

```

14
15 def basic_clean(df):
16     # 1. normalize column names
17     df.columns = [c.strip() for c in df.columns]
18
19     # 2. ensure numeric types for flux and error columns
20     for c in FLUX_COLS + ERR_COLS:
21         if c in df.columns:
22             df[c] = pd.to_numeric(df[c], errors='coerce')
23
24     # 3. replace negative/zero flux values with NaN (physical flux >= 0)
25     for c in FLUX_COLS:
26         if c in df.columns:
27             df.loc[df[c] <= 0, c] = np.nan
28
29     # 4. replace non-positive errors with NaN
30     for c in ERR_COLS:
31         if c in df.columns:
32             df.loc[df[c] <= 0, c] = np.nan
33
34     # 5. drop rows with all fluxes missing (useless rows)
35     if set(FLUX_COLS).issubset(df.columns):
36         df = df.dropna(subset=FLUX_COLS, how='all').reset_index(drop=True)
37
38     # 6. drop exact duplicate rows
39     df = df.drop_duplicates().reset_index(drop=True)
40
41     return df
42
43 def derive_features(df):
44     # a) total_flux (sum of available fluxes)
45     available_flux = [c for c in FLUX_COLS if c in df.columns]
46     df['total_flux'] = df[available_flux].sum(axis=1, skipna=True)
47
48     # b) mean_flux_err
49     available_err = [c for c in ERR_COLS if c in df.columns]
50     if available_err:
51         df['mean_flux_err'] = df[available_err].mean(axis=1, skipna=True)
52     else:
53         df['mean_flux_err'] = np.nan
54
55     # c) flux_count (# of bands with detection)
56     df['flux_count'] = df[available_flux].notna().sum(axis=1)
57
58     # d) brightness class (quantiles)
59     # add 1e-12 to avoid zero problems
60     df['brightness_class'] = pd.qcut(df['total_flux'].fillna(0)+1e-12, q=3, labels=['Low', 'Med',
61
62     # e) spectral_index_proxy: log(freq1/freq9) / log(freq1/freq9)
63     # approximate freq1=30 GHz, freq9=857 GHz for Planck typical bands
64     if available_flux:
65         f1, f9 = available_flux[0], available_flux[-1]
66         df['spectral_index_proxy'] = np.log(df[f1].replace(0, np.nan) / df[f9].replace(0, np.n
67
68     return df

```

```

42
43     def derive_features(df):
44         # a) total_flux (sum of available fluxes)
45         available_flux = [c for c in FLUX_COLS if c in df.columns]
46         df['total_flux'] = df[available_flux].sum(axis=1, skipna=True)
47
48         # b) mean_flux_err
49         available_err = [c for c in ERR_COLS if c in df.columns]
50         if available_err:
51             df['mean_flux_err'] = df[available_err].mean(axis=1, skipna=True)
52         else:
53             df['mean_flux_err'] = np.nan
54
55         # c) flux_count (# of bands with detection)
56         df['flux_count'] = df[available_flux].notna().sum(axis=1)
57
58         # d) brightness class (quantiles)
59         # add 1e-12 to avoid zero problems
60         df['brightness_class'] = pd.qcut(df['total_flux'].fillna(0)+1e-12, q=3, labels=['Low','Medium','High'])
61
62         # e) spectral_index_proxy: log(freq1/freq9) / log(freq1/freq9)
63         # approximate freq1=30 GHz, freq9=857 GHz for Planck typical bands
64         if available_flux:
65             f1, f9 = available_flux[0], available_flux[-1]
66             df['spectral_index_proxy'] = np.log(df[f1].replace(0, np.nan) / df[f9].replace(0, np.nan)) / np.log(30/857)
67
68     return df
69
70
71     def normalize_fluxes(df):
72         # create normalized (0-1) versions for each flux band
73         for c in FLUX_COLS:
74             if c in df.columns:
75                 col = df[c]
76                 minv, maxv = col.min(skipna=True), col.max(skipna=True)
77                 if pd.notna(minv) and pd.notna(maxv) and maxv > minv:
78                     df[f'{c}_norm'] = (col - minv) / (maxv - minv)
79                 else:
80                     df[f'{c}_norm'] = np.nan
81
82         # create z-score columns as well
83         present = [c for c in FLUX_COLS if c in df.columns]
84         if present:
85             scaler = StandardScaler()
86             df_z = df[present].fillna(0)
87             zvals = scaler.fit_transform(df_z)
88             for i, c in enumerate(present):
89                 df[f'{c}_z'] = zvals[:, i]
90
91     return df
92
93     def filter_reliable(df, rel_err_threshold=0.5):
94         # mark reliable if mean_flux_err <= rel_err_threshold * total_flux
95         df['reliable'] = True
96         mask = (df['mean_flux_err'] > rel_err_threshold * df['total_flux'])
97         df.loc[mask, 'reliable'] = False
98         return df
99
100    def iqr_filter(df, col='total_flux', k=1.5):
101        q1 = df[col].quantile(0.25)
102        q3 = df[col].quantile(0.75)
103        iqr = q3 - q1
104        lower, upper = q1 - k*iqr, q3 + k*iqr
105        return df[(df[col] >= lower) & (df[col] <= upper)].reset_index(drop=True)
106
107    def full_pipeline(path="planck.csv"):
108        raw = load_raw(path)
109        cleaned = basic_clean(raw.copy())
110        derived = derive_features(cleaned)
111        normed = normalize_fluxes(derived)
112        normed = filter_reliable(normed)
113        return raw, normed
114

```