

Statistical Performance Indicators: Data Use Text Mining

SPI Team

2020-07-13



WORLD BANK GROUP

Introduction

- Goal is to examine the use of text mining of National Development Plans (NDPs) & National Poverty Reduction Strategies for the SPI Data Use dimension
- In this presentation, will discuss:
 - Source for National Development Plan (FAOLEX)
 - Dictionary based approach for text mining
 - Initial results from text mining NDPs
- Methods and data rely heavily on Paris21 approach to text mining

Background on Data Use Dimension

- Indicator intended to be Proxy for the intensity of use of official statistics by each respective user segment in country.

Data Use dimension contains 5 indicators:

- Data use by national legislature
- Data use by national executive branch
- Data use by civil society
- Data use by academia
- Data use by international organizations

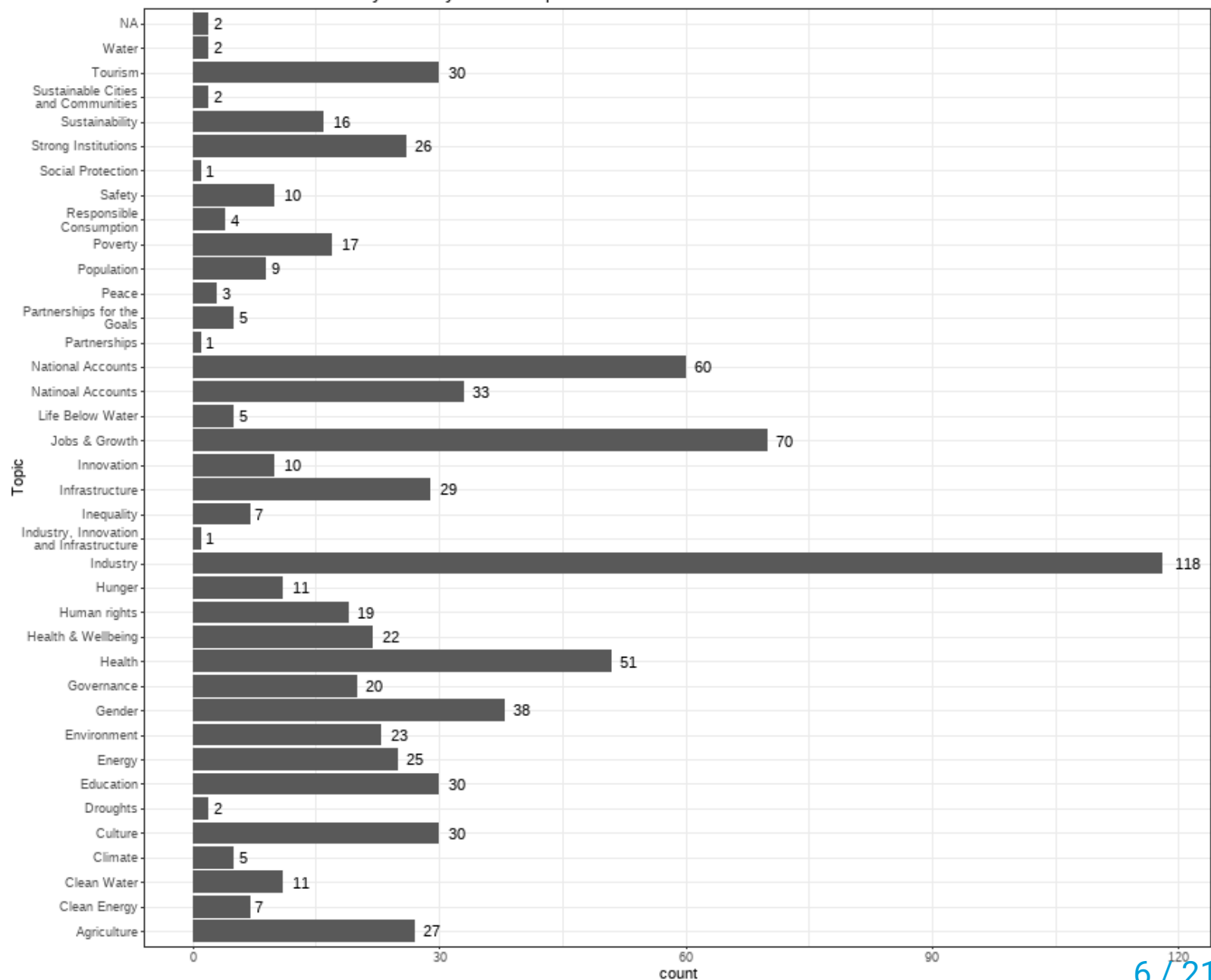
Data Use Indicator Methodology

- Indicator for each of 5 areas is formed by:
 1. Searches the site for references to keywords related to each SDG goal
 2. Searches the site for the subset of references that also include the word statistics
 3. Calculates the ratio between the two
- Algorithm captures density of discussion of each topic in each organization assessed

Dictionary Approach to Text Mining

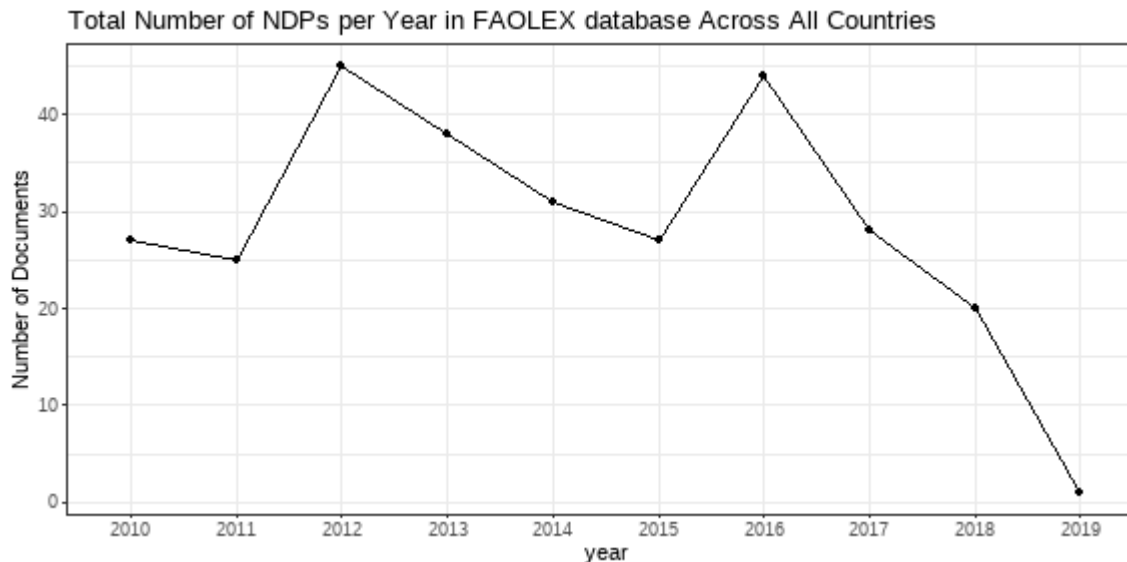
- Paris21 has developed set of text mining tools for analyzing NDPs
- Created dictionary of 782 search terms related to different topic areas
- I categorized keywords into 38 topic areas

Counts of the Number of Keywords by Overall Topic Area

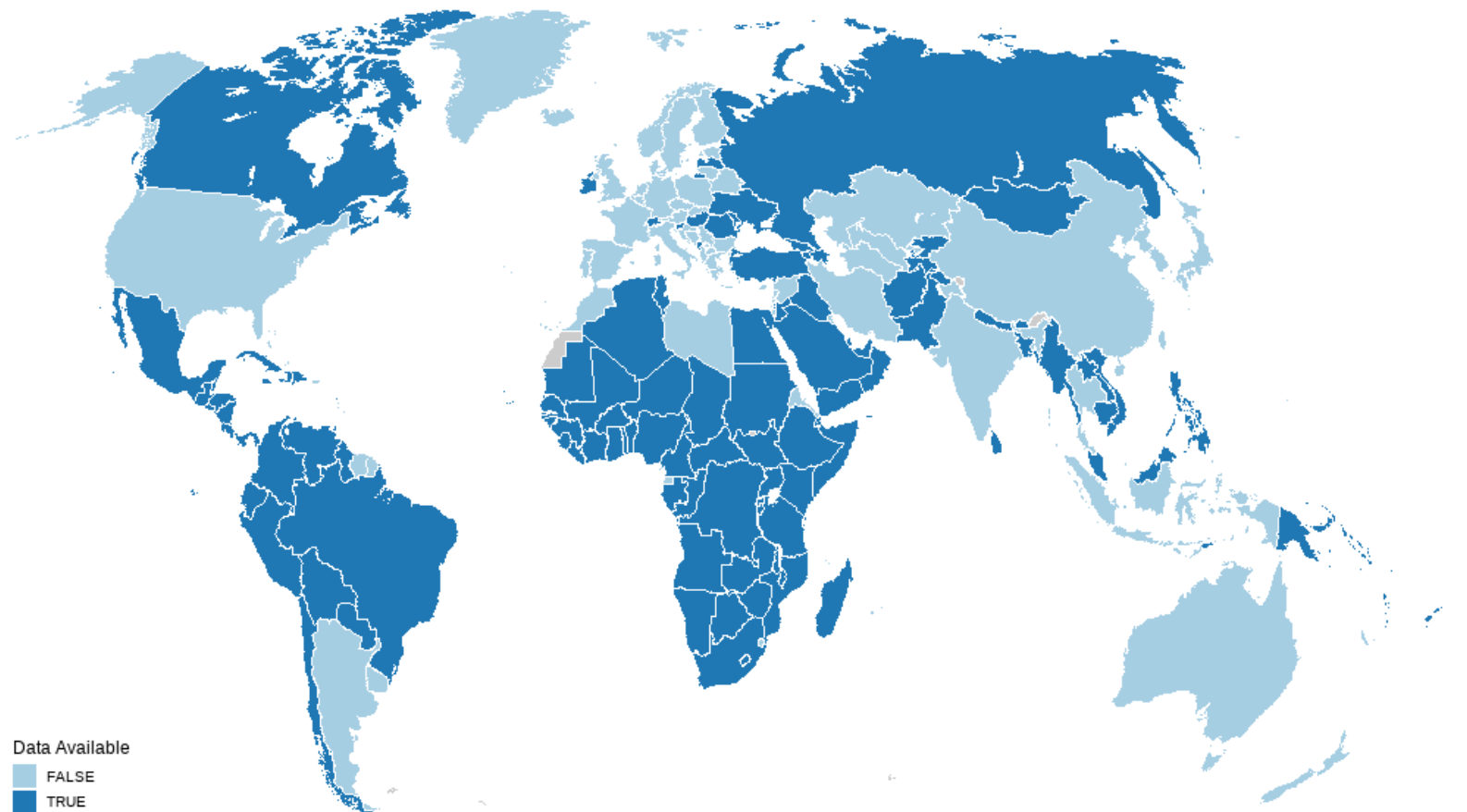


Data Sources for NDPs

- Paris21 used the FAOLEX database:
 - Database of national legislation, policies and other
 - Updated with an average of 8,000 new entries per year
 - **Includes 286 NDPs since 2010 from 123 Countries**
 - Limited number of documents per year



Availability of National Development Plan or Poverty Reduction Plan Documents by Country since 2010 from FAOLEX



Source: FAOLEX

Text Mining

- Using set of dictionary terms, we can look for sentences in NDPs containing these terms
- Then we can check to see if a number is cited in that sentence, as a measure of data use
- As an example going forward, we will consider the "National Strategic Development Plan 2014-2018" from Cambodia

Example sentence 1

- The sentence in the document: "in contrast, malaysia and korea have this proportion in the range of 5% of gdp."
- Contains the keyword: **gdp**
- Did the sentence use data according to model?: **TRUE**

Example sentence 2

- The sentence in the document: "¾ established 14 border liaison offices (blo's) along the border of the country and the mekong river to control and curb drug trafficking and precursor chemicals across the border."
- Contains the keyword: **traffic**
- Did the sentence use data according to model?: **TRUE**

Example sentence 3

- The sentence in the document: "some citizens did not take advantage of employment services, and there was little cooperation between the concerned parties such as trainers, public citizens, and authorities at all levels. x limited cooperation between service providers, ministries, agencies, local authorities, and employers, and the national job and employment agency."
- Contains the keyword: **services**
- Did the sentence use data according to model?: **FALSE**

Example sentence 4

- The sentence in the document: "¾ study the root causes of the problems hindering the waterworks' revenue from covering their expenditures and try all possible solutions to resolve the problems."
- Contains the keyword: **revenue**
- Did the sentence use data according to model?: **FALSE**

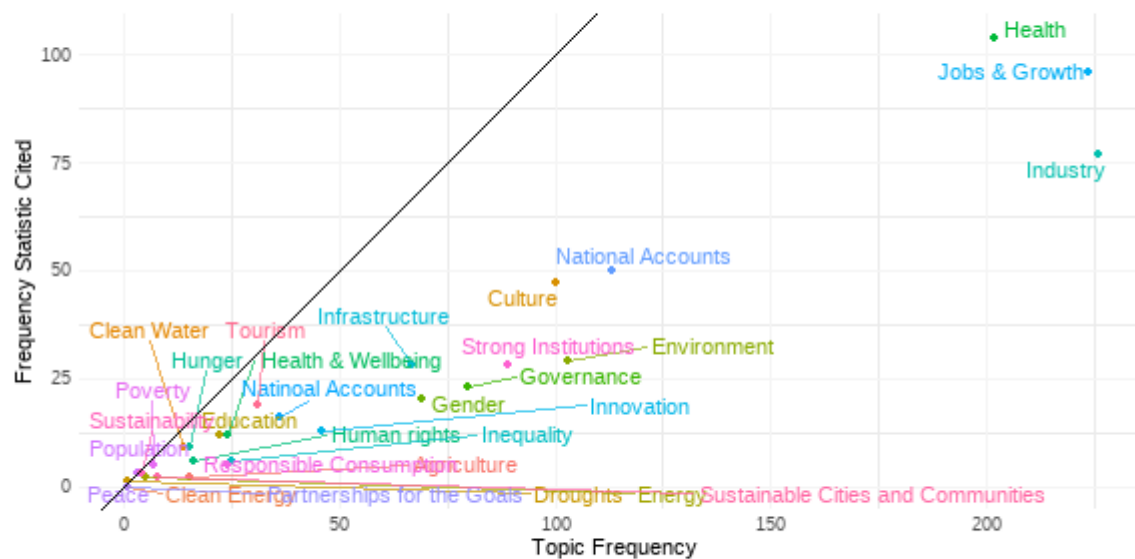
Example sentence 5

- The sentence in the document: "cambodia's gdp per capita is approximated at usd 1,036 in 2013 (estimated)."
- Contains the keyword: **gdp**
- Did the sentence use data according to model?: **TRUE**

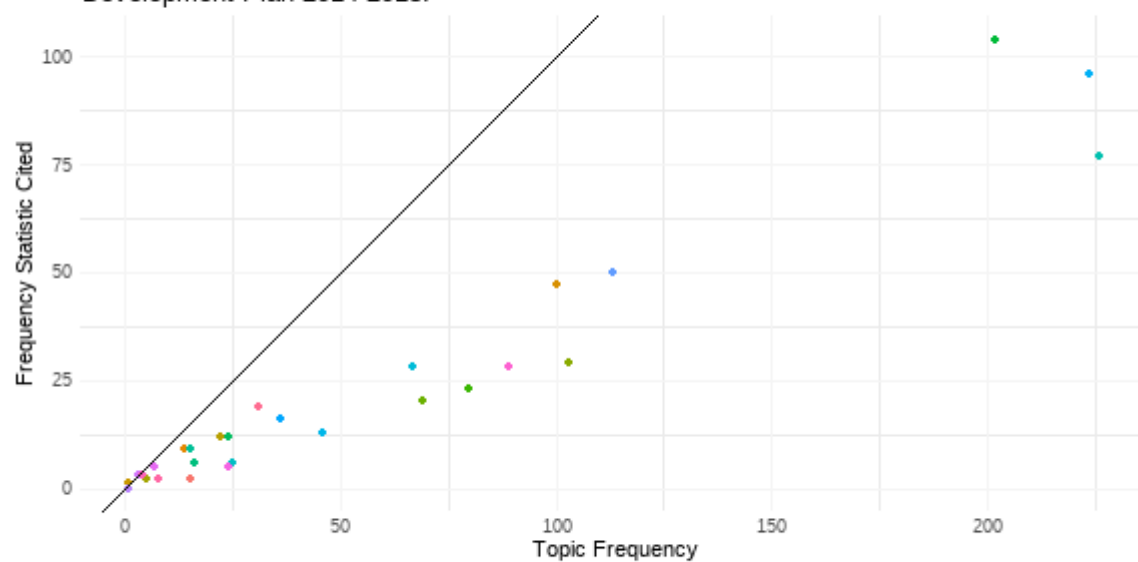
Terms Usage & Data Mentions

- Next, we will plot the number of times a topic was mentioned (on x-axis) against the of times statistics were cited according to our algorithm

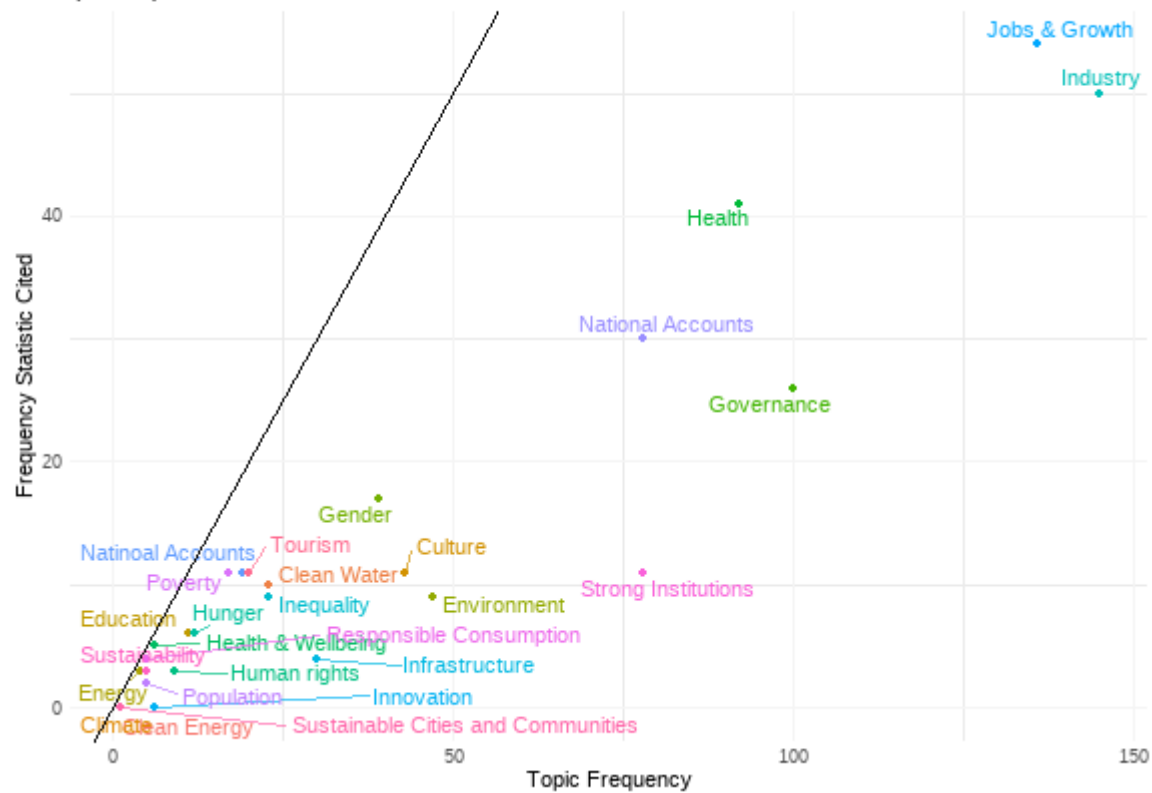
Statistics Cited per Topic Cited - Cambodia - 2014 - National Strategic Development Plan 2014-2018.



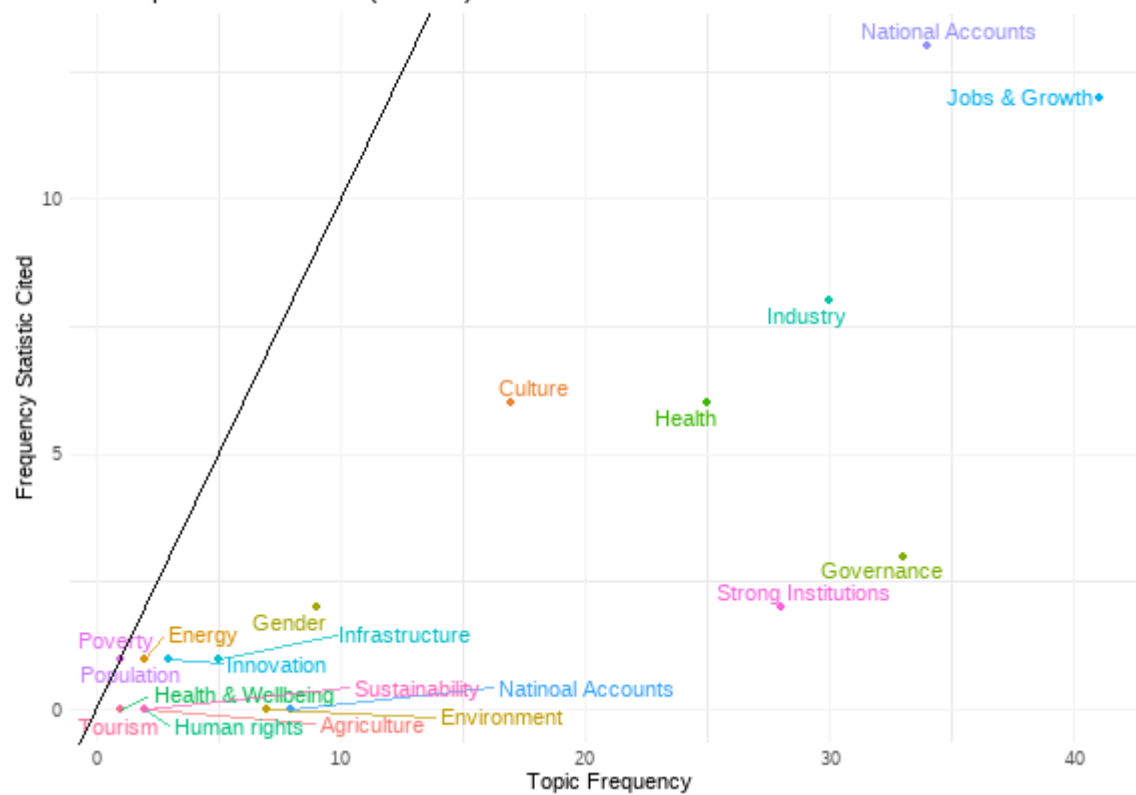
Statistics Cited per Topic Cited (no labels) - Cambodia - 2014 - National Strategic Development Plan 2014-2018.



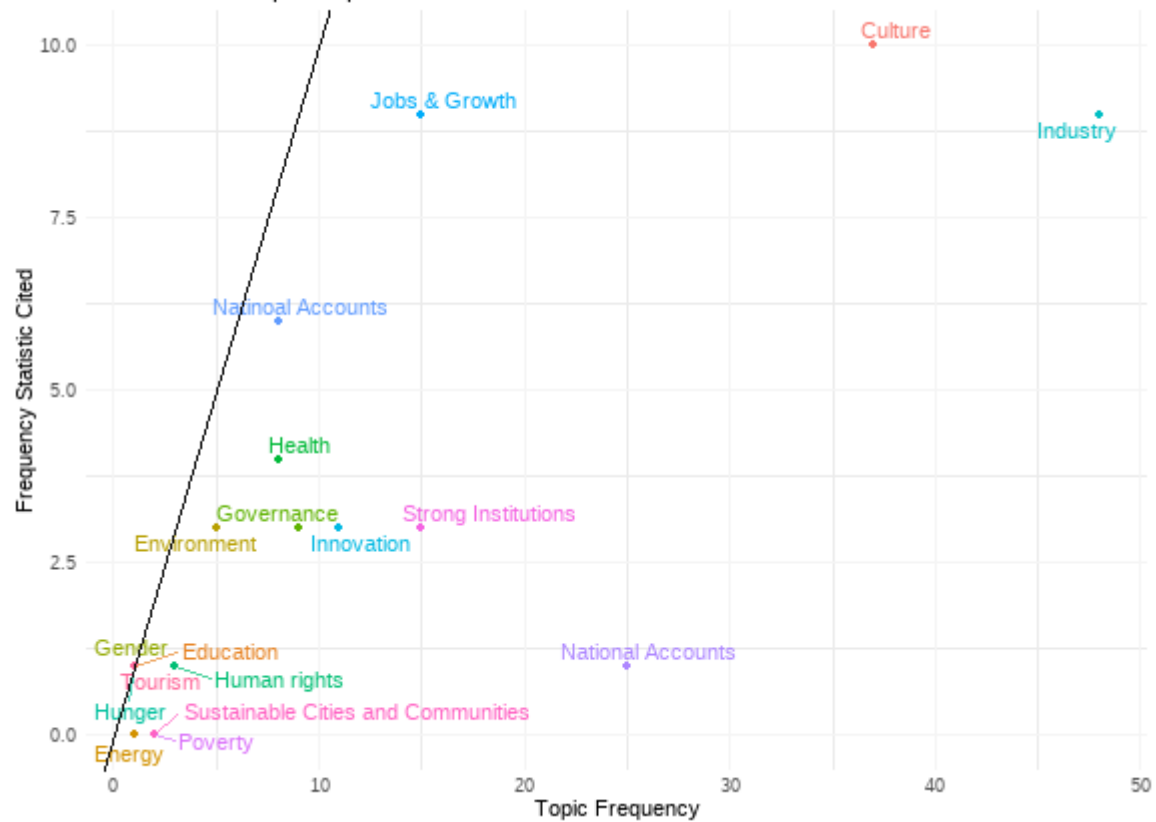
Statistics Cited per Topic Cited - Guinea - 2013 - Poverty Reduction Strategy Paper (PRSP).



Statistics Cited per Topic Cited - Afghanistan - 2017 - Afghanistan National Peace and Development Framework (ANPDF) 2017 to 2021.



Statistics Cited per Topic Cited - Brazil - 2010 - Plan for Brazil 2022.



Statistics Cited per Topic Cited - Tanzania - 2016 - National Five Year Development Plan

