

The spread of information on social media

Lília Perfeito, Pedro Duarte, Joana Gonçalves-Sá

Social Physics and Complexity @ LIP



FARE

Fake News and Real People: Using big data
to understand human behaviour



About me

2004-2008

PhD in Biology

The distribution of fitness effects of spontaneous mutations in bacteria



2008-2012

Collaboration between Institute for genetics and Institute for theoretical physics

Empirical molecular fitness landscapes



University of Cologne

2012-2018

Principal Investigator




Can we predict evolution?



2019-2020

Data Science and Policy



HEALTH		<p>Online vs. Offline Patterns Emergency Now-casting Antibiotic Over-prescription</p>	<p>Google Trends SNS24 Twitter ER acceptance /times SPMS e-prescriptions</p>	<p>Math Modelling ML Epidemiology</p>
POLICY		<p>Political Decisions Gender Differences Agenda Setting Voting vs. Discourse</p>	<p>Media records Twitter Facebook Parliament data</p>	<p>NLP Networks Math Modelling Complex Systems</p>
BEHAVIOU		<p>Cognitive Biases Attitudes Towards Science Tracking Anxiety</p>	<p>Large scale surveys Behavioral experiments Twitter Facebook</p>	<p>Networks Math Modelling Psychology Information</p>

The spread of information on social media: a model based on evolutionary principles

Lília Perfeito, Pedro Duarte, Joana Gonçalves-Sá

Social Physics and Complexity @ LIP



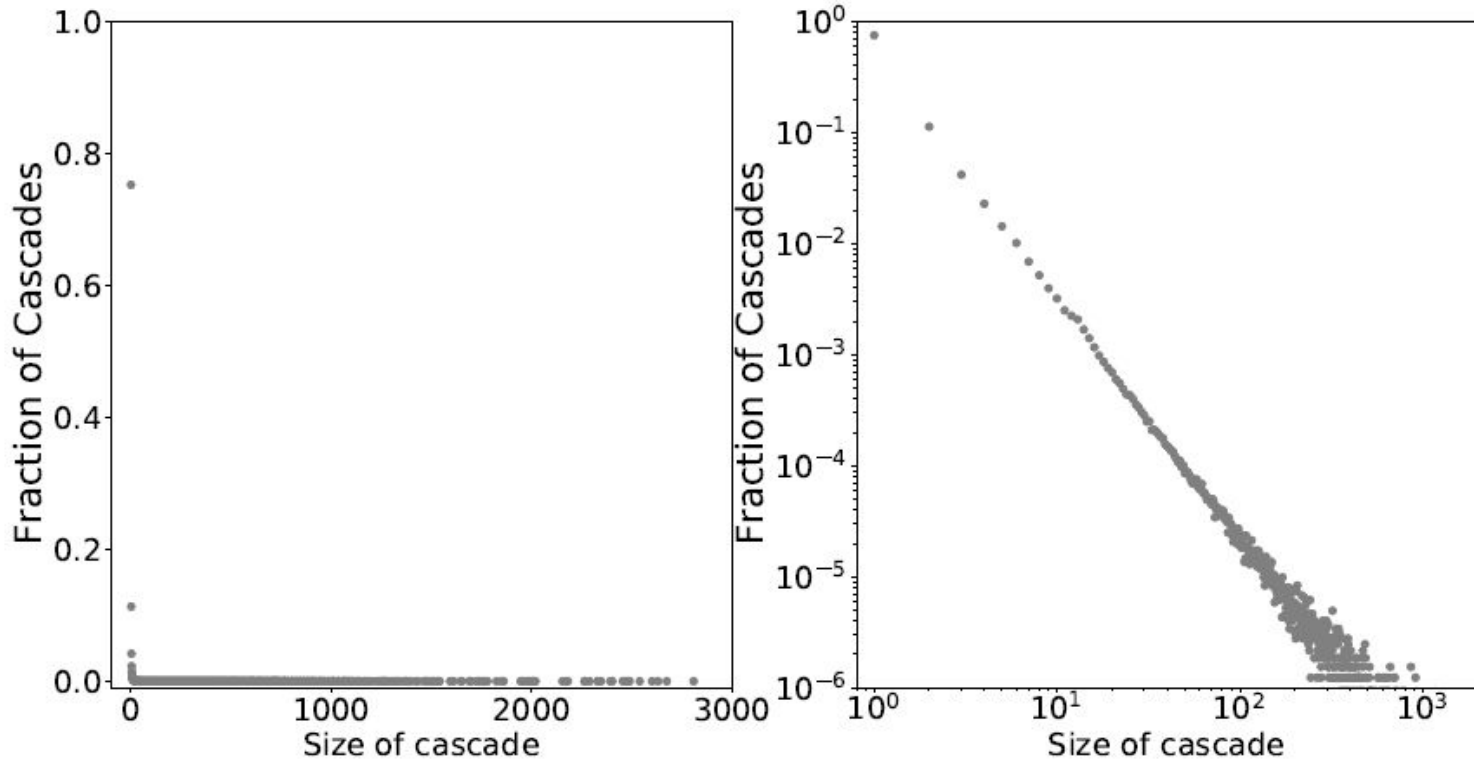
FARE

Fake News and Real People: Using big data
to understand human behaviour

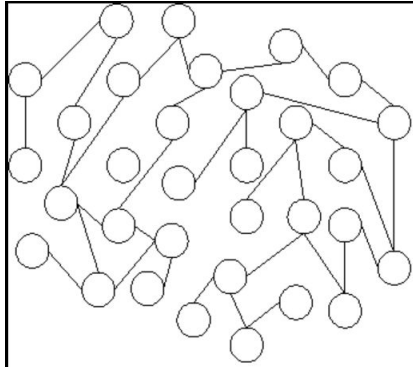




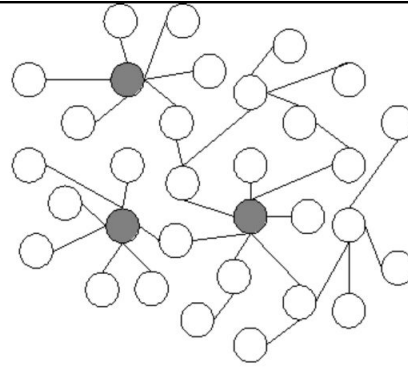
Image credit: all-free-download.com



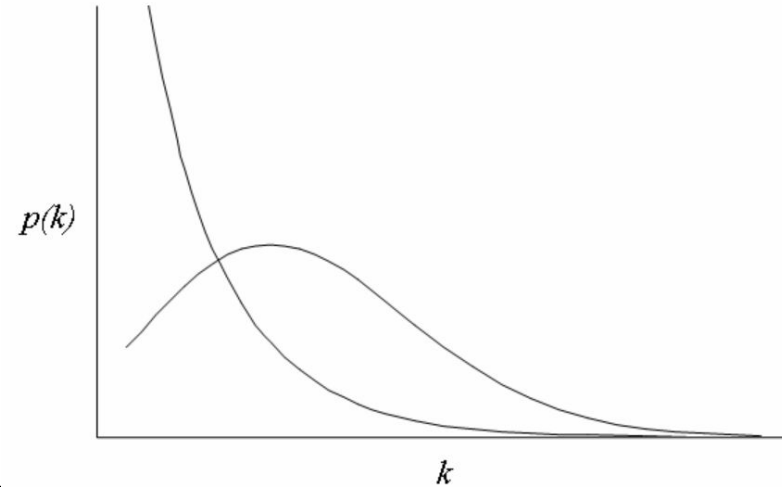
What are the main forces behind the success/failure of memes?

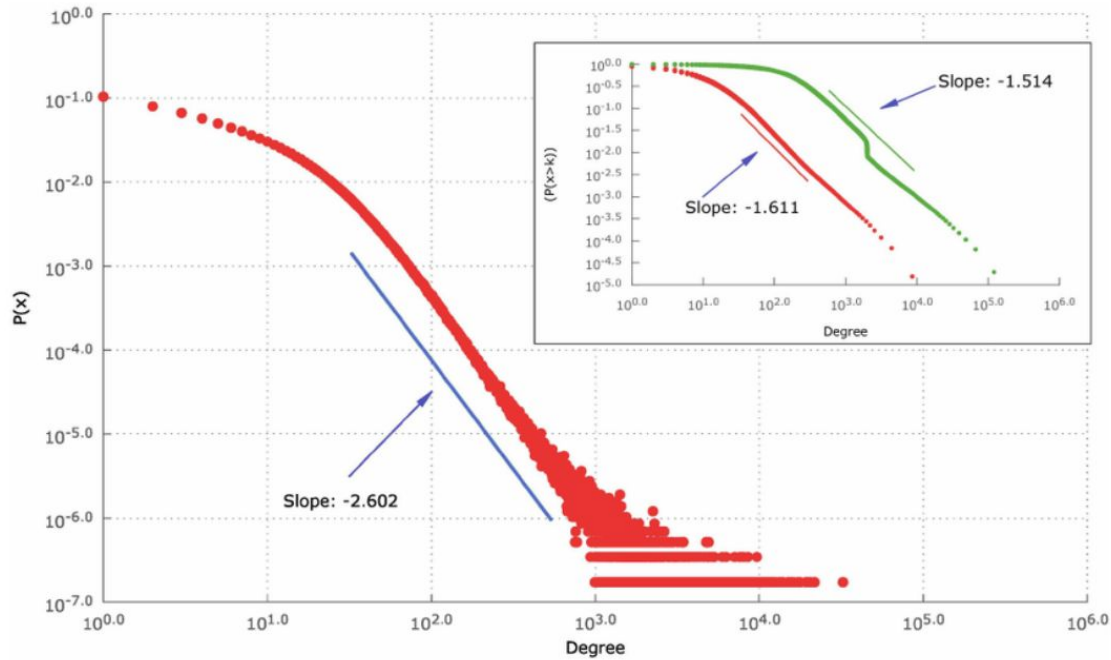


(a) Random network



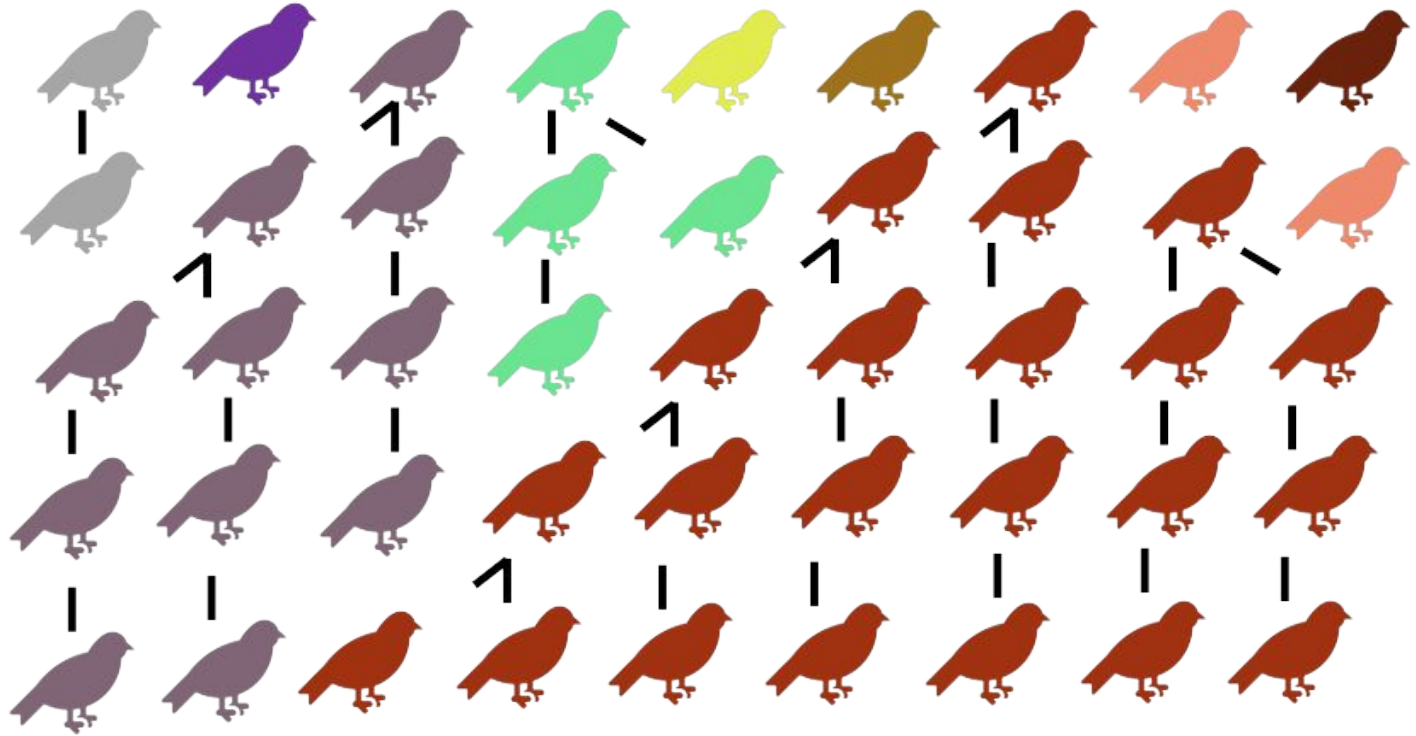
(b) Scale-free network

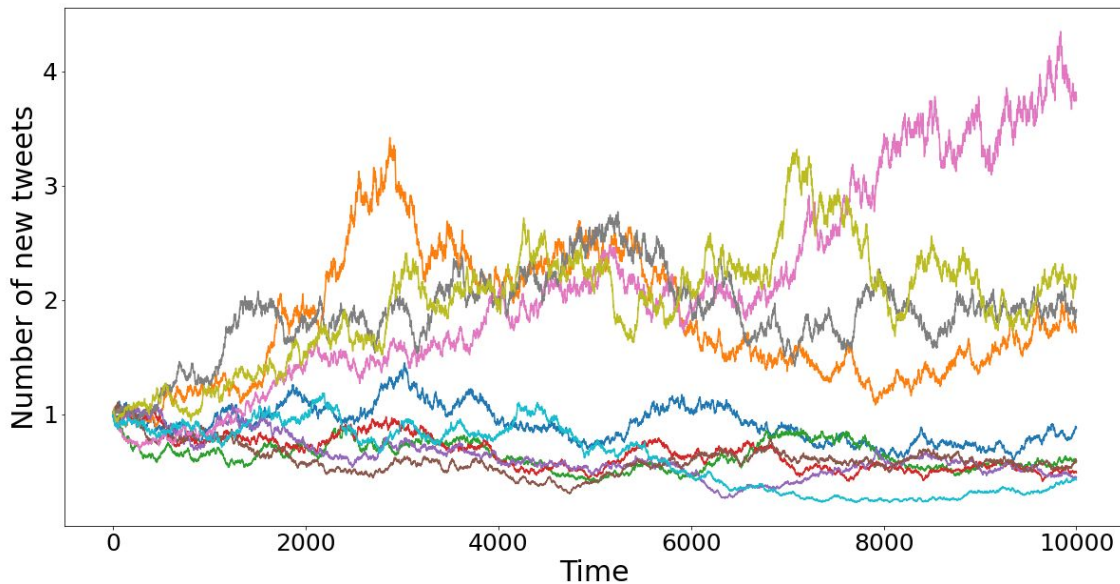




Szüle, J., Kondor, D., Dobos, L., Csabai, I., & Vattay, G. (2014). Lost in the city: revisiting Milgram's experiment in the age of social networks. *PLoS one*, 9(11), e111973.

Time





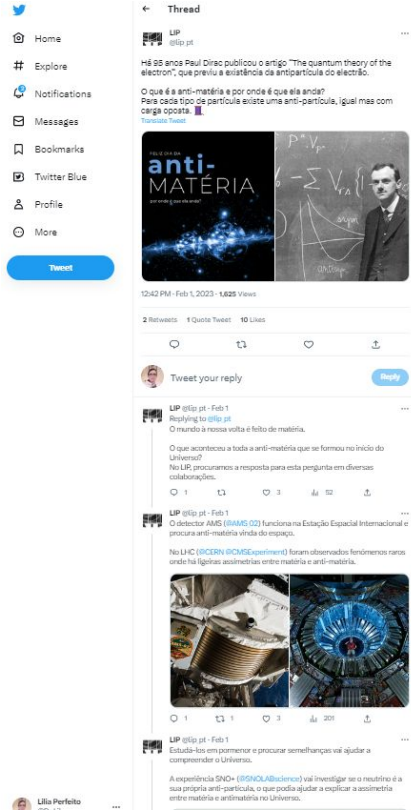
$$\frac{dN}{dt} = x \cdot \epsilon$$

$$\frac{d \log(N)}{dt} = \log(x) + \log(\epsilon)$$

where $\log(\epsilon)$ is a
Gaussian random
variable

Can we tease apart the role of the network and growth?

Can we measure selection acting on information?



Thread

LIP @lip_pt

Há 95 anos Paul Dirac publicou o artigo "The quantum theory of the electron", que previa a existência da antipartícula do electrão.

O que é a anti-matéria e por onde é que ela anda? Para cada tipo de partícula existe uma anti-partícula, igual mas com carga oposta.

anti-MATÉRIA

12:42 PM - Feb 1, 2023 - 1.822 Views

2 Retweets · 1 Quote Tweet · 10 Likes

Tweet your reply

LIP @lip_pt - Feb 1
Replying to lip_pt
O mundo à nossa volta é feito de matéria.

O que aconteceria a toda a anti-matéria que se formou no início do Universo?
No LIP, procuramos a resposta para esta pergunta em diversas colaborações.

LIP @lip_pt - Feb 1
O detector AMS (ALICE) funciona na Estação Espacial Internacional e procura anti-matéria (víde do espaço).

No LHC (CMS e ATLAS) foram observados fenómenos raros onde há ligeiras assimetrias entre matéria e anti-matéria.

LIP @lip_pt - Feb 1
Estudo-los em português e procurar semelhanças vai ajudar a compreender o Universo.

A experiência SND@ (SNS) vai investigar se o neutrino é a sua própria anti-partícula, o que podia ajudar a explicar a assimetria entre matéria e anti-matéria no Universo.

Search Twitter

Relevant people

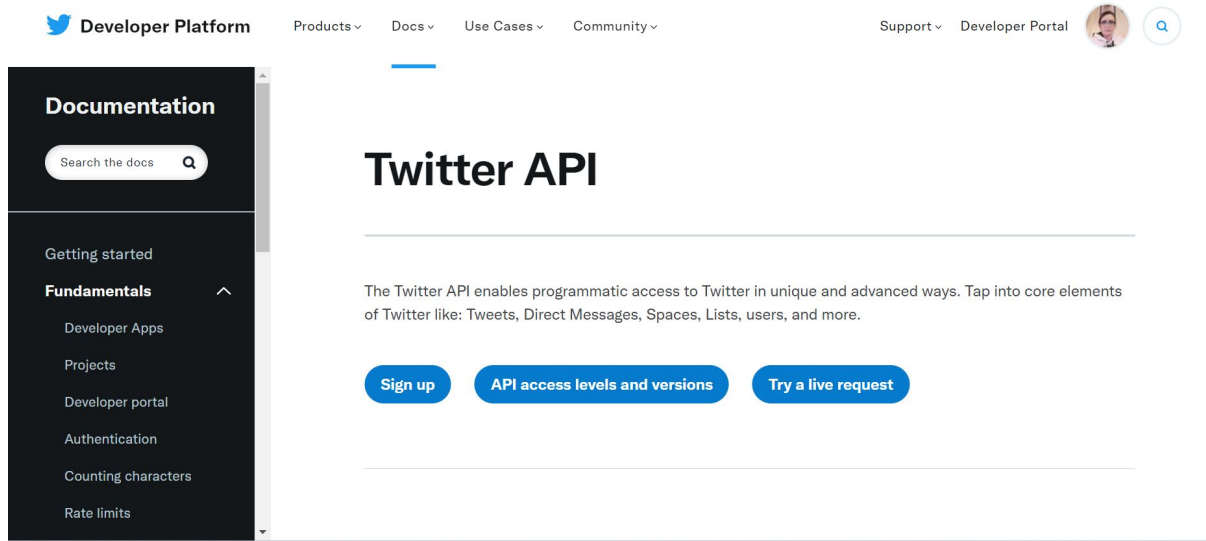
LIP @lip_pt Following

LIP is the research institute for experimental particle physics and associated technologies in Portugal.

Something went wrong. Try reloading.

Retry

Terms of Service Privacy Policy Cookie Policy Accessibility Ad info More... © 2023 Twitter, Inc.



Developer Platform

Products ▾ Docs ▾ Use Cases ▾ Community ▾

Support ▾ Developer Portal

Documentation

Search the docs

Getting started

Fundamentals

- Developer Apps
- Projects
- Developer portal
- Authentication
- Counting characters
- Rate limits

Twitter API

The Twitter API enables programmatic access to Twitter in unique and advanced ways. Tap into core elements of Twitter like: Tweets, Direct Messages, Spaces, Lists, users, and more.

Sign up API access levels and versions Try a live request

Tweets + re-tweets

Language	Covid	Music	Film
Portuguese	66 185 221	24 375 572	2 928 429
Italian	6 738 712	1 016 834	995 971
German	13 566 605	447 831	444 034
Dutch	7 159 327	191 366	186 972

Tweets + re-tweets

Language	Covid	Music	Film
Portuguese	66 185 221	24 375 572	2 928 429
Italian	6 738 712	1 016 834	995 971
German	13 566 605	447 831	444 034
Dutch	7 159 327	191 366	186 972

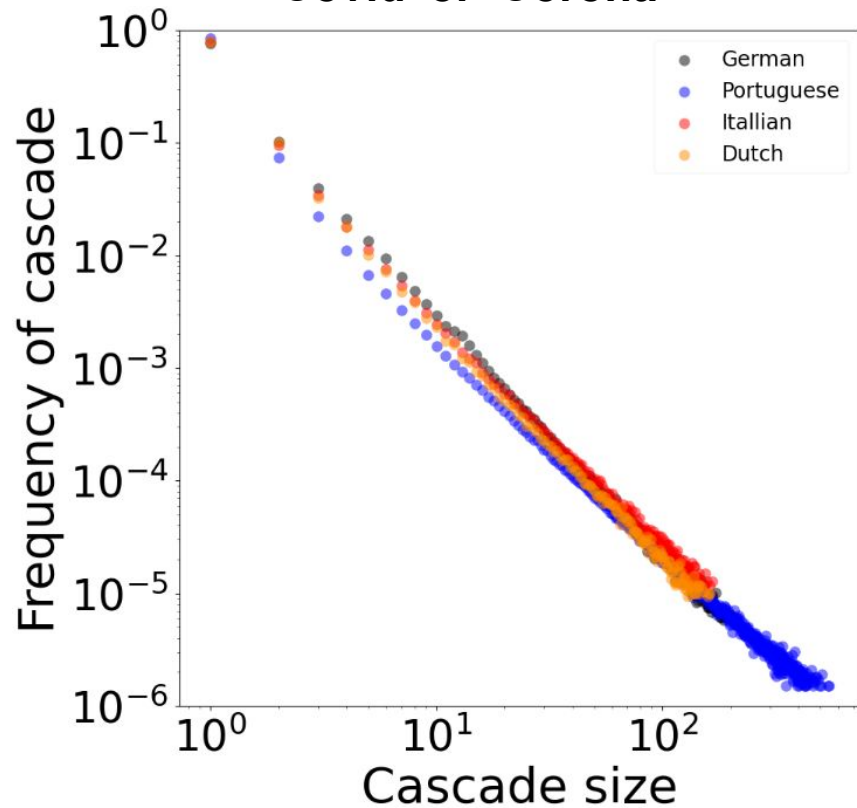
April 2020 to June 2021

Accounts
Original tweets

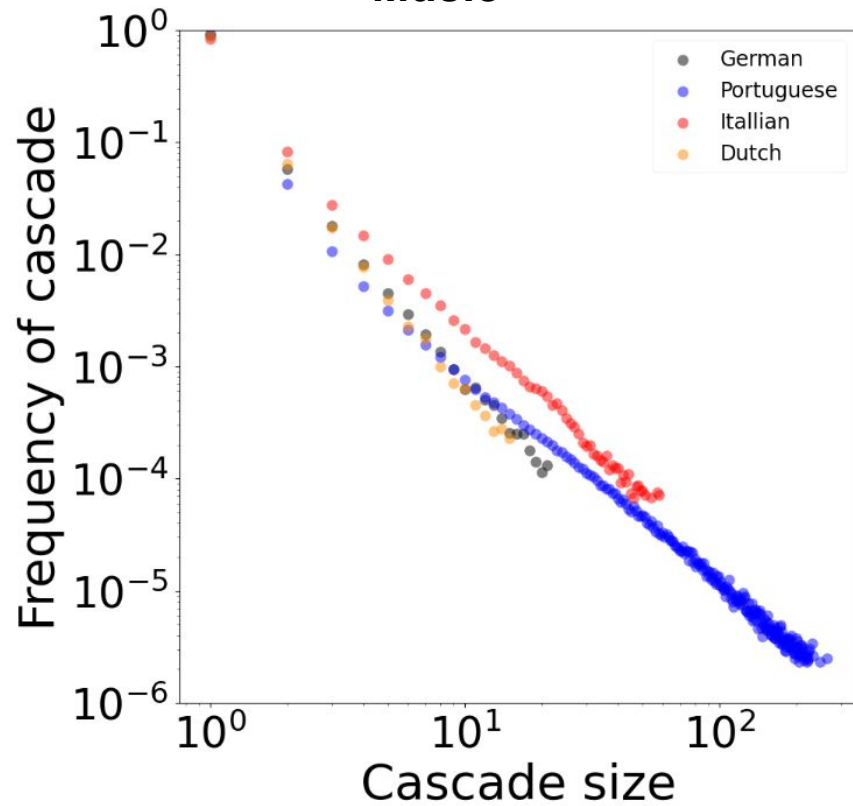
Language	Accounts			Original tweets		
	Covid	Music	Film	Covid	Music	Film
Portuguese	6 914 262	3 905 141	2 928 429	20 595 795	13 251 392	7 888 742
Italian	728 279	728 279	728 279	2 719 486	490 764	638 865
German	808 096	154 463	170 588	5 392 367	324 803	297 891
Dutch	461 353	69 203	95 396	3 293 620	138 851	115 969

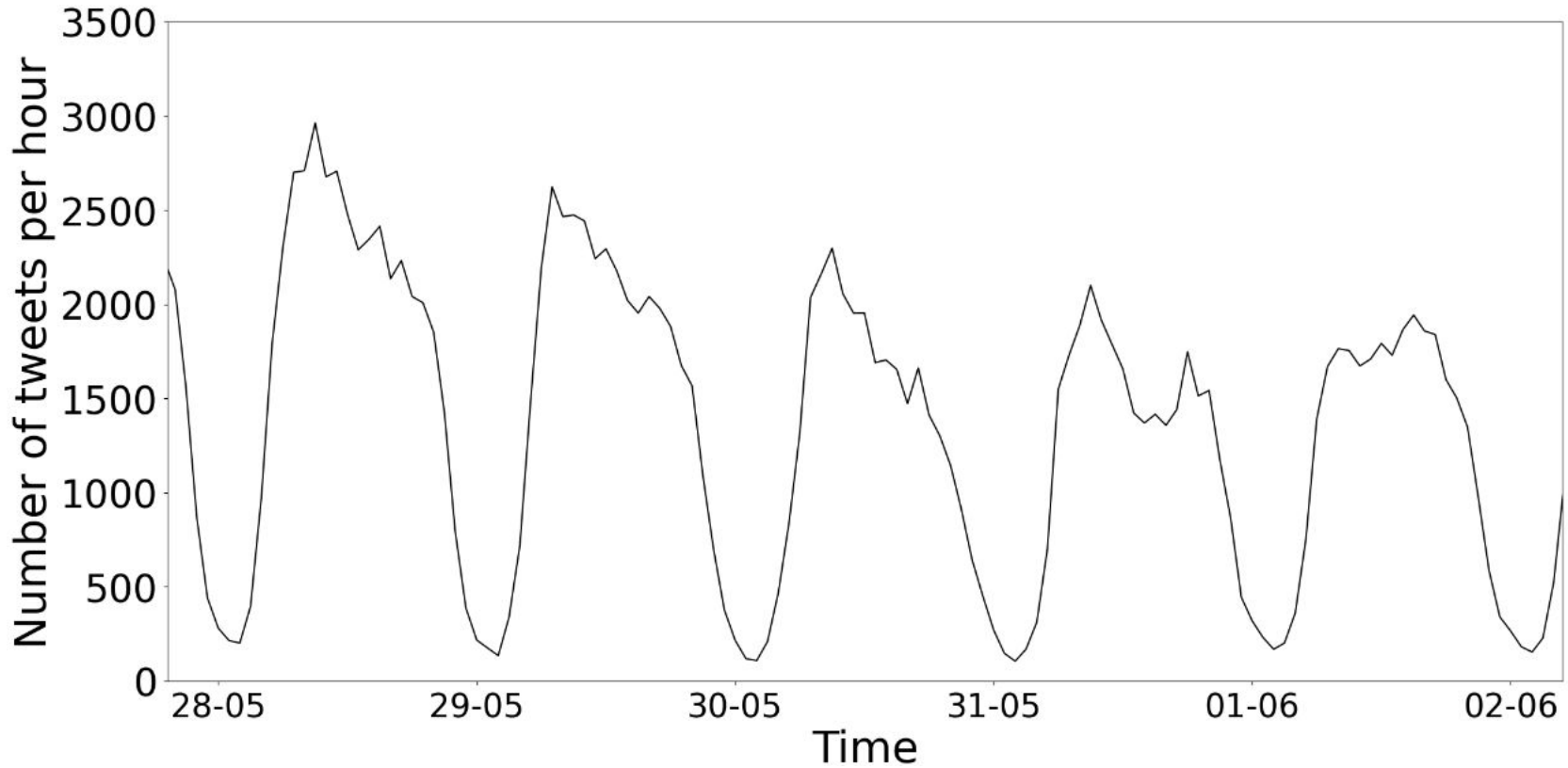
Cascade: Set of identical tweets, probably originating in a single seeding event.

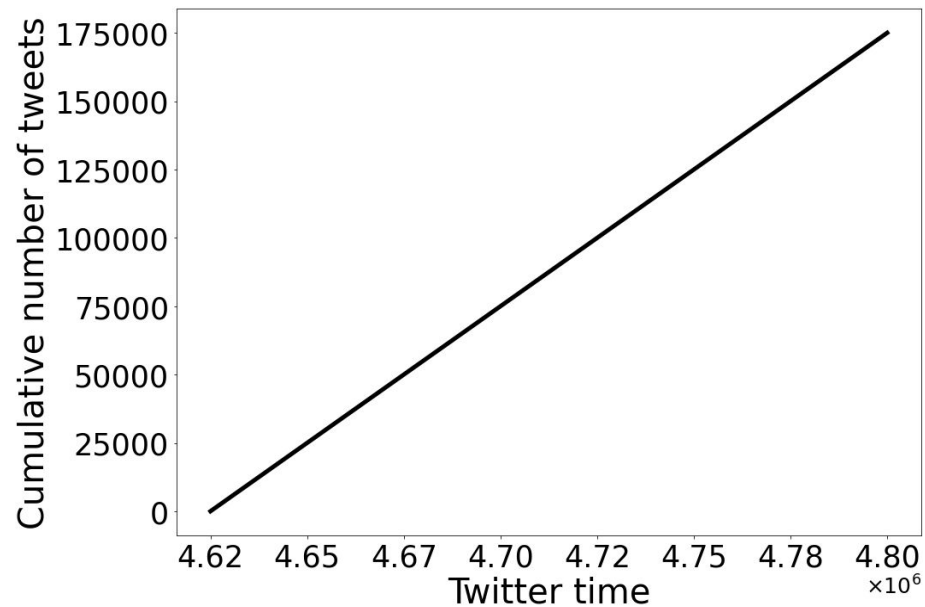
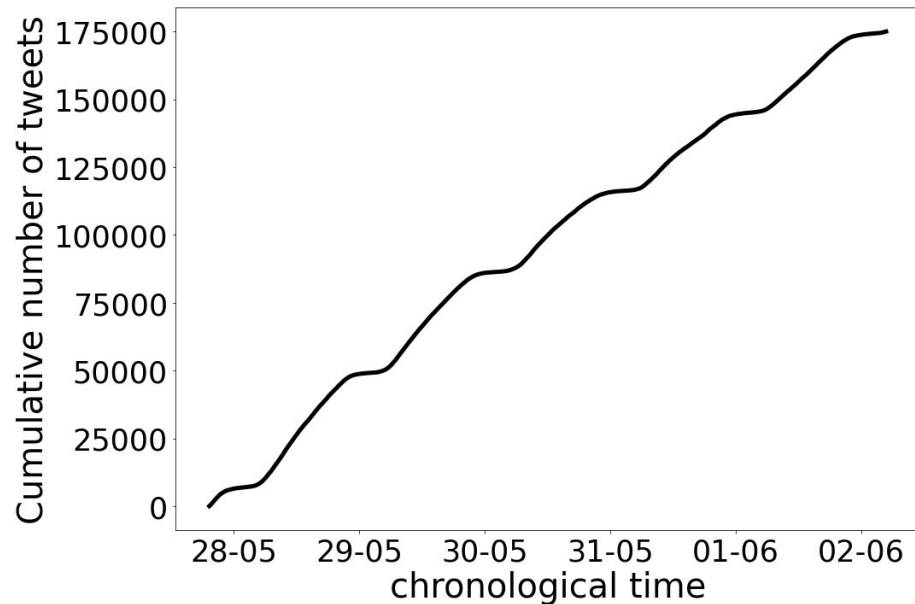
'Covid' or 'Corona'

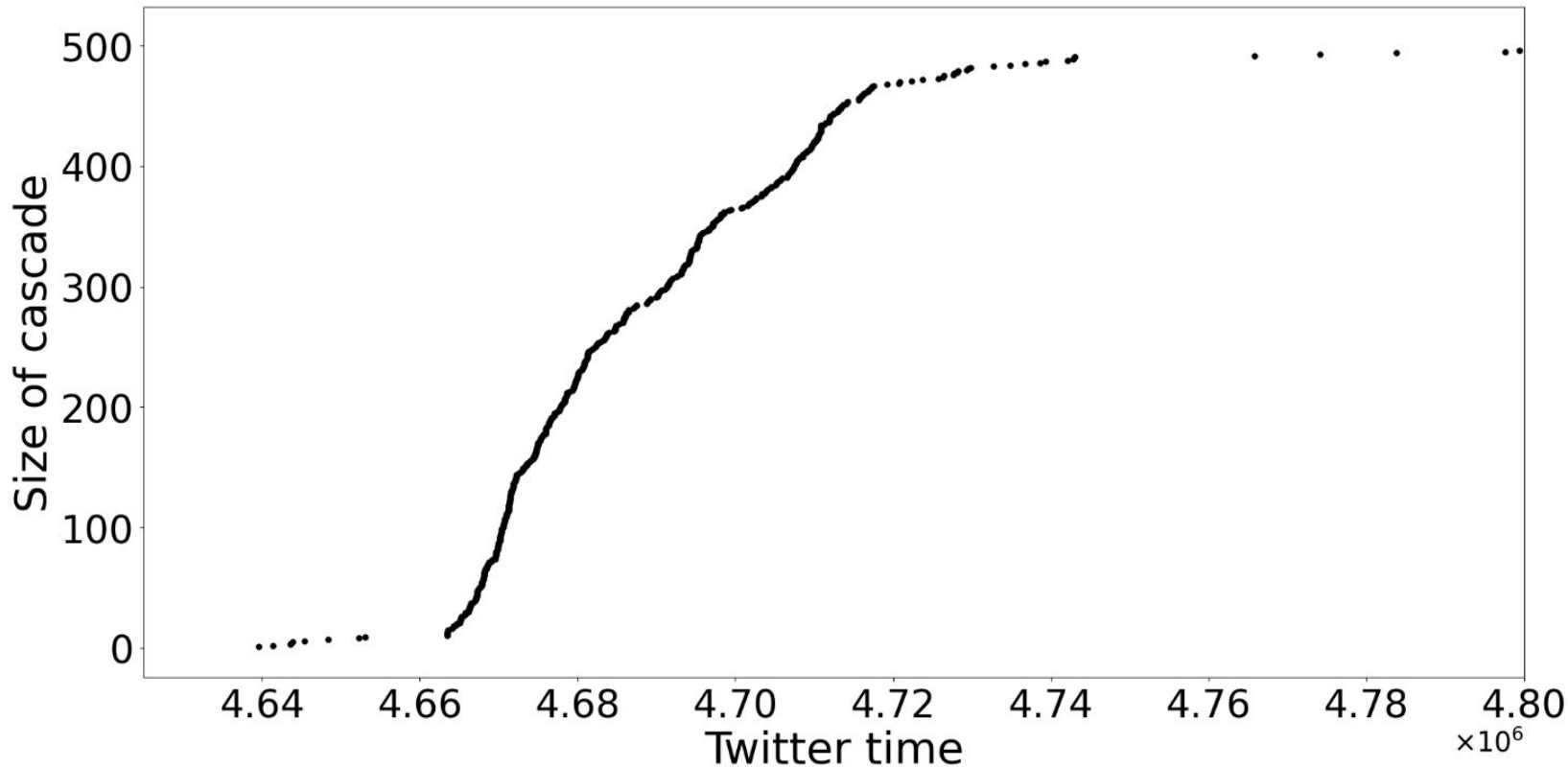


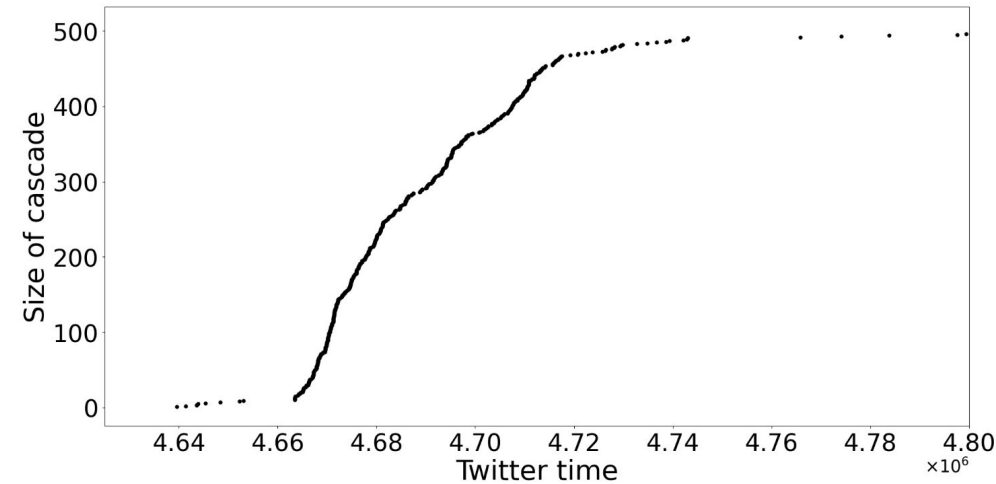
'Music'





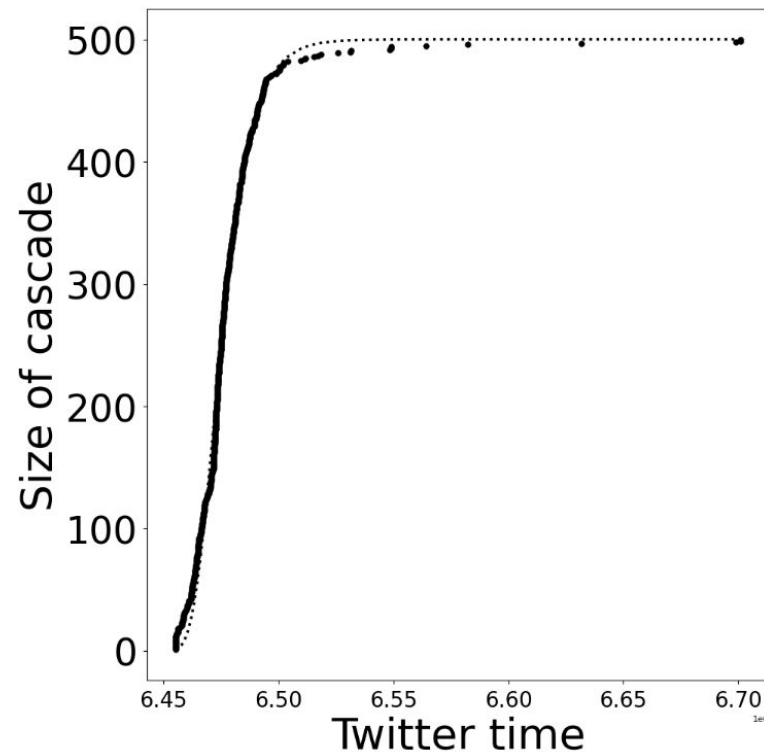
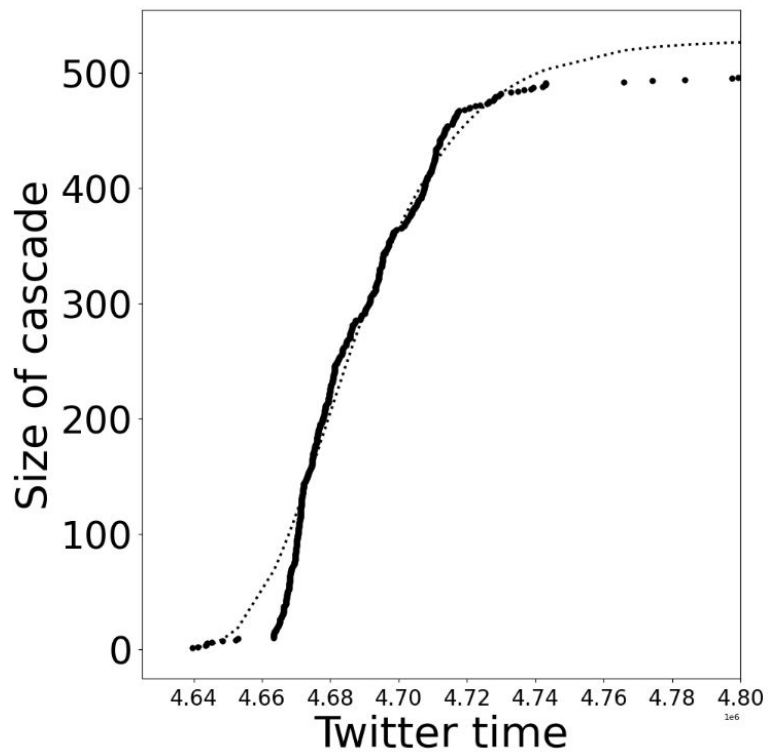






$$\frac{dN}{dt} = a \cdot e^{-g \cdot t}$$

$$N(t) = N(0) \cdot e^{a/g} e^{-a \cdot e^{-g \cdot t}/g}$$



If we let all cascades grow to maximum size, (i.e., $t \rightarrow +\infty$), then:

$$\frac{dN}{dt} = a \cdot e^{-g \cdot t}$$

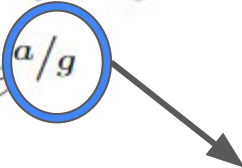
$$N(t) = N(0) \cdot e^{a/g} e^{-a \cdot e^{-g \cdot t}/g}$$

$$\lim_{t \rightarrow \infty} N(t) = e^{a/g}$$

If we let all cascades grow to maximum size, (i.e., $t \rightarrow +\infty$), then:

$$\frac{dN}{dt} = a \cdot e^{-g \cdot t}$$

$$N(t) = N(0) \cdot e^{a/g} e^{-a \cdot e^{-g \cdot t}/g}$$

$$\lim_{t \rightarrow \infty} N(t) = e^{a/g}$$


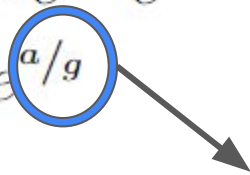
Fitness

If we let all cascades grow to maximum size, (i.e., $t \rightarrow +\infty$), then:

$$\frac{dN}{dt} = a \cdot e^{-g \cdot t}$$

$$N(t) = N(0) \cdot e^{a/g} e^{-a \cdot e^{-g \cdot t}/g}$$

$$\lim_{t \rightarrow \infty} N(t) = e^{a/g}$$



Fitness, ω

What is the distribution of cascade **sizes ($f(N)$)**, given a certain distribution of **fitnesses ($g(\omega)$)**?

$$f(N) = \frac{1}{N} \cdot g(\log N)$$

If we let all cascades grow to maximum size, (i.e., $t \rightarrow +\infty$), then:

$$\frac{dN}{dt} = a \cdot e^{-g \cdot t}$$

$$N(t) = N(0) \cdot e^{a/g} e^{-a \cdot e^{-g \cdot t}/g}$$

$$\lim_{t \rightarrow \infty} N(t) = e^{a/g}$$

Fitness, ω

What is the distribution of cascade **sizes ($f(N)$)**, given a certain distribution of **fitnesses ($g(\omega)$)**?

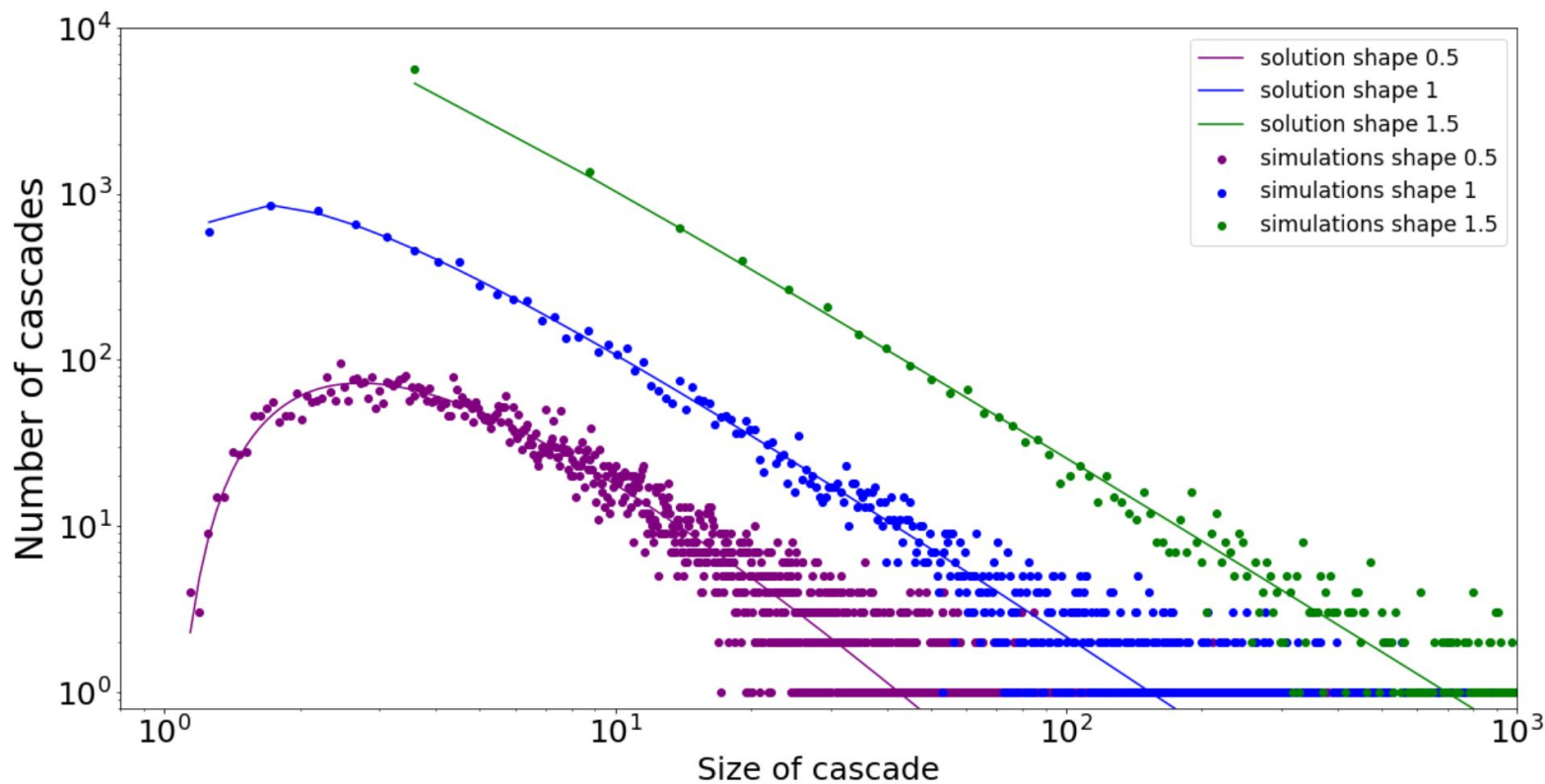
$$f(N) = \frac{1}{N} \cdot g(\log N).$$

Example: **$g(\omega)$ is exponential**

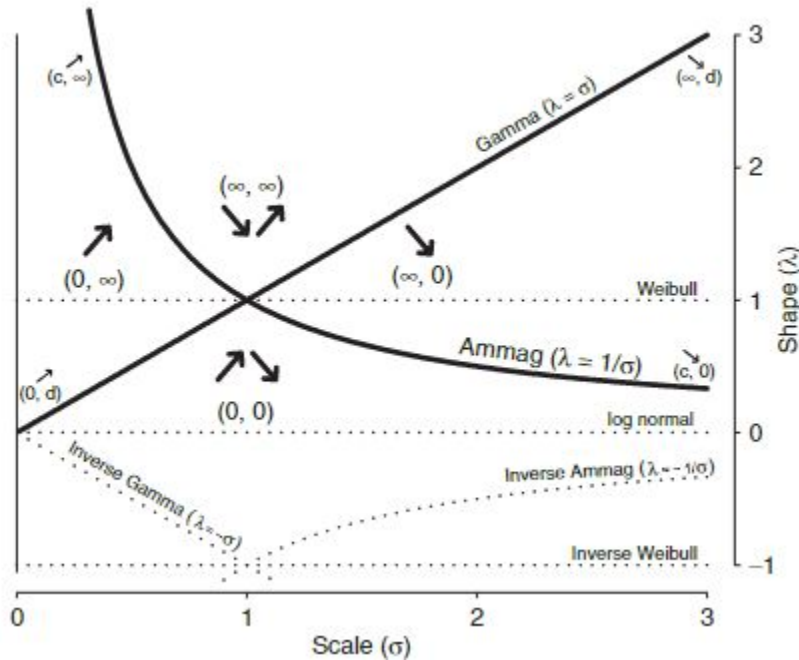
$$g(\omega) = \lambda \cdot e^{-\lambda \cdot \omega}$$

$$f(N) = \frac{\lambda \cdot e^{-\lambda \cdot \log(N)}}{N}$$

$$f(N) = \lambda \cdot N^{-\lambda-1}$$



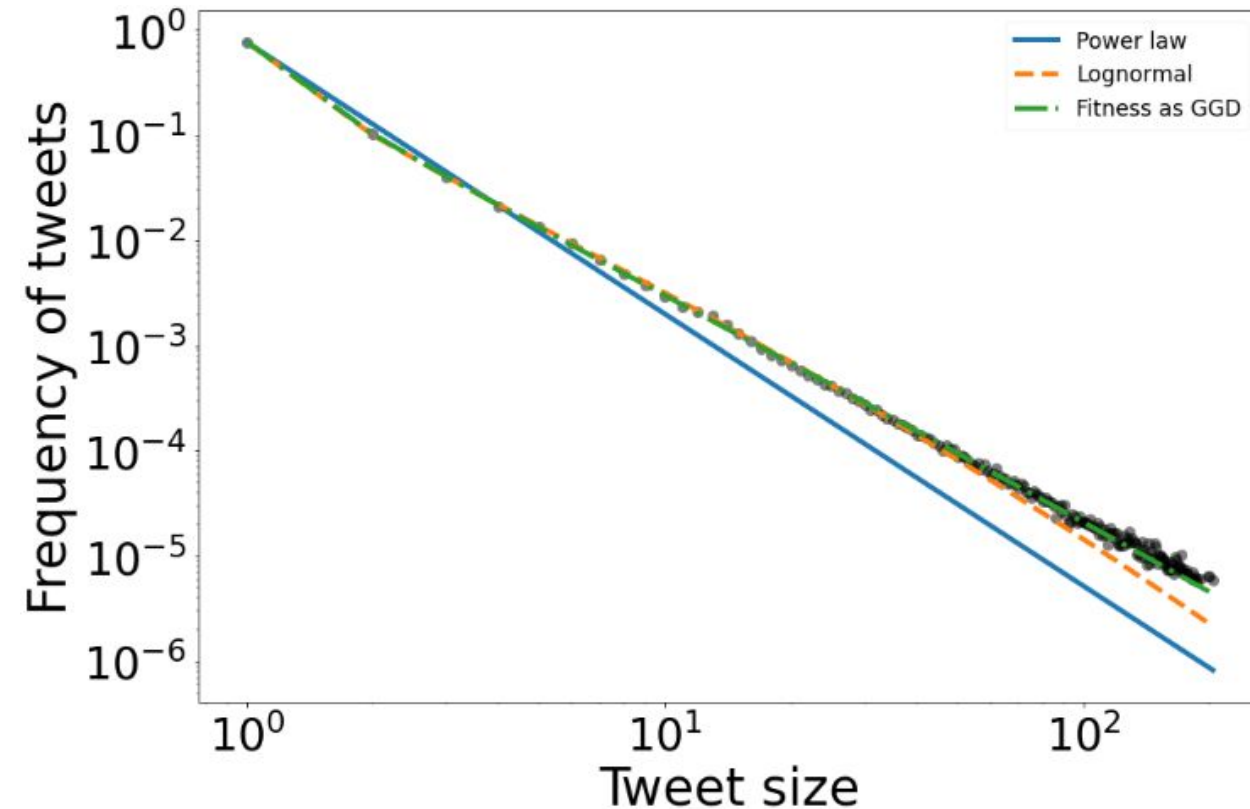
What is a good model for the fitness distribution?



Exponential
Gamma
Lognormal
Powerlaw

Generalized Gamma

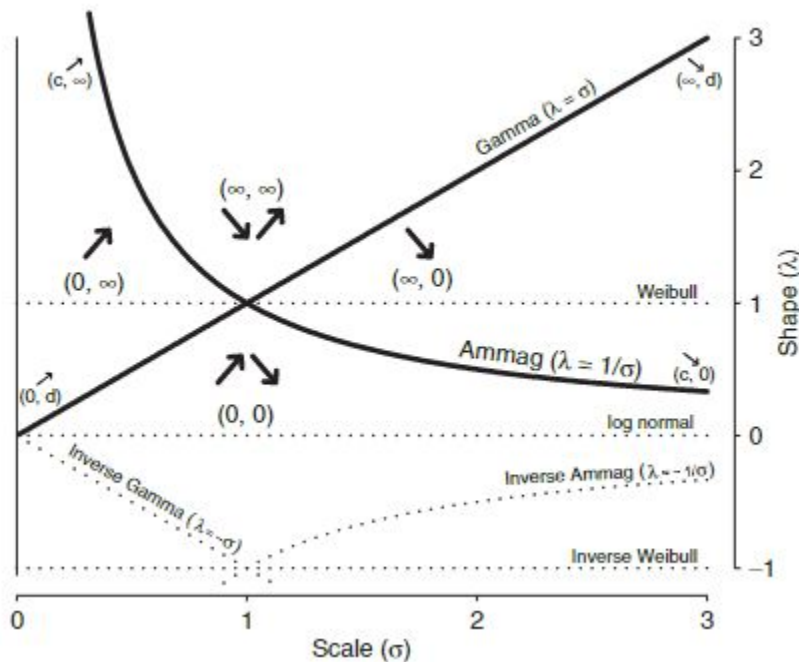
Cox, C., Chu, H., Schneider, M.F., Munoz, A.: Parametric survival analysis and taxonomy of hazard functions for the generalized gamma distribution. *Statistics in medicine* **26**(23), 4352–4374 (2007)



$$f(N) = a \cdot N^{-k}$$

$$f(N) = \text{lognormal}(N, \mu, \sigma, \text{loc})$$

$$f(N) = \text{GGD}(\log(N), \sigma, \beta, \text{loc}) \cdot N^{-1}$$

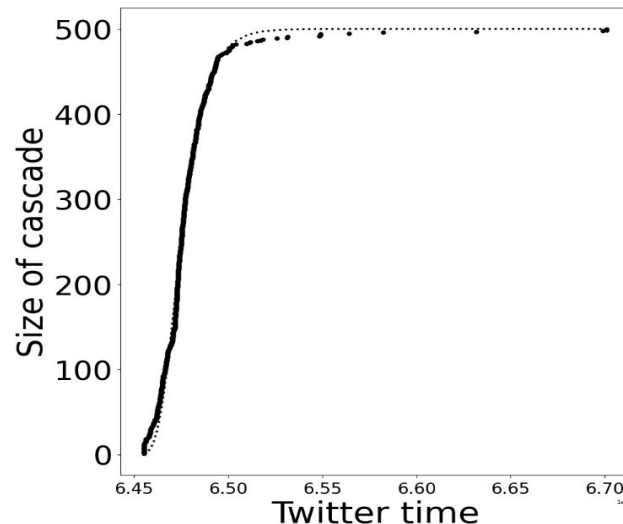
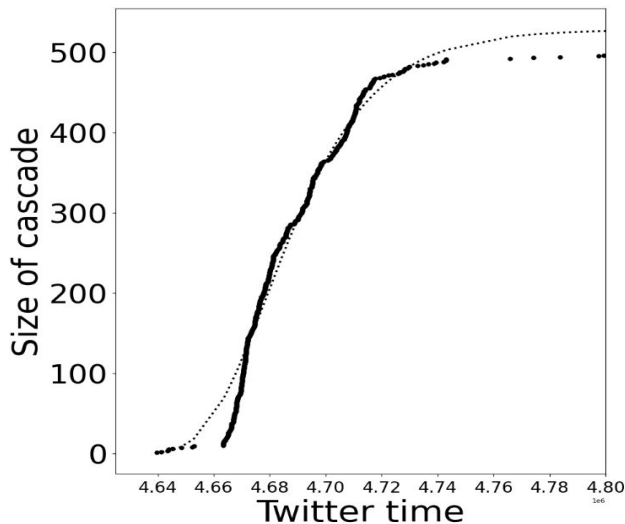


Language	Dataset	Parameters of fitness distribution			
		λ	σ	β	location
Portuguese	Covid	0.23	1.28	20.58	-0.35
	Music	0.29	1.69	19.33	-0.21
	Film	0.29	1.57	18.11	-0.22
Italian	Covid	1.28	1.56	-0.20	-0.27
	Music	1.45	1.79	-0.51	-0.23
	Film	0.36	1.35	11.13	-0.31
German	Covid	1.96	1.53	-1.84	-0.23
	Music	2.22	2.05	-1.57	-0.15
	Film	1.98	2.05	-1.20	-0.17
Dutch	Covid	0.64	1.23	3.51	-0.39
	Music	0.76	1.26	2.80	-0.33
	Film	1.32	1.82	0.46	-0.21

We introduce the concept of twitter time

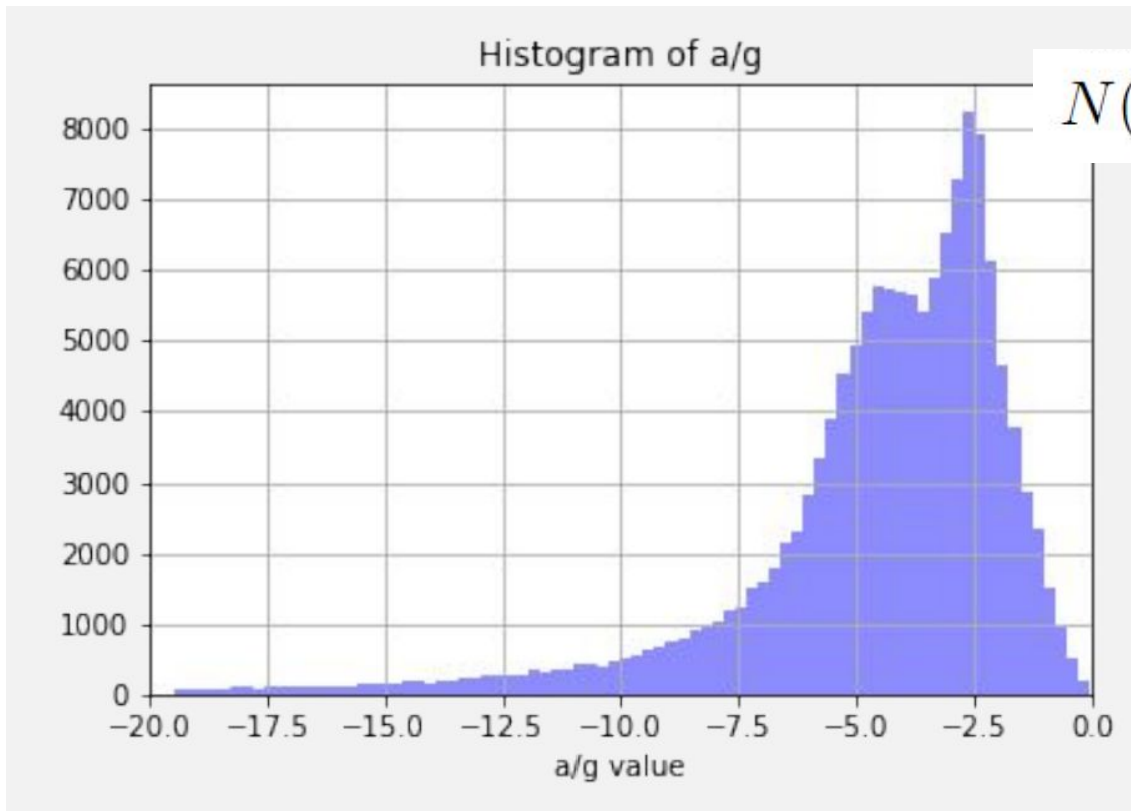
We show the social network is not necessary to reproduce the cascade size distribution

The distribution of fitnesses is well approximated by a generalized Gamma distribution with an exponential-to-heavy tail

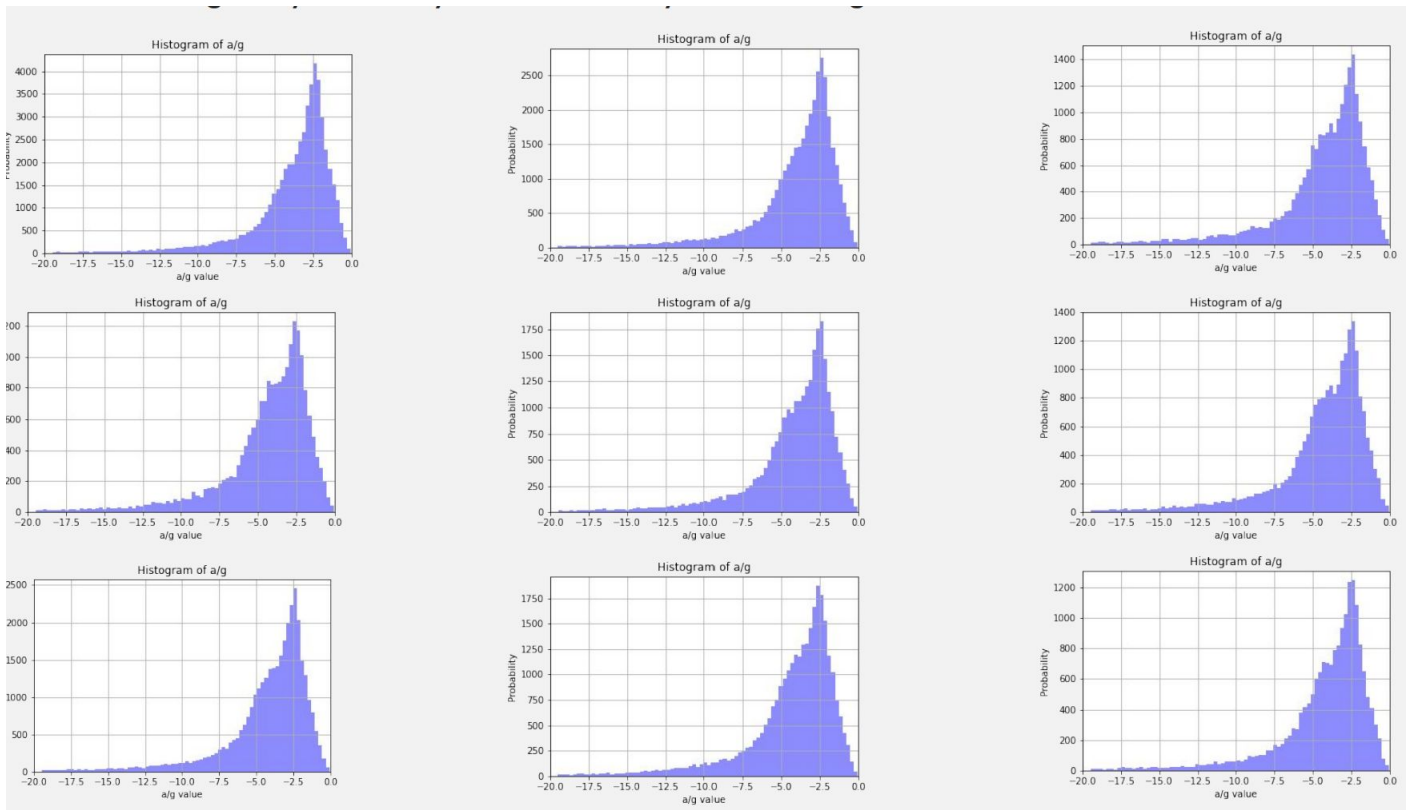


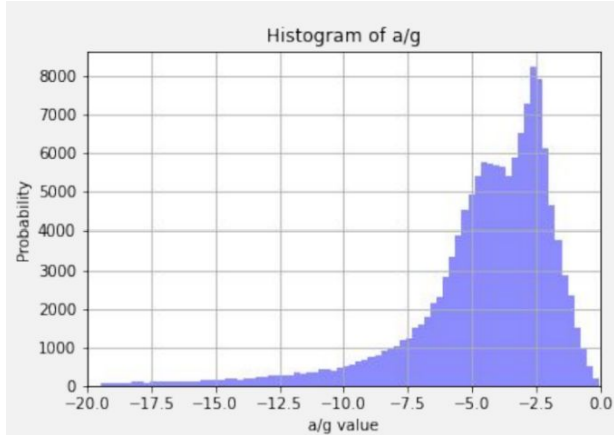
Another way to estimate the distribution of fitness parameters: fit each cascade 1 by 1

Time consuming, not applicable to all cases

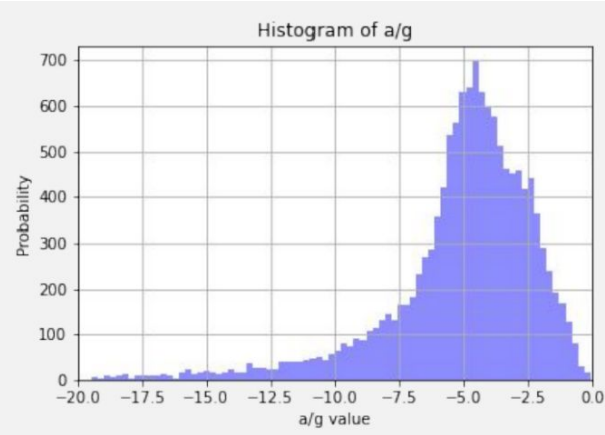


$$N(t) = N(0) \cdot e^{a/g} e^{-a \cdot e^{-g \cdot t}/g}$$

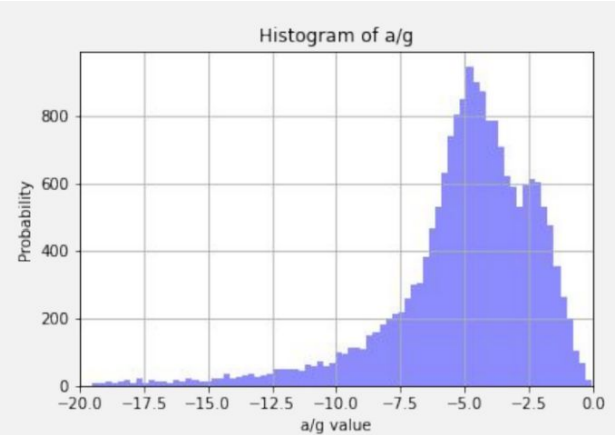




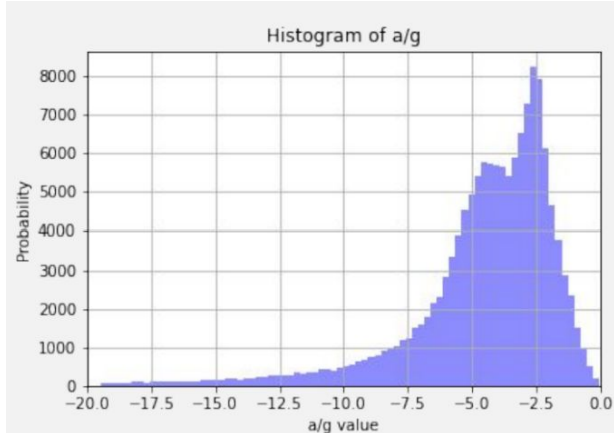
Covid



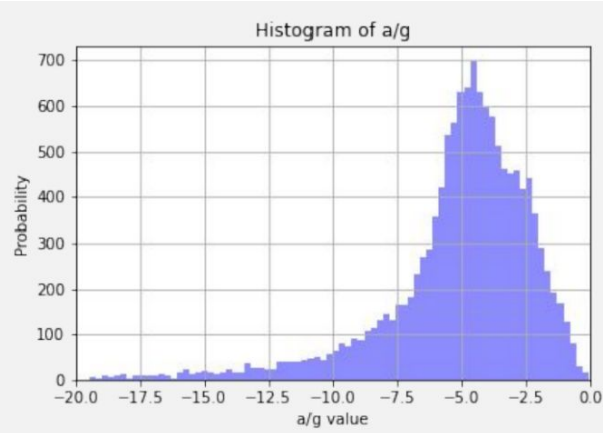
Film



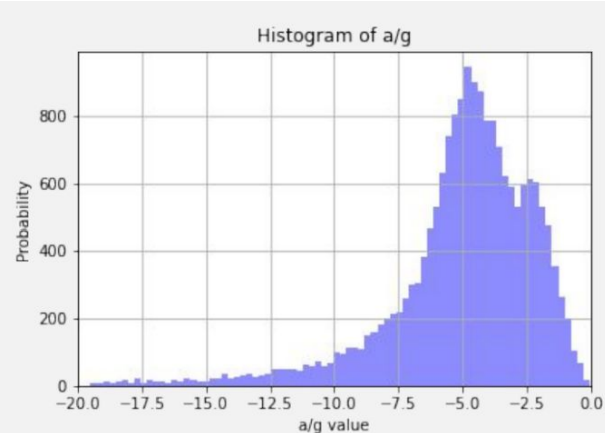
Music



Covid



Film



Music

How can we separate them?

Can we tease apart the role of the network and growth?

Can we measure selection acting on information?

Q1: Can we tease apart the role of the network and growth?

A1: The network is not necessary.

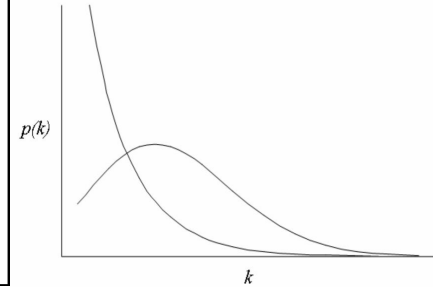
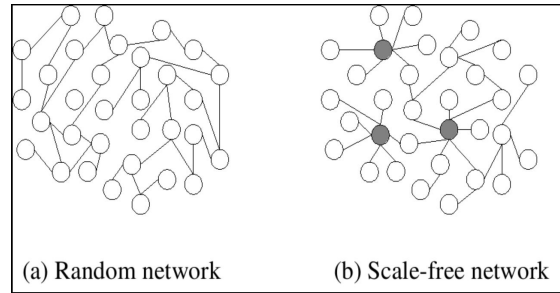
Q1.1: Which one explains more of the variability in the data?

Q2: Can we measure selection acting on information?

A2: Yes

1. Build the network (1000 nodes)

- a. Uniform/Poisson
- b. Realistic, ie, powerlaw



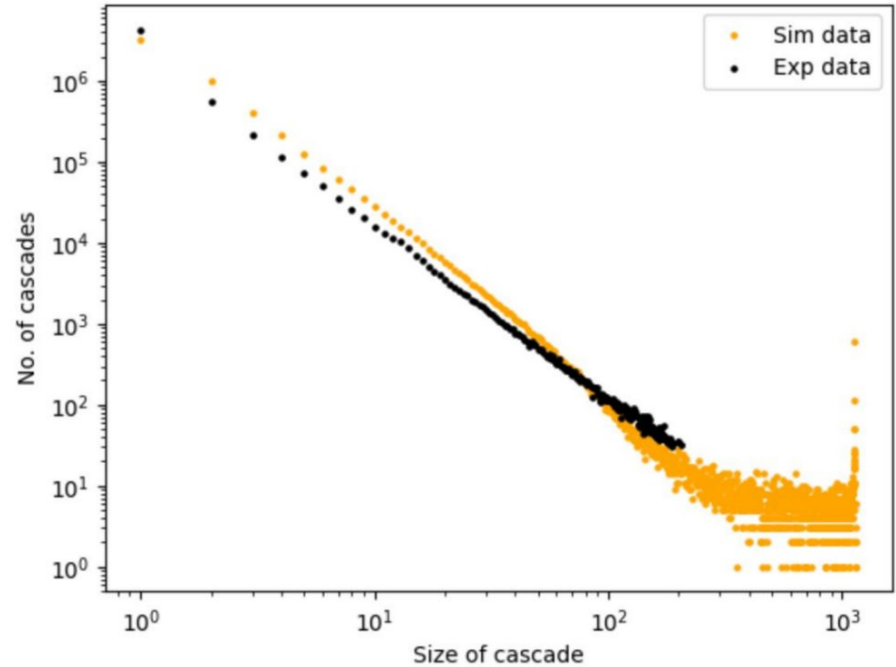
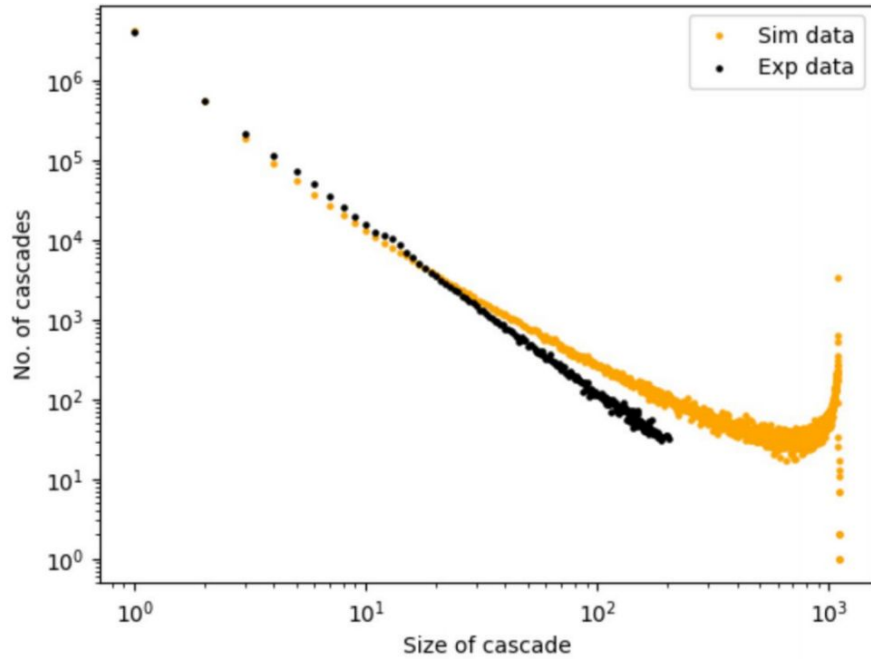
2. Choose a node at random to post the tweet

3. Let tweets spread

- a. Following $p_i(\text{share}) = k_i \cdot a_c \cdot e^{-g_c \cdot t}$
- b. Introducing a set of resistant nodes

4. Repeat thousands of times, find the simulation parameters that best fit the data (Approximate Bayesian Computation)

5. Compare the goodness of fit of the best model without network with the best model with network



- Simple vs complex contagion
- Introducing heterogeneity in nodes (personality types)
- True vs misinformation
- ...