

Documentação Técnica – Algoritmo rank_similarity: Cálculos e Premissas

Objetivo

Este documento descreve, em nível matemático e procedimental, o algoritmo rank_similarity para ranquear candidatos (anos ou ano-mês) pela similaridade em relação a um alvo. Cobre notação, pré-processamentos, agregação temporal, remoção de outliers, métricas, normalizações, cálculo de scores e critérios de cobertura.

Notação

Seja $X_v(t)$ a série da variável v após agregação/normalização na malha temporal T . Seja $Y_v^c(t)$ a curva do candidato c para a variável v , alinhada à mesma malha. O conjunto de variáveis é V ; os candidatos formam o conjunto C .

Construção do alvo (target)

Quando o alvo não é fornecido, é construído por estatística + método de agregação:

- compare_by='year', agg_freq='monthly': target_v(m) = estatística(média/min/máx) por mês $m=1..12$, agregando todos os anos.
- compare_by='year', agg_freq='daily': target_v(doy) = estatística por dia-do-ano.
- compare_by='year', agg_freq='hourly': target_v(k) com $k=(\text{dia-do-ano}, \text{hora})$.
- compare_by='month': restringe ao mês escolhido e agrega por dia-do-mês (e eventualmente hora).

Agregação temporal

Dados originais $s_v(t_{\text{raw}})$ são reamostrados conforme $\text{agg_freq} \in \{\text{monthly}, \text{daily}, \text{hourly}\}$ e agregados por $\text{agg_stat} \in \{\text{mean}, \text{min}, \text{max}\}$. Pré e/ou pós-agregação podem aplicar filtros de outliers.

Remoção de outliers

Os métodos marcam pontos extremos como NaN (não removem linhas fisicamente):

- Z-score: $z_i = (x_i - \mu) / \sigma$; marca $|z_i| > z_{\text{thresh}}$.
- IQR: $Q1, Q3, IQR = Q3 - Q1$; marca fora de $[Q1 - k \cdot IQR, Q3 + k \cdot IQR]$.

- MAD: mediana m ; $MAD = \text{median}(|x_i - m|)$; $z_{\text{rob}} = 0.6745 \cdot (x_i - m) / MAD$; marca $|z_{\text{rob}}| > \text{mad_thresh}$.

Escopos: `pre_agg` (dentro de cada janela antes da estatística), `post_agg` (na curva agregada), `both`.

Normalização (opcional)

Após agregação, cada série pode ser normalizada:

- zscore: $x' = (x - \text{mean}(x)) / \text{std}(x)$
- minmax: $x' = (x - \min(x)) / (\max(x) - \min(x))$
- none: mantém nível original.

Alinhamento e cobertura

Alinhamos alvo $X_v(t)$ e candidato $Y_{v^c}(t)$ via:

- inner: T = interseção dos índices.
- left_target: T = índices do alvo; candidato é reindexado.

Removem-se pares com NaN. Define-se $\text{coverage}_{v^c} = (\# \text{ pares válidos}) / (\# \text{ pontos válidos do alvo})$. O candidato c só é considerado se $\text{coverage}_{v^c} \geq \text{min_coverage}$ para TODAS as variáveis $v \in V$.

Métricas (modo curve)

Para sequências alinhadas $x = X_v(t)$, $y = Y_{v^c}(t)$:

- Correlação de Pearson (corr, ↑ melhor): $\text{corr} = \text{Cov}(x, y) / (\sigma_x \cdot \sigma_y)$.
- RMSE (↓ melhor): $\text{RMSE} = \sqrt{(1/n) \cdot \sum (x_i - y_i)^2}$.
- MAE (↓ melhor): $\text{MAE} = (1/n) \cdot \sum |x_i - y_i|$.
- DTW (↓ melhor): distância de Dynamic Time Warping com janela Sakoe–Chiba (raio r).

Métricas (modo value)

Quando `compare_mode='value'`, comparamos escalares por grupo: $\text{abs_diff} = |\mu_y - \mu_x|$; $\text{pct_diff} = \text{abs_diff} / (|\mu_x| + \epsilon)$.

Normalização das métricas (min–max por variável)

Para cada variável v e métrica m , normaliza-se para $[0,1]$:

- Se ↑ melhor (ex.: corr): $\text{score}_m = (m - m_{\text{min}}) / (m_{\text{max}} - m_{\text{min}})$.

- Se ↓ melhor (ex.: rmse, mae, dtw): $\text{score}_m = (m_{\max} - m) / (m_{\max} - m_{\min})$.

Se $m_{\max} == m_{\min}$ ou não finito, define-se $\text{score}_m = 1.0$.

Score por variável

Com pesos de métrica w_m (distribuição que soma 1):

$\text{score_var}(v,c) = \sum_m w_m \cdot \text{score}_m(v,c)$. No modo value, utiliza-se 0.5 para abs_diff e 0.5 para pct_diff.

Score final do candidato

Com pesos por variável W_v (normalizados):

$\text{score_final}(c) = \sum_{v \in V} W_v \cdot \text{score_var}(v,c)$. O ranking ordena score_final em ordem decrescente.

Tratamento de NaNs

A remoção de outliers insere NaNs; o alinhamento remove pares com NaN. Coverage é calculada antes do scoring; candidatos abaixo do limiar são descartados.

Limitações e extensões

Exemplo numérico (ilustrativo)

Suponha, para uma variável v e três candidatos A,B,C: $\text{corr}=\{0.90,0.60,0.30\}$, $\text{rmse}=\{2,4,1\}$, $\text{mae}=\{1.2,2.2,1.5\}$.

Normalização (min-max por métrica):

- $\text{corr}\uparrow$: A=1.00, B=0.50, C=0.00; $\text{rmse}\downarrow$: A≈0.667, B=0.00, C=1.00; $\text{mae}\downarrow$: A=1.00, B=0.00, C=0.70.

Com $\text{metric_weights} \{\text{corr}:0.5, \text{rmse}:0.3, \text{mae}:0.2\}$: $\text{score_var}(A)=0.90$; B=0.25; C=0.44. Com múltiplas variáveis, o score_final é a média ponderada por var_weights.

Gerado automaticamente – rank_similarity (documentação de cálculos e premissas).