

# Allotaxonomy

Last updated: 2021/10/06, 20:25:28 EDT

Principles of Complex Systems, Vols. 1 & 2  
CSYS/MATH 300 and 303, 2021–2022 | @pocsvox

Prof. Peter Sheridan Dodds | @peterdodds

Computational Story Lab | Vermont Complex Systems Center  
Vermont Advanced Computing Core | University of Vermont



Licensed under the [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 License](https://creativecommons.org/licenses/by-nc-sa/3.0/).



1 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbulence divergence  
Explorations  
References



2 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbulence divergence  
Explorations  
References

## Outline

A plentitude of distances

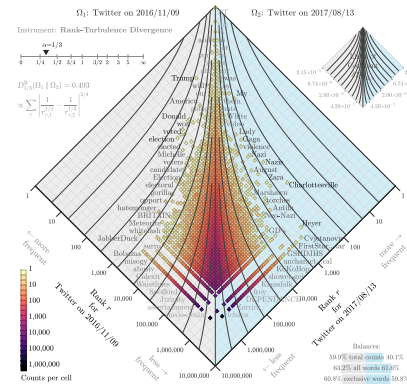
Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

## Goal—Understand this:



Site (papers, examples, code):  
<http://compstorylab.org/allotaxonomy/>

## Foundational papers:

“Allotaxonomy and rank-turbulence divergence: A universal instrument for comparing complex systems”  
Dodds et al., 2020. [5]

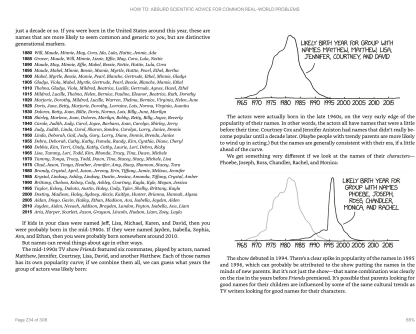
“Probability-turbulence divergence: A tunable allotaxonomic instrument for comparing heavy-tailed categorical distributions”  
Dodds et al., 2020. [6]

## Basic science = Describe + Explain:

- Dashboards of single scale instruments helps us understand, monitor, and control systems.
- Archetype: Cockpit dashboard for flying a plane
- Okay if comprehensible.
- Complex systems present two problems for dashboards:
  - Scale with internal diversity of components: We need meters for every species, every company, every word.
  - Tracking change: We need to re-arrange meters on the fly.
- Goal—Create comprehensible, dynamically-adjusting, differential dashboards showing two pieces:
  - ‘Big picture’ map-like overview,
  - A tunable ranking of components.

1 See the [lexicocalorimeter](#)

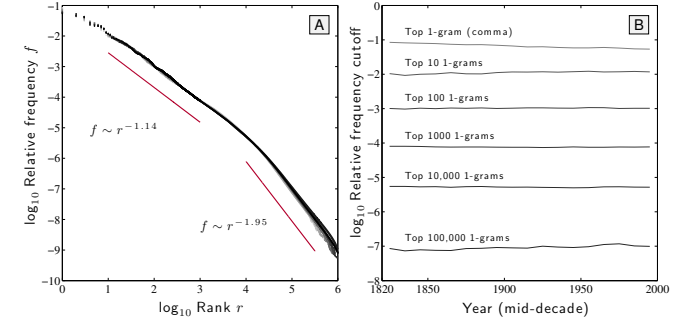
## Baby names, much studied: [12]



How to build a dynamical dashboard that helps sort through a massive number of interconnected time series?



“Is language evolution grinding to a halt? The scaling of lexical turbulence in English fiction suggests it is not”  
Pechenick, Danforth, Dodds, Alshaabi, Adams, Dewhurst, Reagan, Danforth, Reagan, and Danforth.  
Journal of Computational Science, 21, 24–37, 2017. [14]



## For language, Zipf's law has two scaling regimes: [18]

$$f \sim \begin{cases} r^{-\alpha} & \text{for } r \ll r_b \\ r^{-\alpha'} & \text{for } r \gg r_b \end{cases}$$

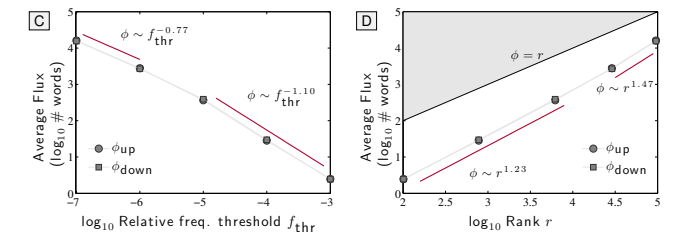
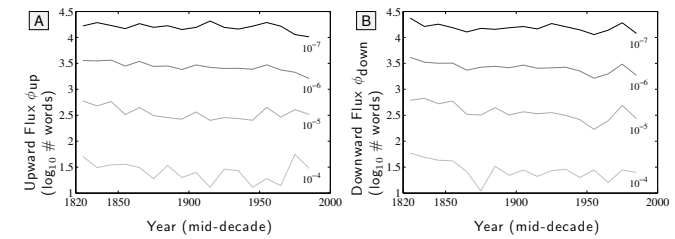
When comparing two texts, define Lexical turbulence as flux of words across a frequency threshold:

$$\phi \sim \begin{cases} f_{thr}^{-\mu} & \text{for } f_{thr} \ll f_b \\ f_{thr}^{-\mu'} & \text{for } f_{thr} \gg f_b \end{cases}$$

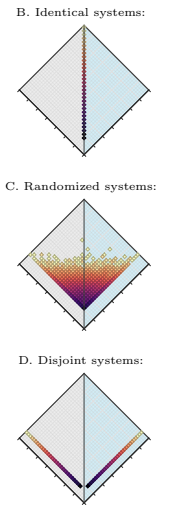
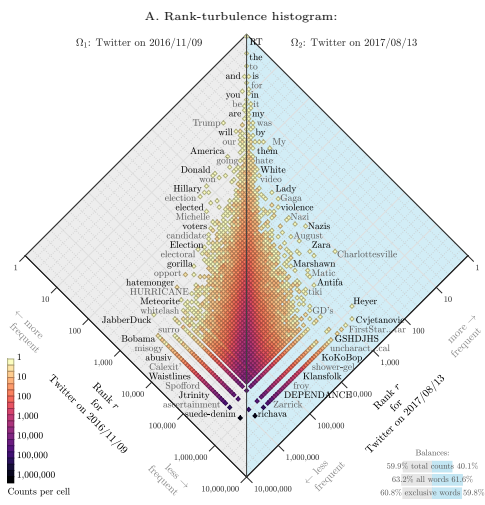
Estimates:  $\mu \approx 0.77$  and  $\mu' \approx 1.10$ , and  $f_b$  is the scaling break point.

$$\phi \sim \begin{cases} r^\nu = r^{\alpha\mu'} & \text{for } r \ll r_b \\ r^{\nu'} = r^{\alpha'\mu} & \text{for } r \gg r_b \end{cases}$$

Estimates: Lower and upper exponents  $\nu \approx 1.23$  and  $\nu' \approx 1.47$ .



7 of 65



## Quite the festival:

Table 1. $L_p$ Minkowski family		
1. Euclidean $L_2$	$d_{Eu} = \sqrt{\sum_{i=1}^n  P_i - Q_i ^2}$	(1)
2. City block $L_1$	$d_{Cb} = \sum_{i=1}^n  P_i - Q_i $	(2)
3. Minkowski $L_p$	$d_{Mk} = \left( \sum_{i=1}^n  P_i - Q_i ^p \right)^{1/p}$	(3)
4. Chebyshev $L_\infty$	$d_{Ck} = \max_i  P_i - Q_i $	(4)

Table 2. $L_1$ family		
5. Sørensen	$d_{Sv} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$	(5)
6. Gower	$d_{Gw} = \frac{1}{d} \frac{\sum_{i=1}^n  P_i - Q_i }{R}$	(6)
7. Soergel	$d_{So} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n \max(P_i, Q_i)}$	(7)
8. Kulczyński $d$	$d_{Kd} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n \min(P_i, Q_i)}$	(8)
9. Canberra	$d_{Cm} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$	(9)
10. Lorentzian	$d_{Lr} = \sum_{i=1}^n \ln(1 +  P_i - Q_i )$	(10)

\*  $L_1$  family  $\supset$  {Intersectoin (13), Wave Hedges (15), Czekanowski (16), Ruzicka (21), Tanimoto (23), etc.}

PoCS @pocsvox Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

14 of 65

### Shannon's Entropy:

$$H(P) = \langle \log_2 \frac{1}{p_\tau} \rangle = \sum_{\tau \in R_{1,2;\alpha}} p_\tau \log_2 \frac{1}{p_\tau} \quad (1)$$

### Kullback-Liebler (KL) divergence:

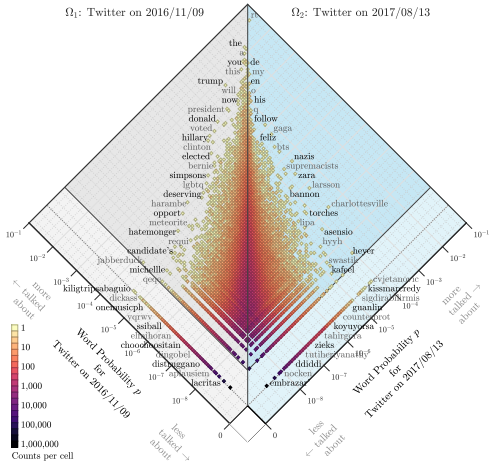
$$D^{KL}(P_2 || P_1) = \left\langle \log_2 \frac{1}{p_{2,\tau}} - \log_2 \frac{1}{p_{1,\tau}} \right\rangle_{P_2}$$

$$= \sum_{\tau \in R_{1,2;\alpha}} p_{2,\tau} \left[ \log_2 \frac{1}{p_{2,\tau}} - \log_2 \frac{1}{p_{1,\tau}} \right]$$

$$= \sum_{\tau \in R_{1,2;\alpha}} p_{2,\tau} \log_2 \frac{p_{1,\tau}}{p_{2,\tau}} \quad (2)$$

- Problem: If just one component type in system 2 is not present in system 1, KL divergence =  $\infty$ .
  - Solution: If we can't compare a spork and a platypus directly, we create a fictional **spork-platypus hybrid**.
  - New problem: Re-read solution.
- 17 of 65

## Zipf-turbulence histogram for probability:



PoCS @pocsvox Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

12 of 65

- We want two main things:
  1. A measure of difference between systems
  2. A way of sorting which types/species/words contribute to that difference
- For sorting, many comparisons give the same ordering.
- A few basic building blocks:
  - $|P_i - Q_i|$  (dominant)
  - $\max(P_i, Q_i)$
  - $\min(P_i, Q_i)$
  - $P_i Q_i$
  - $|P_i^{1/2} - Q_i^{1/2}|$  (Hellinger)

Table 1. $L_p$ Minkowski family		
1. Euclidean $L_2$	$d_{Eu} = \sqrt{\sum_{i=1}^n  P_i - Q_i ^2}$	(1)
2. City block $L_1$	$d_{Cb} = \sum_{i=1}^n  P_i - Q_i $	(2)
3. Minkowski $L_p$	$d_{Mk} = \left( \sum_{i=1}^n  P_i - Q_i ^p \right)^{1/p}$	(3)
4. Chebyshev $L_\infty$	$d_{Ck} = \max_i  P_i - Q_i $	(4)

Table 2. $L_1$ family		
5. Sørensen	$d_{Sv} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$	(5)
6. Gower	$d_{Gw} = \frac{1}{d} \frac{\sum_{i=1}^n  P_i - Q_i }{R}$	(6)
7. Soergel	$d_{So} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n \max(P_i, Q_i)}$	(7)
8. Kulczyński $d$	$d_{Kd} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n \min(P_i, Q_i)}$	(8)
9. Canberra	$d_{Cm} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$	(9)
10. Lorentzian	$d_{Lr} = \sum_{i=1}^n \ln(1 +  P_i - Q_i )$	(10)

\*  $L_1$  family  $\supset$  {Intersectoin (13), Wave Hedges (15), Czekanowski (16), Ruzicka (21), Tanimoto (23), etc.}

PoCS @pocsvox Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

15 of 65

## So, so many ways to compare probability distributions:

- "Families of Alpha- Beta- and Gamma-Divergences: Flexible and Robust Measures of Similarities" Cichocki and Amari, Entropy, 12, 1532-1568, 2010. [2]
- "Comprehensive survey on distance/similarity measures between probability density functions" Sung-Hyuk Cha, International Journal of Mathematical Models and Methods in Applied Sciences, 1, 300-307, 2007. [1]
- Comparisons are distances, divergences, similarities, inner products, fidelities ...
- A worry: Subsampled distributions with very heavy tails
- 60ish kinds of comparisons grouped into 20 families

PoCS @pocsvox Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

13 of 65

- Information theoretic sortings are more opaque
- No tunability

Table 1. $L_p$ Minkowski family		
1. Euclidean $L_2$	$d_{Eu} = \sqrt{\sum_{i=1}^n  P_i - Q_i ^2}$	(1)
2. City block $L_1$	$d_{Cb} = \sum_{i=1}^n  P_i - Q_i $	(2)
3. Minkowski $L_p$	$d_{Mk} = \left( \sum_{i=1}^n  P_i - Q_i ^p \right)^{1/p}$	(3)
4. Chebyshev $L_\infty$	$d_{Ck} = \max_i  P_i - Q_i $	(4)

Table 2. $L_1$ family		
5. Sørensen	$d_{Sv} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$	(5)
6. Gower	$d_{Gw} = \frac{1}{d} \frac{\sum_{i=1}^n  P_i - Q_i }{R}$	(6)
7. Soergel	$d_{So} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n \max(P_i, Q_i)}$	(7)
8. Kulczyński $d$	$d_{Kd} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n \min(P_i, Q_i)}$	(8)
9. Canberra	$d_{Cm} = \frac{\sum_{i=1}^n  P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$	(9)
10. Lorentzian	$d_{Lr} = \sum_{i=1}^n \ln(1 +  P_i - Q_i )$	(10)

\*  $L_1$  family  $\supset$  {Intersectoin (13), Wave Hedges (15), Czekanowski (16), Ruzicka (21), Tanimoto (23), etc.}

PoCS @pocsvox Allotaxonomy

A plentitude of distances

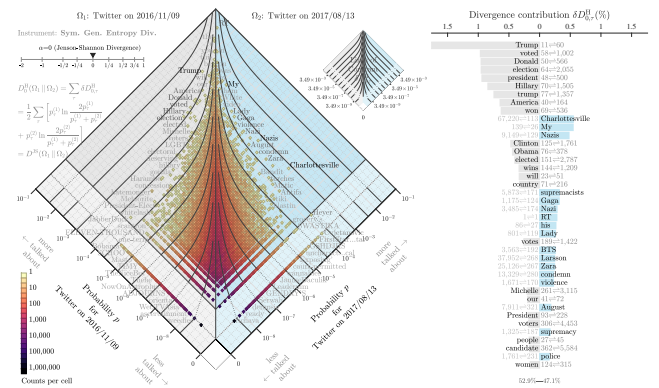
Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

16 of 65



PoCS @pocsvox Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

17 of 65

PoCS @pocsvox Allotaxonomy

A plentitude of distances

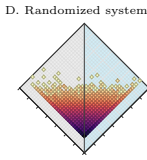
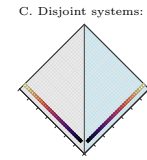
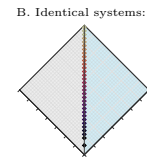
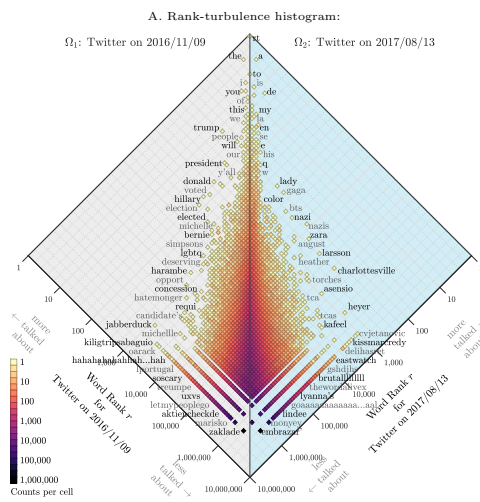
Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

18 of 65



PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbulence divergence  
Explorations  
References

21 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbulence divergence  
Explorations  
References

22 of 65

### Some good things about ranks:

- Working with ranks is intuitive
- Affords some powerful statistics (e.g., Spearman's rank correlation coefficient)
- Can be used to generalize beyond systems with probabilities

### A start:

$$\left| \frac{1}{r_{\tau,1}} - \frac{1}{r_{\tau,2}} \right| \quad (5)$$

- Inverse of rank gives an increasing measure of 'importance'
- High rank means closer to rank 1
- We assign tied ranks for components of equal 'size'
- Issue: Biases toward high rank components

### We introduce a tuning parameter:

$$\left| \frac{1}{[r_{\tau,1}]^\alpha} - \frac{1}{[r_{\tau,2}]^\alpha} \right|^{1/\alpha} \quad (6)$$

- As  $\alpha \rightarrow 0$ , high ranked components are increasingly dampened
- For words in texts, for example, the weight of common words and rare words move increasingly closer together.
- As  $\alpha \rightarrow \infty$ , high rank components will dominate.
- For texts, the contributions of rare words will vanish.

### Trouble:

- The limit of  $\alpha \rightarrow 0$  does not behave well for

$$\left| \frac{1}{[r_{\tau,1}]^\alpha} - \frac{1}{[r_{\tau,2}]^\alpha} \right|^{1/\alpha}$$

- The leading order term is:

$$(1 - \delta_{r_{\tau,1} r_{\tau,2}}) \alpha^{1/\alpha} \left| \ln \frac{r_{\tau,1}}{r_{\tau,2}} \right|^{1/\alpha}, \quad (7)$$

- which heads toward  $\infty$  as  $\alpha \rightarrow 0$ .
- Oops.
- But the insides look nutritious:

$$\left| \ln \frac{r_{\tau,1}}{r_{\tau,2}} \right|$$

is a nicely interpretable log-ratio of ranks.

### Some reworking:

$$\delta D_{\alpha, \tau}^R(R_1 || R_2) \propto \frac{\alpha + 1}{\alpha} \left| \frac{1}{[r_{\tau,1}]^\alpha} - \frac{1}{[r_{\tau,2}]^\alpha} \right|^{1/(\alpha+1)} \quad (8)$$

- Keeps the core structure.
- Large  $\alpha$  limit remains the same.
- $\alpha \rightarrow 0$  limit now returns log-ratio of ranks.
- Next: Sum over  $\tau$  to get divergence.
- Still have an option for normalization.

### Rank-turbulence divergence:

$$D_{\alpha}^R(R_1 || R_2) = \frac{1}{\mathcal{N}_{1,2;\alpha}} \sum_{\tau \in R_{1,2;\alpha}} \delta D_{\alpha, \tau}^R(R_1 || R_2) \quad (9)$$

### Normalization:

- Take a data-driven rather than analytic approach to determining  $\mathcal{N}_{1,2;\alpha}$ .
- Compute  $\mathcal{N}_{1,2;\alpha}$  by taking the two systems to be disjoint while maintaining their underlying Zipf distributions.
- Ensures:  $0 \leq D_{\alpha}^R(R_1 || R_2) \leq 1$
- Limits of 0 and 1 correspond to the two systems having identical and disjoint Zipf distributions.

### Rank-turbulence divergence:

Summing over all types, dividing by a normalization prefactor  $\mathcal{N}_{1,2;\alpha}$  we have our prototype:

$$D_{\alpha}^R(R_1 || R_2) = \frac{1}{\mathcal{N}_{1,2;\alpha}} \frac{\alpha + 1}{\alpha} \sum_{\tau \in R_{1,2;\alpha}} \left| \frac{1}{[r_{\tau,1}]^\alpha} - \frac{1}{[r_{\tau,2}]^\alpha} \right|^{1/(\alpha+1)} \quad (10)$$

### Exclusive types:

- We call types that are present in one system only 'exclusive types'.
- When warranted, we will use expressions of the form  $\Omega^{(1)}$ -exclusive and  $\Omega^{(2)}$ -exclusive to indicate to which system an exclusive type belongs.

### Desirable rank-turbulence divergence features:

- Rank-based.
- Symmetric.
- Semi-positive:  $D_{\alpha}^R(\Omega_1 || \Omega_2) \geq 0$ .
- Linearly separable, for interpretability.
- Subsystem applicable: Ranked lists of any principled subset may be equally well compared (e.g., hashtags on Twitter, stock prices of a certain sector, etc.).
- Zipfophilic: Able to handle systems with rank-ordered component size distribution that are heavy-tailed.
- Scalable: Allow for sensible comparisons across system sizes.
- Tunable.
- Story-finding: Features 1-8 combine to show which component types are most 'important'

## General normalization:

☞ If the Zipf distributions are disjoint, then in  $\Omega^{(1)}$ 's merged ranking, the rank of all  $\Omega^{(2)}$  types will be  $r = N_1 + \frac{1}{2}N_2$ , where  $N_1$  and  $N_2$  are the number of distinct types in each system.

☞ Similarly,  $\Omega^{(2)}$ 's merged ranking will have all of  $\Omega^{(1)}$ 's types in last place with rank  $r = N_2 + \frac{1}{2}N_1$ .

☞ The normalization is then:

$$\mathcal{N}_{1,2;\alpha} = \frac{\alpha+1}{\alpha} \sum_{\tau \in R_1} \left| \frac{1}{[r_{\tau,1}]^\alpha} - \frac{1}{[N_1 + \frac{1}{2}N_2]^\alpha} \right|^{1/(\alpha+1)} + \frac{\alpha+1}{\alpha} \sum_{\tau \in R_2} \left| \frac{1}{[N_2 + \frac{1}{2}N_1]^\alpha} - \frac{1}{[r_{\tau,2}]^\alpha} \right|^{1/(\alpha+1)} \quad (11)$$

## Limit of $\alpha \rightarrow 0$ :

$$D_0^R(R_1 \| R_2) = \sum_{\tau \in R_{1,2;\alpha}} \delta D_{0,\tau}^R = \frac{1}{\mathcal{N}_{1,2;0}} \sum_{\tau \in R_{1,2;\alpha}} \left| \ln \frac{r_{\tau,1}}{r_{\tau,2}} \right|, \quad (12)$$

where

$$\mathcal{N}_{1,2;0} = \sum_{\tau \in R_1} \left| \ln \frac{r_{\tau,1}}{N_1 + \frac{1}{2}N_2} \right| + \sum_{\tau \in R_2} \left| \ln \frac{r_{\tau,2}}{\frac{1}{2}N_1 + N_2} \right|. \quad (13)$$

☞ Largest rank ratios dominate.

## Limit of $\alpha \rightarrow \infty$ :

$$D_\infty^R(R_1 \| R_2) = \sum_{\tau \in R_{1,2;\alpha}} \delta D_{\infty,\tau}^R = \frac{1}{\mathcal{N}_{1,2;\infty}} \sum_{\tau \in R_{1,2;\alpha}} (1 - \delta_{r_{\tau,1} r_{\tau,2}}) \max_\tau \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}. \quad (14)$$

where

$$\mathcal{N}_{1,2;\infty} = \sum_{\tau \in R_1} \frac{1}{r_{\tau,1}} + \sum_{\tau \in R_2} \frac{1}{r_{\tau,2}}. \quad (15)$$

☞ Highest ranks dominate.

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



29 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



30 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



31 of 65

## Probability-turbulence divergence:

$$D_\alpha^P(P_1 \| P_2) = \frac{1}{\mathcal{N}_{1,2;\alpha}^P} \frac{\alpha+1}{\alpha} \sum_{\tau \in R_{1,2;\alpha}} \left| [p_{\tau,1}]^\alpha - [p_{\tau,2}]^\alpha \right|^{1/(\alpha+1)}. \quad (16)$$

☞ For the unnormalized version ( $\mathcal{N}_{1,2;\alpha}^P=1$ ), some troubles return with 0 probabilities and  $\alpha \rightarrow 0$ .

☞ Weep not:  $\mathcal{N}_{1,2;\alpha}^P$  will save the day.

## Normalization:

With no matching types, the probability of a type present in one system is zero in the other, and the sum can be split between the two systems' types:

$$\mathcal{N}_{1,2;\alpha}^P = \frac{\alpha+1}{\alpha} \sum_{\tau \in R_1} [p_{\tau,1}]^{\alpha/(\alpha+1)} + \frac{\alpha+1}{\alpha} \sum_{\tau \in R_2} [p_{\tau,2}]^{\alpha/(\alpha+1)} \quad (17)$$

## Limit of $\alpha=0$ for probability-turbulence divergence

☞ if both  $p_{\tau,1} > 0$  and  $p_{\tau,2} > 0$  then

$$\lim_{\alpha \rightarrow 0} \frac{\alpha+1}{\alpha} \left| [p_{\tau,1}]^\alpha - [p_{\tau,2}]^\alpha \right|^{1/(\alpha+1)} = \left| \ln \frac{p_{\tau,2}}{p_{\tau,1}} \right|. \quad (18)$$

☞ But if  $p_{\tau,1} = 0$  or  $p_{\tau,2} = 0$ , limit diverges as  $1/\alpha$ .

## Limit of $\alpha=0$ for probability-turbulence divergence

☞ Normalization:

$$\mathcal{N}_{1,2;\alpha}^P \rightarrow \frac{1}{\alpha} (N_1 + N_2). \quad (19)$$

☞ Because the normalization also diverges as  $1/\alpha$ , the divergence will be zero when there are no exclusive types and non-zero when there are exclusive types.

## Combine these cases into a single expression:

$$D_0^P(P_1 \| P_2) = \frac{1}{(N_1 + N_2)} \sum_{\tau \in R_{1,2;0}} (\delta_{p_{\tau,1},0} + \delta_{0,p_{\tau,2}}). \quad (20)$$

☞ The term  $(\delta_{p_{\tau,1},0} + \delta_{0,p_{\tau,2}})$  returns 1 if either  $p_{\tau,1} = 0$  or  $p_{\tau,2} = 0$ , and 0 otherwise when both  $p_{\tau,1} > 0$  and  $p_{\tau,2} > 0$ .

☞ Ratio of types that are exclusive to one system relative to the total possible such types,

## Type contribution ordering for the limit of $\alpha=0$

☞ In terms of contribution to the divergence score, all exclusive types supply a weight of  $1/(N_1 + N_2)$ . We can order them by preserving their ordering as  $\alpha \rightarrow 0$ , which amounts to ordering by descending probability in the system in which they appear.

☞ And while types that appear in both systems make no contribution to  $D_0^P(P_1 \| P_2)$ , we can still order them according to the log ratio of their probabilities.

☞ The overall ordering of types by divergence contribution for  $\alpha=0$  is then: (1) exclusive types by descending probability and then (2) types appearing in both systems by descending log ratio.

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



33 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



34 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



35 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



36 of 65

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



37 of 65



# Limit of $\alpha=\infty$ for probability-turbulence divergence

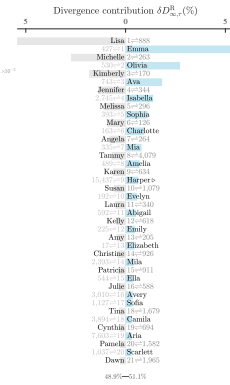
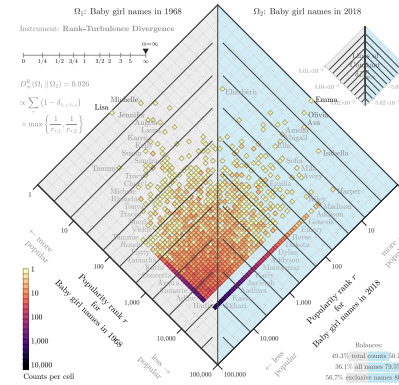
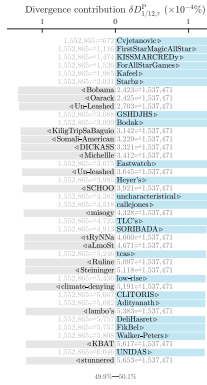
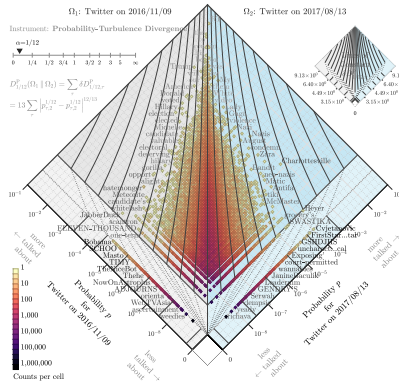
$$D_{\infty}^P(P_1 \| P_2) = \frac{1}{2} \sum_{\tau \in R_{1,2;\infty}} (1 - \delta_{p_{\tau,1}, p_{\tau,2}}) \max(p_{\tau,1}, p_{\tau,2}) \quad (21)$$

where

$$\mathcal{N}_{1,2;\infty}^P = \sum_{\tau \in R_{1,2;\infty}} (p_{\tau,1} + p_{\tau,2}) = 1 + 1 = 2. \quad (22)$$

PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbuler divergence  
Explorations  
References



PoCS  
@pocsvox  
Allotaxonomy

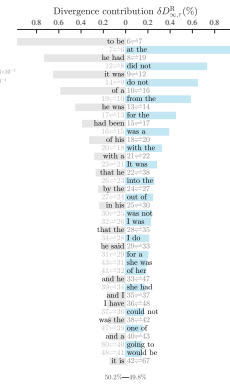
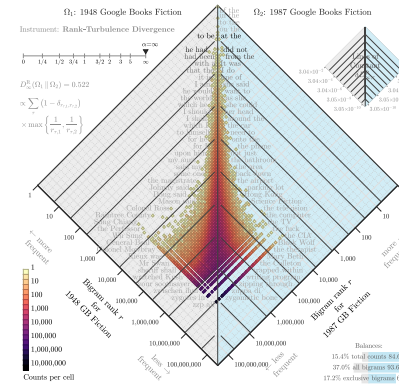
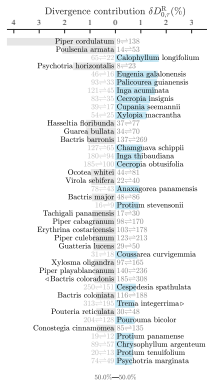
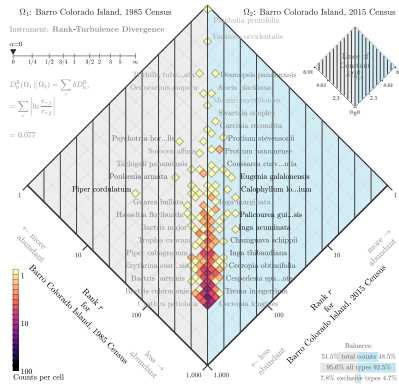
38 of 65

# Connections for PTD:

- $\alpha = 0$ : Similarity measure Sørensen-Dice coefficient [4, 16, 10],  $F_1$  score of a test's accuracy [17, 15].
- $\alpha = 1/2$ : Hellinger distance [8] and Mautusita distance [11].
- $\alpha = 1$ : Many including all  $L^p$ -norm type constructions.
- $\alpha = \infty$ : Motyka distance [3].

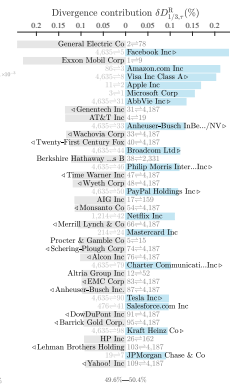
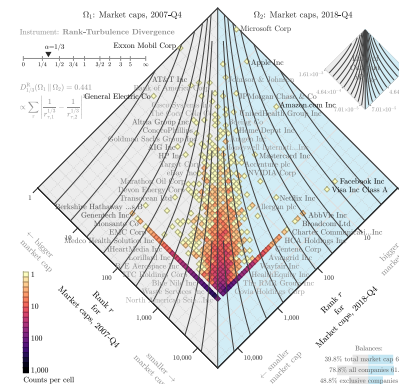
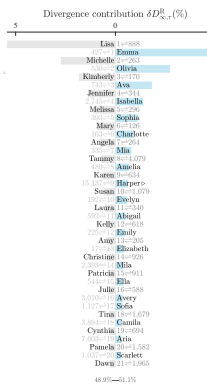
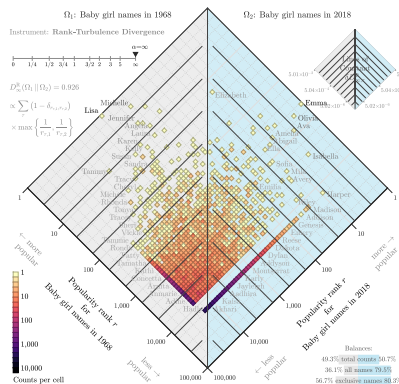
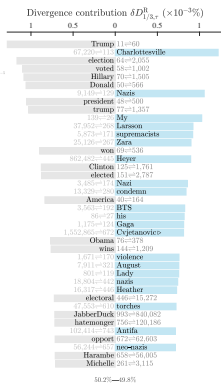
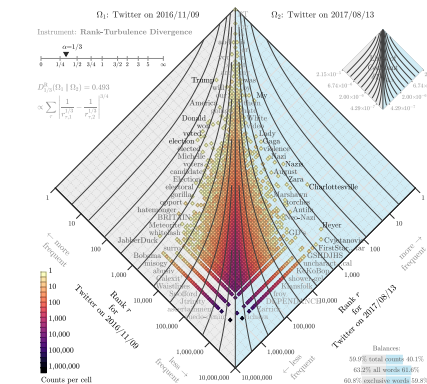
PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbuler divergence  
Explorations  
References

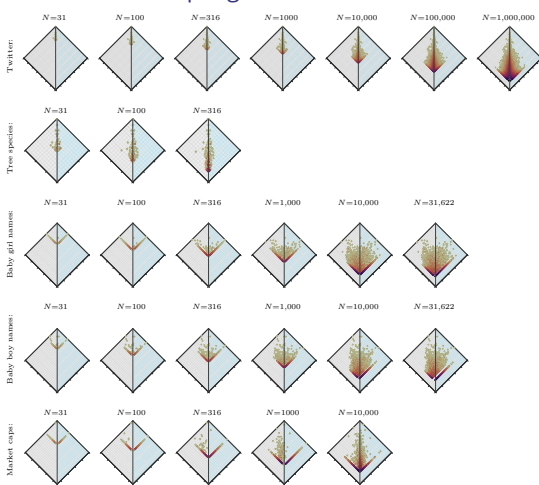


PoCS  
@pocsvox  
Allotaxonomy

39 of 65



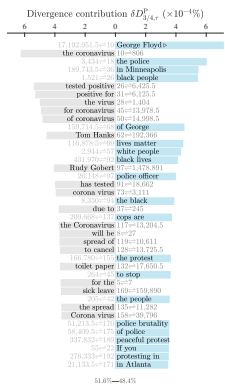
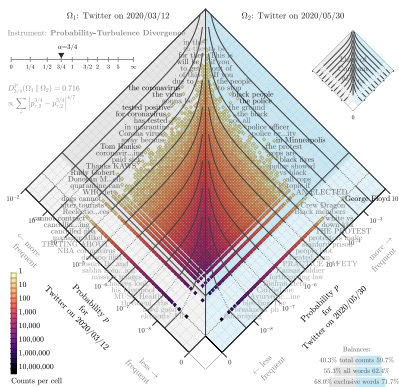
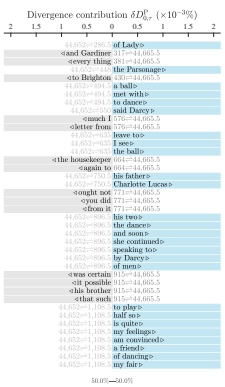
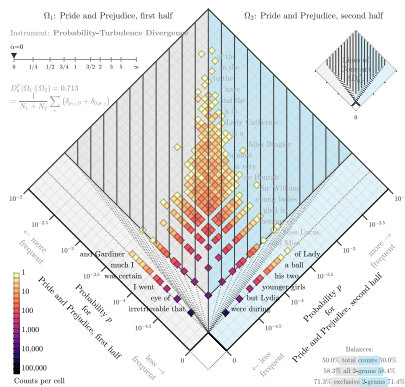
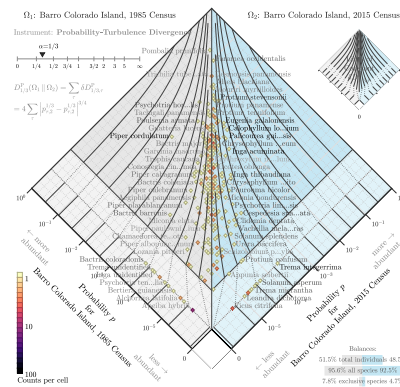
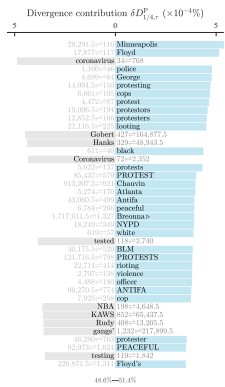
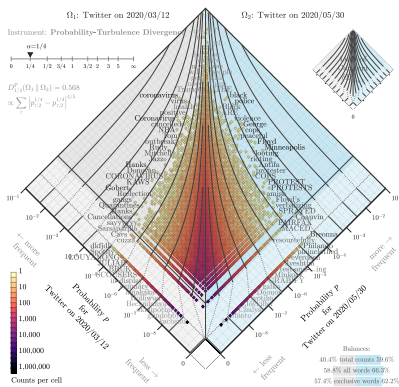
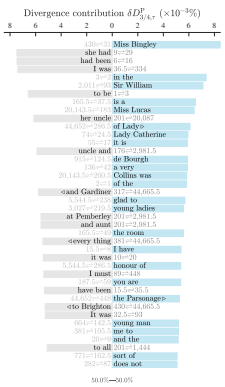
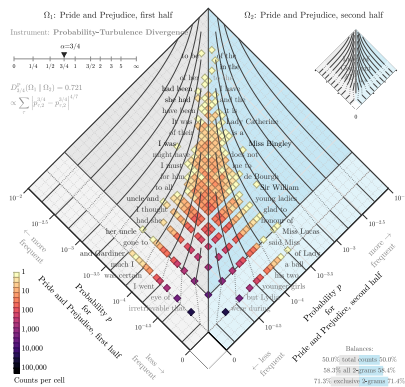
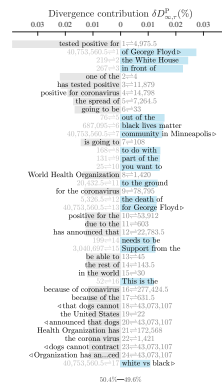
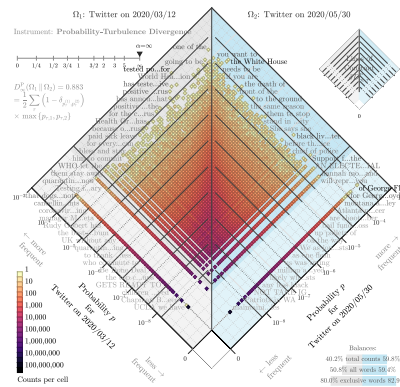
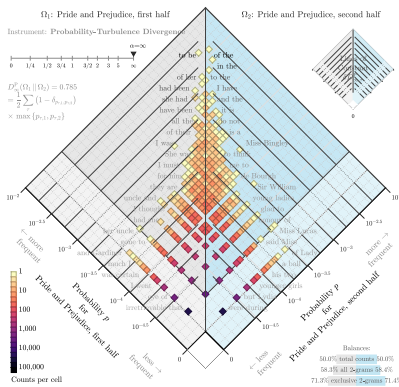
# Effect of subsampling:



PoCS  
@pocsvox  
Allotaxonomy

A plentitude of distances  
Rank-turbulence divergence  
Probability-turbulence divergence  
Explorations  
References

47 of 65



## Flipbooks:

Twitter:

- [instrument-flipbook-1-rank-div.pdf](#)
- [instrument-flipbook-2-probability-div.pdf](#)
- [instrument-flipbook-3-gen-entropy-div.pdf](#)

Market caps:

- [instrument-flipbook-4-marketcaps-6years-rank-div.pdf](#)

Baby names:

- [instrument-flipbook-5-babynames-girls-50years-rank-div.pdf](#)
- [instrument-flipbook-6-babynames-boys-50years-rank-div.pdf](#)

Google books:

- [instrument-flipbook-7-google-books-onegrams-rank-div.pdf](#)
- [instrument-flipbook-8-google-books-bigrams-rank-div.pdf](#)
- [instrument-flipbook-9-google-books-trigrams-rank-div.pdf](#)

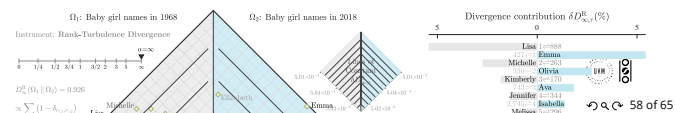
# Flipbooks:

- Pride and Prejudice, 1-grams
- Pride and Prejudice, 2-grams
- Pride and Prejudice, 3-grams
- Twitter, 1-grams
- Twitter, 2-grams
- Twitter, 3-grams
- Barro Colorado Island

Code:  
<https://gitlab.com/compstorylab/allotaxonometer>

## Claims, exaggerations, reminders:

- Needed for comparing large-scale complex systems:  
 Comprehensible, dynamically-adjusting, differential dashboards
- Many measures seem poorly motivated and largely unexamined (e.g., JSD)
- Of value: Combining big-picture maps with ranked lists
- Maybe one day: Online tunable version of rank-turbulence divergence (plus many other instruments)



# References I

[1] S.-H. Cha. Comprehensive survey on distance/similarity measures between probability density functions. *International Journal of Mathematical Models and Methods in Applied Sciences*, 1:300–307, 2007. pdf

[2] A. Cichocki and S.-i. Amari. Families of Alpha- Beta- and Gamma-divergences: Flexible and robust measures of similarities. *Entropy*, 12:1532–1568, 2010. pdf

[3] M.-M. Deza and E. Deza. *Dictionary of Distances*. Elsevier, 2006.

# References II

[4] L. R. Dice. Measures of the amount of ecologic association between species. *Ecology*, 26:297–302, 1945.

[5] P. S. Dodds, J. R. Minot, M. V. Arnold, T. Alshaabi, J. L. Adams, D. R. Dewhurst, T. J. Gray, M. R. Frank, A. J. Reagan, and C. M. Danforth. Allotaxonomy and rank-turbulence divergence: A universal instrument for comparing complex systems, 2020. Available online at <https://arxiv.org/abs/2002.09770>. pdf

# References III

[6] P. S. Dodds, J. R. Minot, M. V. Arnold, T. Alshaabi, J. L. Adams, D. R. Dewhurst, A. J. Reagan, and C. M. Danforth. Probability-turbulence divergence: A tunable allotaxonomic instrument for comparing heavy-tailed categorical distributions, 2020. Available online at <http://arxiv.org/abs/2008.13078>. pdf

[7] D. M. Endres and J. E. Schindelin. A new metric for probability distributions. *IEEE Transactions on Information theory*, 2003. pdf

# References IV

[8] E. Hellinger. Neue begründung der theorie quadratischer formen von unendlichvielen veränderlichen. *Journal für die reine und angewandte Mathematik (Crelles Journal)*, 1909(136):210–271, 1909. pdf

[9] J. Lin. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information theory*, 37(1):145–151, 1991. pdf

[10] J. Looman and J. B. Campbell. Adaptation of Sørensen’s k (1948) for estimating unit affinities in prairie vegetation. *Ecology*, 41(3):409–416, 1960. pdf

# References V

[11] K. Matusita et al. Decision rules, based on the distance, for problems of fit, two samples, and estimation. *The Annals of Mathematical Statistics*, 26(4):631–640, 1955. pdf

[12] R. Munroe. *How To: Absurd Scientific Advice for Common Real-World Problems*. Penguin, 2019.

[13] F. Oesterreicher and I. Vajda. A new class of metric divergences on probability spaces and its applicability in statistics. *Annals of the Institute of Statistical Mathematics*, 55(3):639–653, 2003.

# References VI

[14] E. A. Pechenick, C. M. Danforth, and P. S. Dodds. Is language evolution grinding to a halt? The scaling of lexical turbulence in English fiction suggests it is not. *Journal of Computational Science*, 21:24–37, 2017. pdf

[15] Y. Sasaki. The truth of the  $f$ -measure, 2007.

[16] T. Sorensen. A method of establishing groups of equal amplitude in plant sociology based on similarity of species content and its application to analyses of the vegetation on Danish commons. *Videnskabeligt Selskab Biologiske Skrifter*, 5:1–34, 1948.

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

57 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

58 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

59 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

60 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

61 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

62 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

63 of 65

PoCS  
 @pocsvox  
 Allotaxonomy

A plenitude of distances  
 Rank-turbulence divergence  
 Probability-turbulence divergence  
 Explorations  
 References

64 of 65



## References VII

- [17] C. J. Van Rijsbergen.  
**Information retrieval.**  
Butterworth-Heinemann, 2nd edition, 1979.
- [18] J. R. Williams, J. P. Bagrow, C. M. Danforth, and  
P. S. Dodds.  
Text mixing shapes the anatomy of  
rank-frequency distributions.  
[Physical Review E](#), 91:052811, 2015. [pdf](#)

PoCS  
@pocsvox  
Allotaxonomy

A plenitude of  
distances

Rank-turbulence  
divergence

Probability-  
turbulence  
divergence

Explorations

**References**



65 of 65