

</> Hide code

1 battles

```
1 def compute_mle_elo(df, SCALE=400, BASE=10, INIT_RATING=1000):
2     from sklearn.linear_model import LogisticRegression
3     models = pd.concat([df["model_a"], df["model_b"]]).unique()
4     models = pd.Series(np.arange(len(models)), index=models)
5
6     # duplicate battles
7     df = pd.concat([df, df], ignore_index=True)
8     p = len(models.index)
9     n = df.shape[0]
10
11     X = np.zeros([n, p])
12     X[np.arange(n), models[df["model_a"]]] = +math.log(BASE)
13     X[np.arange(n), models[df["model_b"]]] = -math.log(BASE)
14
15     # one A win => two A win
16     Y = np.zeros(n)
17     Y[df["winner"] == "model_a"] = 1.0
18
19     # one tie => one A win + one B win
20     # find tie + tie (both bad) index
21     tie_idx = (df["winner"] == "tie") | (df["winner"] == "tie
22 (bothbad)")
23     tie_idx[len(tie_idx)//2:] = False
24     Y[tie_idx] = 1.0
25
26     lr = LogisticRegression(fit_intercept=False, penalty=None, tol=1e-
27 8)
28     lr.fit(X, Y)
29
30     elo_scores = SCALE * lr.coef_[0] + INIT_RATING
31
32     # set anchor as mixtral = 1114
33     if "mixtral-8x7b-instruct-v0.1" in models.index:
34         elo_scores += 1114 - elo_scores[models["mixtral-8x7b-instruct-
35 v0.1"]]
```

```
1 def get_bootstrap_result(battles, func_compute_elo, num_round):
2     rows = []
3     for i in tqdm(range(num_round), desc="bootstrap"):
4         rows.append(func_compute_elo(battles.sample(frac=1.0,
5 replace=True)))
6     df = pd.DataFrame(rows)
7     return df[df.median().sort_values(ascending=False).index]
```

```
1 def compute_online_elo(battles, K=4, SCALE=400, BASE=10,
2 INIT_RATING=1000):
```

```

2     rating = defaultdict(lambda: INIT_RATING)
3
4     for rd, model_a, model_b, winner in battles[['model_a',
'model_b', 'winner']].itertuples():
5         ra = rating[model_a]
6         rb = rating[model_b]
7         ea = 1 / (1 + BASE ** ((rb - ra) / SCALE))
8         eb = 1 / (1 + BASE ** ((ra - rb) / SCALE))
9         if winner == "model_a":
10             sa = 1
11         elif winner == "model_b":
12             sa = 0
13         elif winner == "tie" or winner == "tie (bothbad)":
14             sa = 0.5
15         else:
16             raise Exception(f"unexpected vote {winner}")
17         rating[model_a] += K * (sa - ea)
18         rating[model_b] += K * (1 - sa - eb)
19
20     # calibrate llama-13b to 800
21     delta = (800-rating["llama-13b"])
22     for model in battles["model_a"].unique():
23         rating[model] += delta
24
25     return rating

```

	Model	Elo rating
1	### Model A: claude-3-opus-20240229	1256.17
2	### Model B: claude-3-opus-20240229	1226.26
3	### Model B: claude-3-sonnet-20240229	1181.77
4	### Model B: gpt-4-turbo-2024-04-09	1176.08
5	claude-3-opus-20240229	1161.58
...
224	### Model B: gemma-2b-it	829.69
225	### Model B: gemma-7b-it	822.82
226	### Model A: gemma-2b-it	819.47
227	dolly-v2-12b	811.64
228	llama-13b	800.00

228 rows × 2 columns

```

1     def pretty_print_model_ratings(ratings):
2         df = pd.DataFrame([
3             [n, ratings[n]] for n in ratings.keys()
4             ], columns=["Model", "Elo rating"]).sort_values("Elo rating",
ascending=False).reset_index(drop=True)
5         # df["Elo rating"] = (df["Elo rating"] + 0.5).astype(int)
6         df.index = df.index + 1
7         return df
8

```

```

9
10 online_elo_ratings = compute_online_elo(battles)

```

	Model	Elo rating
1	### Model A: claude-3-opus-20240229	1363.86
2	### Model B: claude-3-opus-20240229	1360.53
3	### Model B: gpt-4-turbo-2024-04-09	1340.93
4	### Model A: gpt-4-turbo-2024-04-09	1339.17
5	### Model A: gpt-4-0125-assistants-api	1322.65
...
224	### Model B: gemma-1.1-2b-it	868.82
225	fastchat-t5-3b	866.37
226	stablalm-tuned-alpha-7b	831.28
227	dolly-v2-12b	809.97
228	llama-13b	792.67

228 rows × 2 columns

```

1 elo_mle_ratings = compute_mle_elo(battles)
2 pretty_print_model_ratings(elo_mle_ratings)

```

```

1 BOOTSTRAP_ROUNDS = 100
2
3 np.random.seed(42)
4 bootstrap_elo_lu = get_bootstrap_result(battles, compute_mle_elo,
    BOOTSTRAP_ROUNDS)

```

```

1 v def predict_win_rate(elo_ratings, SCALE=400, BASE=10,
    INIT_RATING=1000):
2     names = sorted(list(elo_ratings.keys()))
3     wins = defaultdict(lambda: defaultdict(lambda: 0))
4 v     for a in names:
5 v         for b in names:
6             ea = 1 / (1 + BASE ** ((elo_ratings[b] - elo_ratings[a])
    / SCALE))
7             wins[a][b] = ea
8             wins[b][a] = 1 - ea
9
10    data = {
11        a: [wins[a][b] if a != b else np.NAN for b in names]
12        for a in names
13    }
14
15    df = pd.DataFrame(data, index=names)
16    df.index.name = "model_a"
17    df.columns.name = "model_b"
18

```



```

4 v fig = px.imshow(win_rate.loc[ordered_models, ordered_models],
5                    color_continuous_scale='RdBu', text_auto=".2f",
6                    title="Predicted Win Rate Using Elo Ratings for Model
A in an A vs. B Battle")
7 v fig.update_layout(xaxis_title="Model B",
8                    yaxis_title="Model A",
9                    xaxis_side="top", height=900, width=900,
10                   title_y=0.07, title_x=0.5)
11 v fig.update_traces(hovtemplate=
12                    "Model A: %{y}<br>Model B: %{x}<br>Win Rate: %{z}
<extra></extra>")
13 fig

```

```

1  from collections import defaultdict
2  import json, math, gdown
3  import numpy as np
4  import pandas as pd
5  import plotly.express as px
6  import sklearn
7  from tqdm import tqdm
8  import requests
9  pd.options.display.float_format = '{:.2f}'.format

```

```

1  url =
    "https://storage.googleapis.com/arena_external_data/public/clean_battle_20240419.json"
2  response = requests.get(url)
3
4 v with open('local_file_name.json', 'wb') as file:
5     file.write(response.content)
6
7  # load the JSON data from the local file
8 v with open('local_file_name.json', 'r') as file:
9     battles = pd.read_json(file).sort_values(ascending=True, by=
    ["tstamp"])

```

```

1  import marimo as mo

```

