

# APRENDIZAGEM COMPUTACIONAL

up202208296 - Lucas Greco do Espírito Santo Jorge

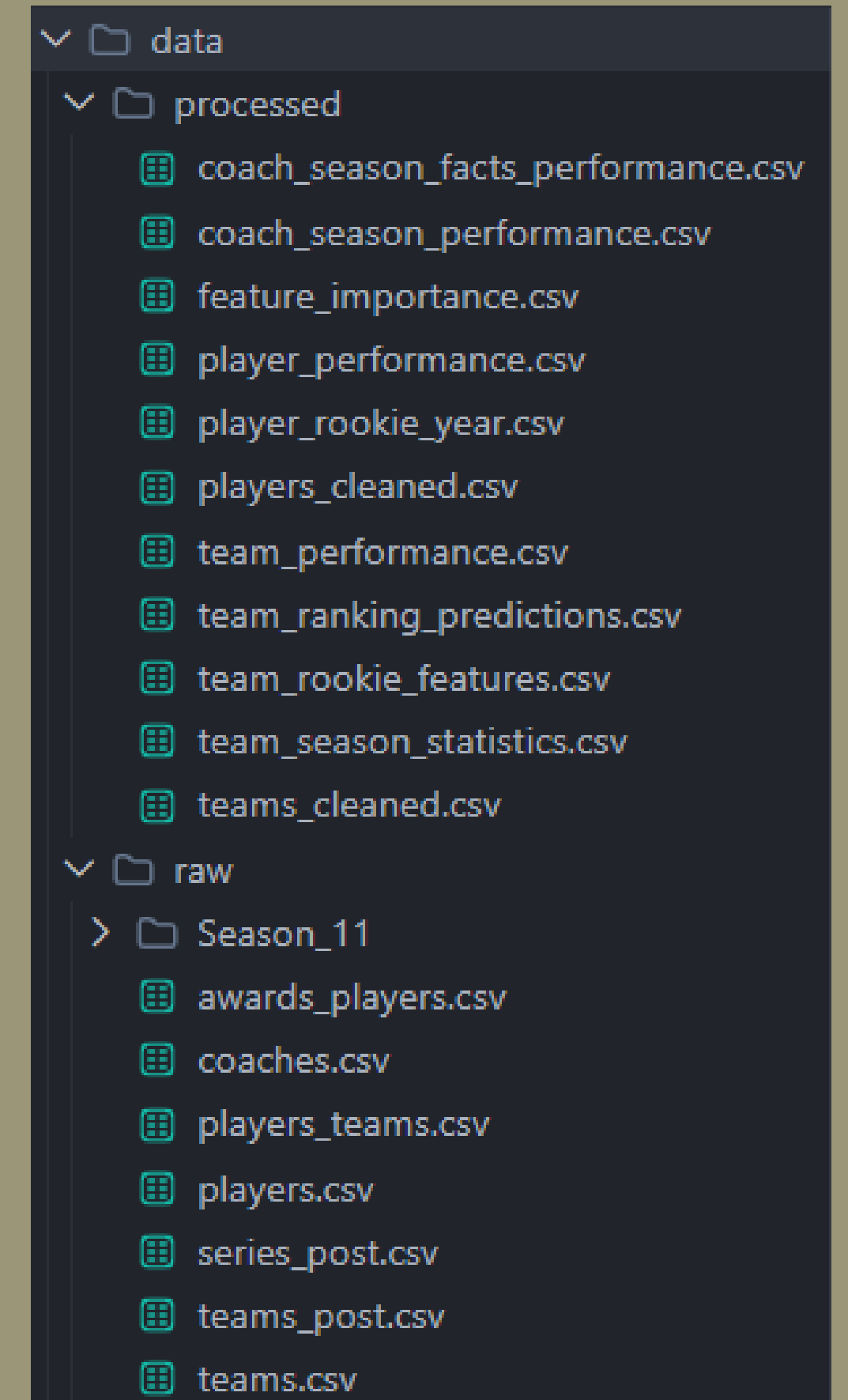
**Grupo 71**

# DESCRIÇÃO DE DOMÍNIO

PREVER O RANKEAMENTO FINAL DAS EQUIPES DE CADA CONFERÊNCIA DE UM ANO FUTURO.

VARIAVEL ALVO: "RANK"

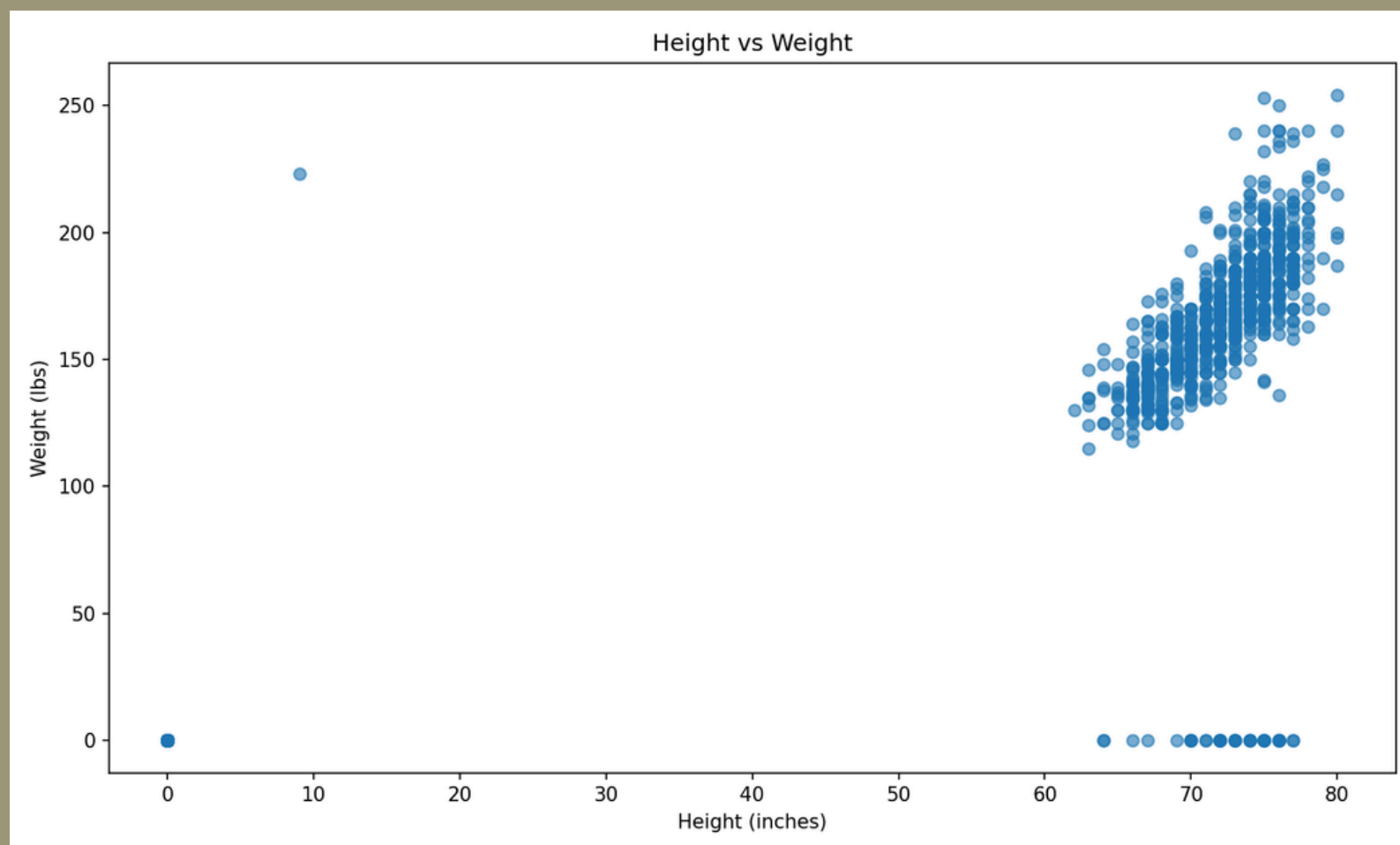
- DATASET:
  - RAW: CSVS QUE RECEBEMOS
    - INCLUSIVE A PASTA SEASON\_11, QUE CONTÉM OS DADOS DE TESTE, COMO O TIME, OS JOGADORES E OS TÉCNICOS.
  - PROCESSED: CSVS LIMPOS E PROCESSADOS



# EXPLORAÇÃO DE DADOS

- VÁRIOS JOGADORES APRESENTAM POSIÇÕES INDEFINIDAS.
- EXISTEM REGISTROS COM ALTURAS E PESOS INCOMPATÍVEIS.
- ALGUMAS COLUNAS DE TEAMS, COMO “DIV ID”, ESTÃO VAZIAS.

- HÁ COLUNAS COM VALORES CONSTANTES ENTRE TODOS OS TIMES.
- FORAM IDENTIFICADAS DUAS CONFERÊNCIAS DISTINTAS NO MESMO ANO.
- AS POSIÇÕES NO BASQUETE VARIAM AMPLAMENTE ENTRE ATLETAS.
- O TEMPO DE JOGO NÃO É UNIFORME ENTRE OS JOGADORES.
- ROOKIES NÃO POSSUEM HISTÓRICO PARA COMPARAÇÃO ESTATÍSTICA.
- TREINADORES PODEM TROCAR DE EQUIPE DURANTE A TEMPORADA.



- TABELA QUE RELACIONA A ALTURA E O PESO DAS JOGADORAS ANTES DA LIMPEZA

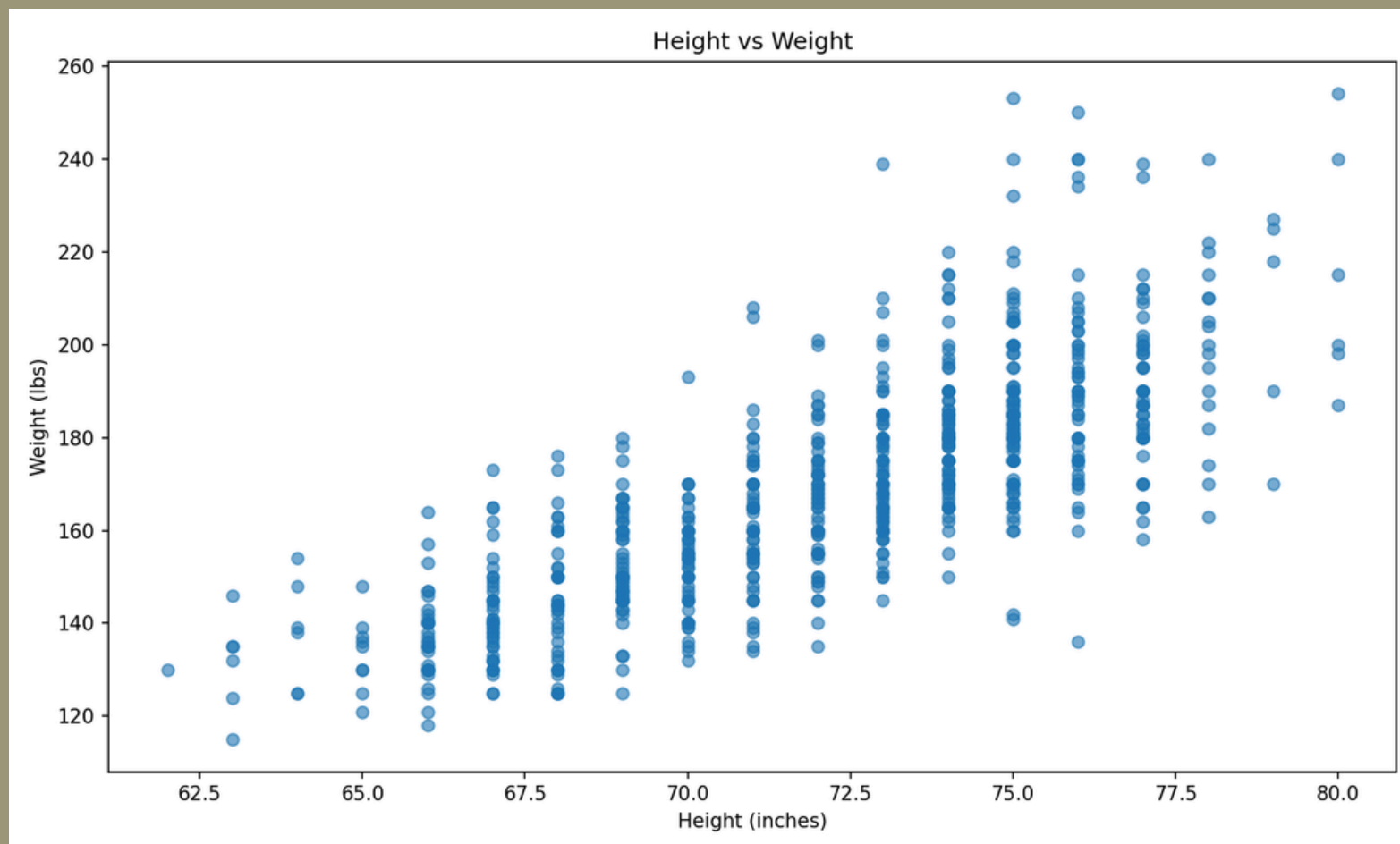
# EXPLORAÇÃO DE DADOS E LIMPEZA

- LIMPEZA DA TABELA PLAYERS

- REMOÇÃO DE REGISTROS DUPLICADOS.
- POSIÇÕES VAZIAS SÃO CLASSIFICADAS COMO UNKNOWN.
- APLICAÇÃO DE LIMITES MÍNIMOS: PESO  $\geq$  60 KG E ALTURA  $\geq$  24 CM.

- LIMPEZA DA TABELA TEAMS

- VALIDAÇÃO DA CONSISTÊNCIA:  $WON + LOST = GAMES\ PLAYED$ .
- CAMPOS CATEGÓRICOS VAZIOS SÃO PREENCHIDOS COM UNKNOWN (DIVID, CONFID, ARENA).
- ELIMINAÇÃO DE COLUNAS SEM UTILIDADE, ISTO É, COM VALORES IDÊNTICOS EM TODAS AS LINHAS.



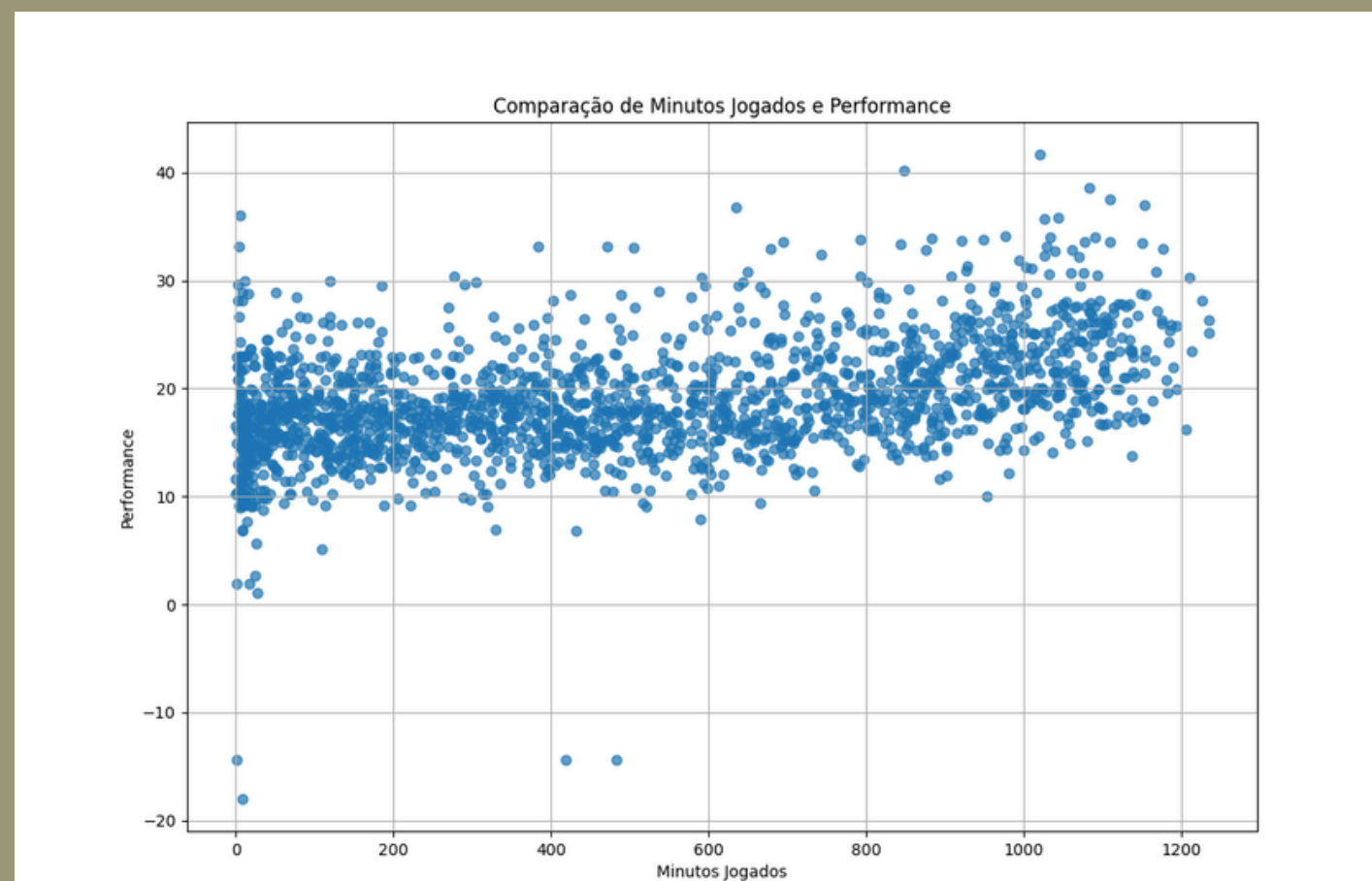
- TABELA QUE RELACIONA A ALTURA E O PESO DAS JOGADORAS APÓS A LIMPEZA

# PREVISÃO DO RANKING

- DEFINIÇÃO DO PROBLEMA
  - DESENVOLVER UM MODELO PREDITIVO PARA ESTIMAR O RANKING DAS EQUIPES EM SUAS CONFERÊNCIAS, UTILIZANDO DADOS HISTÓRICOS E PADRÕES DE DESEMPENHO.
- VARIÁVEL-ALVO PARA PREVISÃO: **"RANK"**- DEFINIDA EM TEAMS.CSV.
  - DEVIDO À ALTA VOLATILIDADE ENTRE TEMPORADAS, O MODELO UTILIZA COMO FEATURES:
  - TENDÊNCIA RECENTE DA EQUIPE, DESEMPENHO DA EQUIPE, DESEMPENHO DO TÉCNICO E RANKINGS ANTERIORES.
- PIPELINE DO MODELO
  - LIMPEZA DE DADOS → CLEAN\_PLAYERS.PY, CLEAN\_TEAMS.PY
  - PERFORMANCE DOS JOGADORES → PLAYER\_PERFORMANCE.PY
  - PERFORMANCE DAS EQUIPES → TEAM\_PERFORMANCE.PY
  - PERFORMANCE DOS TÉCNICOS → COACH\_SEASON\_FACTS\_PERFORMANCE.CSV
  - MODELO DE PREVISÃO DE RANKING → TEAM\_RANKING\_MODEL.PY
  - GERAÇÃO DE GRÁFICOS → MODEL\_GRAPHICS.PY

# PREVISÃO DO RANKING

- PREPARAÇÃO DOS DADOS
  - CÁLCULO DA PERFORMANCE DOS JOGADORES BASEADO APENAS EM ESTATÍSTICAS INDIVIDUAIS DA TEMPORADA, INDEPENDENTEMENTE DO DESEMPENHO DA EQUIPE.
  - A PERFORMANCE É ESTIMADA POR MÉTRICAS **PER-36**, COMBINADAS POR MÉDIA PONDERADA CONFORME A POSIÇÃO.
  - PARA JOGADORES COM TEMPO DE QUADRA ABAIXO DO MÍNIMO ESTABELECIDO, APLICA-SE UM SISTEMA DE FALLBACK HIERÁRQUICO PARA SUBSTITUIR SUA PERFORMANCE.



- TABELA QUE MOSTRA A RELAÇÃO DA PERFORMANCE COM OS MINUTOS JOGADOS, COMO PODEMOS VER, ESTÁ NORMALIZADA GRAÇAS AO PER 36, NÃO TEM UMA INCLINAÇÃO MUITO ACENTUADA.

# PREVISÃO DO RANKING

- PREPARAÇÃO DOS DADOS:

- TIME:

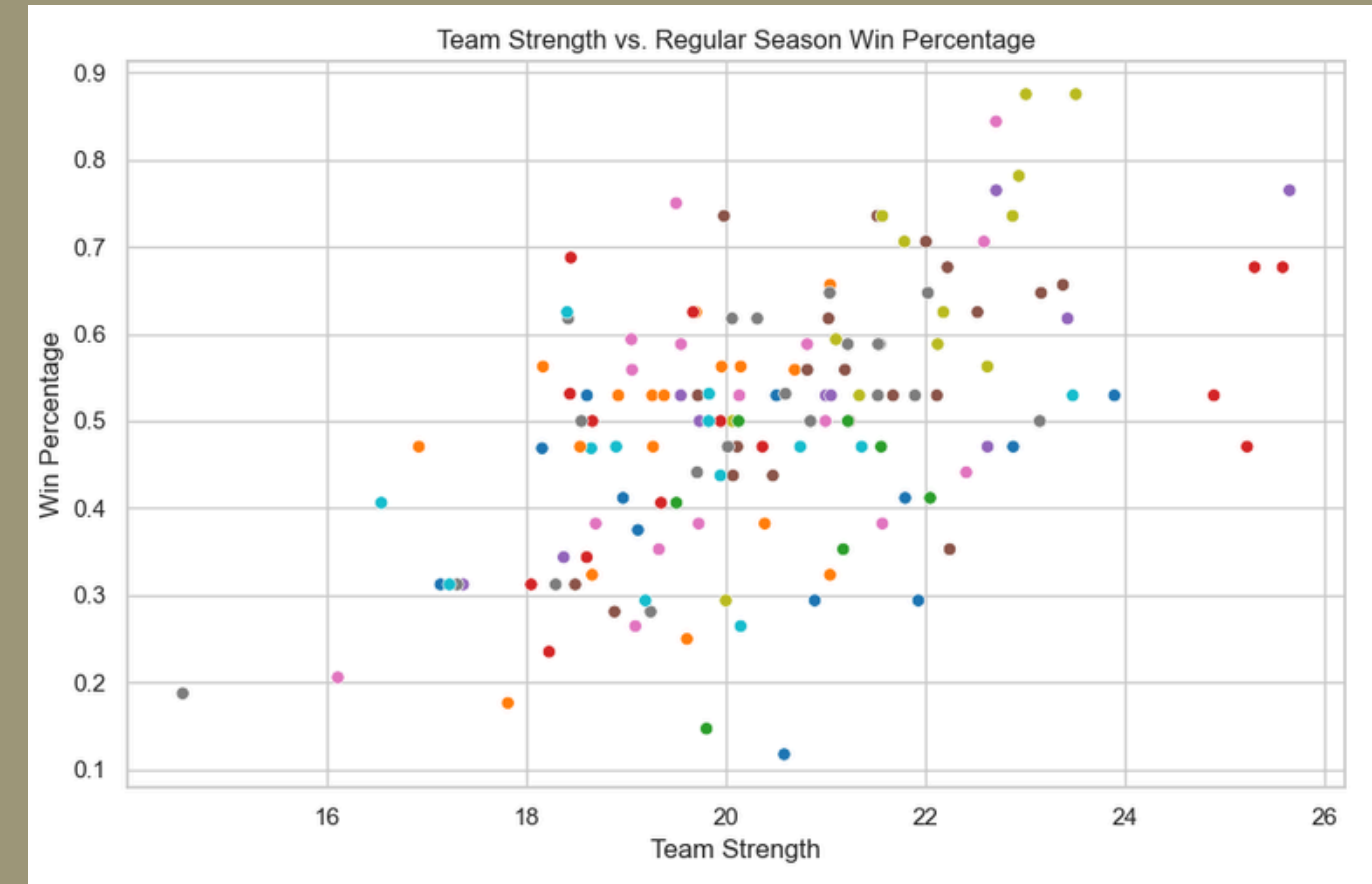
- **FORÇA DO TIME:** MÉDIA PONDERADA DA PERFORMANCE DOS JOGADORES PELO TEMPO JOGADO, ESTIMANDO A QUALIDADE GLOBAL DO ELENCO.
  - **ESTATÍSTICA PITAGÓRICA:**PREVISÃO DA TAXA DE VITÓRIAS A PARTIR DOS PONTOS MARCADOS E SOFRIDOS

- $(\text{PYTHAG\_WIN\_PCT} = \text{PF}^{\text{EXP}} / (\text{PF}^{\text{EXP}} + \text{PA}^{\text{EXP}}))$ .

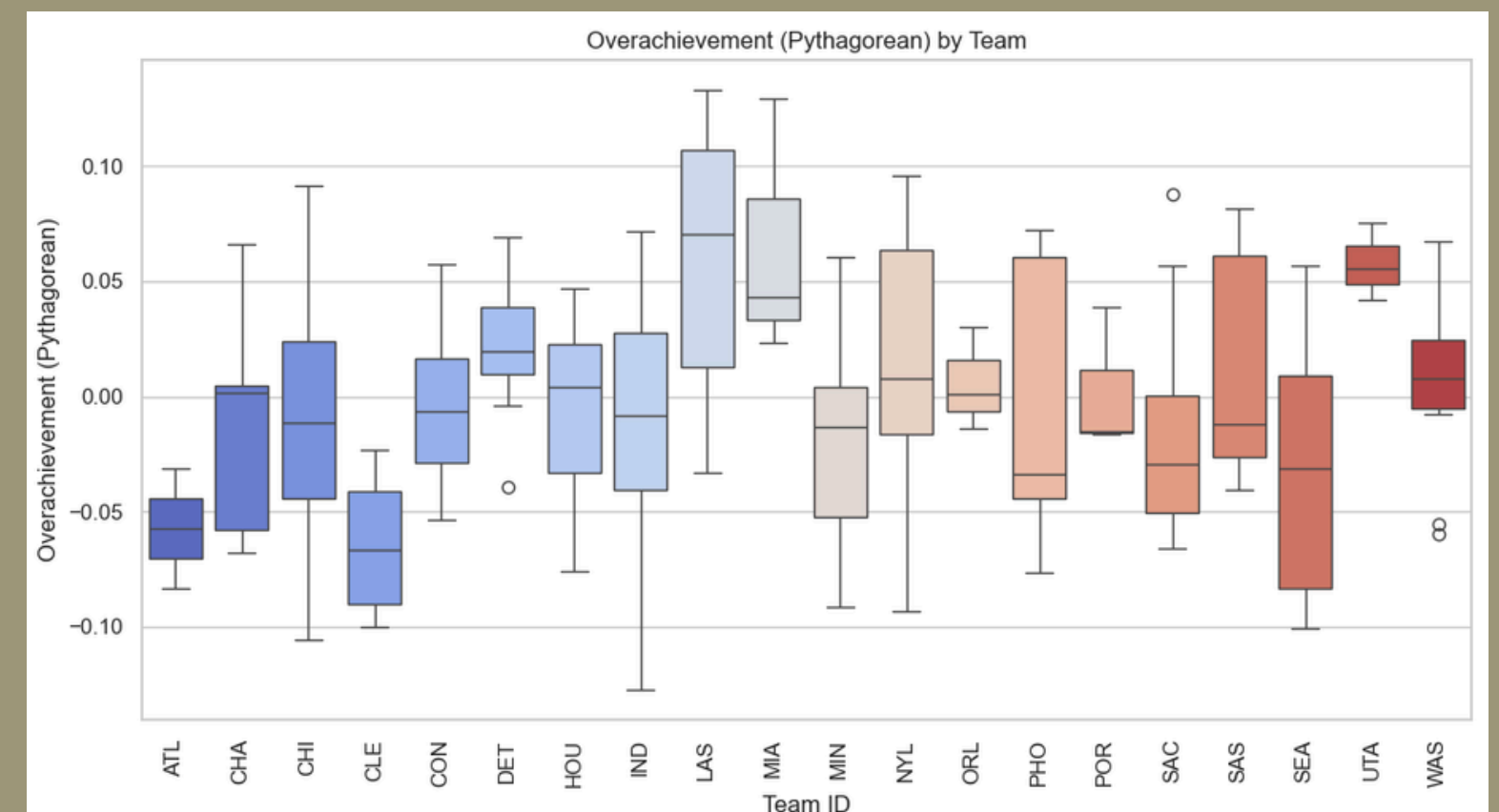
- PF: PONTOS A FAVOR

- PA: PONTOS CONTRA

- **REGRESSÃO LINEAR:** MODELO QUE RELACIONA A FORÇA DO ELENCO À EXPECTATIVA DE VITÓRIAS, ESTIMANDO QUANTAS VITÓRIAS O TIME “DEVERIA” ALCANÇAR.



- A RELAÇÃO ENTRE A "FORÇA DO ELENCO" CALCULADA E A % DE VITÓRIAS REAL NA TEMPORADA.

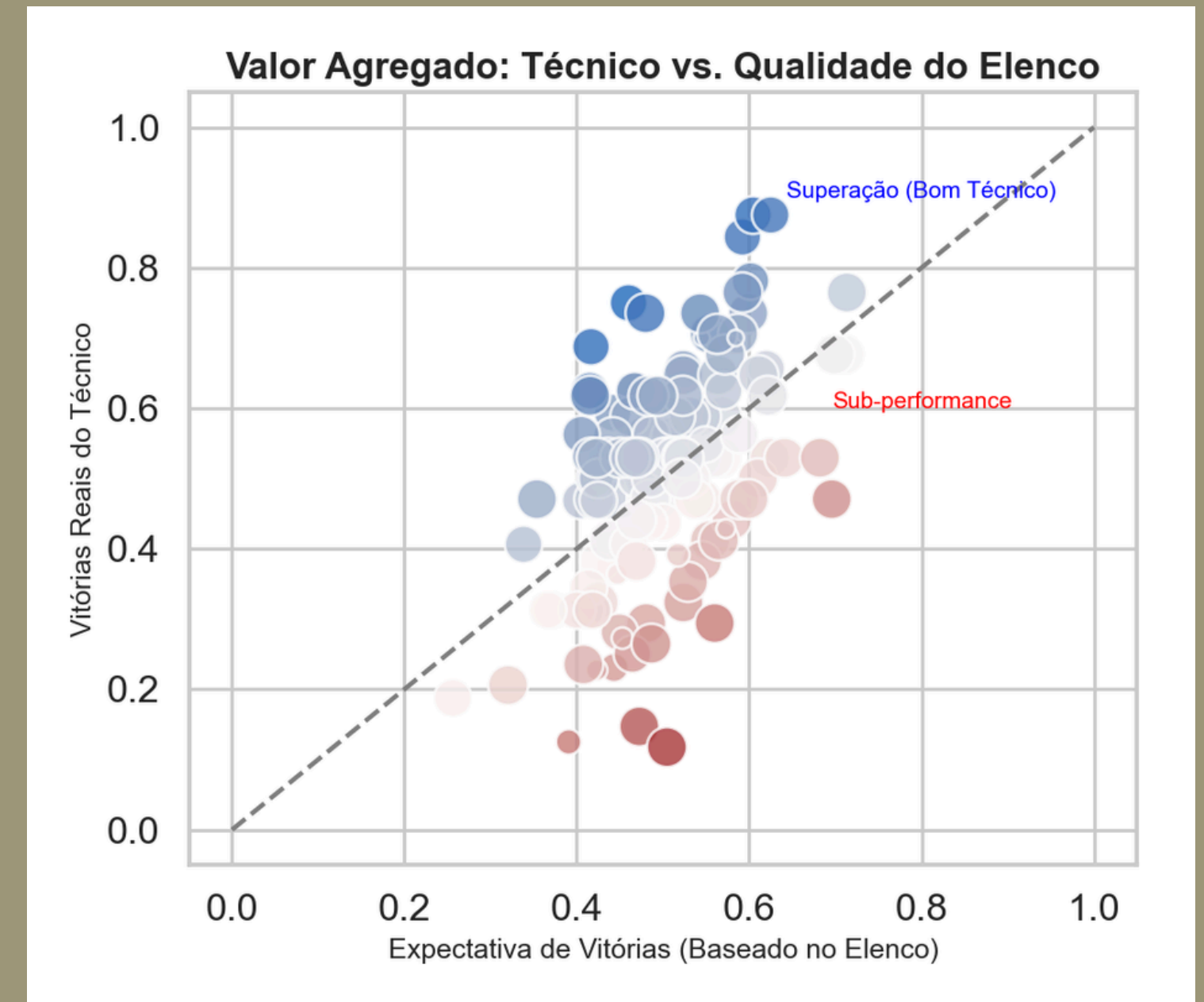


- A DIFERENÇA ENTRE QUANTAS VITÓRIAS O TIME TEVE E QUANTAS ELE DEVERIA TER TIDO BASEADO NO SALDO DE PONTOS



# PREVISÃO DO RANKING

- PREPARAÇÃO DOS DADOS:
  - TÉCNICO
    - CÁLCULO DE MÉTRICAS BÁSICAS, COMO O PERCENTUAL DE VITÓRIAS (**RS\_WIN\_PCT\_COACH**).
    - USO DE CÓDIGO **STINT-AWARE** PARA TRATAR SEPARADAMENTE PERÍODOS COM TROCA DE TÉCNICO.
    - APLICAÇÃO DE **SUAVIZAÇÃO BAYESIANA EMPÍRICA** PARA CORRIGIR AMOSTRAS PEQUENAS, AJUSTANDO O DESEMPENHO EM DIREÇÃO À MÉDIA DA LIGA.
    - AVALIAÇÃO DE VALOR AGREGADO: COMPARA-SE A TAXA DE VITÓRIAS DO TÉCNICO COM A FORÇA DO ELENCO (**TEAM\_STRENGTH**) PARA IDENTIFICAR SE O TREINADOR ENTREGA RESULTADOS ACIMA DO ESPERADO.



- RELAÇÃO DAS VITÓRIAS DOS TÉCNICOS COM A QUANTIDADE DE VITÓRIAS ESPERADAS DO ELENCO
- TÉCNICOS BONS TRAZEM MAIS VITÓRIAS DO QUE ESPERADO



# PREVISÃO DO RANKING

- PREVISÃO DE RANKING
- MODELO UTILIZADO: **CLASSIFICATIVO**, AO INVÉS DE REGRESSIVO.
- MÉTODO: PAIRWISE LEARNING-TO-RANK COM GRADIENT BOOSTING CLASSIFIER.
  - PREVÊ RANKINGS DE EQUIPES USANDO DADOS HISTÓRICOS E CARACTERÍSTICAS TEMPORAIS (MÉDIAS MÓVEIS E TENDÊNCIAS).
- SETUP EXPERIMENTAL:
  - ITERAÇÃO 1: TREINO (ANOS 1-6) | VALIDAÇÃO (ANO 7)
  - ITERAÇÃO 2: TREINO (ANOS 1-7) | VALIDAÇÃO (ANO 8)
  - TESTE FINAL: TREINO (ANOS 1-8) | TESTE (ANOS 9-10)

# PREVISÃO DO RANKING

- FEATURES DO HISTÓRICO, CONSISTÊNCIA E QUALIDADE ESTRUTURAL
  - **COACH\_CAREER\_RS\_WIN\_PCT\_MA5**: MÉDIA DE VITÓRIAS DO TÉCNICO NA CARREIRA (5 ANOS), REFLETINDO EXPERIÊNCIA.
  - **TEAM\_STRENGTH\_MA5**: FORÇA MÉDIA DO ELENCO NOS ÚLTIMOS 5 ANOS, INDICANDO POTENCIAL DO GRUPO.
  - **CONF\_WIN\_PCT\_MA5**: MÉDIA DE VITÓRIAS NA CONFERÊNCIA (5 ANOS), MEDINDO COMPETITIVIDADE REGIONAL.
  - **OFF\_EFF\_MA5**: EFICIÊNCIA OFENSIVA MÉDIA EM 5 ANOS, CAPTURANDO QUALIDADE ESTRUTURAL.
  - **DEF\_EFF\_MA5**: EFICIÊNCIA DEFENSIVA MÉDIA EM 5 ANOS.
  - **PYTHAG\_WIN\_PCT\_MA5**: PORCENTAGEM DE VITÓRIAS ESPERADAS (PYTHAGOREAN) MÉDIA EM 5 ANOS, AJUSTADA POR PONTOS.
  - **POINT\_DIFF\_MA5**: SALDO DE PONTOS MÉDIO DOS ÚLTIMOS 5 ANOS, REFLETINDO ATAQUE/DEFESA.

# PREVISÃO DO RANKING

- FEATURES DE TENDÊNCIAS, ESTABILIDADE E RECÊNCIA
  - **POINT\_DIFF\_TREND5**: TENDÊNCIA DO SALDO DE PONTOS EM 5 ANOS, INDICANDO EVOLUÇÃO.
  - **PYTHAG\_WIN\_PCT\_TREND5**: TENDÊNCIA DA PORCENTAGEM DE VITÓRIAS ESPERADAS EM 5 ANOS.
  - **WIN\_PCT\_CONSISTENCY**: ESTABILIDADE DA TAXA DE VITÓRIAS AO LONGO DO TEMPO.
  - **PREV\_RANK\_MA5**: POSIÇÃO MÉDIA DOS ÚLTIMOS 5 ANOS (CONSISTÊNCIA HISTÓRICA).
  - **PREV\_RANK\_MA3**: POSIÇÃO MÉDIA DOS ÚLTIMOS 3 ANOS (JANELA MAIS CURTA).
  - **COACH\_TENURE\_PREV**: TEMPO DO TREINADOR NO TIME NA TEMPORADA ANTERIOR (ESTABILIDADE TÉCNICA).
  - **PREV\_RANK\_1**: POSIÇÃO DO TIME NO ANO ANTERIOR (EFEITO IMEDIATO).

# FUNCIONAMENTO DA PREVISÃO DO RANKING

- A FUNÇÃO `GENERATE_PAIRWISE_DATA` CONSTRÓI PARES DE EQUIPAS E CALCULA DIFERENÇAS VETORIAIS ENTRE AS SUAS FEATURES ESTATÍSTICAS:

Equipa	Feature 1 ( <code>win_pct</code> )	Feature 2 ( <code>strength</code> )	Feature 3 ( <code>tenure</code> )
A (Lakers)	0.60 (60%)	15.5	1 ano
B (Celtics)	0.55 (55%)	18.0	4 anos

- EXEMPLO DE VETORES DE FEATURES
- (E.G., EQUIPA A - EQUIPA B = [+0.05, -2.5, -3])

- O MODELO `GRADIENTBOOSTINGCLASSIFIER` APRENDE PADRÕES NESSAS DIFERENÇAS RELATIVAS, ESTIMANDO PROBABILIDADES DE SUPERIORIDADE  $P(A > B)$ .
- NA FUNÇÃO `PREDICT_RANKS_PAIRWISE`, O MODELO REALIZA COMPARAÇÕES EXAUSTIVAS ENTRE TODAS AS EQUIPAS DE UMA CONFERÊNCIA. PARA CADA PAR, PREVÊ A PROBABILIDADE DE VITÓRIA (E.G.,  $P(A > B) = 0.8$ ).
- CÁLCULO DO SCORE FINAL:
  - $SCORE(A) = \sum P(A > I)$  PARA TODAS AS EQUIPAS I
  - $SCORE(B) = \sum P(B > I)$  PARA TODAS AS EQUIPAS I
- EXEMPLO:  $SCORE(A) = 1.7 \rightarrow RANK\ 1$ ;  $SCORE(B) = 0.8 \rightarrow RANK\ 2$
- VANTAGEM: ESTA ABORDAGEM PAIRWISE CAPTURA DIFERENÇAS ORDINAIS ENTRE EQUIPAS, ALINHANDO-SE COM O OBJETIVO DE CLASSIFICAÇÃO RELATIVA EM VEZ DE PREVISÃO ABSOLUTA DE DESEMPENHO.

# FUNCIONAMENTO DA PREVISÃO DO RANKING

- PROCESSO DE PREVISÃO E RANKING
  - A FUNÇÃO **PREDICT\_RANKS\_PAIRWISE** EXECUTA COMPARAÇÕES EXAUSTIVAS ENTRE TODAS AS EQUIPAS DE UMA CONFERÊNCIA, CALCULANDO PROBABILIDADES DE SUPERIORIDADE PARA CADA PAR.
  - ALGORITMO DE SCORING:
  - PARA CADA EQUIPA I, O SCORE FINAL É DADO POR:
    - **$\text{SCORE}(I) = \sum_j P(I > J)$**
  - ONDE J REPRESENTA TODAS AS OUTRAS EQUIPAS DA CONFERÊNCIA.
- EXEMPLO ILUSTRATIVO:
  - EQUIPA A:  $P(A > B) + P(A > C) = 0.8 + 0.9 = 1.7 \rightarrow \text{RANK } 1$
  - EQUIPA B:  $P(B > A) + P(B > C) = 0.2 + 0.6 = 0.8 \rightarrow \text{RANK } 2$
  - EQUIPA C:  $P(C > A) + P(C > B) = 0.1 + 0.4 = 0.5 \rightarrow \text{RANK } 3$
- FUNDAMENTAÇÃO: ESTA ABORDAGEM PRIORIZA A ORDENAÇÃO RELATIVA SOBRE MÉTRICAS ABSOLUTAS. EM CONTEXTO COMPETITIVO, A CLASSIFICAÇÃO BASEIA-SE NA SUPERIORIDADE COMPARATIVA ENTRE EQUIPAS, NÃO EM VALORES NUMÉRICOS ISOLADOS DE DESEMPENHO.

# PREVISÃO DO RANKING

- MÉTRICAS USADAS:
- **MAE** (MEAN ABSOLUTE ERROR): QUANTIFICA, EM TERMOS LINEARES, O DESVIO MÉDIO DAS POSIÇÕES ATRIBUÍDAS PELO MODELO.
- **SPEARMAN** (SPEARMAN’S RANK CORRELATION COEFFICIENT): INDICA O GRAU EM QUE A ORDEM RELATIVA DOS ELEMENTOS É PRESERVADA ENTRE O RANKING VERDADEIRO E O RANKING PREVISTO.
- **NDCG** (NORMALIZED DISCOUNTED CUMULATIVE GAIN): UTILIZA UM FATOR DE DESCONTO LOGARÍTMICO PARA PENALIZAR ERROS EM POSIÇÕES MAIS BAIXAS, NORMALIZANDO O RESULTADO DE MODO A FACILITAR COMPARAÇÕES ENTRE LISTAS.
- **WITHIN-K ACCURACY**: MEDE A PROPORÇÃO DE INSTÂNCIAS EM QUE O ITEM CORRETO APARECE ENTRE AS K PRIMEIRAS POSIÇÕES DO RANKING PRODUZIDO PELO MODELO. QUANTIFICA A CAPACIDADE DO MÉTODO EM POSICIONAR O ELEMENTO RELEVANTE PRÓXIMO AO TOPO, MESMO QUANDO NÃO ACERTA A POSIÇÃO EXATA.

TRAIN METRICS (AGGREGATED)			TEST METRICS		
MAE_rank: 1.5826					
Mean_Spearman: 0.4516					
Mean_NDCG@10: 0.9164					
n_groups: 16					
VALIDATION METRICS			Test Year 9:		
Validation Year 7:			Conference	EA	WE
MAE_rank	0.5714	0.2857	MAE_rank	1.1429	1.7143
Spearman	0.8929	0.9643	Spearman	0.7857	0.1786
NDCG@10	0.9894	0.9960	NDCG@10	0.9651	0.9002
Within-1 Acc	85.71%	100.00%	Within-1 Acc	71.43%	57.14%
Within-3 Acc	100.00%	100.00%	Within-3 Acc	100.00%	85.71%
n_teams	7	7	n_teams	7	7
Validation Year 8:			Test Year 10:		
Conference	EA	WE	Conference	EA	WE
MAE_rank	1.3333	2.2857	MAE_rank	2.5714	0.6667
Spearman	0.6000	-0.1071	Spearman	-0.1429	0.8286
NDCG@10	0.9174	0.7634	NDCG@10	0.7992	0.9913
Within-1 Acc	50.00%	42.86%	Within-1 Acc	28.57%	83.33%
Within-3 Acc	100.00%	71.43%	Within-3 Acc	71.43%	100.00%
n_teams	6	7	n_teams	7	6
OVERFITTING DIAGNOSIS					
Train-Test MAE gap: 0.0271					
Val-Test MAE gap: -0.4444					



# PREVISÃO DO RANKING

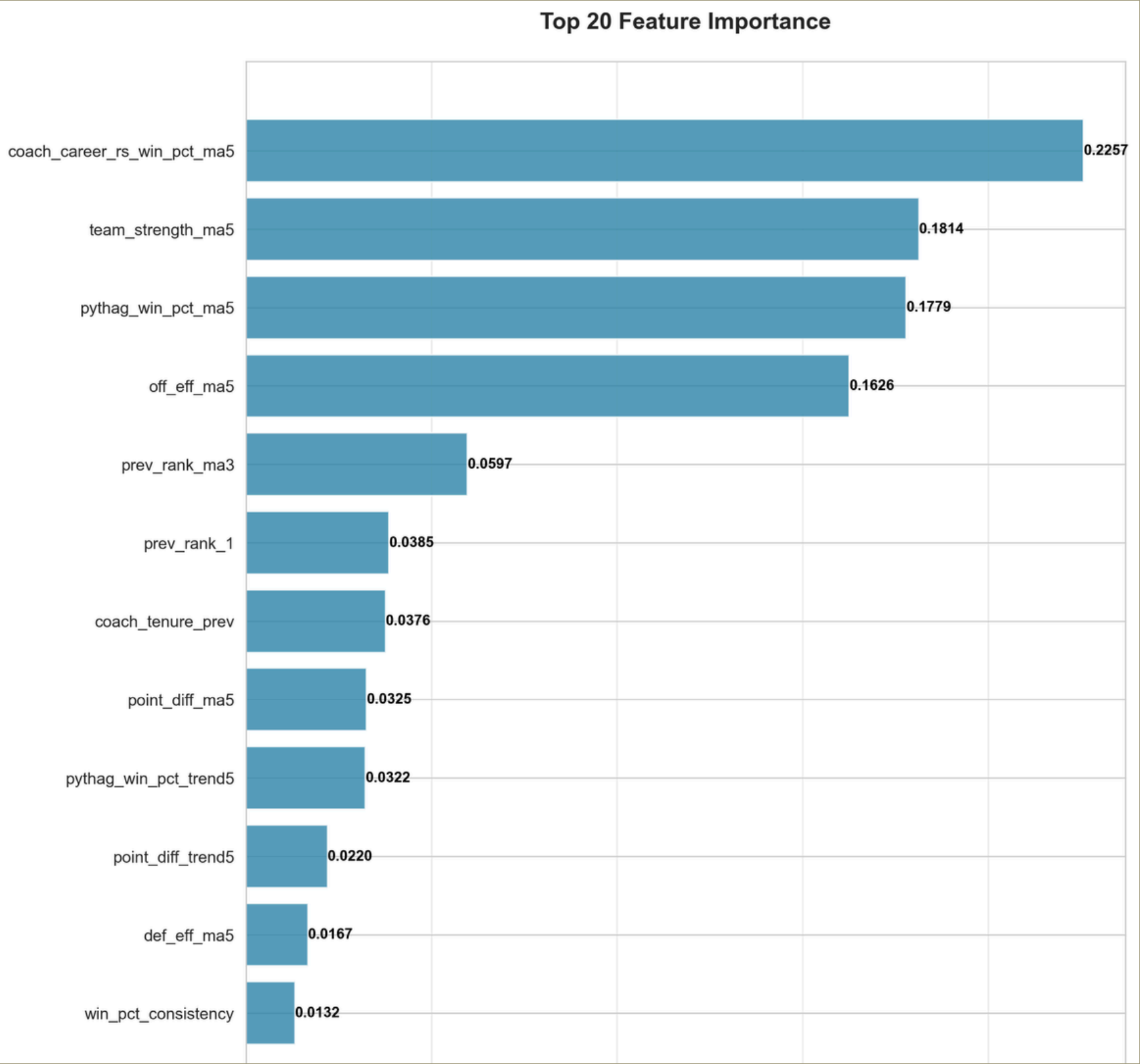
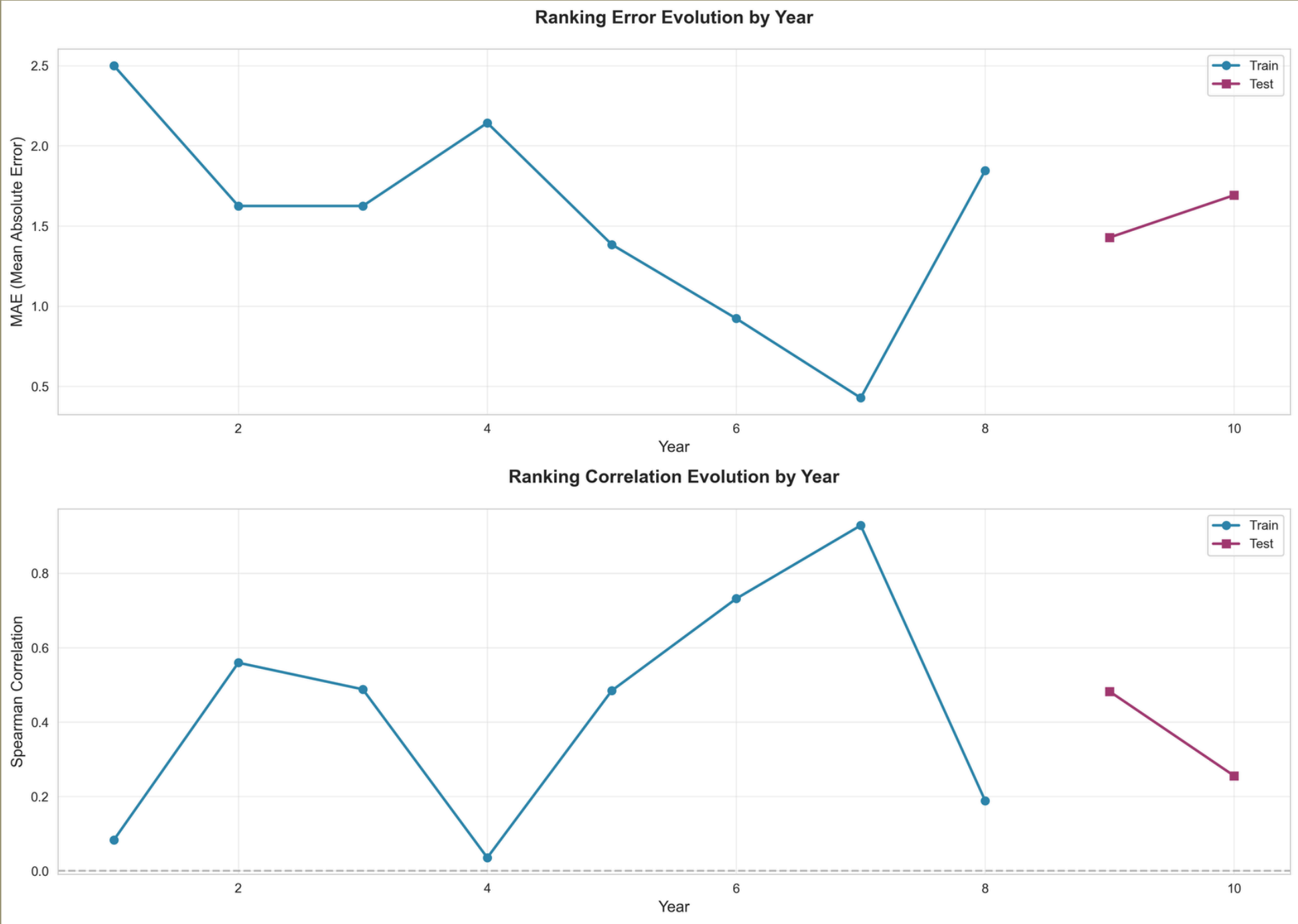
- **RESULTADOS:**
  - A DIFERENÇA DE DESEMPENHO ENTRE TREINO E TESTE É BAIXA, INDICANDO QUE O MODELO APRENDEU PADRÕES RELEVANTES SEM INCORRER EM OVERFITTING.
  - OBSERVA-SE VARIAÇÃO SIGNIFICATIVA ENTRE DIFERENTES CONFERÊNCIAS.
  - EMBORA O MODELO NÃO ACERTE EXATAMENTE O RANKING DE TODOS OS TIMES, APRESENTA ALTA PRECISÃO EM PROXIMIDADE, ALCANÇANDO WITHIN-3 ACCURACY DE 100% EM DIVERSOS CASOS.
  - DESTACAM-SE COMO OUTLIERS NEGATIVOS A CONFERÊNCIA EA NO ANO 10 E A CONFERÊNCIA WE NO ANO 8.

=====			=====		
TRAIN METRICS (AGGREGATED)			TEST METRICS		
=====			=====		
MAE_rank: 1.5826					
Mean_Spearman: 0.4516					
Mean_NDCG@10: 0.9164					
n_groups: 16					
			Test Year 9:		
=====			Conference EA WE		
VALIDATION METRICS			MAE_rank 1.1429 1.7143		
=====			Spearman 0.7857 0.1786		
Validation Year 7:			NDCG@10 0.9651 0.9002		
Conference EA WE			Within-1 Acc 71.43% 57.14%		
MAE_rank 0.5714 0.2857			Within-3 Acc 100.00% 85.71%		
Spearman 0.8929 0.9643			n_teams 7 7		
NDCG@10 0.9894 0.9960					
Within-1 Acc 85.71% 100.00%			Test Year 10:		
Within-3 Acc 100.00% 100.00%			Conference EA WE		
n_teams 7 7			MAE_rank 2.5714 0.6667		
			Spearman -0.1429 0.8286		
Validation Year 8:			NDCG@10 0.7992 0.9913		
Conference EA WE			Within-1 Acc 28.57% 83.33%		
MAE_rank 1.3333 2.2857			Within-3 Acc 71.43% 100.00%		
Spearman 0.6000 -0.1071			n_teams 7 6		
NDCG@10 0.9174 0.7634					
Within-1 Acc 50.00% 42.86%			=====		
Within-3 Acc 100.00% 71.43%			OVERFITTING DIAGNOSIS		
n_teams 6 7			=====		
			Train-Test MAE gap: 0.0271		
			Val-Test MAE gap: -0.4444		



# PREVISÃO DO RANKING

• GRÁFICOS:



# PREVISÃO DO RANKING

- ANO 11:
  - TREINAR DO ANO 1 AO 10:
  - TESTAR ANO 11:

