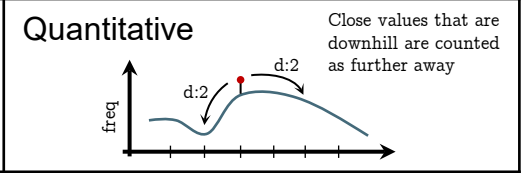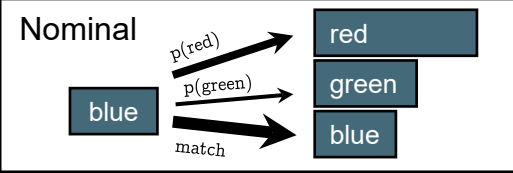# Similarity Based on Resample Exposure

Anton D. Lautrup, Hafiz Saud Arshad, Tobias Hyrup,
Muhammad Rajabinasab, Arthur Zimek, and Peter Schneider-Kamp
*University of Southern Denmark, Odense, Denmark*

**Problem:** Existing measures of proximity in heterogeneous data discards most data insights beyond mere matching and therefore causes unstable results due to no granularity.

**Goal:** Develop a data-driven metric that accounts for frequency of categorical levels and employ distribution shape of the numerical variables.

**Solution:** Resample exposure (REX) similarity

$$\text{REX}(q,t) = \sum_k^m \begin{cases} 0 & \text{if } q_k, \text{or } t_k \text{ is missing} \\ 1 & \text{if } q_k = t_k \\ \Pr(t_k) & \text{if } k\text{th variable is nominal} \\ \left(1 - \dfrac{|q_k - t_k|}{\text{rng}_k}\right)\left(1 - \dfrac{b(q_k t_k)}{\text{trav}_k}\right) & \text{if } k\text{th variable is quantitative} \end{cases}$$



Nominal — p(red), p(green), match → red / green / blue

Quantitative — Close values that are downhill are counted as further away; d:2



| Binary | | | |
|---|---|---|---|
| q | t | GOW | REX |
| A | T | T | 1 | 1 |
| B | T | T | 1 | 1 |
| C | F | F | 1 | 1 |
| D | F | T | 0 | 0.8 |
| E | T | F | 0 | 0.2 |
| F | T | F | 0 | 0.2 |

| Categorical | | | |
|---|---|---|---|
| q | t | GOW | REX |
| A | r | r | 1 | 1 |
| B | g | r | 0 | 0.7 |
| C | b | r | 0 | 0.7 |
| D | r | g | 0 | 0.2 |
| E | b | g | 0 | 0.2 |
| F | r | b | 0 | 0.1 |

| Binary | |
|---|---|
| T | 0.8 |
| F | 0.2 |

| Categorical | |
|---|---|
| r | 0.7 |
| g | 0.2 |
| b | 0.1 |

Uniform  Multimodal

| Uniform | | | |
|---|---|---|---|
| q | t | GOW | REX |
| A | 1.0 | 1.0 | 1 | 1 |
| B | 2.0 | 1.0 | 0.66 | 0.66 |
| C | 1.0 | 2.0 | 0.66 | 0.66 |
| D | 1.0 | 3.0 | 0.33 | 0.33 |
| E | 2.0 | 4.0 | 0.33 | 0.33 |
| F | 1.0 | 4.0 | 0.0 | 0.0 |

| Multimodal | | | |
|---|---|---|---|
| q | t | GOW | REX |
| A | 1.0 | 1.0 | 1 | 1 |
| B | 1.0 | 2.0 | 0.66 | 0.66 |
| C | 2.0 | 1.0 | 0.66 | 0.5 |
| D | 1.0 | 3.0 | 0.33 | 0.22 |
| E | 2.0 | 4.0 | 0.33 | 0.22 |
| F | 4.0 | 1.0 | 0.0 | 0.0 |

| Overall Score | |
|---|---|
| GOW | REX |
| A | 4 | 4 |
| B | 2.33 | 3.03 |
| C | 2.33 | 2.86 |
| D | 0.66 | 1.55 |
| E | 0.66 | 0.95 |
| F | 0.0 | 0.3 |



Euclidean Distance    Gower Distance    Resample Exposure