

CMSN process

The provided code implements a mathematical analysis of trajectories in the Collatz-Matthews process. Below is its mathematical abstraction:

1. Collatz-Matthews Process

For a starting integer n , iteratively apply the transformation:

$$a_{k+1} = \begin{cases} \frac{a_k}{2} & \text{if } a_k \text{ is even} \\ 3a_k + 1 & \text{if } a_k \text{ is odd} \end{cases}$$

until $a_k = 1$. During this process, track the following metrics for n :

- $b_x(n)$: Total steps (even + odd).
- $b_z(n)$: Count of odd steps.
- $b_y(n)$: Cumulative sum of $\log_2(a_k)$ over odd steps.
- $G(n)$: Cumulative sum of $\log_2\left(3 + \frac{1}{a_k}\right)$ over odd steps.
- $\max_a(n)$: Maximum a_k encountered.

2. Metric Definitions

For a trajectory starting at n :

$$\begin{aligned} b_x(n) &= \text{Total iterations until } a_k = 1, \\ b_z(n) &= \sum_{k \text{ odd step}} 1, \\ b_y(n) &= \sum_{k \text{ odd step}} \log_2(a_k), \\ G(n) &= \sum_{k \text{ odd step}} \log_2\left(3 + \frac{1}{a_k}\right), \\ \max_a(n) &= \max\{a_1, a_2, \dots, a_{b_x(n)}\}. \end{aligned}$$

3. Derived Ratios and Conditions

For analysis, define:

$$\begin{aligned} \text{bz_bx_ratio}(n) &= \frac{b_z(n)}{b_x(n)}, \\ \text{by_bx_ratio}(n) &= \frac{b_y(n)}{b_x(n)}, \\ \text{bx_minus_bz}(n) &= b_x(n) - b_z(n), \\ \text{net_log_balance}(n) &= (b_x(n) - b_z(n)) - G(n) - \log_2(n). \end{aligned}$$

Key conditions tested:

$$\begin{aligned} b_y(n) &< b_x(n) - b_z(n), \\ b_x(n) - b_z(n) &> G(n). \end{aligned}$$

4. Statistical Analysis

For $n = 1, 2, \dots, N$, compute:

- Descriptive Statistics:** Mean, standard deviation, quantiles for $b_x, b_y, b_z, G, \text{max_a}$.
- Correlations:** Matrix of Pearson correlations between metrics.
- Linear Regression:** Fit $b_z(n)$ vs. $b_x(n)$:

$$b_z = \beta \cdot b_x + \epsilon, \quad \text{with slope } \beta \text{ and } R^2.$$

- Distributions:** Histograms of metrics and ratios (e.g., b_z/b_x , max_a in log scale).

5. Key Equations

- Growth Relationship:**

$$\text{net_log_balance}(n) \approx 0 \quad (\text{empirical observation}).$$

- Regression:**

$$b_z(n) \propto b_x(n) \quad (\text{slope } \beta \approx 0.388 \text{ observed}).$$

This abstraction captures the dynamical system, metric definitions, and statistical relationships analyzed in the code. The Collatz conjecture's unresolved nature implies these relationships remain empirical observations rather than proven theorems.

```

1
2 import sys
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt

```

```

6 from scipy import stats
7 import seaborn as sns
8 from multiprocessing import Pool
9
10 # CMSN Generation
11 def collatz_matthews_step(a, bx, by, bz, g, max_a):
12     if a % 2 == 0:
13         return a // 2, bx + 1, by, bz, g, max_a
14     new_a = 3 * a + 1
15     return new_a, bx + 1, by + np.log2(a), bz + 1, g + np.log2(3 + 1/a),
16     max(new_a, max_a)
17
18 def generate_cmsn(start):
19     a, bx, by, bz, g, max_a = start, 0, 0, 0, 0, start
20     while a != 1:
21         a, bx, by, bz, g, max_a = collatz_matthews_step(a, bx, by, bz, g,
22         max_a)
23     return (start, bx, by, bz, g, max_a)
24
25 def generate_chunk(start_end):
26     start, end = start_end
27     return [generate_cmsn(n) for n in range(start, end + 1)]
28
29 def generate_data(max_n, num_processes=16):
30     print(f"Generating CMSN data for n=1 to {max_n} with {num_processes}
31     processes...")
32     chunk_size = max_n // num_processes
33     ranges = [(i * chunk_size + 1, (i + 1) * chunk_size if i < num_processes -
34     1 else max_n)
35               for i in range(num_processes)]
36
37     with Pool(num_processes) as pool:
38         results = pool.map(generate_chunk, ranges)
39
40     flat_results = [item for sublist in results for item in sublist]
41     df = pd.DataFrame(flat_results, columns=['n', 'bx', 'by', 'bz', 'g',
42     'max_a'])
43     df.to_csv("cmsn_data_30M_with_max_a.csv.gz", compression='gzip',
44     index=False)
45     print(f"Data saved to 'cmsn_data_30M_with_max_a.csv.gz'")
46     return df
47
48 # Analysis Functions (unchanged from prior, but adjusted file names)
49 def analyze_data(df):
50     df['bz_bx_ratio'] = df['bz'] / df['bx']
51     df['by_bx_ratio'] = df['by'] / df['bx']
52     df['by_bz_ratio'] = df['by'] / df['bz'].replace(0, np.nan)
53     df['g_bz_ratio'] = df['g'] / df['bz'].replace(0, np.nan)
54     df['bx_minus_bz'] = df['bx'] - df['bz']
55     df['by_lt_bx_minus_bz'] = df['by'] < df['bx_minus_bz']
56     df['bx_minus_bz_gt_g'] = df['bx_minus_bz'] > df['g']
57     df['log_n'] = np.log2(df['n'])
58     df['net_log_balance'] = df['bx_minus_bz'] - df['g'] - df['log_n']

```

```

53
54     print("\nBasic Statistics:")
55     stats_summary = df[['bx', 'by', 'bz', 'g', 'max_a']].describe()
56     print(stats_summary.round(2))
57     stats_summary.to_csv("cmsn_stats_summary_30M.csv")
58
59     print("\nKey Conditions:")
60     print(f"Mean b_z / b_x: {df['bz_bx_ratio'].mean():.3f}, Max:
61 {df['bz_bx_ratio'].max():.3f}")
62     print(f"Sequences where b_y < b_x - b_z:
63 {df['by_lt_bx_minus_bz'].mean()*100:.2f}%")
64     print(f"Sequences where b_x - b_z > G:
65 {df['bx_minus_bz_gt_g'].mean()*100:.2f}%")
66     print(f"Mean net log balance (should ≈ 0 without b_y):
67 {df['net_log_balance'].mean():.3f}")
68
69     for col, label in [('bx', 'Steps'), ('by', 'Log Odd Sum'), ('bz', 'Odd
70 Count'),
71                        ('g', 'Growth G'), ('max_a', 'Max Odd Value')]:
72         top5 = df.nlargest(5, col)[['n', col]]
73         print(f"\nTop 5 by {label}:")
74         print(top5)
75         top5.to_csv(f"cmsn_top5_{col}_30M.csv", index=False)
76
77     plt.figure(figsize=(12, 12))
78     for i, (col, label) in enumerate([
79         ('bx', 'Steps (b_x)'), ('by', 'Log Odd Sum (b_y)'),
80         ('bz', 'Odd Count (b_z)'), ('g', 'Odd Growth (G)'), ('max_a', 'Max Odd
81 Value (log scale)'),
82     ], 1):
83         plt.subplot(5, 1, i)
84         sns.histplot(df[col], bins=50, kde=True, log_scale=(col == 'max_a',
85 True))
86         plt.title(f'Distribution of {label}')
87         plt.xlabel(label)
88     plt.tight_layout()
89     plt.savefig('cmsn_distributions_30M.png')
90     plt.close()
91
92     plt.figure(figsize=(10, 6))
93     sns.histplot(df['bz_bx_ratio'].dropna(), bins=50, kde=True)
94     plt.axvline(0.388, color='r', linestyle='--', label='Threshold 0.388')
95     plt.title('Distribution of b_z / b_x')
96     plt.xlabel('b_z / b_x')
97     plt.legend()
98     plt.savefig('bz_bx_distribution_30M.png')
99     plt.close()
100
101     sample = df.sample(min(10000, len(df)), random_state=42)
102     plt.figure(figsize=(10, 6))
103     plt.scatter(sample['bx_minus_bz'], sample['g'], s=1, alpha=0.5)
104     plt.plot([0, max(sample['bx_minus_bz'])], [0, max(sample['bx_minus_bz'])],
105 'r--', label='y = x')

```

```

98     plt.xlabel('b_x - b_z (Even Steps)')
99     plt.ylabel('G (Odd Growth)')
100    plt.title('b_x - b_z vs G')
101    plt.legend()
102    plt.savefig('bx_minus_bz_vs_g_30M.png')
103    plt.close()
104
105    corr_matrix = df[['bx', 'by', 'bz', 'g', 'max_a']].corr()
106    print("\nCorrelation Matrix:")
107    print(corr_matrix.round(3))
108    plt.figure(figsize=(8, 6))
109    sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', vmin=-1, vmax=1)
110    plt.savefig('cmsn_correlation_heatmap_30M.png')
111    plt.close()
112
113    slope, _, r_value, _, _ = stats.linregress(df['bx'], df['bz'])
114    print(f"\nLinear Regression (b_z vs b_x): slope={slope:.3f}, R²=
115    {r_value**2:.3f}")
116
117    def main(max_n):
118        df = generate_data(max_n)
119        analyze_data(df)
120        print(f"\nAnalysis complete. Data saved to
121        'cmsn_data_30M_with_max_a.csv.gz' and related files.")
122
123    if __name__ == "__main__":
124        if len(sys.argv) != 2:
125            print("Usage: python cmsn_full_30M.py <max_n>")
126            sys.exit(1)
127        try:
128            max_n = int(sys.argv[1])
129            if max_n < 1:
130                raise ValueError("max_n must be positive")
131        except ValueError as e:
132            print(f"Error: {e}")
133            sys.exit(1)
134    main(max_n)
135    ...

```