# Discovery, Big Data, web analytics, and Bayesian A/B testing at Wikipedia

**Mikhail Popov** (@bearloga)
Data Scientist // Discovery // Wikimedia Foundation
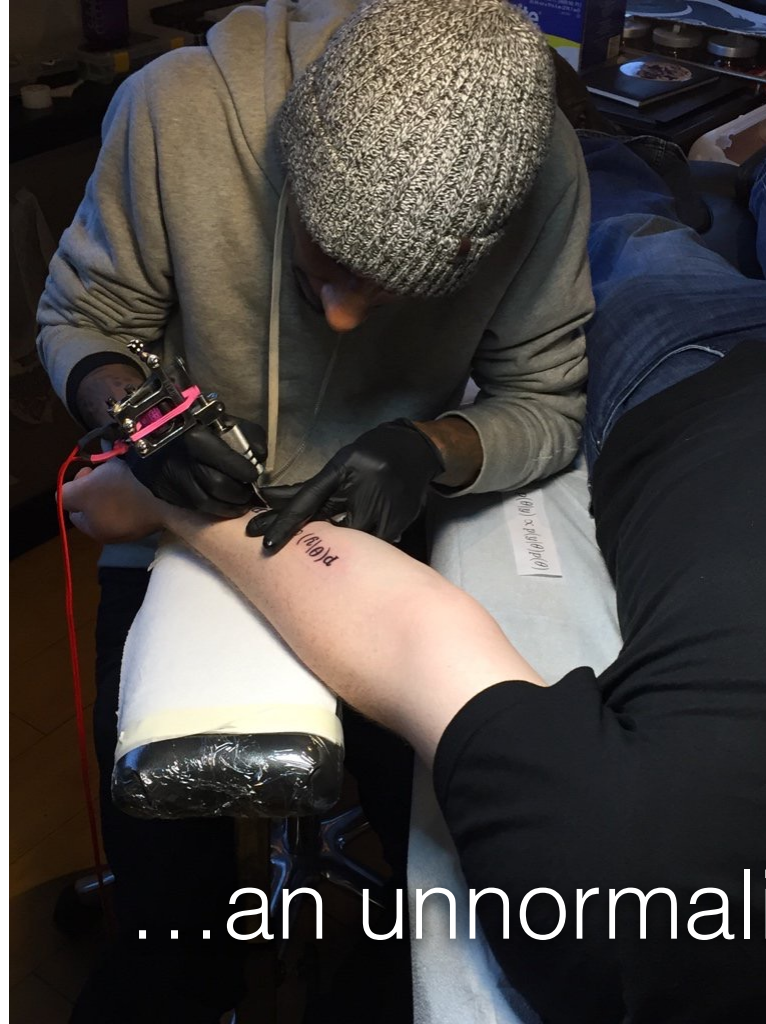
# I was/am…

🍕 Bachelor in Mathematics (CSUF NSM, 2012)

🍕 Master in Statistical Practice (CMU, 2013)

🍕 Statistician for a neuropsychology research program at University of Pittsburgh / UPMC for 2 years

🍕 Data Scientist / Data Analyst / Quantitative Analyst / Statistician / why do we have so many names??? 🤔 at Wikimedia Foundation
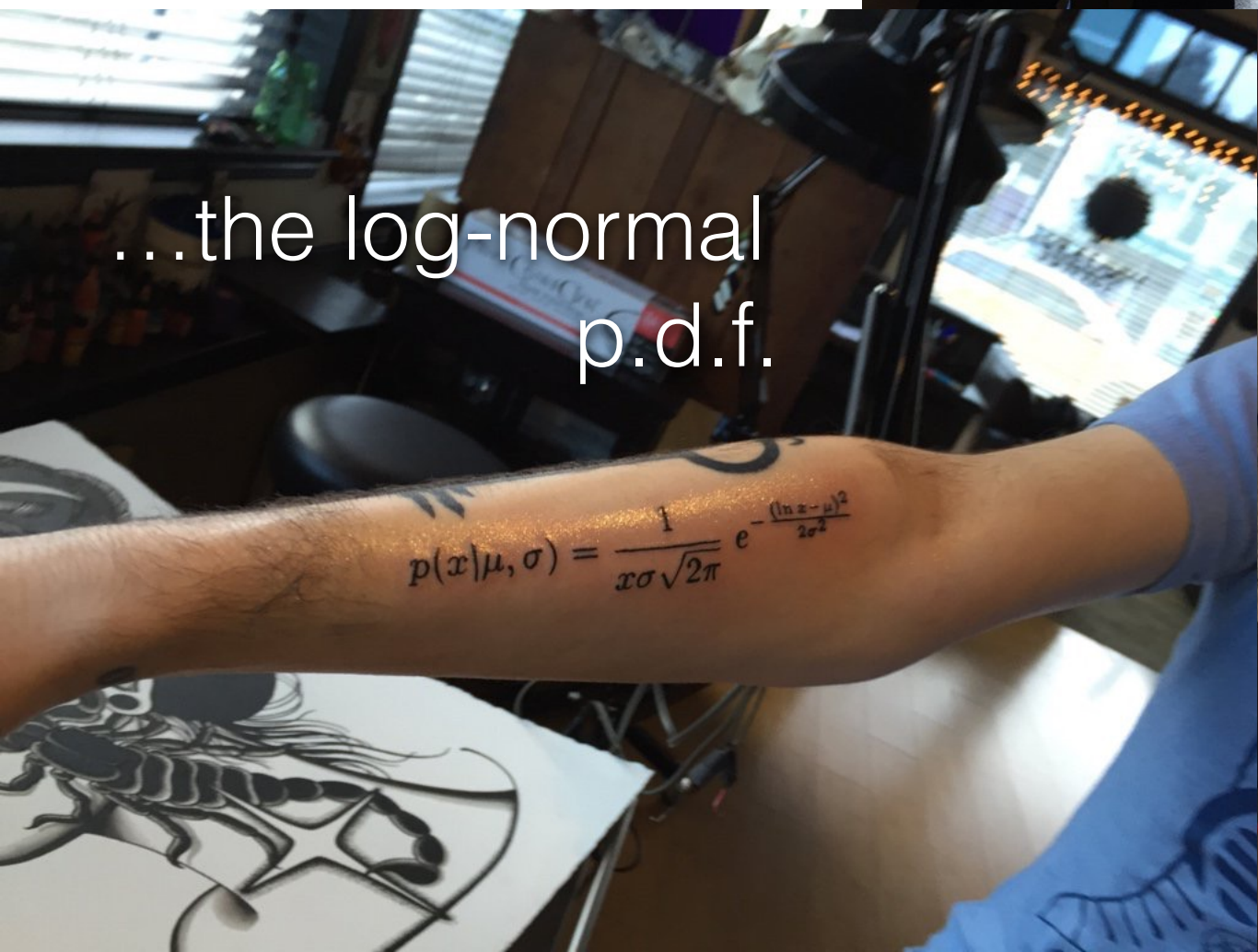
🍕 Tattoo designer-typesetter (Wait, what?)

9 January 2016:

Three data science dorks walk into a tattoo parlor in Berkeley, holding LaTeX-typeset formulas for…

…an unnormalized posterior density!

$$p(\theta|y) \propto p(y|\theta)\, p(\theta)$$

…the log-normal p.d.f.

$$p(x|\mu,\sigma) = \frac{1}{x\sigma\sqrt{2\pi}}\, e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}$$

…the Poisson p.m.f.

$$p(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

# The Wikimedia Foundation (WMF)

🍏 A nonprofit charitable organization.

🍎 Committed to creating a world in which every single human being can freely share in the sum of all knowledge (e.g. Wikipedia, Wikipedia Zero).

🍏 Collaborates with users around the world.

🍎 Develops MediaWiki software used by many organizations and companies (e.g. Intel, Wikia, NASA).

🍏 Fights for people's privacy (see *Wikimedia Foundation, et al. v. National Security Agency, et al.* – https://en.wikipedia.org/wiki/Wikimedia_Foundation_v._NSA)

I'M THE TOP 6TH WEBSITE!

Wikipedia

Wikimania

Wikibooks

wiki voyage

Community

Commons

Wikiquote

Wikisource

Wikispecies

WIKIMEDIA

Wikinews

MediaWiki

Wiktionary

Incubator

Wikiversity

WIKIDATA

# Discovery Department

*Building the anonymous path of discovery to a trusted and relevant source of knowledge.*

**Projects**:

🔍 Search feature & APIs

🔍 Wikipedia Portal (wikipedia.org)

🔍 Maps, in collaboration with OpenStreetMap

🔍 Wikidata Query Service

# Discovery's Analysts

📊 Integrated within department instead of separate dep't

📊 Provide ad-hoc analyses and reports as needed

📊 Build and maintain dashboards for tracking *key performance indicators* (KPIs) and other metrics

📊 Consult with teams in design of experiments – "A/B tests" – e.g. narrowing down affected population

📊 Work with engineers to design event logging schemas

# Ad-hoc analyses include…



**Clickthrough Rates by Top 10 Primary Accept-Languages**

Overall Clickthrough Rate — Search Clickthrough Rate

Legend:
- Dutch
- English
- French
- German
- Italian
- Korean
- Portuguese
- Russian
- Spanish

Overall Clickthrough Rate values: +7.8%, +10.6%, +26.3%, +11.0%, +20.9%, -2.9%, +6.2%, +11.4%, +5.7%

Search Clickthrough Rate values: +4.9%, +8.7%, +11.6%, -1.8%, +11.6%, +5.3%, -0.4%, +3.4%

Y-axis: Clickthrough Rate (80%, 70%, 60%, 50%, 40%, 30%, 20%)

X-axis: Accept-Language — Does NOT include English / Includes English

**38 countries where IE 7 is one of the top 10 browsers**

↑ Global browser usage and JavaScript support.

← Engagement with wikipedia.org based on whether your list of preferred languages includes English or not.

discovery.wmflabs.org

The Internet! (So cloudy!)

Kafka, Camus, etc.

Web requests and event logs on Hadoop/MySQL

Data acquisition & aggregation R scripts

Published to

Public aggregated datasets (datasets.wikimedia.org)

# Discovery Dashboards

R/Shiny-powered dashboards

(discovery.wmflabs.org)

KPI summary for Nov 27th-Oct 29th , and % change from Dec 28th-Nov 28th :

449ms    29.0%    122M    26.6%

Key Performance Indicators

Additional information

Questions, bug reports, and feature suggestions

Title usage

# Database-to-Dashboard Pipeline

# Data & Technologies

## Web & Search Requests

- *MapReduce* with Hadoop Distributed File System (HDFS)

- Kafka (log buffer) → HDFS via LinkedIn's Camus pipeline

- Includes IP addresses, request referrers, user agents, queries

- Retrieve and aggregate data with HiveQL and User-defined functions (UDFs) written in Java

## Event Logging (User Actions)

- JavaScript event triggers (e.g. clickthroughs, edits, pings)

- Kafka → MariaDB (MySQL)

- Varying degrees of privacy (anonymized → containing PII) & sampling (all vs. 1 in n users)

- Retrieve and aggregate using Structured Query Language

# Event Logging Volume

- Only 0.5% of sessions (users) get selected for anonymous tracking

- Monday, 2 May 2016, on Desktop: 23K sessions (users) recorded

| event | total events | average per user |
|---|---|---|
| search results page | 226K | 9.93 |
| opened a result | 2.8K | 1.37 |
| "I'm alive!" check-in | 14K | 8.04 |
| clicked | 22K | 1.73 |

Mon, 2 May 2016, Desktop:
- ☑ ~237.43M unique users*
- ☑ ~7.73B http requests total
- ☑ Query execution took:
  - ☑ 3.52 days of CPU time
  - ☑ 27 minutes of real time

# My team only uses R

**Development**: devtools, testthat, roxygen2, Rcpp, RStudio + GitHub

**Data Manipulation & Analysis**

🛠 Data wrangling: readr, dplyr, tidyr, data.table, lubridate, xts

🛠 Web analytics: urltools, rgeolocate, iptools, uaparser

🛠 Bayesian analysis of categorical data: BCDA

**Data Visualization & Reporting**

🖍 ggplot2 + ggthemes + ggally

🖍 RMarkdown + knitr for reproducible reporting

🖍 shiny + shinydashboard + dygraphs (for time series)

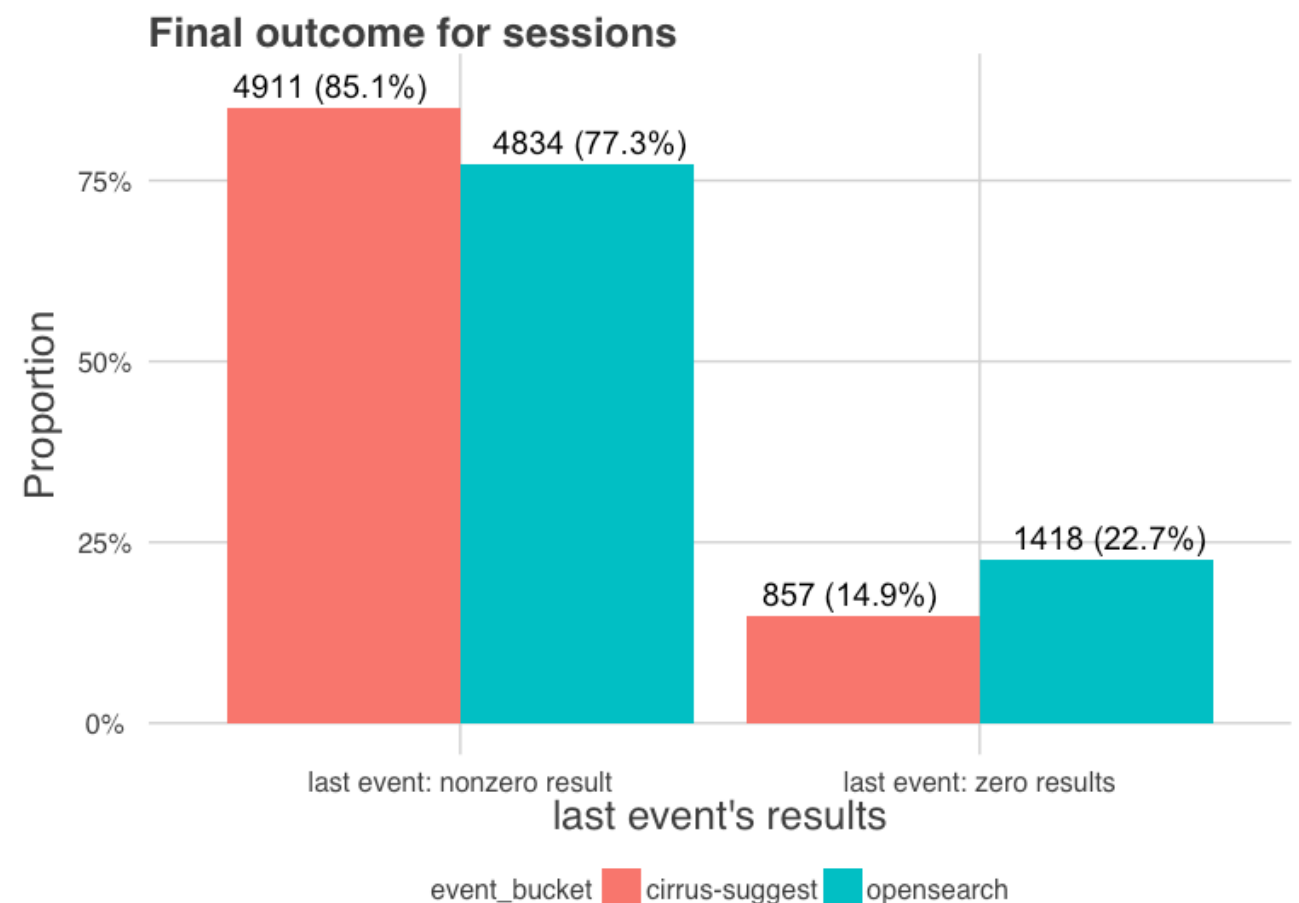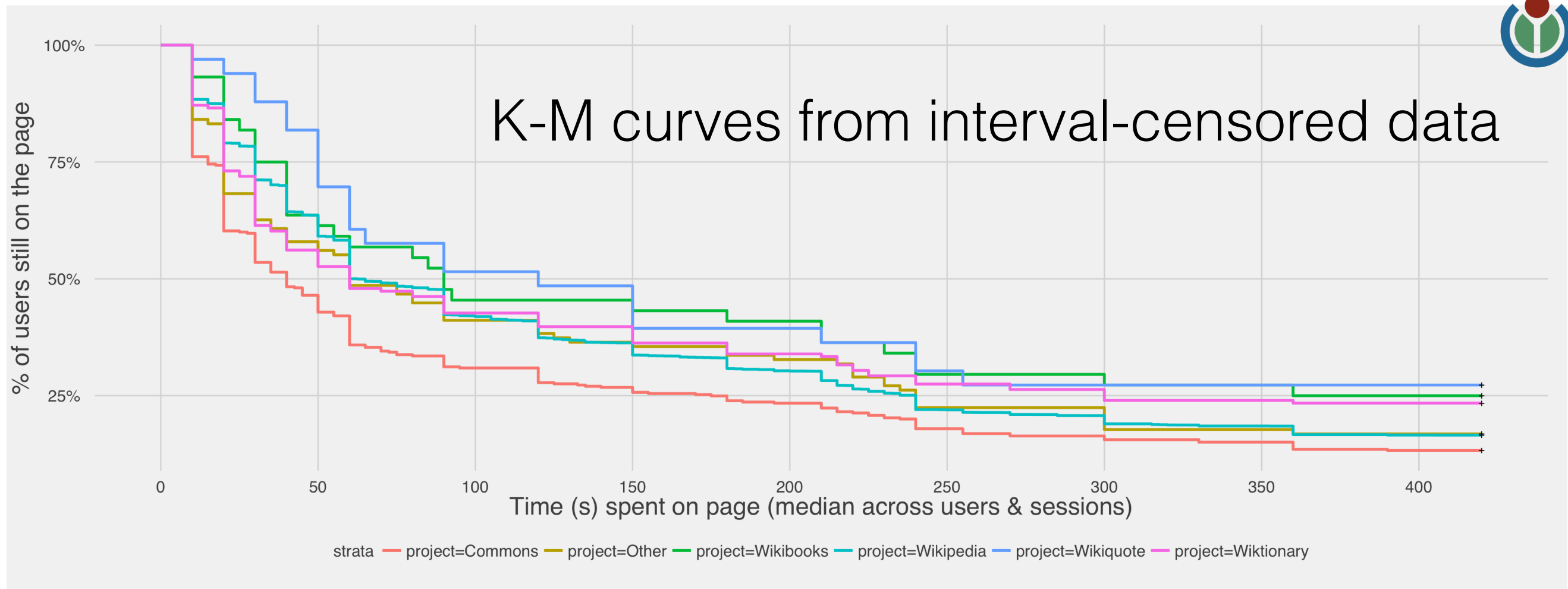🖍 All code (but not data) is public at github.com/wikimedia-research

# Testing Completion Suggester's Impact on Zero Results Rate

| foofigters                                      🔍 |
|---------------------------------------------------|
| Foo Fighters                                      |
| Foo Fighters discography                          |
| Foo Fighters (album)                              |
| Foo Fighters live performances                    |
| Foo Fighters: Sonic Highways                      |
| Foo Fighters: Back and Forth                      |
| Foo Fighters Tour                                 |
| Foo Fighters Greatest Hits                        |
| containing...                                     |
| *foofigters*                                      |

🔍 Tolerance for typos and stop words

🔍 Fewer typos suggested (e.g. searching for *Airton Senna* will properly suggest *Ayrton Senna*)

🔍 Used Bayesian methodology (Beta-Binomial model) to compute difference, relative risk, and credible intervals.

🔍 Deployed as opt-in beta → deployed to production!

**Final outcome for sessions**

4911 (85.1%)  
4834 (77.3%)  
857 (14.9%)  
1418 (22.7%)

Proportion — 75%, 50%, 25%, 0%

last event: nonzero result   last event: zero results  
last event's results

event_bucket ▮ cirrus-suggest ▮ opensearch

K-M curves from interval-censored data

% of users still on the page

Time (s) spent on page (median across users & sessions)

strata — project=Commons — project=Other — project=Wikibooks — project=Wikipedia — project=Wikiquote — project=Wiktionary

# Human-Computer Interaction Research: User Satisfaction

**Scenario 1**: Find out Barack Obama's birth date.

**Scenario 2**: Learn about Black Lives Matter movement.

**Scenario 3**: I search for Lin-Manuel Miranda's *Hamilton*, stay there for 5s, and then decide I want to read about Aaron Burr.

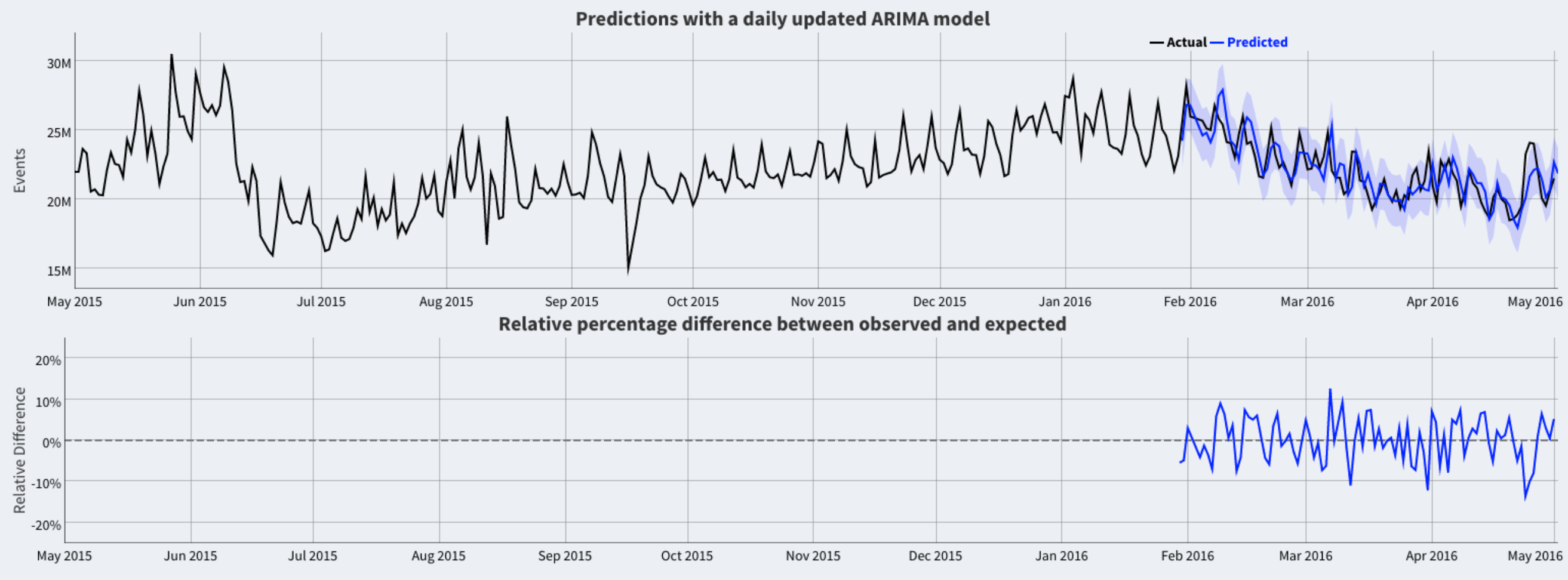**Scenario 4**: I search a vague term with no specific task in mind.

# **Project**: daily forecasts to detect anomalies & abuse of service



Prototype using ARIMA: http://discovery-experimental.wmflabs.org/forecast/

Future work will probably focus on Bayesian Structural Time Series Modeling

# Thank you!

Got additional questions?

**Work email**:
mpopov@wikimedia.org

**Personal email**:
mikhail@mpopov.com

**Twitter**: @bearloga

We're hiring data analysts!

See the jobs page at WMF:
https://wikimediafoundation.org/wiki/Work_with_us