

Allotaxonomy

Last updated: 2021/10/06, 20:25:47 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



PoCS
@pocsvox
Allotaxonomy

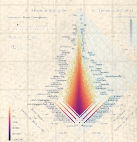
A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References



Licensed under the *Creative Commons Attribution-NonCommercial-ShareAlike 3.0 License*.



These slides are brought to you by:

PoCS
@pocsvox
Allotaxonomy

Sealie & Lambie
Productions



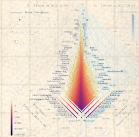
A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

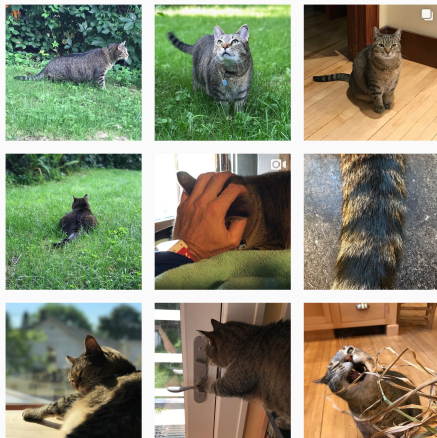
References



These slides are also brought to you by:

PoCS
@pocsvox
Allotaxonomy

Special Guest Executive Producer



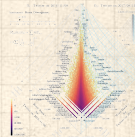
A plentitude of distances



Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



 On Instagram at [pratchett_the_cat](https://www.instagram.com/pratchett_the_cat) 



Outline

PoCS
@pocsvox
Allotaxonomy

A plenitude of distances

A plenitude of
distances

Rank-turbulence divergence

Rank-turbulence
divergence

Probability-turbulence divergence

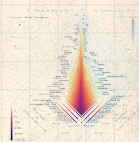
Probability-
turbulence
divergence

Explorations

Explorations

References

References



The Boggoracle Speaks:

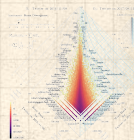
A plenitude of
distances

Rank-turbulence
divergence

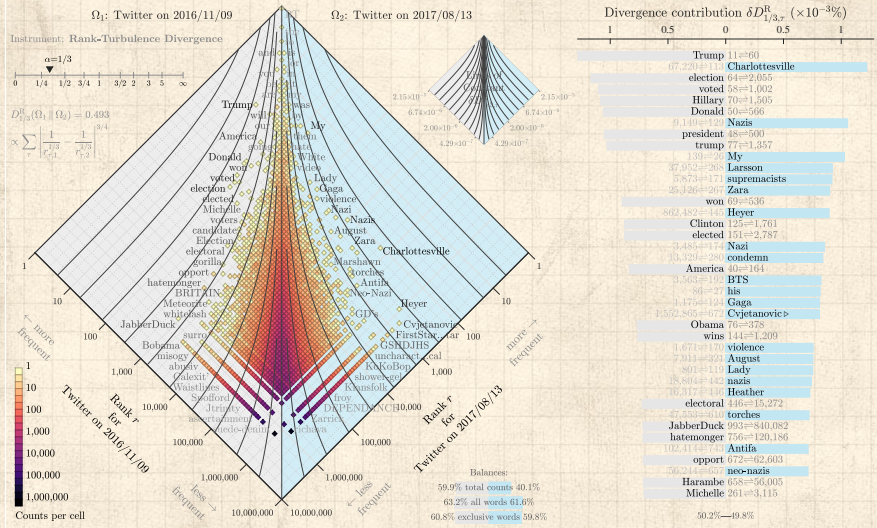
Probability-
turbulence
divergence

Explorations

References



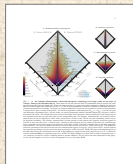
Goal—Understand this:



Site (papers, examples, code):

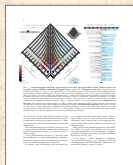
<http://compstorylab.org/allotaxonomy/>

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:

1. Scale with internal diversity of components: We need meters for every species, every company, every word.
2. Tracking change: We need to re-arrange meters on the fly.

 Goal—Create comprehensible, dynamically-adjusting, differential dashboards showing two pieces:¹

1. 'Big picture' map-like overview,
2. A tunable ranking of components.

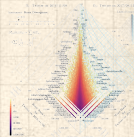
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



¹See the [lexicocalorimeter](#) 

Baby names, much studied:^[12]

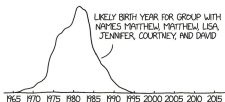
HOW TO: ABSURD SCIENTIFIC ADVICE FOR COMMON REAL-WORLD PROBLEMS

just a decade or so. If you were born in the United States around this year, these are names that are more likely to seem common and generic to you, but are distinctive generational markers.

1890 Will, Maudie, Minnie, May, Cora, Ida, Lela, Hattie, Annie, Ada
1885 Gracey, Maudie, Will, Minnie, Lela, Effie, May, Cora, Lela, Nellie
1900 Maudie, May, Minnie, Effie, Mabel, Bessie, Nellie, Hattie, Lela, Cora
1895 Maudie, Mabel, Minnie, Bessie, Minnie, Myrtle, Hattie, Pearl, Ethel, Bertha
1900 Mabel, Myrtle, Bessie, Minnie, Pearl, Blanche, Gertrude, Ethel, Minnie, Gladys
1905 Gladys, Vada, Mabel, Myrtle, Gertrude, Pearl, Bessie, Blanche, Marnie, Ethel
1910 Thelma, Gladys, Vada, Mildred, Beatrice, Lucille, Gertrude, Agnes, Hazel, Ethel
1915 Mildred, Lucille, Thelma, Helen, Bernice, Pauline, Eleanor, Beatrice, Ruth, Dorothy
1920 Marjorie, Dorothy, Mildred, Lucille, Myrtle, Thelma, Bernice, Virginia, Helen, June
1925 Doris, Jane, Betty, Marjorie, Dorothy, Lorraine, Lela, Susan, Virginia, Beverly
1930 Dolores, Betty, Joan, Ethel, Doris, Norma, Lela, Billy, Jane, Marilyn
1935 Shirley, Marlene, Joan, Dolores, Marilyn, Bobby, Betty, Billy, Joyce, Beverly
1940 Corde, Judith, Judy, Carol, Joyce, Barbara, Joan, Carolyn, Shirley, Jerry
1945 Judy, Judith, Linda, Carol, Sharon, Sandra, Carolyn, Larry, Anita, Dennis
1950 Linda, Deborah, Gail, Andy, Gary, Larry, Diane, Dennis, Brenda, Janice
1955 Debra, Deborah, Cathy, Kathy, Pamela, Randy, Kim, Cynthia, Diane, Cheryl
1960 Debbie, Kim, Tami, Cindy, Kathy, Cathy, Laverie, Lori, Debra, Ricky
1965 Lisa, Tammy, Lori, Tiffni, Kim, Alexandra, Tracy, Tina, Dana, Michele
1970 Tammy, Tanya, Tracy, Todd, Dana, Tina, Sherry, Stacy, Michele, Lisa
1975 Cheri, Susan, Tanya, Heather, Jennifer, Amy, Stacy, Shannon, Sherry, Tary
1980 Brenda, Crystal, April, Susan, Arlene, Kim, Tiffany, Janice, Melissa, Jennifer
1985 Crystal, Lashay, Ashley, Lindsey, Doreen, Jessica, Amanda, Tiffany, Amber
1990 Britany, Chelsea, Kelsey, Cody, Ashley, Courtney, Ryan, Kyle, Megan, Jessica
1995 Taylor, Kelley, Dakota, Austin, Haley, Cody, Tyler, Shelby, Brittany, Kayla
2000 Destiny, Madison, Hailey, Sydney, Alexis, Kaitlyn, Hunter, Brianna, Hannah, Alyssa
2005 Aislin, Dkya, Guisli, Hailey, Ethan, Madison, Ava, Isabella, Jayden, Aiden
2010 Jayden, Aislin, Noelle, Addison, Braxton, London, Peyton, Isabella, Ava, Liam
2015 Arias, Harper, Scarlett, Jason, Grayson, Alexander, Hudson, Liam, Diego, Layla

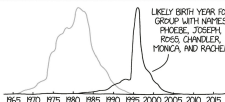
If kids in your class were named Jeff, Lisa, Michael, Karan, and David, then you were probably born in the mid-1960s. If they were named Jayden, Isabella, Sophia, Ava, and Ethan, then you were probably born somewhere around 2010. But names can reveal things about age in other ways.

The mid-1990s TV show *Friends* featured six roommates, played by actors, named Matthew, Jennifer, Courtney, Lisa, David, and another Matthew. Each of those names has its own popularity curve. If we combine them all, we can guess what years the group of actors was likely born:



The actors were actually born in the late 1960s, on the very early edge of the popularity of their names. In other words, the actors all have names that were a little before their time. Courtney Cox and Jennifer Aniston had names that didn't really become popular until a decade later. (Maybe people with trendy parents are more likely to wind up in acting.) But the names are generally consistent with their era, if a little ahead of the curve.

We get something very different if we look at the names of their characters—Phoebe, Joseph, Ross, Chandler, Rachel, and Monica:



The show debuted in 1994. There's a clear spike in popularity of the names in 1995 and 1996, which can probably be attributed to the show putting the names in the minds of new parents. But it's not just the show—that name combination was clearly on the rise in the years before *Friends* premiered. It's possible that parents looking for good names for their children are influenced by some of the same cultural trends as TV writers looking for good names for their characters.

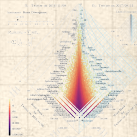
A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



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

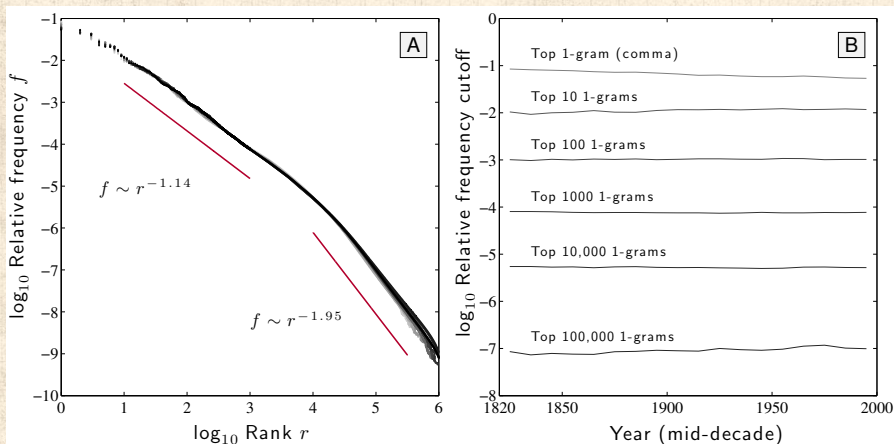




"Is language evolution grinding to a halt? The scaling of lexical turbulence in English fiction suggests it is not" ↗

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:

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

Estimates: $\mu \simeq 0.77$ and $\mu' \simeq 1.10$, 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 $\nu \simeq 1.23$ and $\nu' \simeq 1.47$.

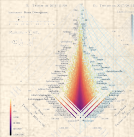
A plentitude of distances

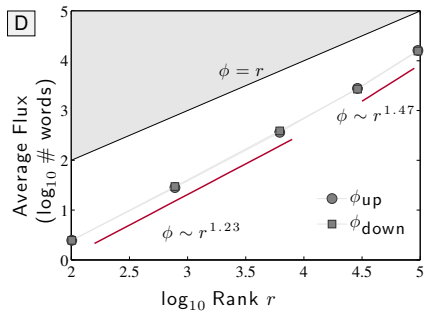
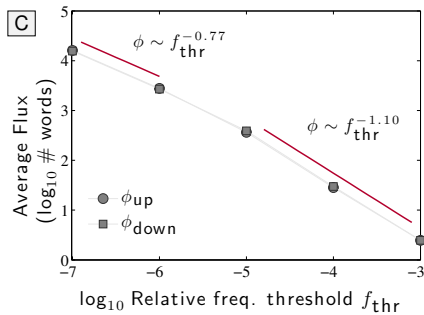
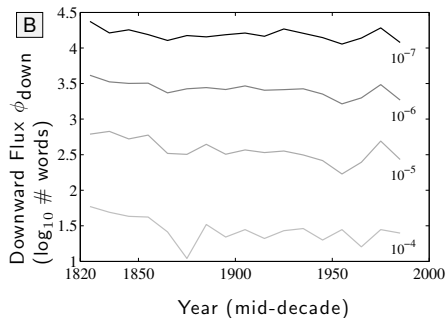
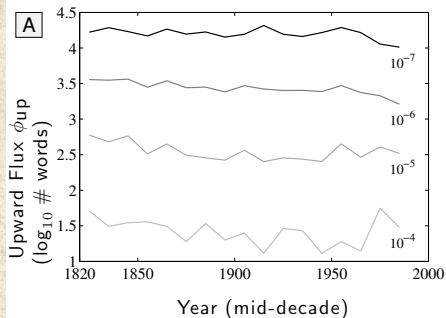
Rank-turbulence divergence

Probability-turbulence divergence

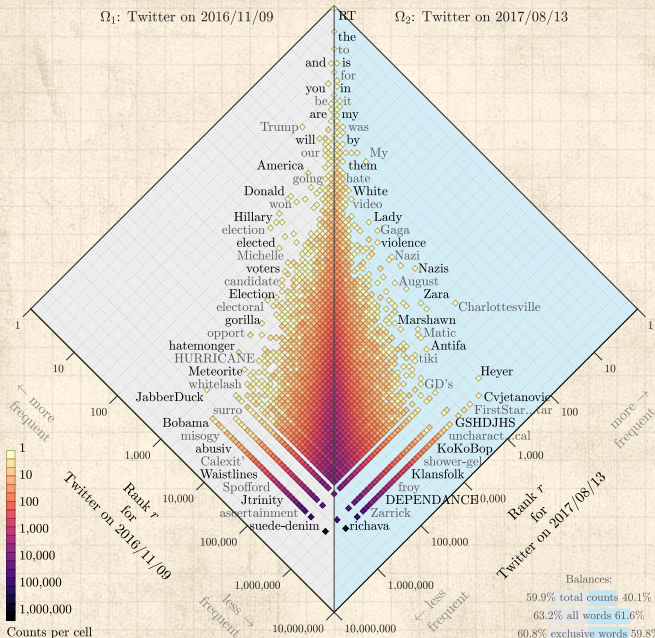
Explorations

References

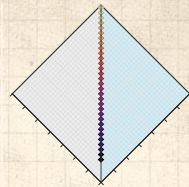




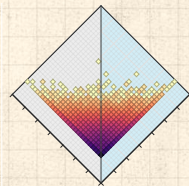
A. Rank-turbulence histogram:



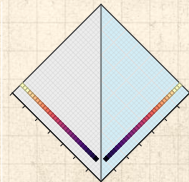
B. Identical systems:



C. Randomized systems:



D. Disjoint systems:



Balances:
59.9% total counts 40.1%
63.2% all words 61.6%
60.8% exclusive words 59.8%

So, so many ways to compare probability distributions:

PoCS
@pocsvox
Allotaxonomy



“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]

A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



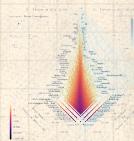
Comparisons are distances, divergences, similarities, inner products, fidelities ...



A worry: Subsampled distributions with very heavy tails



60ish kinds of comparisons grouped into 10 families



Quite the festival:

Table 1. L₁ Minkowski family

1. Euclidean L ₁	$d_{min} = \sqrt{\sum_{i=1}^n P_i - Q_i ^2}$ (1)
2. City block L ₁	$d_{min} = \sum_{i=1}^n P_i - Q_i $ (2)
3. Minkowski L _p	$d_{min} = \sqrt[p]{\sum_{i=1}^n P_i - Q_i ^p}$ (3)
4. Chebyshev L _∞	$d_{min} = \max_i P_i - Q_i $ (4)

Table 2. L₁ family

5. Sorenson	$d_{min} = \frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$ (5)
-------------	---

6. Gower	$d_{min} = \frac{1}{d} \sqrt{\sum_{i=1}^n \frac{ P_i - Q_i }{R_i}}$ (6)
	$\frac{1}{d} \sum_{i=1}^n P_i - Q_i $ (7)

7. Soregol	$d_{min} = \frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n \min(P_i, Q_i)}$ (8)
------------	--

8. Kulczyński d	$d_{min} = \frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n \min(P_i, Q_i)}$ (9)
-----------------	--

9. Canberra	$d_{min} = \sqrt[n]{\frac{\sum_{i=1}^n P_i - Q_i }{P_i + Q_i}}$ (10)
-------------	---

10. Lovaszian	$d_{min} = \sum_{i=1}^n \min(P_i - Q_i)$ (11)
---------------	---

* L₁ family ⇒ Intersection (13), Wave Hedges (15), Czekanowski (16), Ruszka (21), Tanimoto (23), etc.

Table 3. Intersection family

11. Intersection	$s_{in} = \sum_{i=1}^n \min(P_i, Q_i)$ (12)
	$d_{min} = 1 - s_{in} = \frac{1}{n} \sum_{i=1}^n P_i - Q_i $ (13)
12. Wave Hedges	$d_{min} = \sum_{i=1}^n \frac{\min(P_i, Q_i)}{\max(P_i, Q_i)}$ (14)
	$\frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n \max(P_i, Q_i)}$ (15)
13. Czekanowski	$s_{in} = \frac{\sum_{i=1}^n \min(P_i, Q_i)}{\sum_{i=1}^n (P_i + Q_i)}$ (16)
	$d_{min} = 1 - s_{in} = \frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$ (17)

14. Motyka	$s_{in} = \frac{\sum_{i=1}^n \min(P_i, Q_i)}{\sum_{i=1}^n (P_i + Q_i)}$ (18)
------------	--

	$d_{min} = 1 - s_{in} = \frac{\sum_{i=1}^n P_i - Q_i }{\sum_{i=1}^n (P_i + Q_i)}$ (19)
--	---

15. Kulczyński s	$s_{in} = \frac{1}{d_{min}} \frac{\sum_{i=1}^n \min(P_i, Q_i)}{\sum_{i=1}^n (P_i + Q_i)}$ (20)
------------------	--

16. Ruszka	$s_{in} = \frac{\sum_{i=1}^n \min(P_i, Q_i)}{\sum_{i=1}^n \max(P_i, Q_i)}$ (21)
------------	---

17. Tanimoto	$d_{min} = \frac{\sum_{i=1}^n P_i - Q_i + 2 \sum_{i=1}^n \min(P_i, Q_i)}{\sum_{i=1}^n P_i - Q_i + 2 \sum_{i=1}^n \max(P_i, Q_i)}$ (22)
--------------	--

	$\frac{\sum_{i=1}^n \min(P_i, Q_i)}{\sum_{i=1}^n \max(P_i, Q_i)}$ (23)
--	--

Table 4. Inner Product family

18. Inner Product	$s_{in} = P \cdot Q = \sum_{i=1}^n P_i Q_i$ (24)
-------------------	--

19. Harmonic mean	$s_{in} = \frac{\sum_{i=1}^n 2P_i Q_i}{\sum_{i=1}^n (P_i + Q_i)}$ (25)
-------------------	--

20. Cosine	$s_{in} = \frac{\sum_{i=1}^n P_i Q_i}{\sqrt{\sum_{i=1}^n P_i^2} \sqrt{\sum_{i=1}^n Q_i^2}}$ (26)
------------	--

21. Kumar-Hauschok (PCE)	$s_{in} = \frac{\sum_{i=1}^n P_i Q_i}{\sum_{i=1}^n P_i^2 + \sum_{i=1}^n Q_i^2 - \sum_{i=1}^n P_i Q_i}$ (27)
--------------------------	---

22. Jaccard	$s_{in} = \frac{\sum_{i=1}^n P_i Q_i}{\sum_{i=1}^n P_i^2 + \sum_{i=1}^n Q_i^2 - \sum_{i=1}^n P_i Q_i}$ (28)
-------------	---

	$d_{min} = 1 - s_{in} = \frac{\sum_{i=1}^n P_i - Q_i ^2}{\sum_{i=1}^n P_i^2 + \sum_{i=1}^n Q_i^2 - \sum_{i=1}^n P_i Q_i}$ (29)
--	---

23. Dice	$s_{in} = \frac{2 \sum_{i=1}^n P_i Q_i}{\sum_{i=1}^n P_i^2 + \sum_{i=1}^n Q_i^2}$ (30)
----------	--

	$d_{min} = 1 - s_{in} = \frac{\sum_{i=1}^n P_i - Q_i ^2}{\sum_{i=1}^n P_i^2 + \sum_{i=1}^n Q_i^2}$ (31)
--	--

Table 5. Fidelity family or Squared-chord family

24. Fidelity	$s_{in} = \sum_{i=1}^n \sqrt{P_i Q_i}$ (32)
--------------	---

25. Bhattacharyya	$d_{in} = -\ln \sum_{i=1}^n \sqrt{P_i Q_i}$ (33)
-------------------	--

26. Hellinger	$d_{in} = \sqrt{\sum_{i=1}^n (\sqrt{P_i} - \sqrt{Q_i})^2}$ (34)
	$-2 \sqrt{\sum_{i=1}^n P_i Q_i}$ (35)

27. Matusita	$d_{in} = \sqrt{\sum_{i=1}^n (\sqrt{P_i} - \sqrt{Q_i})^2}$ (36)
	$= 2 \sqrt{\sum_{i=1}^n P_i Q_i}$ (37)

28. Squared-chord	$d_{in} = \sum_{i=1}^n \sqrt{ P_i - Q_i }$ (38)
$s_{in} = 1 - d_{in}$	$s_{in} = 2 \sum_{i=1}^n \sqrt{P_i Q_i} - 1$ (39)

Table 6. Squared L₁ family or χ² family

29. Squared Euclidean	$d_{in} = \sum_{i=1}^n P_i - Q_i ^2$ (40)
-----------------------	--

30. Pearson χ ²	$d_{in}(P, Q) = \sum_{i=1}^n \frac{(P_i - Q_i)^2}{Q_i}$ (41)
----------------------------	--

31. Neyman χ ²	$d_{in}(P, Q) = \sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i}$ (42)
---------------------------	--

32. Squared χ ²	$d_{in} = \sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i + Q_i}$ (43)
----------------------------	--

33. Probabilistic Symmetric χ ²	$d_{in} = \sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i + Q_i}$ (44)
--	--

34. Divergence	$d_{in} = 2 \sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i + Q_i}$ (45)
----------------	--

35. Clark	$d_{in} = \sqrt{\sum_{i=1}^n \frac{ P_i - Q_i }{P_i + Q_i}}$ (46)
-----------	---

36. Additive Symmetric χ ²	$d_{in} = \sum_{i=1}^n \frac{(P_i - Q_i)^2 (P_i + Q_i)}{P_i Q_i}$ (47)
---------------------------------------	--

* Squared L₁ family ⇒ Jaccard (29), Dice (31)

37. Kullback-Leibler	$d_{in} = \sum_{i=1}^n P_i \ln \frac{P_i}{Q_i}$ (48)
----------------------	--

38. Jeffreys	$d_{in} = \sum_{i=1}^n (P_i - Q_i) \ln \frac{P_i}{Q_i}$ (49)
--------------	--

39. K. divergence	$d_{in} = \sum_{i=1}^n P_i \ln \frac{2P_i}{P_i + Q_i}$ (50)
-------------------	---

40. Topoc	$d_{in} = \sum_{i=1}^n P_i \left[\ln \left(\frac{2P_i}{P_i + Q_i} \right) - Q_i \ln \left(\frac{2Q_i}{P_i + Q_i} \right) \right]$ (51)
-----------	---

41. Jensen-Shannon	$d_{in} = \frac{1}{2} \sum_{i=1}^n P_i \left[\ln \left(\frac{2P_i}{P_i + Q_i} \right) \frac{Q_i}{P_i} + \ln \left(\frac{2Q_i}{P_i + Q_i} \right) \frac{P_i}{Q_i} \right]$ (52)
--------------------	---

42. Jensen divergence	$d_{in} = \sum_{i=1}^n \left[\frac{P_i \ln P_i + Q_i \ln Q_i}{2} - \left(\frac{P_i + Q_i}{2} \right) \ln \left(\frac{P_i + Q_i}{2} \right) \right]$ (53)
-----------------------	---

Table 8. Combinations

43. Taneja	$d_{in} = \sum_{i=1}^n \frac{ P_i - Q_i }{2} \ln \left \frac{P_i + Q_i}{2 \sqrt{P_i Q_i}} \right $ (54)
------------	--

44. Kumar-Johnson	$d_{in} = \sum_{i=1}^n \left[\frac{(P_i - Q_i)^2}{2(P_i Q_i)^2} \right]$ (55)
-------------------	--

45. Avg(L ₁ , L _∞)	$d_{in} = \frac{\sum_{i=1}^n P_i - Q_i + \max_i P_i - Q_i }{2}$ (56)
---	---

Table 10. Vicissitude

Vicis-Wave Hedges	$d_{min} = \sum_{i=1}^n \frac{ P_i - Q_i }{\max(P_i, Q_i)}$ (60)
-------------------	--

Vicis-Symmetric χ ²	$d_{min} = \sum_{i=1}^n \frac{ P_i - Q_i ^2}{\max(P_i, Q_i)^2}$ (61)
--------------------------------	--

Vicis-Symmetric χ ²	$d_{min} = \sum_{i=1}^n \frac{ P_i - Q_i ^2}{\max(P_i, Q_i)}$ (62)
--------------------------------	--

Vicis-Symmetric χ ²	$d_{min} = \sum_{i=1}^n \frac{ P_i - Q_i ^2}{\max(P_i, Q_i)}$ (63)
--------------------------------	--

max-Symmetric	$d_{in} = \max \left(\sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i}, \sum_{i=1}^n \frac{(P_i - Q_i)^2}{Q_i} \right)$ (64)
---------------	--

min-symmetric	$d_{in} = \min \left(\sum_{i=1}^n \frac{(P_i - Q_i)^2}{P_i}, \sum_{i=1}^n \frac{(P_i - Q_i)^2}{Q_i} \right)$ (65)
---------------	--

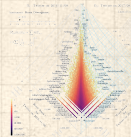
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References





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	$d_{Euc} = \sqrt{\sum_{i=1}^d P_i - Q_i ^2}$	(1)
--------------------	---	-----

2. City block L_1	$d_{CB} = \sum_{i=1}^d P_i - Q_i $	(2)
---------------------	-------------------------------------	-----

3. Minkowski L_p	$d_{Mk} = \sqrt[p]{\sum_{i=1}^d P_i - Q_i ^p}$	(3)
--------------------	---	-----

4. Chebyshev L_∞	$d_{Cheb} = \max_i P_i - Q_i $	(4)
-------------------------	---------------------------------	-----

Table 2. L_1 family

5. Sørensen	$d_{sor} = \frac{\sum_{i=1}^d P_i - Q_i }{\sum_{i=1}^d (P_i + Q_i)}$	(5)
-------------	---	-----

6. Gower	$d_{gow} = \frac{1}{d} \sum_{i=1}^d \frac{ P_i - Q_i }{R_i}$	(6)
----------	--	-----

	$= \frac{1}{d} \sum_{i=1}^d P_i - Q_i $	(7)
--	--	-----

7. Soergel	$d_{sg} = \frac{\sum_{i=1}^d P_i - Q_i }{\sum_{i=1}^d \max(P_i, Q_i)}$	(8)
------------	---	-----

8. Kulczynski d	$d_{kul} = \frac{\sum_{i=1}^d P_i - Q_i }{\sum_{i=1}^d \min(P_i, Q_i)}$	(9)
-------------------	--	-----

9. Canberra	$d_{can} = \sum_{i=1}^d \frac{ P_i - Q_i }{P_i + Q_i}$	(10)
-------------	--	------

10. Lorentzian	$d_{lor} = \sum_{i=1}^d \ln(1 + P_i - Q_i)$	(11)
----------------	---	------

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

PoCS

@pocsvox

Allotaxonomy

A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References

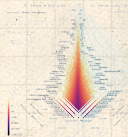


Table 1. L_p Minkowski family

1. Euclidean L_2	$d_{Euc} = \sqrt{\sum_{i=1}^d P_i - Q_i ^2}$	(1)
--------------------	---	-----

2. City block L_1	$d_{CB} = \sum_{i=1}^d P_i - Q_i $	(2)
---------------------	-------------------------------------	-----

3. Minkowski L_p	$d_{Mk} = \sqrt[p]{\sum_{i=1}^d P_i - Q_i ^p}$	(3)
--------------------	---	-----

4. Chebyshev L_{∞}	$d_{Cheb} = \max_i P_i - Q_i $	(4)
---------------------------	---------------------------------	-----

Table 2. L_1 family

5. Sørensen	$d_{sor} = \frac{\sum_{i=1}^d P_i - Q_i }{\sum_{i=1}^d (P_i + Q_i)}$	(5)
-------------	---	-----

6. Gower	$d_{gow} = \frac{1}{d} \sum_{i=1}^d \frac{ P_i - Q_i }{R_i}$	(6)
----------	--	-----

	$= \frac{1}{d} \sum_{i=1}^d P_i - Q_i $	(7)
--	--	-----

7. Soergel	$d_{sg} = \frac{\sum_{i=1}^d P_i - Q_i }{\sum_{i=1}^d \max(P_i, Q_i)}$	(8)
------------	---	-----

8. Kulczynski d	$d_{kul} = \frac{\sum_{i=1}^d P_i - Q_i }{\sum_{i=1}^d \min(P_i, Q_i)}$	(9)
-------------------	--	-----

9. Canberra	$d_{can} = \sum_{i=1}^d \frac{ P_i - Q_i }{P_i + Q_i}$	(10)
-------------	--	------

10. Lorentzian	$d_{Lor} = \sum_{i=1}^d \ln(1 + P_i - Q_i)$	(11)
----------------	---	------

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

PoCS
@pocsvox
Allotaxonomy

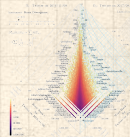
A plentitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



Information theoretic sortings are more opaque


No tunability


Shannon's Entropy:


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

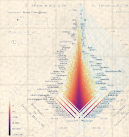
Kullback-Liebler (KL) divergence:

$$\begin{aligned} D^{\text{KL}}(P_2 \parallel P_1) &= \left\langle \log_2 \frac{1}{p_{2,\tau}} - \log_2 \frac{1}{p_{1,\tau}} \right\rangle_{P_2} \\ &= \sum_{\tau \in R_{1,2;\alpha}} p_{2,\tau} \left[\log_2 \frac{1}{p_{2,\tau}} - \log_2 \frac{1}{p_{1,\tau}} \right] \\ &= \sum_{\tau \in R_{1,2;\alpha}} p_{2,\tau} \log_2 \frac{p_{1,\tau}}{p_{2,\tau}}. \end{aligned} \quad (2)$$

 Problem: If just one component type in system 2 is not present in system 1, KL divergence = ∞ .

 Solution: If we can't compare a spork and a platypus directly, we create a fictional **spork-platypus hybrid**.

 New problem: Re-read solution.



🗑️ Jensen-Shannon divergence (JSD): [9, 7, 13, 1]

$$\begin{aligned} D^{\text{JS}}(P_1 \parallel P_2) &= \frac{1}{2} D^{\text{KL}}\left(P_1 \parallel \frac{1}{2}[P_1 + P_2]\right) + \frac{1}{2} D^{\text{KL}}\left(P_2 \parallel \frac{1}{2}[P_1 + P_2]\right) \\ &= \frac{1}{2} \sum_{\tau \in R_{1,2;\alpha}} \left(p_{1,\tau} \log_2 \frac{p_{1,\tau}}{\frac{1}{2}[p_{1,\tau} + p_{2,\tau}]} + p_{2,\tau} \log_2 \frac{p_{2,\tau}}{\frac{1}{2}[p_{1,\tau} + p_{2,\tau}]} \right). \end{aligned} \quad (3)$$

A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

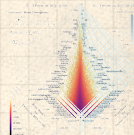
Explorations

References

🗑️ Involving a third intermediate averaged system means JSD is now finite: $0 \leq D^{\text{JS}}(P_1 \parallel P_2) \leq 1$.

🗑️ Generalized entropy divergence: [2]

$$\begin{aligned} D_{\alpha}^{\text{AS2}}(P_1 \parallel P_2) &= \\ \frac{1}{\alpha(\alpha-1)} \sum_{\tau \in R_{1,2;\alpha}} \left[(p_{\tau,1}^{1-\alpha} + p_{\tau,2}^{1-\alpha}) \left(\frac{p_{\tau,1} + p_{\tau,2}}{2} \right)^{\alpha} - (p_{\tau,1} + p_{\tau,2}) \right]. \end{aligned} \quad (4)$$



Produces JSD when $\alpha \rightarrow 0$.



Ω_1 : Twitter on 2016/11/09

Ω_2 : Twitter on 2017/08/13

Divergence contribution $\delta D_{0,r}^H$ (%)

Instrument: Sym. Gen. Entropy Div.

$\alpha=0$ (Jenson-Shannon Divergence)

$-2 \quad -1 \quad -1/2 \quad -1/4 \quad 0 \quad 1/4 \quad 1/2 \quad 3/4 \quad 1$

$$D_{0,r}^H(\Omega_1 || \Omega_2) = \sum \delta D_{0,r}^H$$

$$= \frac{1}{2} \sum_r \left[p_r^{(1)} \ln \frac{2p_r^{(1)}}{p_r^{(1)} + p_r^{(2)}} \right]$$

$$+ p_r^{(2)} \ln \frac{2p_r^{(2)}}{p_r^{(1)} + p_r^{(2)}}$$

$$= D^{JS}(\Omega_1 || \Omega_2)$$

10^{-1}

10^{-2}

10^{-3}

10^{-4}

10^{-5}

10^{-6}

10^{-7}

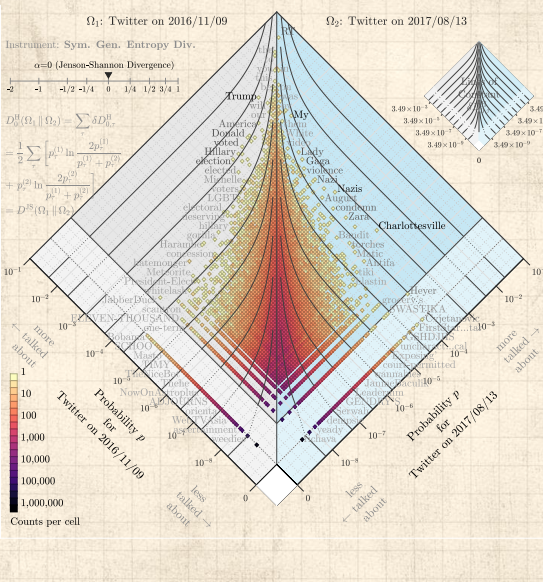
10^{-8}

10^{-9}

10^{-10}

10^{-11}

10^{-12}

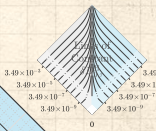


Counts per cell

more talked about
less talked about

Probability p for Twitter on 2016/11/09

Probability p for Twitter on 2017/08/13



Counts per cell

more talked about
less talked about

Probability p for Twitter on 2016/11/09

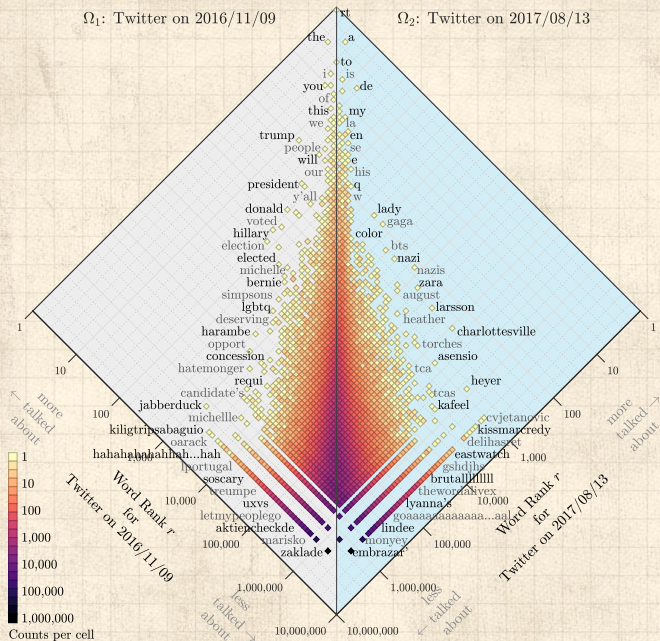
Probability p for Twitter on 2017/08/13

1.5 1 0.5 0 0.5 1 1.5

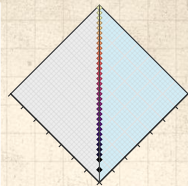
- Trump 11=60
- voted 58=1,002
- Donald 50=566
- election 64=2,055
- president 48=500
- Hillary 70=1,505
- trump 77=1,357
- America 40=164
- won 69=536
- 67,220=113 Charlotteville
- 139=20 My
- 9,149=129 Nazis
- Clinton 125=1,761
- Obama 76=378
- elected 151=2,787
- wins 144=1,209
- will 23=51
- country 71=216
- 5,873=171 supremacists
- 1,175=124 Gaga
- 3,485=174 Nazi
- 1=1 RT
- 86=27 his
- 801=119 Lady
- votes 180=1,422
- 3,563=192 BTS
- 37,952=268 Larsson
- 25,126=267 Zara
- 13,329=280 condemn
- 1,671=170 violence
- Michelle 261=3,115
- our 41=72
- 7,911=321 August
- President 93=228
- voters 306=4,453
- 1,325=187 supremacy
- people 27=45
- candidate 362=5,584
- 1,761=231 police
- women 124=315

52.9%—47.1%

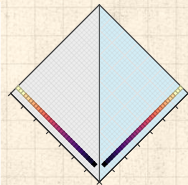
A. Rank-turbulence histogram:



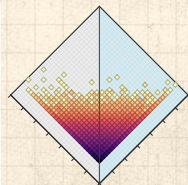
B. Identical systems:



C. Disjoint systems:



D. Randomized systems:



A plenitude of
distances

Rank-turbulence
divergence

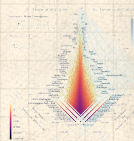
Probability-
turbulence
divergence

Explorations

References

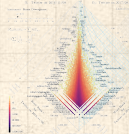
Exclusive types:

- 🧩 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:

1. Rank-based.
2. Symmetric.
3. Semi-positive: $D_{\alpha}^R(\Omega_1 || \Omega_2) \geq 0$.
4. Linearly separable, for interpretability.
5. 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.).
6. Zipfophilic: Able to handle systems with rank-ordered component size distribution that are heavy-tailed.
7. Scalable: Allow for sensible comparisons across system sizes.
8. Tunable.
9. Story-finding: Features 1–8 combine to show which component types are most 'important'



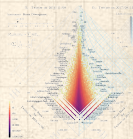
Some good things about ranks:

- 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.

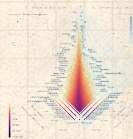
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



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 ∞ 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.

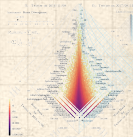
A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References



Some reworking:

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

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

Rank-turbulence divergence:

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

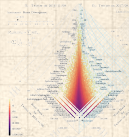
A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References



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: $0 \leq D_{\alpha}^R(R_1 \parallel R_2) \leq 1$
- Limits of 0 and 1 correspond to the two systems having identical and disjoint Zipf distributions.

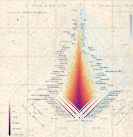
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



Rank-turbulence divergence:

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

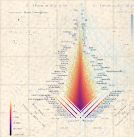
$$D_{\alpha}^R(R_1 || R_2) = \frac{1}{\mathcal{N}_{1,2;\alpha}} \frac{\alpha + 1}{\alpha} \sum_{\tau \in R_{1,2;\alpha}} \left| \frac{1}{[r_{\tau,1}]^{\alpha}} - \frac{1}{[r_{\tau,2}]^{\alpha}} \right|^{1/(\alpha+1)} \quad (10)$$

A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

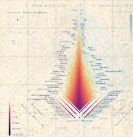
Explorations
References



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:

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




Limit of $\alpha \rightarrow 0$:

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

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.

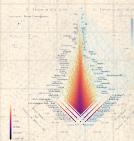
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References




Limit of $\alpha \rightarrow \infty$:

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

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.

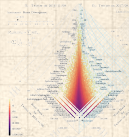
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence



Explorations

References



Probability-turbulence divergence:

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

-  For the unnormalized version ($\mathcal{N}_{1,2;\alpha}^{\text{P}}=1$), some troubles return with 0 probabilities and $\alpha \rightarrow 0$.
-  Weep not: $\mathcal{N}_{1,2;\alpha}^{\text{P}}$ 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:

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

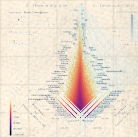
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



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 $1/\alpha$.

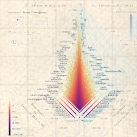
A plenitude of
distances

Rank-turbulence
divergence


Probability-turbulence
divergence

Explorations


References



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.

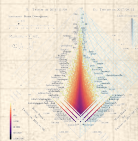
A plenitude of
distances

Rank-turbulence
divergence

Probability-turbulence
divergence



Explorations

References



Combine these cases into a single expression:

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

-  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,

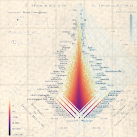
A plenitude of distances

Rank-turbulence divergence




Probability-turbulence divergence

Explorations

References



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 $D_0^P(P_1 \parallel P_2)$, 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.

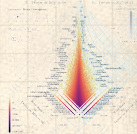
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



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

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

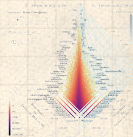
A plenitude of distances

Rank-turbulence divergence





Probability-turbulence divergence

Explorations

References



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].

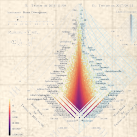
A plenitude of distances

Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



Ω_1 : Twitter on 2016/11/09

Ω_2 : Twitter on 2017/08/13

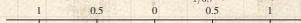
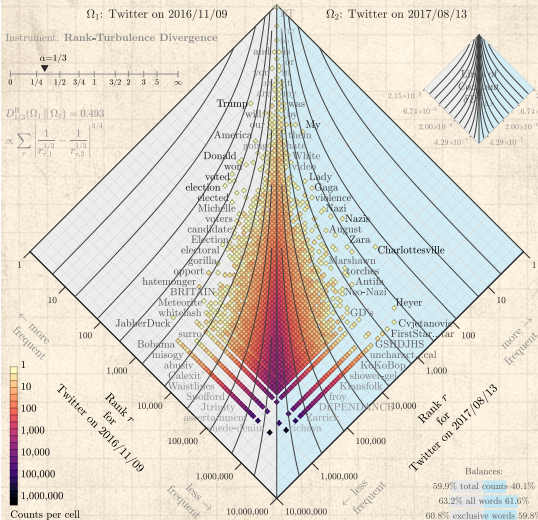
Divergence contribution $\delta D_{1/3,7}^R$ ($\times 10^{-3}\%$)

Instrument: Rank-Turbulence Divergence

$\alpha=1/3$

$$D_{1/3}^R(\Omega_1 || \Omega_2) = 0.493$$

$$\propto \sum_r \left| \frac{1}{r_{-1/3}} - \frac{1}{r_{+1/3}} \right|^{3/4}$$



Trump	11=60
17,220=113	Charlotteville
election	64=2,055
voted	58=1,002
Hillary	70=1,505
Donald	50=566
9,149=129	Nazis
president	48=500
trump	77=1,357
139=20	My
37,952=268	Larsson
5,873=171	supremacists
25,126=267	Zara
won	69=536
862,482=443	Heyer
Clinton	125=1,761
elected	151=2,787
3,485=174	Nazi
13,329=280	condemn
America	40=164
3,503=192	BTS
86=27	his
1,175=124	gaga
1,562,865=673	Cvjetanovic
Obama	76=378
wins	144=1,209
1,671=170	violence
7,911=321	August
801=110	Lady
18,804=442	nazis
16,317=140	Heather
electoral	446=15,272
47,558=610	torches
JabberDuck	993=840,082
hatemonger	756=120,186
102,414=743	Antifa
opport	672=62,603
56,244=657	neo-nazis
Harambe	658=56,005
Michelle	261=3,115
	50.2%—49.8%

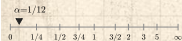
Balances:
 59.9% total counts 40.1%
 63.2% all words 61.6%
 60.8% exclusive words 59.8%

Ω_1 : Twitter on 2016/11/09

Ω_2 : Twitter on 2017/08/13

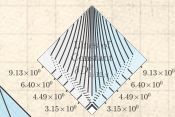
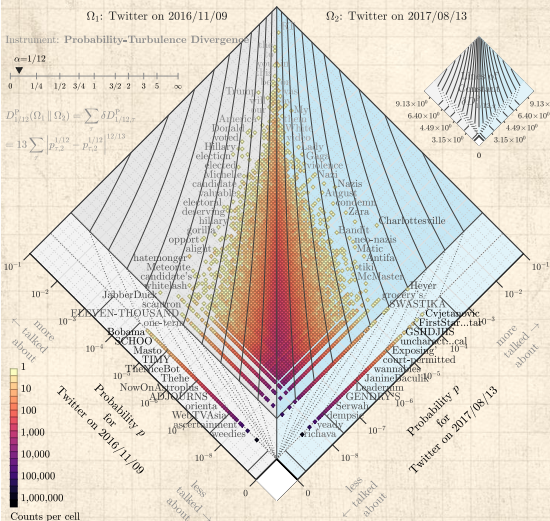
Divergence contribution $\delta D_{1/12,r}^D (\times 10^{-4}\%)$

Instrument: Probability-Turbulence Divergence



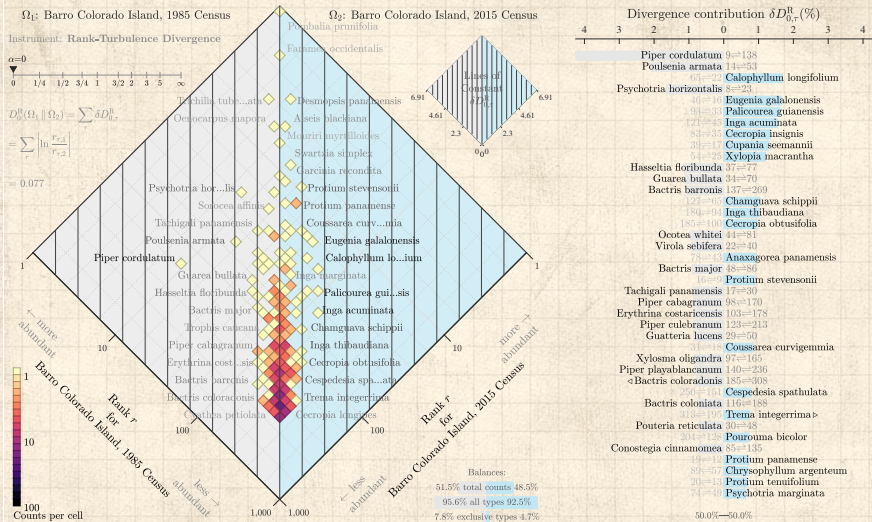
$$D_{1/12}^D(\Omega_1 \parallel \Omega_2) = \sum \delta D_{1/12,r}^D$$

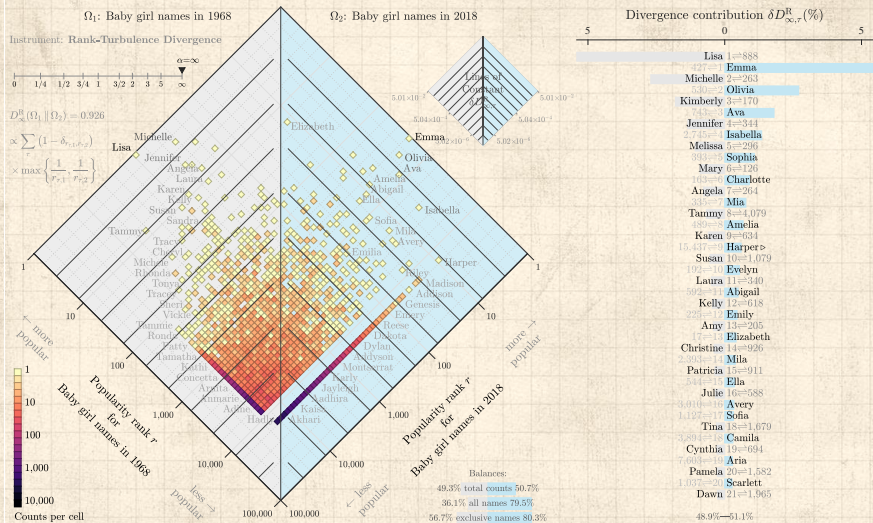
$$= 13 \sum \left[\frac{1/12}{P_{r,2}} - \frac{1/12}{P_{r,1}} \right]$$

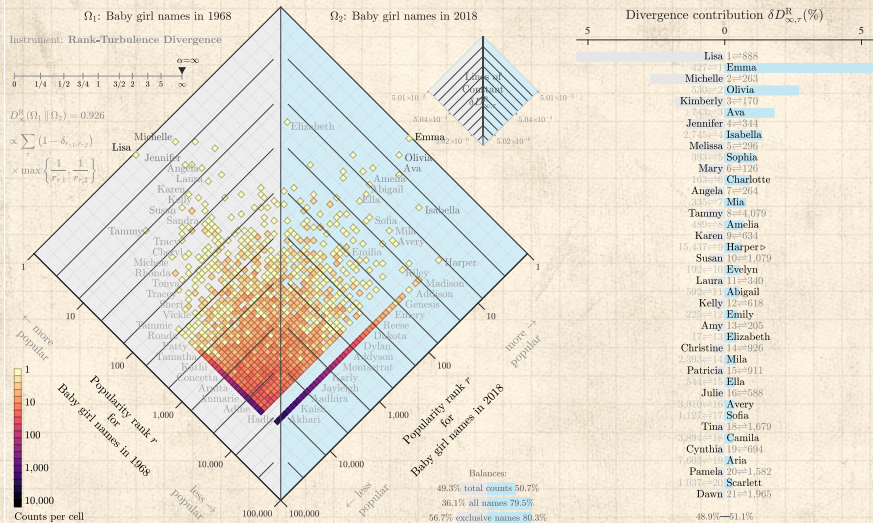


1	0	1
1.552,865=6.73		Cvjetanovic >
1.552,865=1.116		FirstStarMagicAllStar >
1.552,865=1.47		KISSMARCHED >
1.552,865=1.520		ForAllStarGames >
1.552,865=1.985		Kafeel >
1.552,865=2.021		Starbz >
		< Bobama 2,423=1,537,471
		< Oarack 2,425=1,537,471
		< Un-Leashed 2,703=1,537,471
1.552,865=3.088		GSHDJHS >
1.552,865=3.099		Bodak >
< KiligTripSaBagnio 3,142=1,537,471		
< Somali-American 3,229=1,537,471		
< DICKASS 3,321=1,537,471		
< Michelle 3,412=1,537,471		
1.552,865=3.673		Eastwatch >
< Un-leashed 3,645=1,537,471		
1.552,865=3.983		Heyer's >
< SCHOO 3,921=1,537,471		
1.552,865=4.382		uncharacteristical >
1.552,865=4.518		callejones >
		< misogy 4,328=1,537,471
1.552,865=4.723		TLC >
1.552,865=4.913		SORIBADA >
< tRyNna 4,660=1,537,471		
< aLmoSt 4,671=1,537,471		
1.552,865=5.240		tcas >
< Ruline 5,097=1,537,471		
< Steinger 5,118=1,537,471		
1.552,865=5.436		low-rise >
1.552,865=5.662		climate-denying 5,191=1,537,471
1.552,865=5.682		CLITORIS >
1.552,865=5.682		Adityanath >
< lambo's 5,383=1,537,471		
1.552,865=5.755		DelHiHasret >
1.552,865=5.755		FikBel >
1.552,865=5.808		Walker-Peters >
< KBAT 5,617=1,537,471		
1.552,865=6.040		UNIDAS >
< stammered 5,653=1,537,471		

49.9%—50.1%







Ω_1 : 1948 Google Books Fiction

Ω_2 : 1987 Google Books Fiction

Instrument: Rank-Turbulence Divergence

$$D_{\infty}^R(\Omega_1, \Omega_2) = 0.522$$

$$\infty \sum_{\tau} (1 - \delta_{\tau,1} \delta_{\tau,2})$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$D_{\infty}^R(\Omega_1, \Omega_2) = 0.522$$

$$\infty \sum_{\tau} (1 - \delta_{\tau,1} \delta_{\tau,2})$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

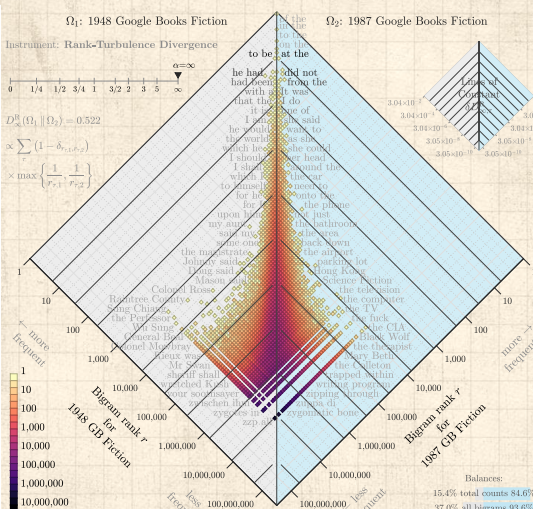
$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

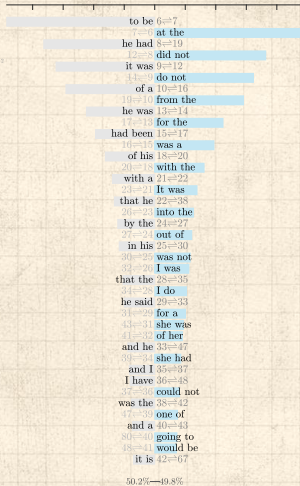
$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$

$$\times \max \left\{ \frac{1}{r_{\tau,1}}, \frac{1}{r_{\tau,2}} \right\}$$



Divergence contribution $\delta_{\infty,r}^R$ (%)

0.8 0.6 0.4 0.2 0 0.2 0.4 0.6 0.8



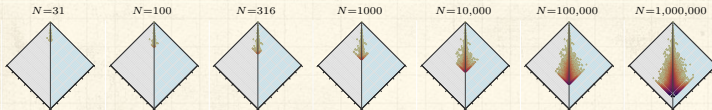
Balances:
 15.4% total counts 84.6%
 37.0% all bigrams 93.6%
 17.2% exclusive bigrams 67.3%

Effect of subsampling:

PoCS
@pocsvox

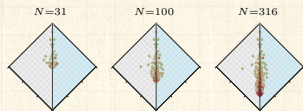
Allotaxonomy

Twitter:



A plenty of distances

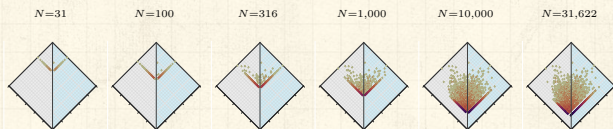
Tree species:



Rank-turbulence divergence

Probability-turbulence divergence

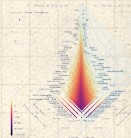
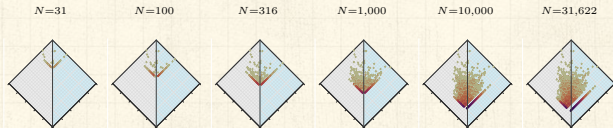
Baby girl names:



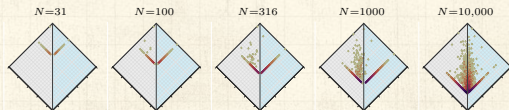
Explorations

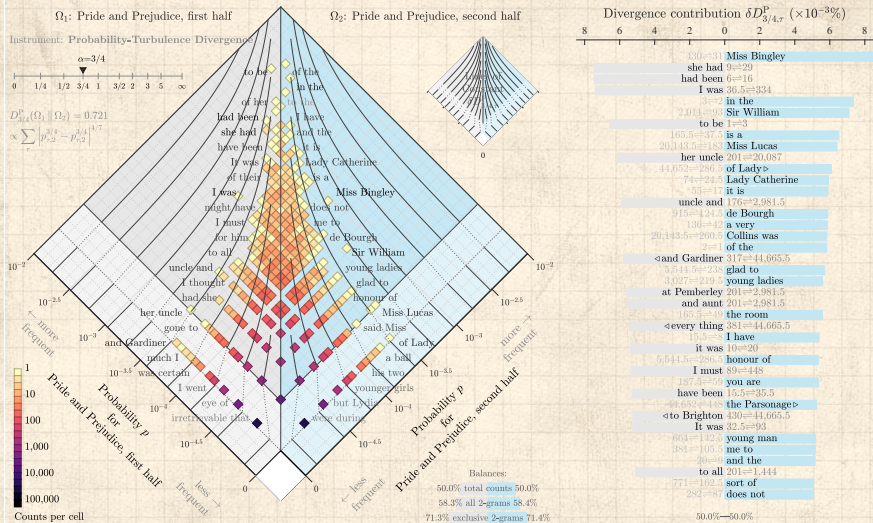
References

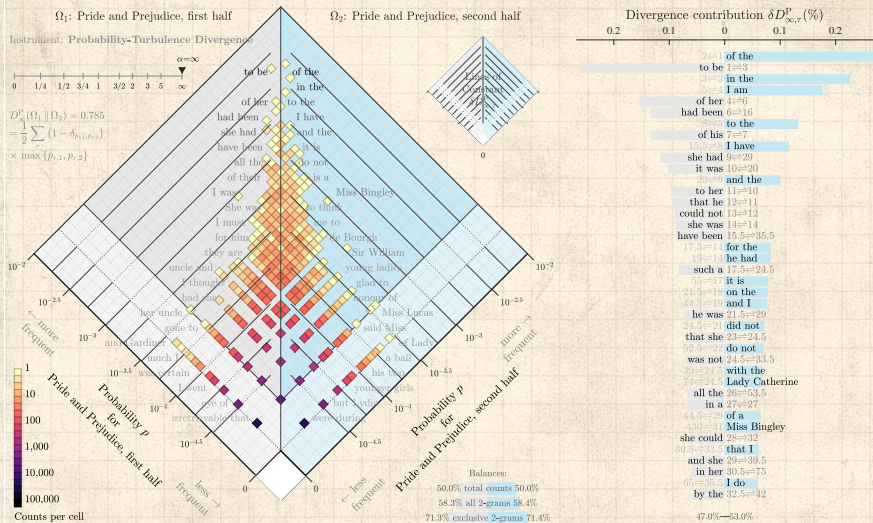
Baby boy names:



Market caps:



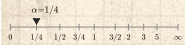




Ω_1 : Twitter on 2020/03/12

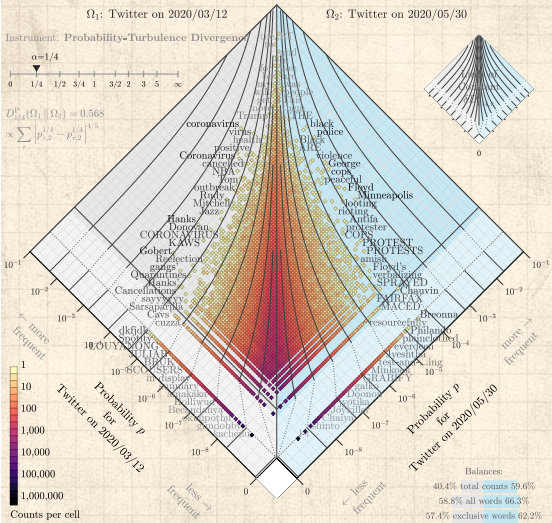
Ω_2 : Twitter on 2020/05/30

Instrument: Probability-Turbulence Divergence

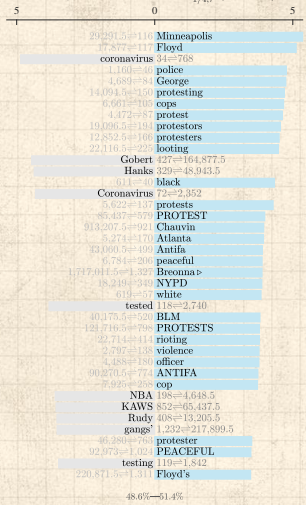


$$D_{1/4}^P(\Omega_1 || \Omega_2) = 0.568$$

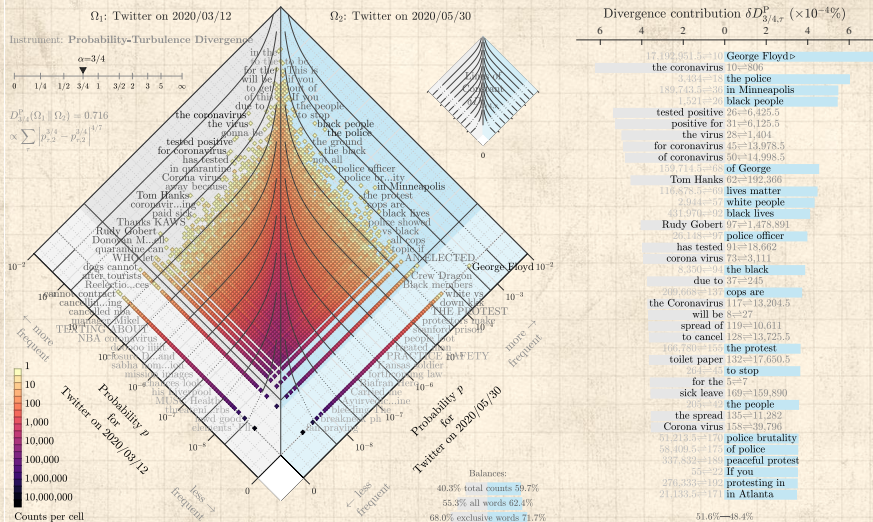
$$\propto \sum_p |p_{\Omega_1}^{1/4} - p_{\Omega_2}^{1/4}|^{1/5}$$



Divergence contribution $\delta D_{1/4,7}^P (\times 10^{-4}\%)$



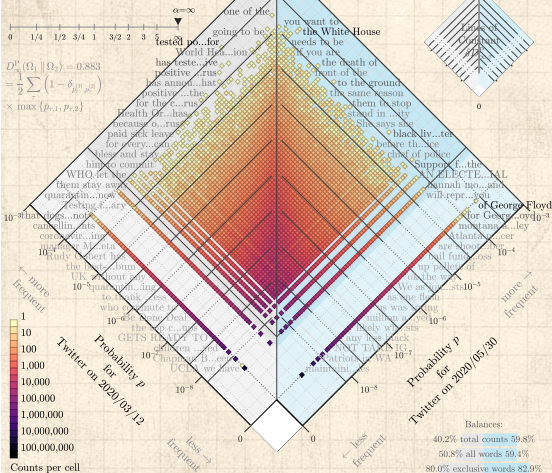
Balances:
 40.4% total counts 59.6%
 58.8% all words 66.3%
 57.4% exclusive words 62.2%



Ω_1 : Twitter on 2020/03/12

Ω_2 : Twitter on 2020/05/30

Instrument: Probability-Turbulence Divergence



$$D_{\infty}^p(\Omega_1, \Omega_2) = 0.883$$

$$= \frac{1}{2} \sum (1 - \delta_{\Omega_1, \Omega_2}^p)$$

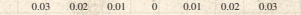
$$\times \max\{p_{r,1}, p_{r,2}\}$$

Counts per cell

1
10
100
1,000
10,000
100,000
1,000,000
10,000,000
100,000,000



Divergence contribution $\delta D_{\infty, r}^p$ (%)



- 1=4,975.5
- 2=219.0
- 3=11,879
- 4=14,798
- 5=7,264.5
- 6=33
- 7=108
- 8=1,420
- 9=78,795
- 10=53,912
- 11=603
- 12=22,783.5
- 13=45
- 14=143.5
- 15=30
- 16=277,424.5
- 17=631.5
- 18=43,073,107
- 19=22
- 20=43,073,107
- 21=172,568
- 22=1,421
- 23=43,073,107
- 24=43,073,107
- 25=10
- 26=7
- 27=108
- 28=131
- 29=9
- 30=11
- 31=12
- 32=13
- 33=10
- 34=11
- 35=12
- 36=13
- 37=10
- 38=11
- 39=12
- 40=13
- 41=10
- 42=11
- 43=12
- 44=13
- 45=10
- 46=11
- 47=12
- 48=13
- 49=10
- 50=11
- 51=12
- 52=13
- 53=10
- 54=11
- 55=12
- 56=13
- 57=10
- 58=11
- 59=12
- 60=13
- 61=10
- 62=11
- 63=12
- 64=13
- 65=10
- 66=11
- 67=12
- 68=13
- 69=10
- 70=11
- 71=12
- 72=13
- 73=10
- 74=11
- 75=12
- 76=13
- 77=10
- 78=11
- 79=12
- 80=13
- 81=10
- 82=11
- 83=12
- 84=13
- 85=10
- 86=11
- 87=12
- 88=13
- 89=10
- 90=11
- 91=12
- 92=13
- 93=10
- 94=11
- 95=12
- 96=13
- 97=10
- 98=11
- 99=12
- 100=13

Balances:
40.2% total counts 59.8%
50.8% all words 59.4%
80.0% exclusive words 82.9%

50.4%—49.6%

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-onigrams-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>

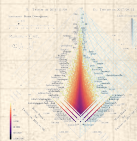
A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

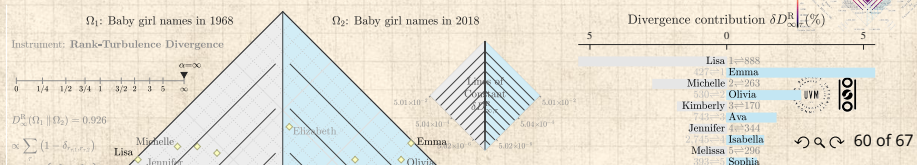
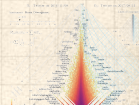
References



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)

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



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](#) 

A plenitude of
distances


Rank-turbulence
divergence

Probability-
turbulence
divergence

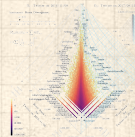
Explorations


References

- [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.



- [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 

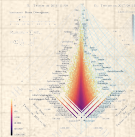
A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References



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 ↗

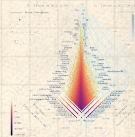
A plenitude of distances




Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



- [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](#) 

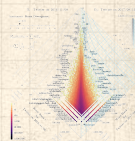
A plenitude of
distances


Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References



- [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. Osterreicher 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.

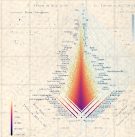
A plenitude of distances


Rank-turbulence divergence

Probability-turbulence divergence

Explorations

References



- [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.
[Videnski Selskab Biologiske Skrifter](#), 5:1–34, 1948.

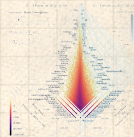
A plenitude of
distances


Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References



- [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](#) 

A plenitude of
distances

Rank-turbulence
divergence

Probability-
turbulence
divergence

Explorations

References

