

NBA MVP Prediction


HOML Project

Eloi Beurtheret, Raphael Leonardi, Baptiste Pras

Université Paris-Saclay

Jeudi 11 Décembre 2025

Introduction

 VOTING RESULTS						
2022-23 KIA NBA MOST VALUABLE PLAYER AWARD						
Player, Team	1 st Place Votes (10 Points)	2 nd Place Votes (7 Points)	3 rd Place Votes (5 Points)	4 th Place Votes (3 Points)	5 th Place Votes (1 Point)	Total Points
Joel Embiid, Philadelphia	73	25	2	0	0	915
Nikola Jokic, Denver	15	52	32	0	0	674
Giannis Antetokounmpo, Milwaukee	12	25	95	0	0	806
Jayson Tatum, Boston	0	0	1	89	8	280
Shai Gilgeous-Alexander, Oklahoma City	0	0	0	6	28	48
Donovan Mitchell, Cleveland	0	0	0	1	27	30
Dominantas Sabonis, Sacramento	0	0	0	1	24	27
Luka Dončić, Dallas	0	0	0	2	4	10
Stephen Curry, Golden State	0	0	0	1	2	5
Jimmy Butler, Miami	0	0	0	0	3	3
De'Aaron Fox, Sacramento	0	0	0	0	2	2
Jalen Brunson, New York	0	0	0	0	1	1
Ja Morant, Memphis	0	0	0	0	1	1

Per Game ▾ Upgraded ▾ Share & Export ▾ ☒ When table is sorted, hide non-qualifiers for rate stats [Geography](#) [Hide Partial Rows](#)

Rk	Player	Age	Team	Pos	G	GS	MP	FG	FGA	FG%	3P	3PA	3P%	2P	2PA	2P%	eFG%	FT	FTA	FT%	ORB	DRB	TRB	AST	STL	BLK	TOV	PF	Pts	Awards
1	Luka Dončić	26	JAL	PG	17	17	37.1	18.8	42.0	24	18.0	7.7	13.2	29.1	35.0	9.8	18.1	10.0	0.8	6.4	9.2	9.1	1.5	0.4	4.4	2.5	25.6			
2	Shai Gilgeous-Alexander	27	CAC	PG	23	23	33.3	10.8	19.4	55.6	2.2	5.3	44.3	8.6	14.4	59.5	45.3	9.0	19.2	46.0	0.5	4.3	4.7	6.4	1.4	0.7	17.2	28		
3	Travis Mays	29	PHI	PG	23	23	38.6	10.7	22.9	46.7	5.7	9.9	57.1	13.6	19.6	69.4	7.3	28.0	0.3	6.6	4.7	7.2	1.7	0.9	2.7	2.3	31.0			
4	Donovan Mitchell	29	IND	SG	23	23	34.3	10.6	21.1	58.2	4.0	18.4	28.0	6.6	10.7	61.5	28.7	5.3	6.3	83.6	1.8	3.0	4.8	5.5	1.4	0.3	3.4	24	35	
5	Nikola Jokic	30	DEN	C/PG	C	23	23	24.9	16.5	27.1	60.2	2.1	5.9	34.4	8.4	12.1	69.4	27.5	6.1	7.3	85.5	3.0	9.8	12.8	14.0	1.3	0.8	3.5	28	2
6	Jayson Tatum	29	BOS	SG	23	23	35.7	10.7	21.4	49.6	2.1	5.9	36.3	8.6	18.7	46.0	18.6	5.6	7.1	78.0	1.2	5.0	6.2	4.9	1.1	0.4	3.6	31	21	
7	Giannis Antetokounmpo	31	MIL	PF	17	17	28.1	11.1	17.3	63.9	0.4	1.4	43.0	10.9	48.7	22.6	6.2	8.8	73.3	6.8	10.1	6.1	8.8	5.9	3.3	2.8	28.9			
8	Anthony Edwards	24	MIN	SG	28	28	24.4	9.7	19.3	50.0	3.4	8.1	41.6	6.3	11.3	56.0	28.7	6.0	7.2	83.3	0.8	4.2	4.9	3.8	1.2	0.8	3.0	21	28	
9	Russell Westbrook	37	OKC	PG	28	28	36.8	8.6	16.9	50.9	2.9	7.7	37.0	5.8	9.2	63.5	28.3	0.4	8.8	27.4	0.8	4.7	5.5	6.7	1.1	0.2	3.4	22	28	
10	Jalen Brunson	29	NYK	SG	15	15	36.7	9.8	19.0	47.1	2.7	7.4	36.9	7.1	13.6	51.9	33.9	5.7	6.7	86.7	0.3	2.6	3.1	6.4	0.8	0.0	2.5	23	28.0	
11	Stephen Curry	37	SAS	PG/SG	PG	18	18	35.1	9.1	19.4	47.1	4.7	12.8	39.1	4.4	7.4	60.2	28.2	4.9	6.4	44.9	0.3	3.4	3.7	4.8	1.3	0.8	1.1	22	27.6
12	LeBron James	39	LAC	PF/SG	PF	22	22	35.1	9.2	19.5	47.2	3.8	8.3	46.1	6.2	11.2	55.5	54.0	6.2	6.8	90.7	2.2	4.4	6.5	2.8	8.8	0.4	12.7	27.6	
13	Cade Cunningham	24	DET	SG	PG	21	21	36.4	9.5	21.0	45.2	1.9	6.4	29.6	7.6	14.7	51.0	48.0	5.5	7.8	70.4	1.5	4.9	6.4	0.3	1.5	0.8	3.9	23	27.5
14	Jacques Nijig	36	MEM	PG	23	23	18.3	7.7	17.1	44.5	0.5	1.6	38.7	4.2	7.9	53.6	38.6	7.8	8.7	89.9	0.3	4.9	5.4	8.3	1.2	0.4	4.0	21	28.4	
15	Victor Oladipo	32	SAS	C/SG	C	12	12	34.7	9.3	18.4	50.2	1.7	4.8	34.6	7.6	13.6	56.0	7.6	7.9	76.7	2.8	10.9	13.9	4.9	1.1	0.6	3.6	27	28.2	
16	Denzel Austin	25	MEM	SG	PF	24	24	24.3	7.8	16.3	47.8	2.5	6.6	38.0	5.3	9.7	54.5	25.5	5.8	8.0	73.6	1.3	6.0	7.2	6.3	0.7	0.6	3.6	33	
17	Michael Porter Jr.	27	BKN	SG/PF	PF	19	19	32.8	9.3	18.7	49.6	3.5	8.6	40.2	5.7	9.7	59.2	29.0	3.8	4.8	79.8	1.2	6.4	7.6	3.2	8.9	0.3	2.2	22	18
18	Karl-Anthony Towns	34	MIN	PG	14	14	31.7	9.1	18.4	49.2	2.3	6.9	39.0	4.8	12.6	38.7	38.6	4.9	6.1	82.2	0.4	4.9	5.3	2.9	0.4	1.9	1.5	25.4		
19	Kevin Durant	37	SEA	SG/PF	PF	18	18	35.3	8.8	17.4	50.5	1.7	4.5	39.4	7.1	12.9	54.7	35.4	5.9	6.7	88.6	0.3	4.5	4.8	2.8	1.6	0.6	2.6	23	2
20	Dennis Rodman	29	PHI	SG/PF	SG	22	22	34.7	8.1	17.9	45.1	1.8	5.6	31.5	6.4	12.2	52.2	28.6	7.0	8.8	79.4	1.9	3.3	4.3	6.7	8.9	0.4	3.7	23	25.0
21	Jamal Murray	29	DEN	SG/PF	PG	22	22	35.0	9.1	18.9	58.6	3.3	7.3	44.7	5.8	10.6	54.7	29.7	3.6	4.8	75.0	0.5	4.0	4.5	6.8	1.1	0.3	2.4	18	25.0
22	Nikola Vucevic	32	SEA	PG	19	19	38.5	8.2	16.2	50.6	3.9	8.7	44.8	3.2	8.4	34.7	38.8	5.3	6.3	86.2	0.4	3.2	3.6	2.4	1.1	0.3	1.9	23	24.8	
23	Dariusz Sklar	31	PHI	SG	PF	13	13	34.2	8.0	18.4	48.1	1.8	4.9	36.6	7.2	13.7	52.2	33.9	4.8	7.0	68.6	1.4	5.4	7.0	4.1	1.3	0.5	2.2	28	24.5
24	De'Aaron Fox	28	IND	SG	PG	15	15	32.7	8.4	17.3	48.6	2.4	6.3	37.9	6.0	10.9	54.9	25.6	5.1	6.8	74.4	0.9	2.9	3.7	6.5	1.3	0.3	3.5	24	24.3
25	Tyler Herro	26	MIA	SG	SG	5	5	33.4	8.6	16.4	52.4	3.8	8.2	46.4	5.6	10.2	54.9	40.6	5.8	3.8	94.7	0.6	4.0	4.6	2.2	1.2	0.8	1.2	14	23.8

Scrapped player statistics and standings from <https://www.basketball-reference.com/>

```
print(f"[INFO] Downloading per-season raw data")
save_dir = os.path.join(raw_data_dir, save_dir)
os.makedirs(save_dir, exist_ok=True)

for year in range(start_year, end_year + 1):
    url = f"https://www.basketball-reference.com/leagues/NBA_{year}_per_game.html"
    print(f"[INFO] Downloading season {year} from {url}...")

    output_path = os.path.join(save_dir, f"NBA_{year}_per_game.csv")
    if os.path.exists(output_path):
        print(f"[SKIP] Skipping season {year}, file already exists.")
        continue

    try:
        response = requests.get(url, impersonate="chromell0")
        response.raise_for_status()

        # Parse the HTML table using pandas
        dfs = pd.read_html(response.text)
        df = dfs[0]

        # Remove header rows duplicated in the table
        df = df[df["Player"] != "Player"]
        df["Season"] = year

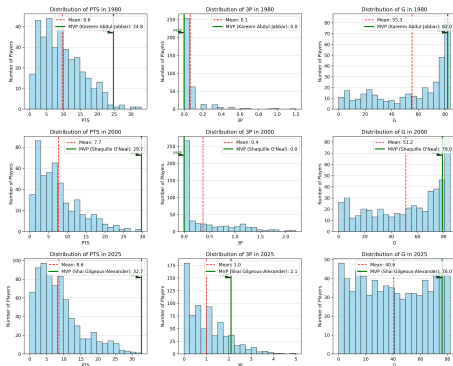
        # Save to CSV
        df.to_csv(output_path, index=False)
        print(f"[OK] Saved to {output_path}")

    except Exception as e:
        print(f"[WARN] Failed for season {year}: {e}, skipping")

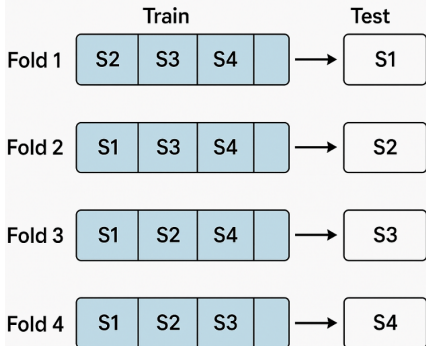
print(f"[DONE] Saved raw data to {save_dir}")
```

Preprocessing

- Encoded player's position into a one-hot-vector
- Transformed Win-Lose ratio into an integer
- Incorporated team standing
- Normalized statistics per season



Leave-One-Season-Out Strategy



```
# Process all years for this fold
for other_year in range(year_start, year_end + 1):
    other_year_str = str(other_year)
    year_folder = os.path.join(data_dir, other_year_str)

    # Check files exist
    file_paths = {file: os.path.join(year_folder, file) for file in files_to_process}
    if not all(os.path.isfile(path) for path in file_paths.values()):
        print(f"[WARN] Missing files for year {other_year}, skipping")
        continue

    # Load data
    df_X = pd.read_csv(file_paths["Data.csv"]).astype(np.float32)
    X = df_X.values
    y_top1 = pd.read_csv(file_paths["y_top1.csv"]).values.squeeze().astype(np.int64)
    y_top10 = pd.read_csv(file_paths["y_top10.csv"]).values.squeeze().astype(np.int64)
    df_names = pd.read_csv(file_paths["Name.csv"])

    if other_year == year:
        # Save test split
        np.savez_compressed(os.path.join(test_dir, "test.npz"),
                           X=X, y_top1=y_top1, y_top10=y_top10)
        df_names.to_csv(os.path.join(test_dir, "Name.csv"), index=False)
    else:
        # Accumulate for train split
        X_train_list.append(X)
        y_top1_train_list.append(y_top1)
        y_top10_train_list.append(y_top10)
        names_train.append(df_names)

# Save train split
if X_train_list:
    X_train = np.concatenate(X_train_list, axis=0)
    y_top1_train = np.concatenate(y_top1_train_list, axis=0)
    y_top10_train = np.concatenate(y_top10_train_list, axis=0)
    df_names_train = pd.concat(names_train, axis=0, ignore_index=True)

    np.savez_compressed(os.path.join(train_dir, "train.npz"),
                       X=X_train, y_top1=y_top1_train, y_top10=y_top10_train)
    df_names_train.to_csv(os.path.join(train_dir, "Name.csv"), index=False)
```

Feature Selection (1)

- Greedy Forward Selection
- Greedy Backward Selection
- Exhaustive Selection (unused because it's over 1 billion possibilities)

```
def greedy_forward_selection(model_class, dataset_dir, pipeline_name, fixed_params,
                             output_dir, year_start, year_end, patience=3):
    """
    Greedy forward selection using Recall@1 as the scoring metric with early stopping.
    """
    print(f"[INFO] Running greedy forward selection (Recall@1) on dataset: {dataset_dir}")

    data = load_dataset(dataset_dir, year_start)
    num_features = data["X_train"].shape[1]
    available_features = list(range(num_features))
    selected_features = []
    best_score = -1
    no_improve_count = 0
    recall_scores = []

    for i in range(1, num_features + 1):
        current_best_score = -1
        feature_to_add = None

        for feat in tqdm(available_features, desc=f"Testing additions (current={len(selected_features)})", leave=False):
            candidate = selected_features + [feat]
            score = compute_recall_at_1_avg(model_class, fixed_params, dataset_dir, year_start, year_end, candidate)
            if score > current_best_score or (score == current_best_score and feat < (feature_to_add or float("inf"))):
                current_best_score = score
                feature_to_add = feat

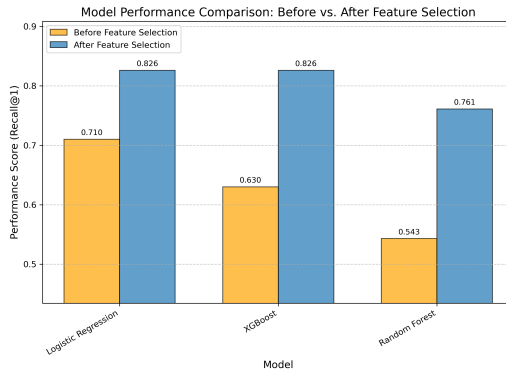
        if current_best_score > best_score:
            best_score = current_best_score
            best_combination = selected_features + [feature_to_add]
            no_improve_count = 0
        else:
            no_improve_count += 1

        recall_scores.append(current_best_score)
        selected_features.append(feature_to_add)
        available_features.remove(feature_to_add)

    print(f"[INFO] Added feature {feature_to_add} | Total: {len(selected_features)} | Recall@1: {current_best_score:.3f}")

    if no_improve_count >= patience:
        print(f"[EARLY STOP] No improvement for {patience} consecutive steps. Stopping.")
        break
```

Feature Selection (2)



- Logistic Regression: ["Team Overall", "G", "MP", "FGA", "FG%", "3P", "FT", "FT%", "ORB", "DRB", "BLK", "PF", "AST"]
- XGBoost: ["Team Overall", "PTS", "TOV", "GS", "MP", "2PA", "3PA", "FT"]
- Random Forest: ["2P", "Team Overall", "PTS", "G", "TOV", "FT%", "PF"]

We focused on three models :

- Logistic Regression ; simple model, originally intended as a baseline
- eXtreme Gradient Boosting ; chosen as our main gradient boosting model because it outperformed the other models in preliminary testing
- Random Forest ; as trees handle tabular data very well

For each player, we predict a probability of being the MVP. The one having the highest probability is then predicted as MVP (allows for a ranking)

Hyperparameters Tuning (1)

```
# Loop over seasons
global_results = []
for year in tqdm(range(year_start, year_end + 1), desc=f"{param_name}-{param_value}", file=sys.stdout):
    # Load dataset
    data = load_dataset(datasets_dir, year)
    X_train = data["X_train"]
    y_train = data["y_top1_train"]
    X_test = data["X_test"]
    y_test = data["y_top1_test"]
    y10_test = data["y_top10_test"]
    player_names_test = data["Name_test"]

    if selected_feature_names is not None:
        all_feature_names = get_feature_names(dataset_name, year)
        selected_indices = [i for i, name in enumerate(all_feature_names) if name in selected_feature_names]
        X_train = X_train[:, selected_indices]
        X_test = X_test[:, selected_indices]

    # Train model
    model = model_class(**model_init_params)
    model.fit(X_train, y_train)

    # Evaluate
    results = evaluate_model(model, X_test, y_test, y10_test, player_names_test, top_ks=top_ks)
    global_results.append((
        "year": year,
        **results
    ))

# Aggregate
df_results = pd.DataFrame(global_results)

# Summarize - same logic as summarize_global_results
summary_metrics = {}
```

Hyperparameters Tuning (2)

Pipeline: allStats_from1980

min_samples_leaf	Recall@1	Recall@3	Recall@5	Recall@10	Precision@1 (top-k)	Precision@3 (top-k)	Precision@5 (top-k)	Precision@10 (top-k)	Precision@1 (top-10)	Precision@3 (top-10)	Precision@5 (top-10)	Precision@10 (top-10)	Mean Abs Rank Error@1	Mean Abs Rank Error@3	Mean Abs Rank Error@5	Mean Abs Rank Error@10	True MVP avg rank	True MVP min rank	True MVP max rank	True MVP correct (rank=1)
1	0.500	0.870	0.913	0.978	0.500	0.652	0.604	0.611	1.000	0.928	0.835	0.611	1.630	2.217	2.822	2.891	2.261	1.000	12.000	0.500
2	0.543	0.826	0.935	1.000	0.543	0.630	0.622	0.644	1.000	0.928	0.830	0.644	1.065	2.246	2.900	2.980	2.109	1.000	8.000	0.543
5	0.478	0.848	0.913	1.000	0.478	0.652	0.626	0.663	1.000	0.942	0.839	0.663	1.326	2.058	2.817	3.024	2.217	1.000	8.000	0.478
10	0.500	0.848	0.957	1.000	0.500	0.630	0.630	0.677	1.000	0.928	0.852	0.677	1.457	2.283	2.848	2.941	2.174	1.000	8.000	0.500
20	0.478	0.761	0.913	1.000	0.478	0.616	0.639	0.690	1.000	0.928	0.870	0.690	1.413	2.457	2.761	2.893	2.435	1.000	9.000	0.478

- Logistic Regression: {"solver": "saga", "penalty": "l2", "class_weight": None, "max_iter": 50000, "C": 5.0}
- XGBoost: {"n_estimators": 100, "learning_rate": 0.1, "max_depth": 4, "scale_pos_weight": 5, "eval_metric": "logloss", "random_state": 42}
- Random Forest: {"n_estimators": 750, "max_depth": 15, "class_weight": None, "min_samples_leaf": 2, "random_state": 42}

Evaluation (1)

```
# Top-k
for k in top_ks:
    # Recall@k: whether the true MVP is present in the top-k predictions (1 = yes, 0 = no)
    top_k_true = sorted_true[:k]
    hit = 1 if top_k_true
    results["top_{k}_hit"] = hit
    vprint(f"[INFO] Top-{k} hit: {hit}")

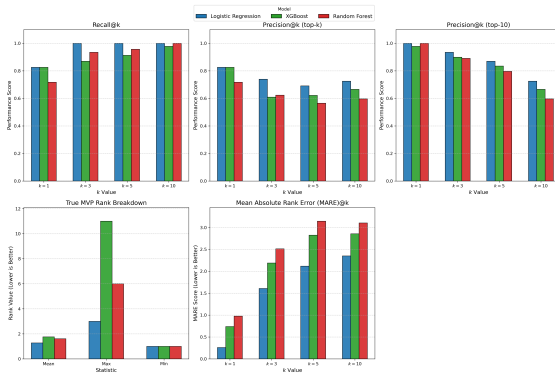
    pred_top_k_names = sorted_names[:k]
    pred_top_k_real_ranks = sorted_y10[:k]

    # Precision@k vs real top-k: proportion of top-k predicted players who are ranked in the real top-k (according to available top-k annotations)
    pred_top_k_real_ranks = sorted_y10[:k]
    real_topk_mask = (y_10 >= 1) & (y_10 <= k)
    nb_real_topk = real_topk_mask.sum()
    n_in_real_topk = sum(1 for r in pred_top_k_real_ranks if 1 <= r <= k)
    denom_topk = min(k, nb_real_topk) if nb_real_topk > 0 else 1
    precision_at_k_exact = n_in_real_topk / denom_topk
    results["precision_at_{k}_topk"] = precision_at_k_exact
    vprint(f"[METRIC] Precision@{k} vs real top-{k}: {precision_at_k_exact:.3f}")

    # Precision@k vs real top-10: proportion of top-k predicted players who are ranked in the real top-10 (according to available top-10 annotations)
    real_top10_mask = (y_10 >= 1) & (y_10 <= 10)
    nb_real_top10 = real_top10_mask.sum()
    n_in_real_top10 = sum(1 for r in pred_top_k_real_ranks if 1 <= r <= 10)
    denom_top10 = min(k, nb_real_top10) if nb_real_top10 > 0 else 1
    precision_at_k_top10 = n_in_real_top10 / denom_top10
    results["precision_at_{k}_top10"] = precision_at_k_top10
    vprint(f"[METRIC] Precision@{k} vs real top-10: {precision_at_k_top10:.3f}")

# Mean absolute rank errors: distance between predicted rank table and real rank table
abs_errors = []
for pred_rank_idx, real_rank in enumerate(pred_top_k_real_ranks, 1):
    if real_rank == -1:
        assumed_rank = 11 # If not in real top-10, treat as 11th
    else:
        assumed_rank = real_rank
    abs_error = abs(pred_rank_idx - assumed_rank)
    abs_errors.append(abs_error)
mean_abs_error = np.mean(abs_errors)
```

Evaluation (2)



Conclusion

Predicted MVP for year 2026 : Shai Gilgeous-Alexander (prob=0.8047)

Predicted top-10 for year 2026 :

- 1 Shai Gilgeous-Alexander (prob=0.8047)
 - 2 Nikola Jokić (prob=0.7473)
 - 3 Luka Dončić (prob=0.0355)
- 4 Cade Cunningham (prob=0.0344)
 - 5 Jalen Duren (prob=0.0037)
- 6 Isaiah Hartenstein (prob=0.0031)
 - 7 Austin Reaves (prob=0.0029)
 - 8 Jamal Murray (prob=0.0028)
 - 9 Ajay Mitchell (prob=0.0021)
- 10 Alperen Şengün (prob=0.0018)