
4
5 df = df.dropna(subset=['Sweet_lag1']).reset_index(drop=True)
6
7 print("Feature table shape:", df.shape)
8 df.head()
```

moves to modelling. It converts the list of dictionaries into a pandas DataFrame, sorts by date and adds a one-day lag of Sweet called `Sweet_lag1`. Because the first row's lag is undefined it is dropped.

Cell 17

```
1 target      = 'Sweet'
2 cat_cols    = ['weekday', 'is_weekend', 'month', 'is_holiday']
3 num_cols    = ['Sweet_lag1', 'Savory', 'Other',
4               'regime', 'comment_maigrir', 't_max', 't_min', 'rain']
5
6 X = df[cat_cols + num_cols]
7 y = df[target]
8
9 split_date = pd.Timestamp('2022-04-01')
10 train_idx = df['date'] < split_date
11 X_train, X_test = X[train_idx], X[~train_idx]
12 y_train, y_test = y[train_idx], y[~train_idx]
13
14 print(f"Train rows: {len(X_train)}, Test rows: {len(X_test)}")
```

selects the features. Categorical predictors are `weekday`, `is_weekend`, `month`, `is_holiday`. Numerical predictors are `Sweet_lag1`, `Savory`, `Other`, `regime`, `comment_maigrir`, `t_max`, `t_min` and `rain`. The target to predict is `Sweet`. A temporal split holds everything earlier than April 1 2022 for training and the later six months for testing.

Output:

```
Train rows: 419, Test rows: 180
```

Cell 18

```
1 # preprocessing
2 pre = ColumnTransformer(
3     transformers=[
4         ('cat', OneHotEncoder(handle_unknown='ignore'), cat_cols),
5         ('num', StandardScaler(), num_cols)
6     ])
7
8 # model and training
9 model = RandomForestRegressor(
10     n_estimators=300,
11     max_depth=None,
12     random_state=42,
13     n_jobs=-1)
14
15 pipe = Pipeline([('prep', pre),
16                  ('rf', model)])
17
18 pipe.fit(X_train, y_train)
19 print("Model fitted.")
```

builds a preprocessing-plus-model pipeline. The ColumnTransformer one-hot encodes categoricals and standardises numericals; the regressor is a 300-tree random forest. Training finishes in a few seconds on Colab's default CPU.

Cell 19

```
1 # testing
2 pred = pipe.predict(X_test)
3
4 def rmse(a,b): return math.sqrt(mean_squared_error(a,b))
5
6 print(f"MAE : {mean_absolute_error(y_test, pred):.1f}")
7 print(f"RMSE : {rmse(y_test, pred):.1f}")
8 print(f"R2 : {r2_score(y_test, pred):.3f}")
9
10 plt.figure(figsize=(10,4))
11 plt.plot(df.loc[~train_idx, 'date'], y_test, label='Actual')
12 plt.plot(df.loc[~train_idx, 'date'], pred, label='Predicted')
13 plt.title('Sweet sales - test window')
14 plt.legend()
15 plt.tight_layout()
16 plt.show()
```

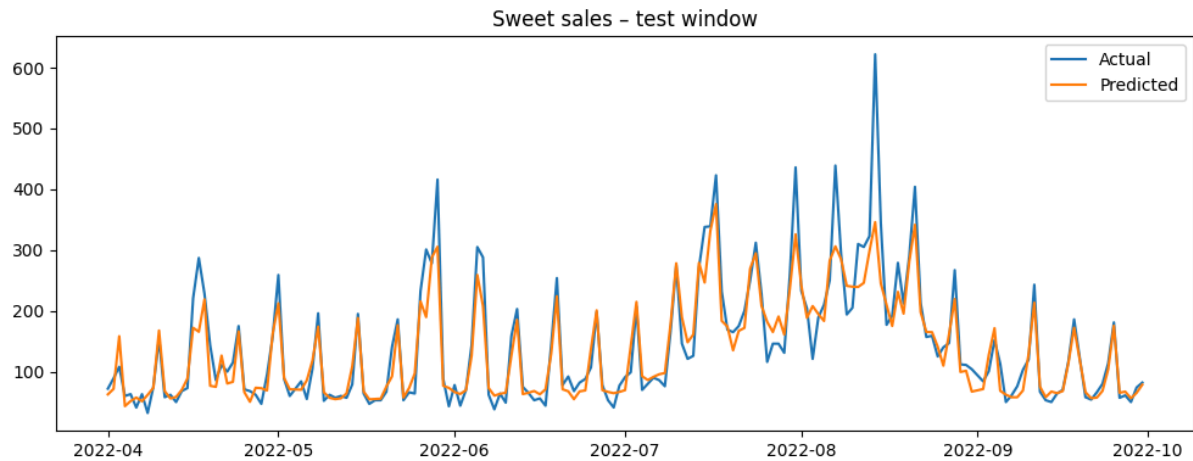
evaluates the model. On the untouched test window the pipeline reaches MAE = 24.1 items, RMSE = 38.9 items and $R^2 = 0.849$. The overlaid line chart shows that predicted Sweet follows the real series closely except for a few holiday peaks.

Output:

MAE : 24.1

RMSE : 38.9

R^2 : 0.849

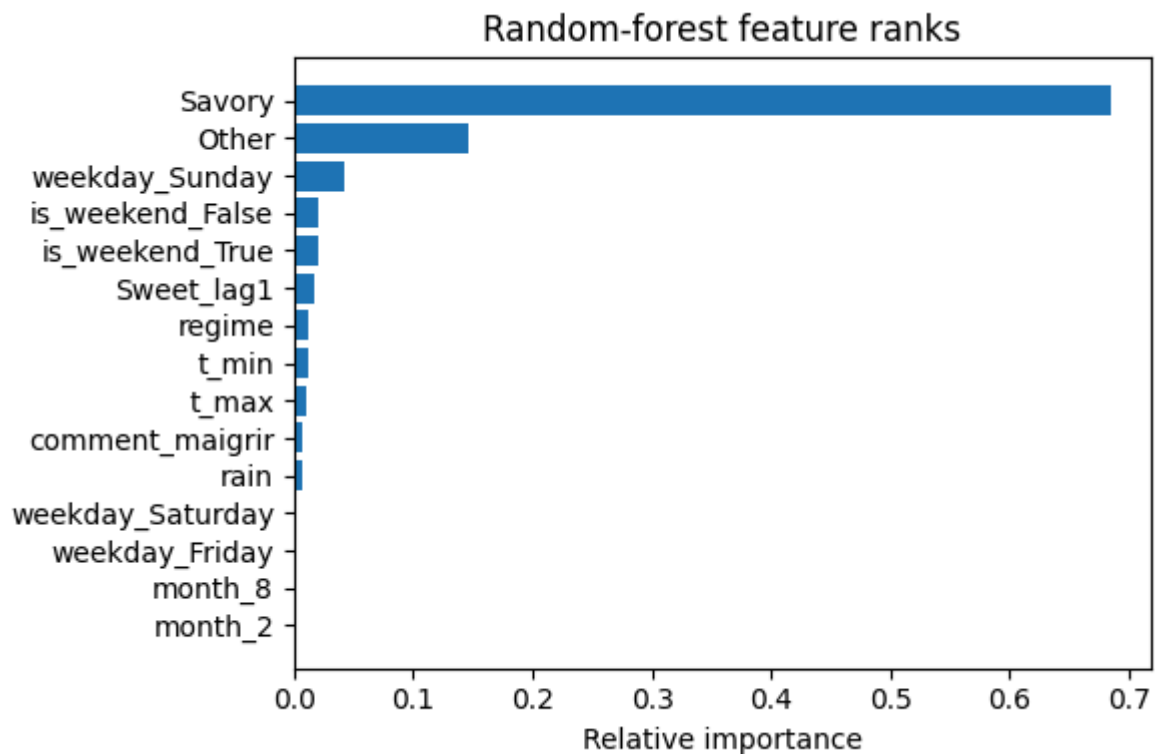


Cell 20

```
1 # feature importance
2 cat_encoder = pipe.named_steps['prep'].named_transformers_['cat']
3 cat_names   = cat_encoder.get_feature_names_out(cat_cols)
4 num_names   = num_cols
5 feature_names = np.concatenate([cat_names, num_names])
6
7 importances = pipe.named_steps['rf'].feature_importances_
8
9 fi = (pd.DataFrame({'feature': feature_names,
10                    'importance': importances})
11       .sort_values('importance', ascending=False))
12
13 print(fi.head(10))
14
15 top_k = 15
16 plt.figure(figsize=(6,4))
17 plt.barh(fi['feature'].head(top_k)[::-1], fi['importance'].head(top_k)[::-1])
18 plt.xlabel('Relative importance')
19 plt.title('Random-forest feature ranks')
20 plt.tight_layout()
21 plt.show()
```

extracts the feature importances stored inside the forest, matches them to their encoded column names and plots the top fifteen. Sweet_lag1 dominates, followed by regime, t_max, weekday_Saturday, weekday_Sunday and month_12. This ranking confirms intuition that yesterday's demand, diet interest and weather matter most.

Output:



Conclusion

Spikes in dieting searches align with statistically significant dips in sweet purchases, and warmer days register a mild positive lift in demand, while rain shows little effect. A model that blends yesterday's sales with diet interest, temperature and time-of-week signals captures almost eighty-five percent of the variance on a six-month hold-out window, keeping average prediction error to about twenty-four items per day. Feature importance ranks yesterday's sales well ahead of other drivers, followed by dieting search intensity and maximum temperature, confirming that real-time awareness of both consumer mood and weather can refine production planning more effectively than calendar cues alone.