

데이터 시각화 ggplot 뛰어넘기

김영진

SK텔레콤

[발표 자료 웹 버전 보기](#)

데이터 시각화



ggplot2

뛰어넘기

김영진

김영진

- 사회학 박사 (2015)
- SKTelecom
Data Analytics Group



데이터 시각화?

이 자료는 RUCK에 공개될 예정입니다.

정보를 보기 좋게 표현하는 것?



MS 4551

Account of grain products, bread, beer, butter oil. Sumer, 32nd c. BC



John Tukey (1915-2000)

LIFE

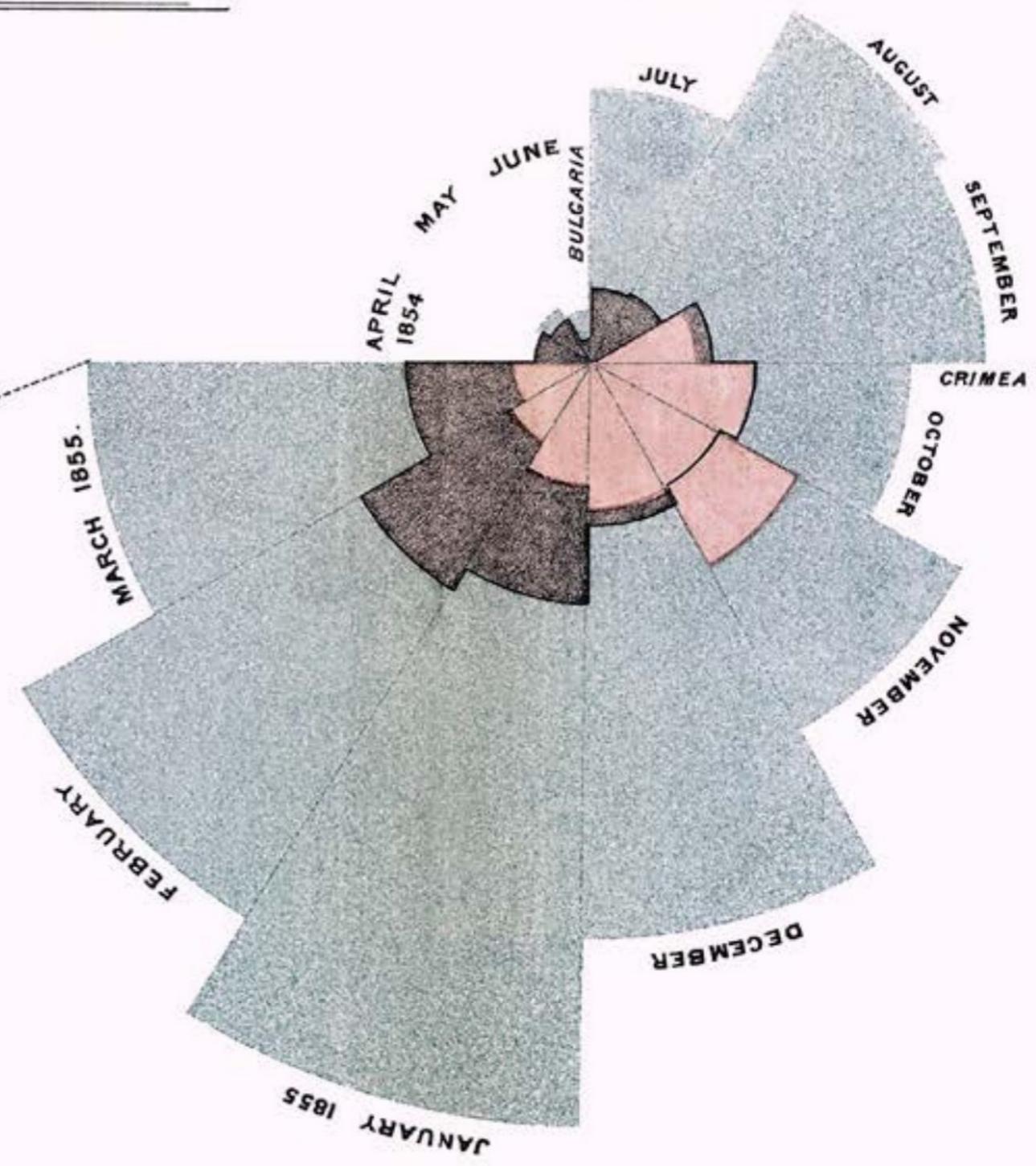
Exploratory Data Analysis 탐색적 데이터 분석

“데이터를 **탐색**하면서,
질문과 가설을 만들고,
분석의 **뼈대**를 만드는 **과정**”

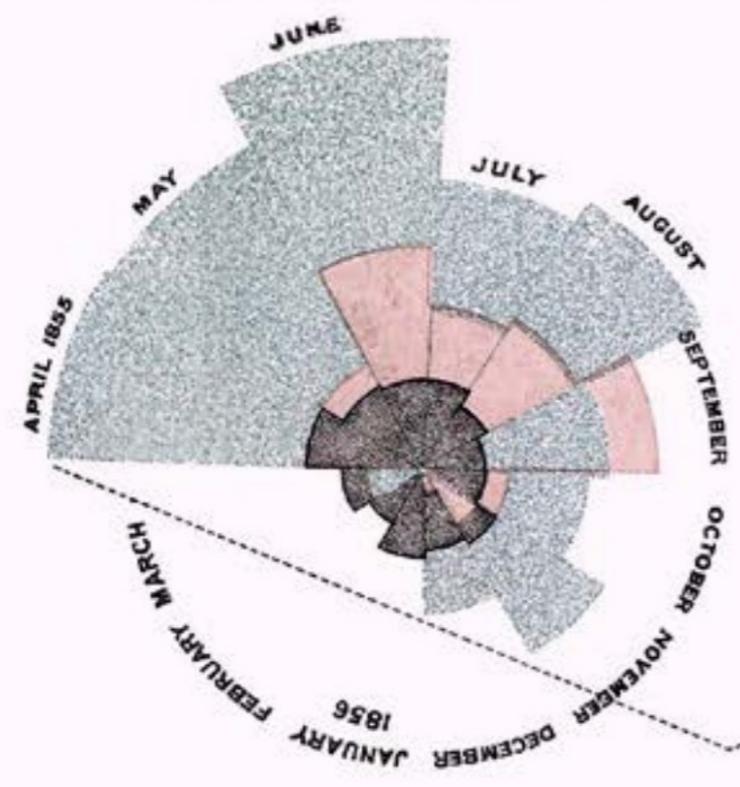
“형사가 **수사**를 할 때
여러 가지 **단서**들로
가설을 세우고
증거를 확립하는 것과 유사”

DIAGRAM OF THE CAUSES OF MORTALITY IN THE ARMY IN THE EAST.

1.
APRIL 1854 TO MARCH 1855.



2.
APRIL 1855 TO MARCH 1856.



The Areas of the blue, red, & black wedges are each measured from the centre as the common vertex.

The blue wedges measured from the centre of the circle represent area for area the deaths from Preventible or Mitigable Zymotic diseases, the red wedges measured from the centre the deaths from wounds, & the black wedges measured from the centre the deaths from all other causes.

The black line across the red triangle in Nov^r 1854 marks the boundary of the deaths from all other causes during the month.

In October 1854, & April 1855, the black area coincides with the red; in January & February 1856, the blue coincides with the black.

The entire areas may be compared by following the blue, the red & the black lines enclosing them.



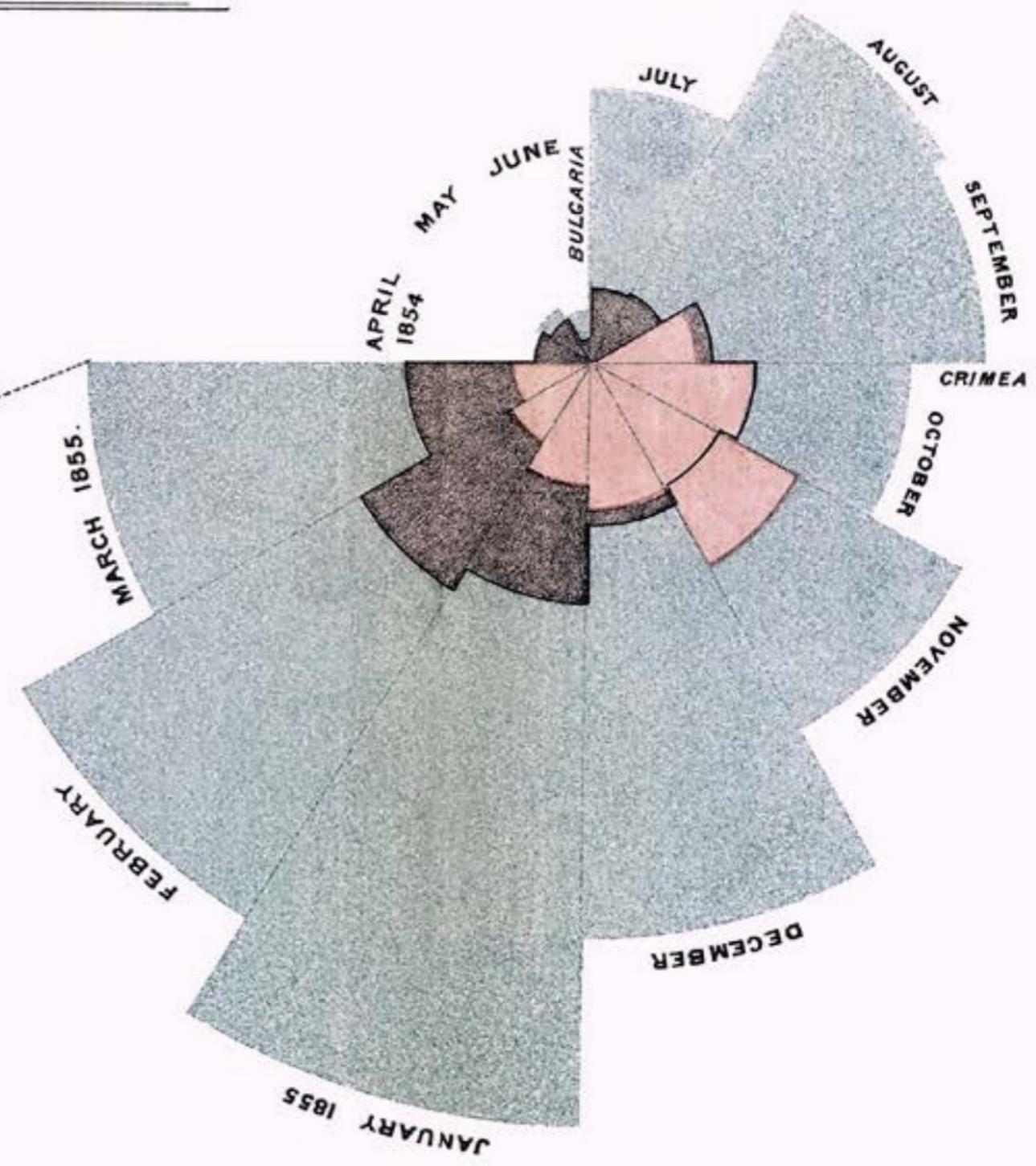
Florence Nightingale

플로렌스 나이팅게일
(1820~1910)

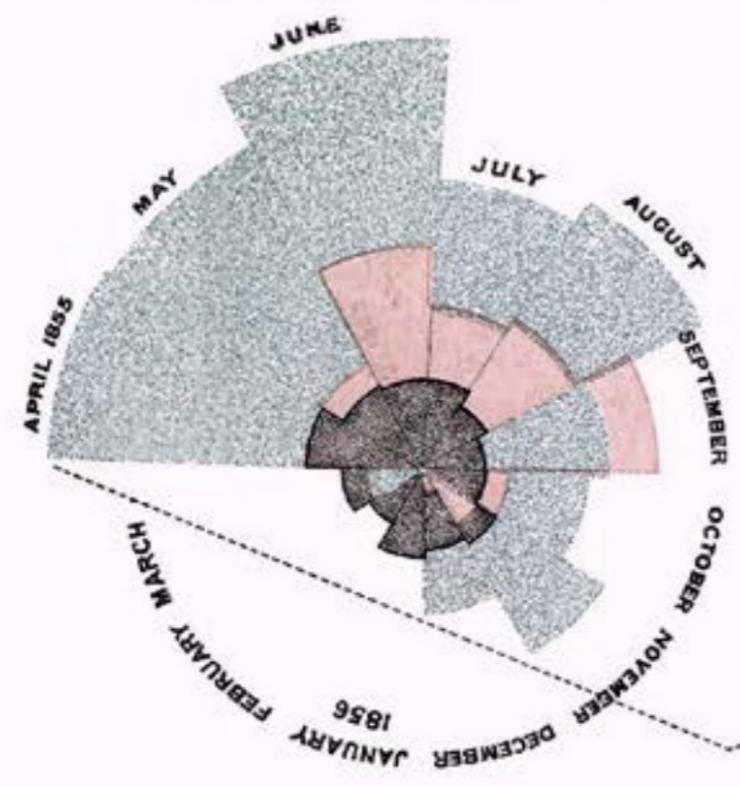
“나는 간호를 받는 사람들의
안녕을 위해 헌신하겠습니다.”

DIAGRAM OF THE CAUSES OF MORTALITY IN THE ARMY IN THE EAST.

1.
APRIL 1854 TO MARCH 1855.



2.
APRIL 1855 TO MARCH 1856.



The Areas of the blue, red, & black wedges are each measured from the centre as the common vertex.

The blue wedges measured from the centre of the circle represent area for area the deaths from Preventible or Mitigable Zymotic diseases, the red wedges measured from the centre the deaths from wounds, & the black wedges measured from the centre the deaths from all other causes.

The black line across the red triangle in Nov^r 1854 marks the boundary of the deaths from all other causes during the month.

In October 1854, & April 1855, the black area coincides with the red; in January & February 1856, the blue coincides with the black.

The entire areas may be compared by following the blue, the red & the black lines enclosing them.

이 도표 한 장으로,
영국 군대는 전투력을 위해 전장에서
위생을 최우선으로 변경하였음

百聞而不如一見

십부문이십이일견



March on Moscow - Charles Minard

Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite. Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Ségur, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk et Mohilow et ont rejoint vers Orscha et Witebsk, avaient toujours marché avec l'armée.

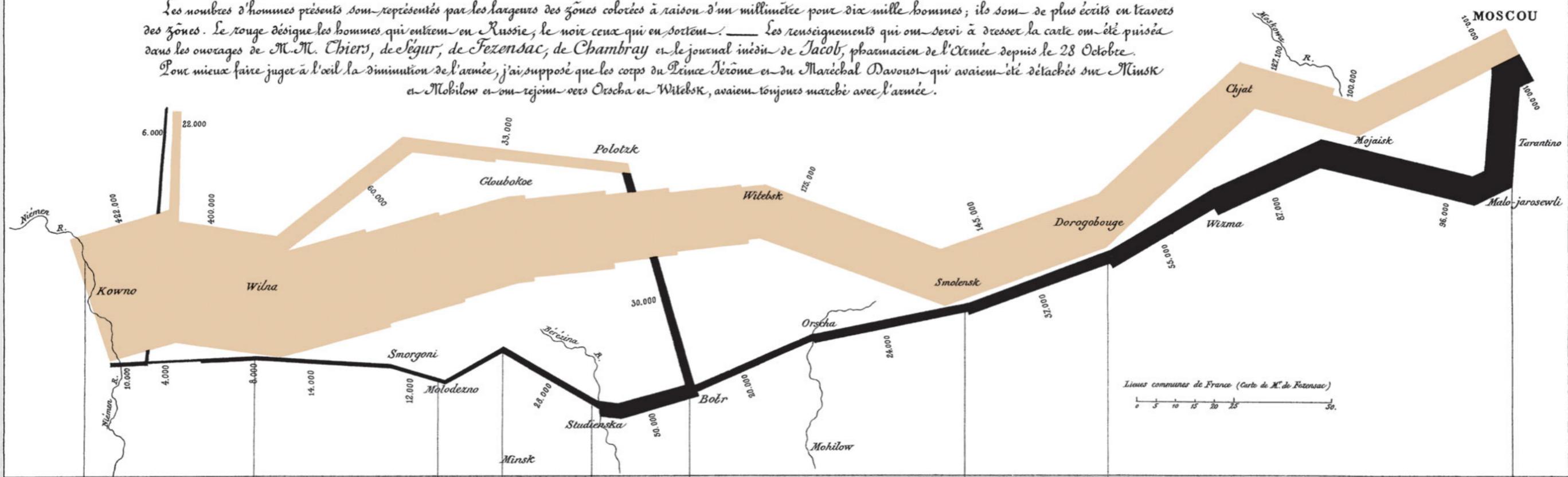
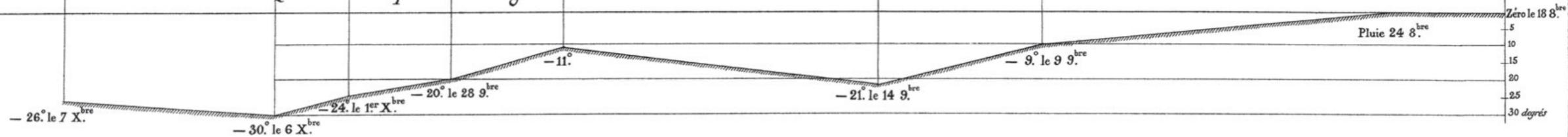


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.



Les Cosaques passent au galop le Niémen gelé.

Autog. par Regnier, 8. Par. S^{te} Marie S^t G^{ermain} à Paris.

Imp. Lith. Regnier et Dourdet.

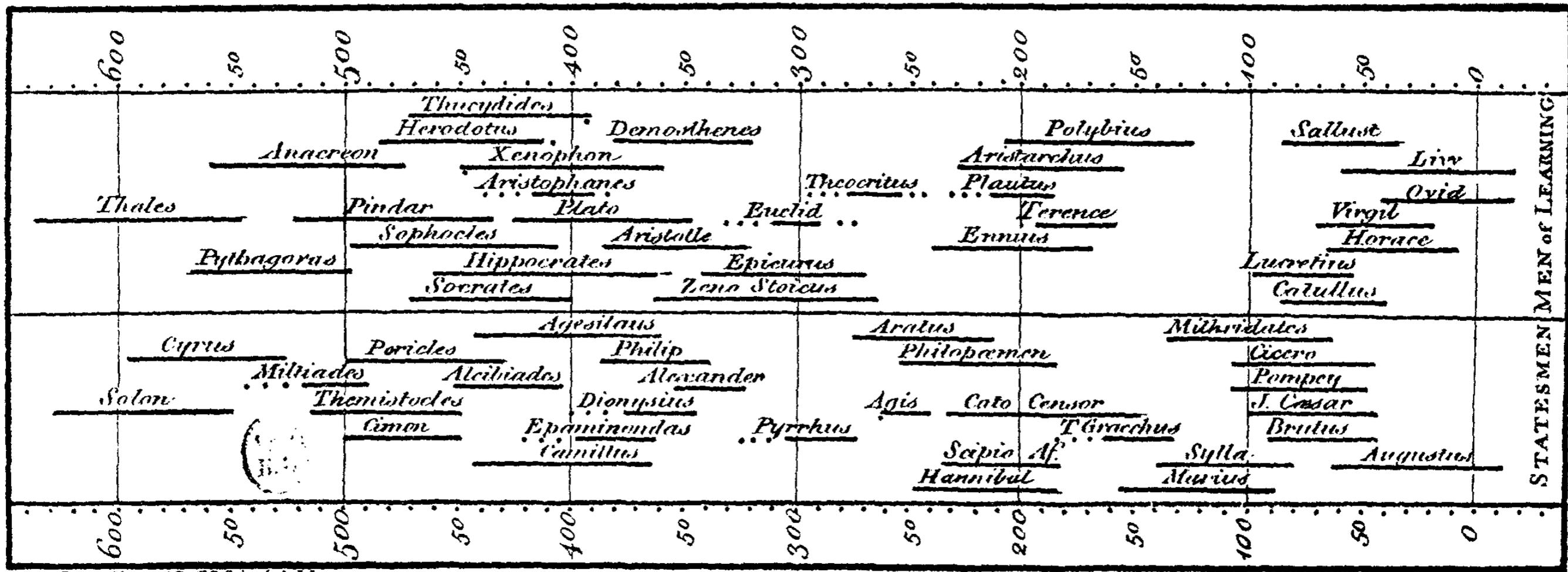
<http://ecowest.org/tag/edward-tufte/>
<https://www.edwardtufte.com/tufte/posters>

하나의 **그래프**가 한 권의 책보다
더 나은 **설명**과 **이해**

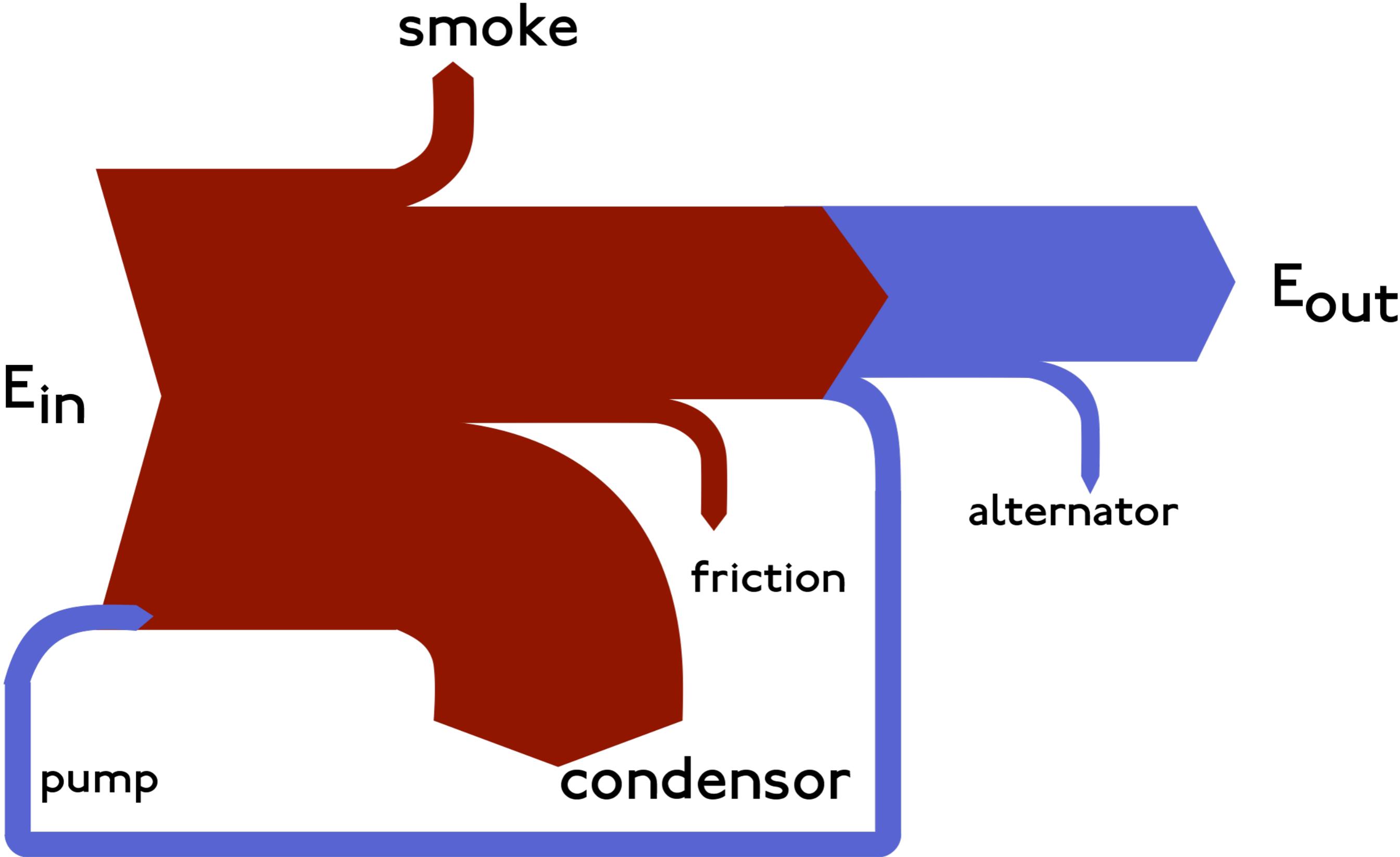
나폴레옹의 행진은 독창적 디자인?

A Chart of Biography (1765)

A Specimen of a Chart of Biography.



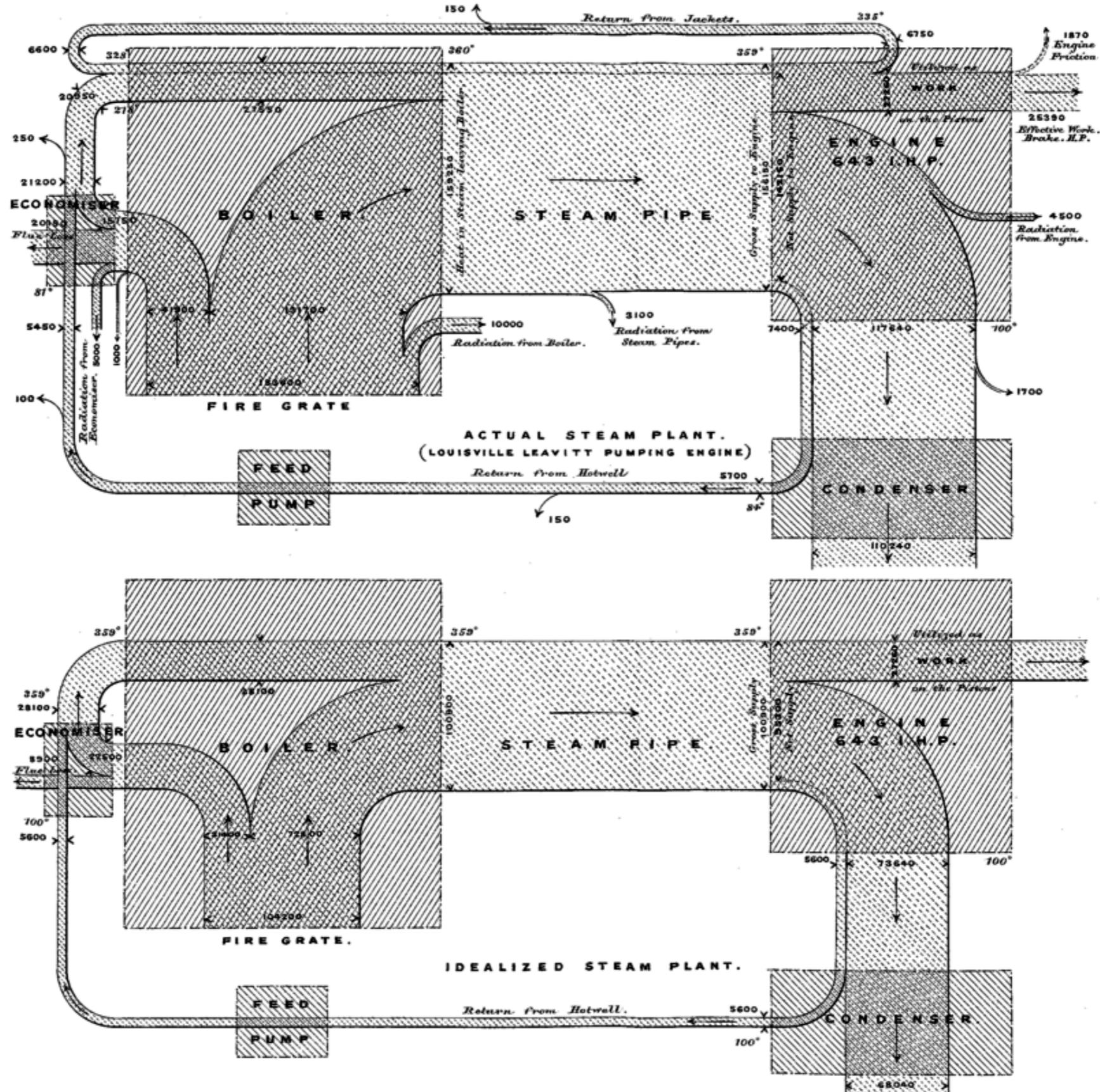
Sankey diagram



Sankey's original 1898

THE THERMAL EFFICIENCY OF STEAM-ENGINES.

PLATE 5.



March on Moscow (1869)- Charles Minard

Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813.
 Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite. Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. — Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M. M. Chiers, de Ségur, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre. Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davout qui avaient été détachés sur Minsk et Mohilow et ont rejoint vers Orscha et Witebsk, avaient toujours marché avec l'armée.

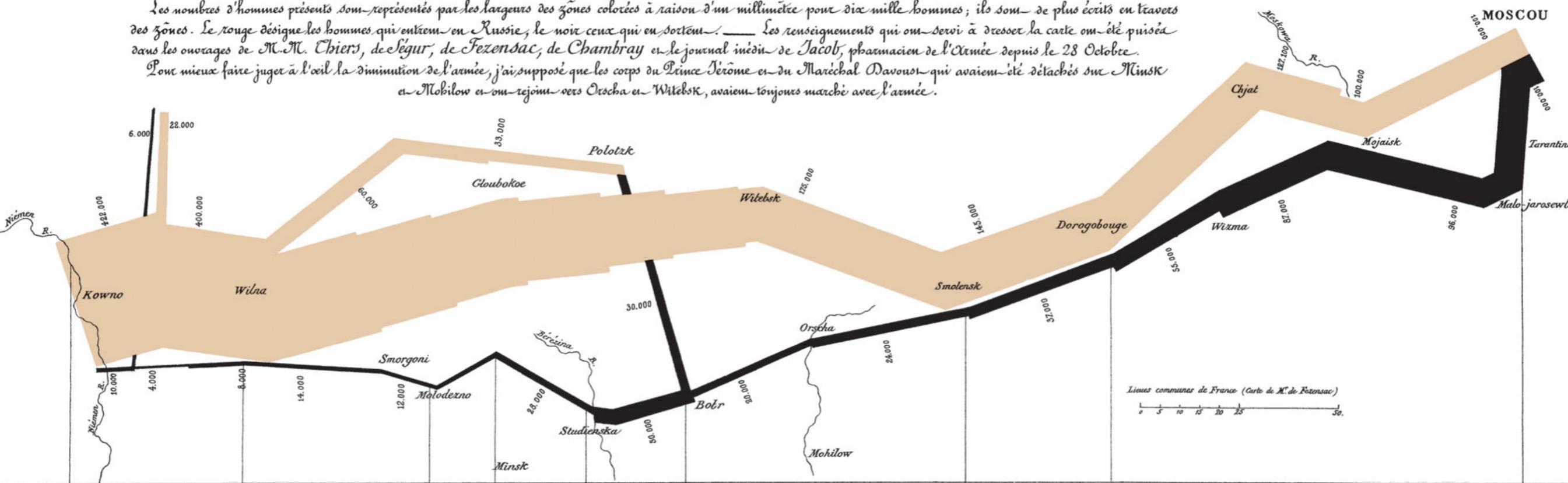
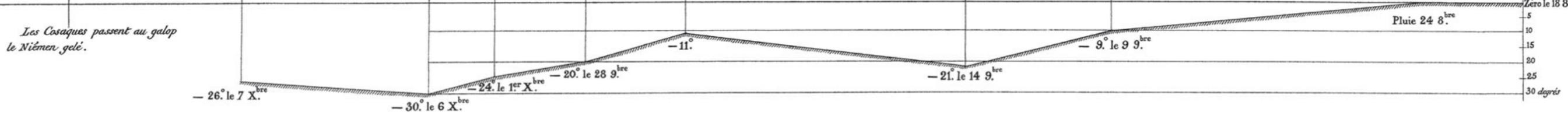


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.



Autog. par Regnier, 8. Par. S^{te} Marie S^t G^{ermain} à Paris.

Imp. Lith. Regnier et Douvret.

<http://ecowest.org/tag/edward-tufte/>
<https://www.edwardtufte.com/tufte/posters>

시간의 흐름과 이벤트의 강도를 표현하고자 한 그래프

Streamgraph

February 23, 2008

FEEDBACK

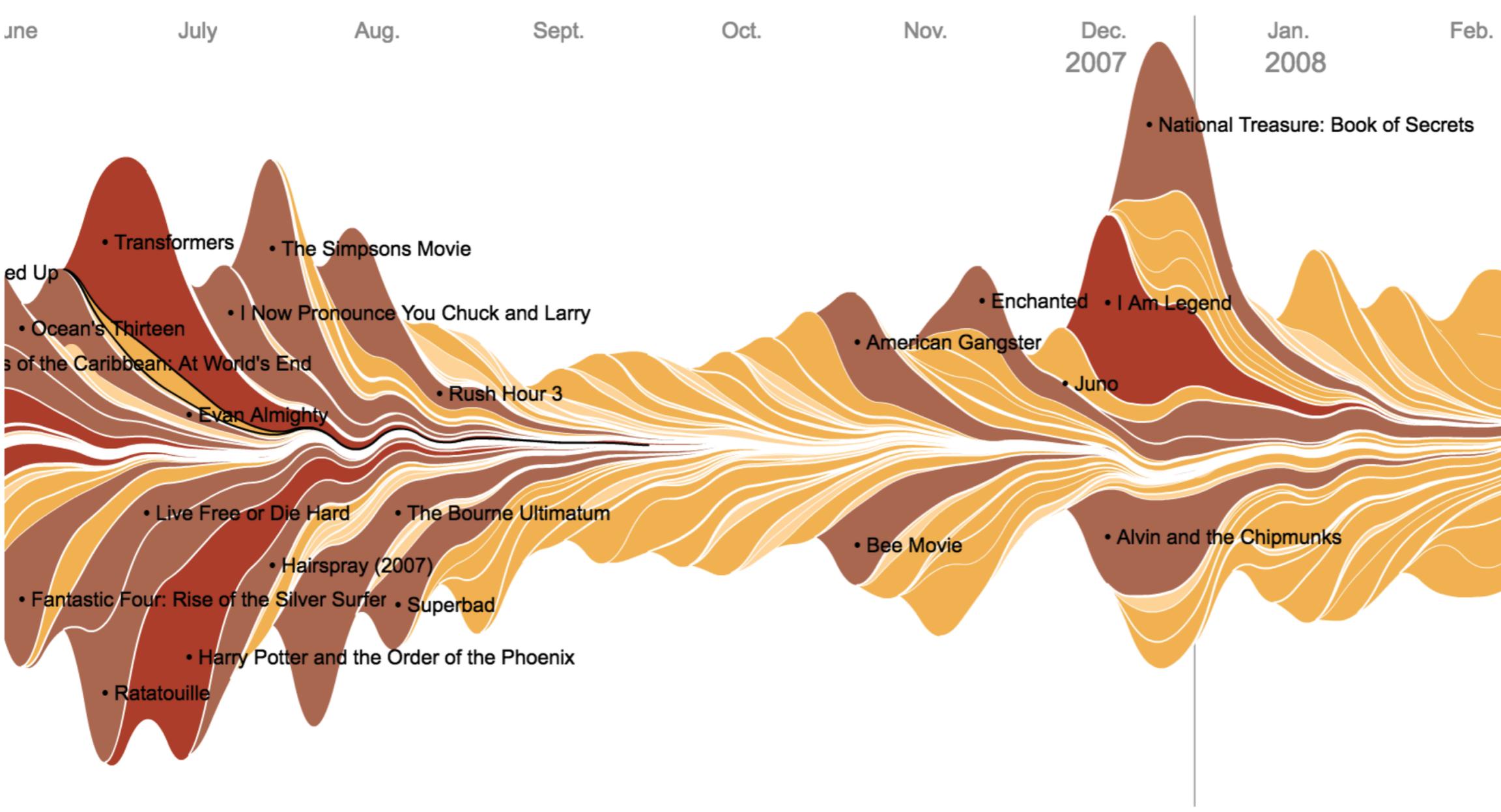
The Ebb and Flow of Movies: Box Office Receipts 1986 – 2008

Summer blockbusters and holiday hits make up the bulk of box office revenue each year, while contenders for the Oscars tend to attract smaller audiences that build over time. Here's a look at how movies have fared at the box office, after adjusting for inflation.

Find Movie

Go

June July Aug. Sept. Oct. Nov. Dec. 2007 Jan. 2008 Feb.

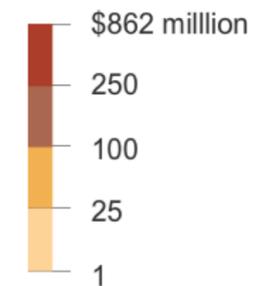


Each shape shows how one film did at the box office.



← **Width** →
shows longevity

The **area** of the shape (and its **color**) corresponds to the film's total domestic gross, through Feb. 21

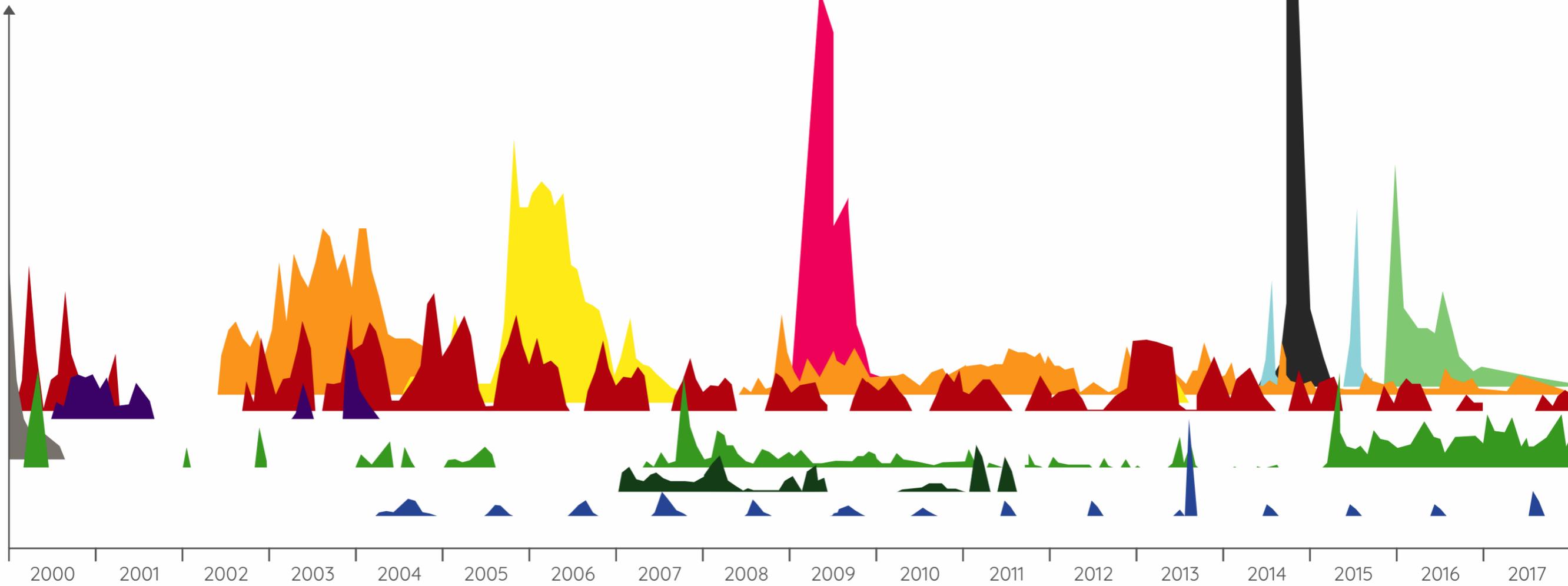


layered density graph

Mountains Out of Molehills

A timeline of media-inflamed fears

INTENSITY
(no. of news media mentions)

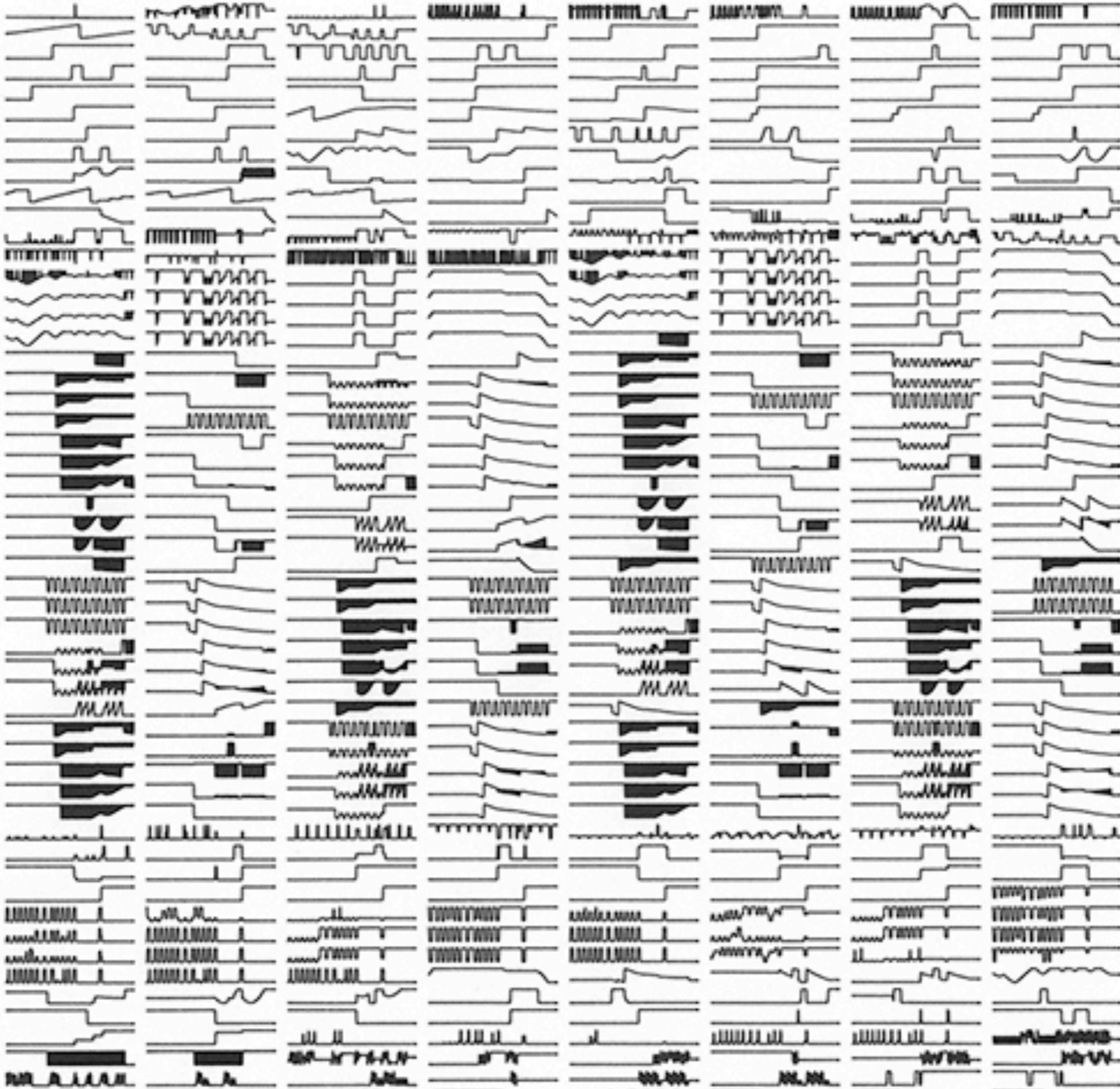


ASTEROIDS BIRD FLU EBOLA KILLER WASPS MAD COW DISEASE MERS MILLENNIUM BUG SARS SWINE FLU VACCINES & AUTISM VIOLENT VIDEO GAMES ZIKA

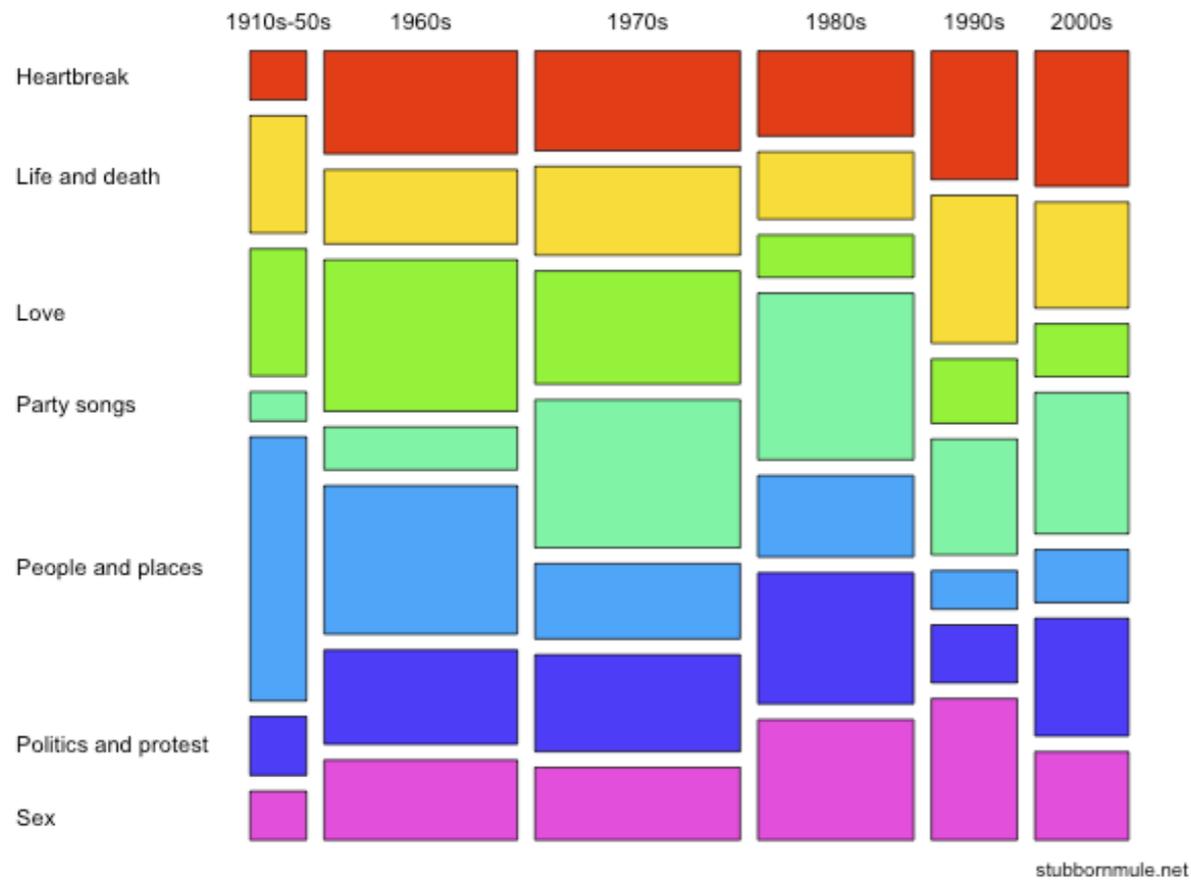
Concept & Design: [David McCandless](#) // Design & Code: [Fabio Bergamaschi](#)

align to baseline scale to fit Ebola scale to deaths

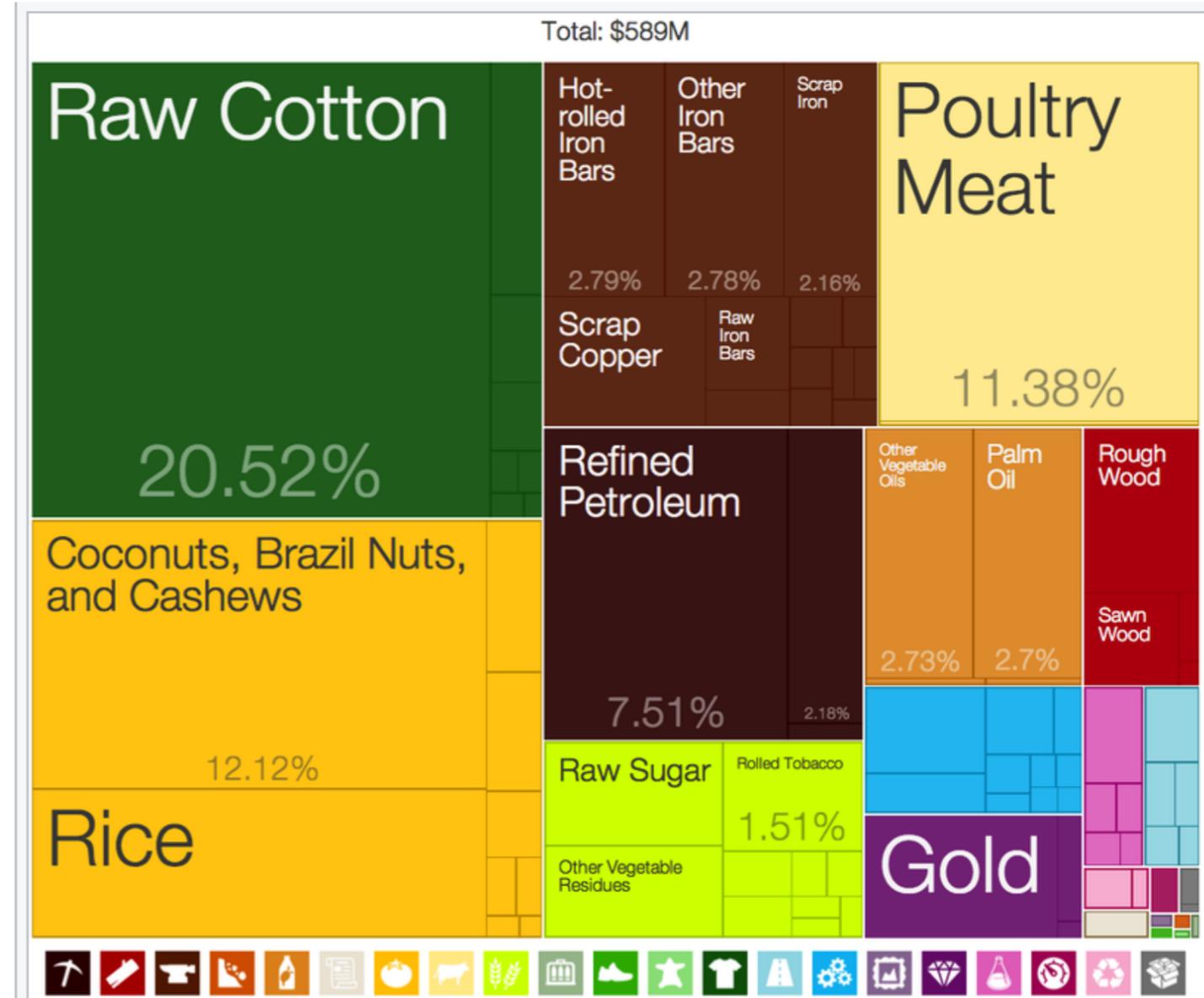
Spikelines



Mosaic plot



TreeMap



Treemap of Benin's exports by product category, 2009. The Product Exports Treemaps are one of the most recent applications of these kind of visualizations, developed by the Harvard-MIT [Observatory of Economic Complexity](#)



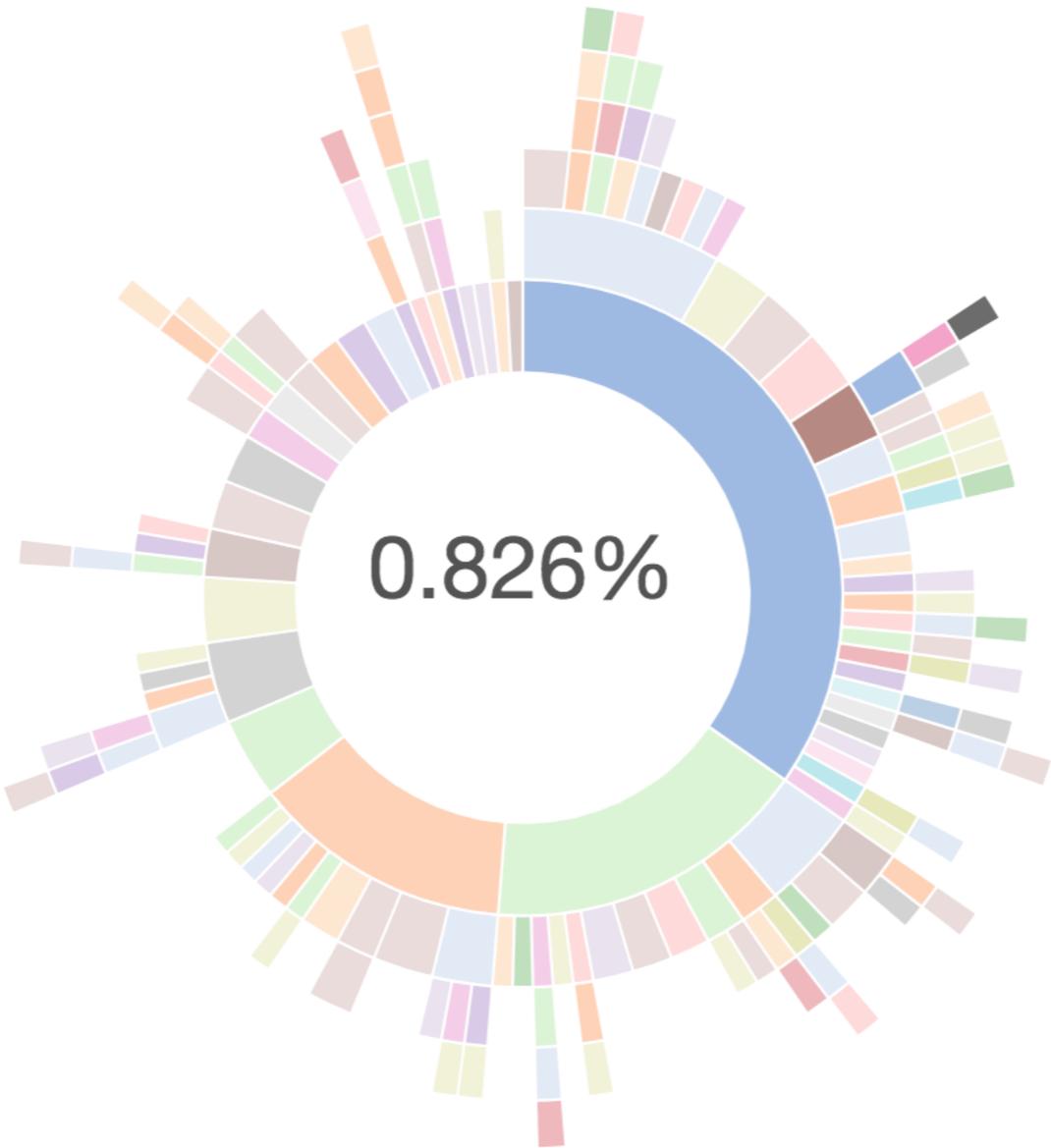
Marimekko

Sunburst Chart

Multi-level Pie Chart, Belt Chart, Radial Treemap.



Legend



진짜 데이터로 R을 이용하여 위에서 이야기한 Graph 그려보기

아파트 실거래가 데이터

<http://rtdown.molit.go.kr/>

The screenshot shows the homepage of the '실거래가 공개시스템' (Real Estate Transaction Price Disclosure System) on the Ministry of Land, Urban Planning and Construction's website. The page features a search section with various filters for transaction dates, property types, and locations. A '다운로드' (Download) button is visible at the bottom of the search area.

국토교통부 실거래가 공개시스템 실거래가 다운로드 일자리가 성장이고 복지입니다.

실거래가 다운로드 > 조건별 자료제공

HOME 실거래가 다운로드 > 조건별 검색

조건별 검색 국토교통부 실거래가 공개시스템을 이용하시면 쉽고 편리하게 이용하실 수 있습니다.

<조건별 자료제공 이용시 유의사항>

□ 본 서비스에서 제공하는 정보는 법적인 효력이 없으므로 참고용으로만 활용하시기 바라며, 외부 공개시에는 반드시 신고일 기준으로 집계되는 공식통계를 이용하여 주시기 바랍니다.

□ 신고정보가 실시간 변경, 해제되어 제공시점에 따라 공개건수 및 내용이 상이할 수 있는 점 참고하시기 바랍니다.

□ 본 자료는 계약일 기준입니다. (※ 7월 계약, 8월 신고건 → 7월 거래건으로 제공)

* 주택매매 거래는 부동산 거래신고 등에 관한 법률 제3조에 따라 계약일로부터 60일 이내 신고토록 하고 있습니다.

계약일자 20180901 ~ 20180930 파일구분 EXCEL

실거래가구분 아파트(매매) 주소구분 지번주소 도로명주소

시도 전체 시군구 전체 읍면동 전체 전체

면적 전체 금액선택 (만원) ~ (만원)

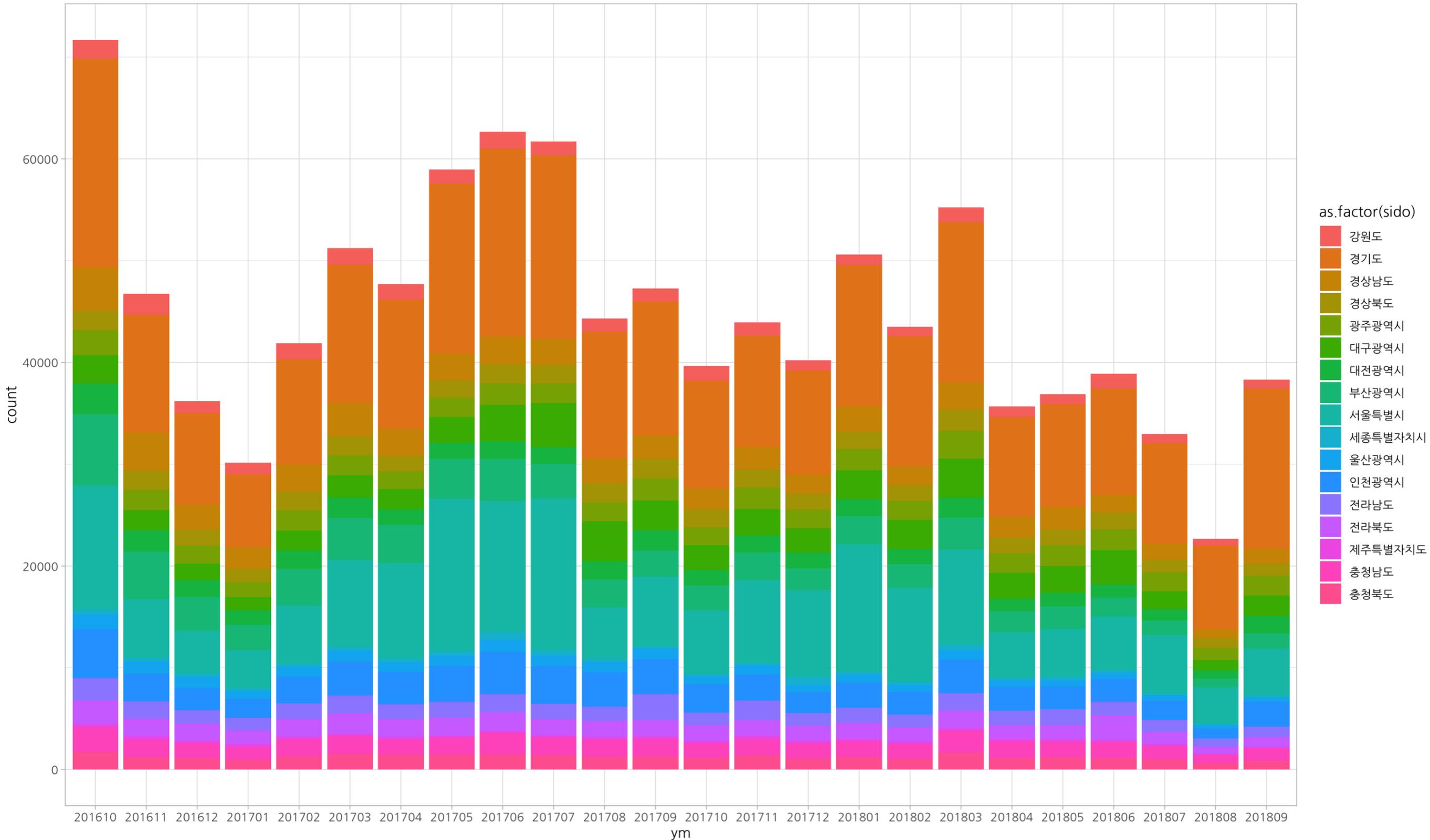
다운로드

2016. 10 ~ 2018. 9

1,078,890 rows



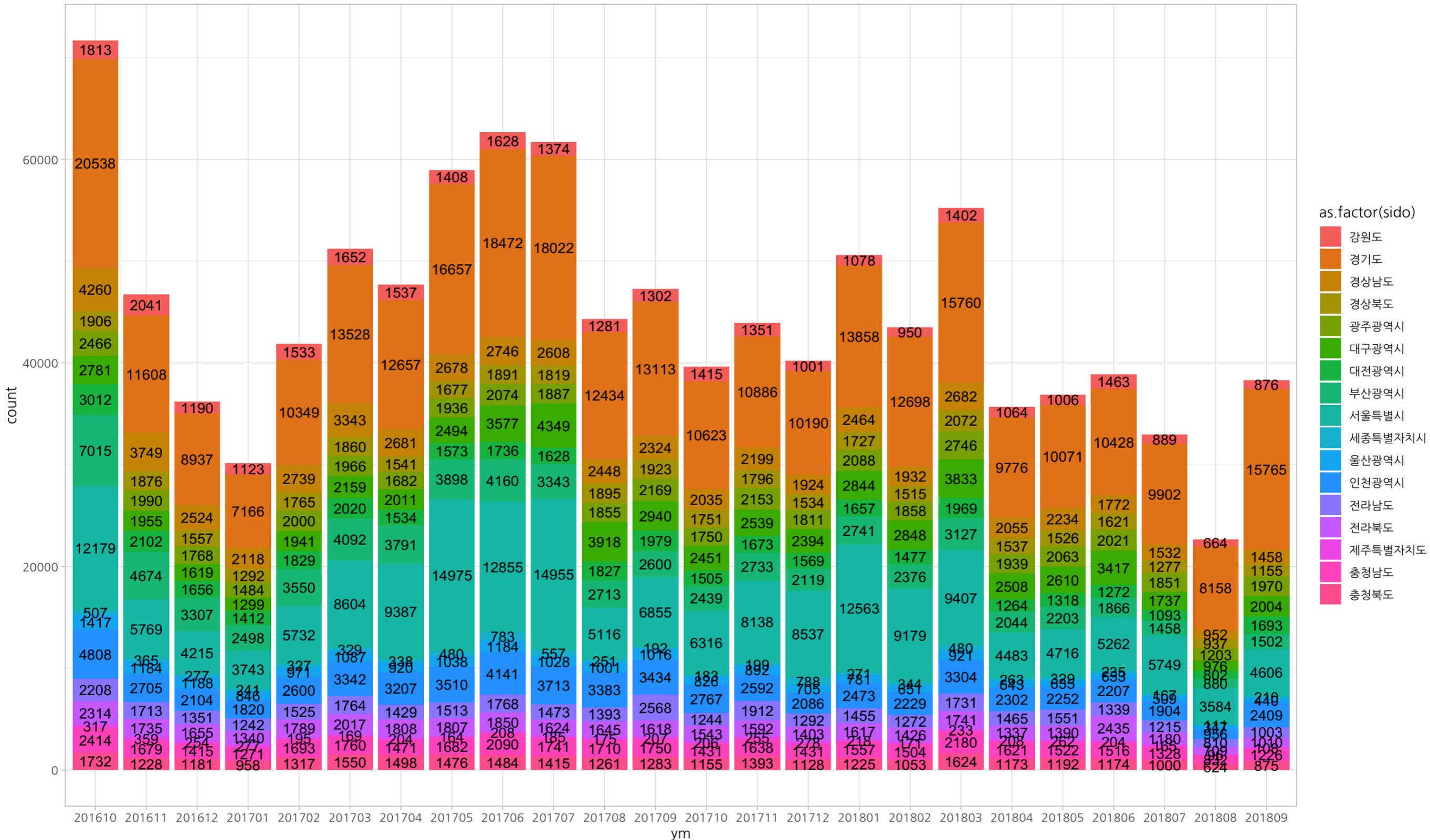
ggplot2 geom_bar()



ggplot(apt, aes(x= ym, fill = as.factor(sido))) + geom_bar()



ggplot2 geom_bar() + geom_text()

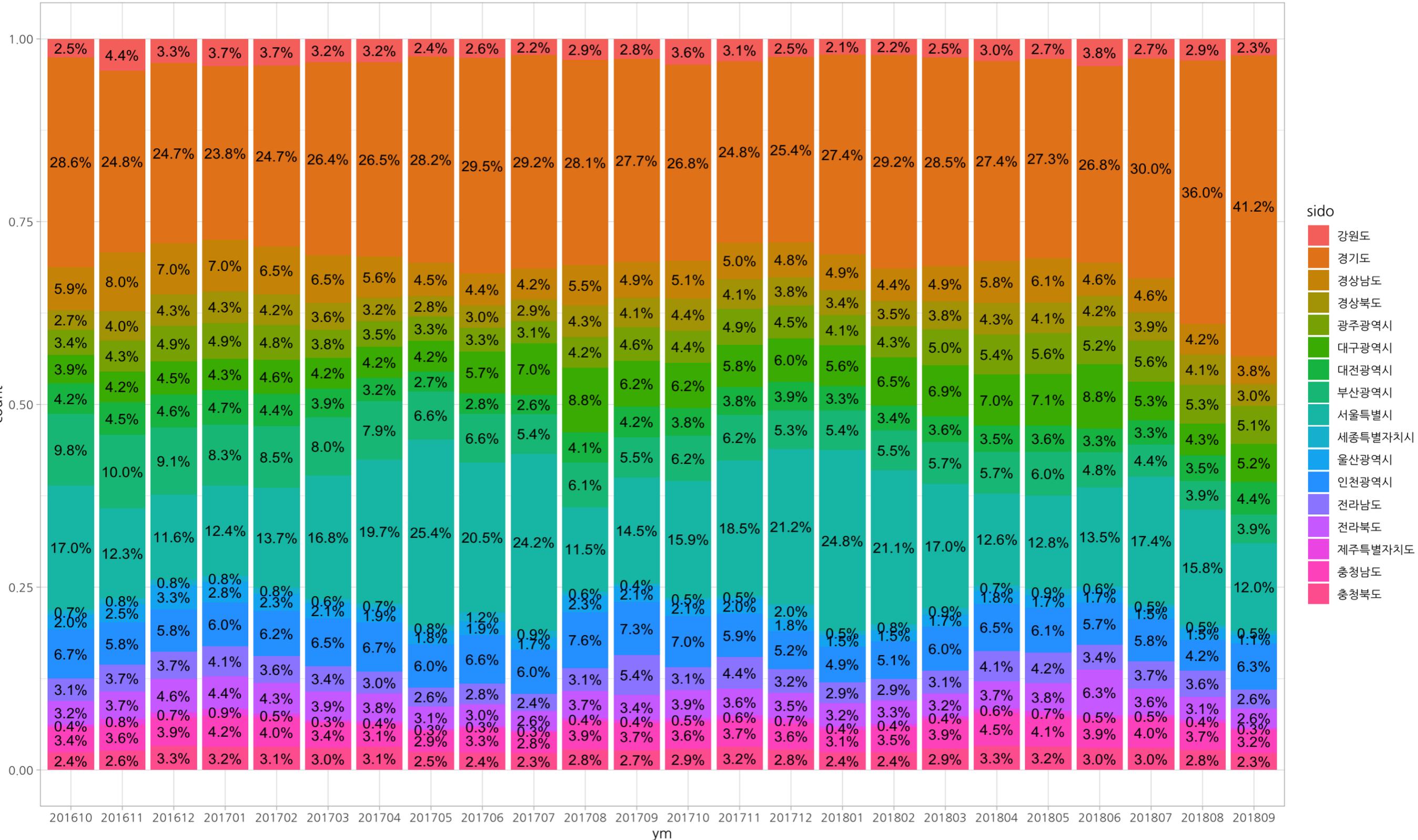


**ggplot(apt, aes(x= ym, fill = as.factor(sido))) + geom_bar() +
geom_text(aes(label=..count..),stat="count",position=position_stack(0.5))**



ggplot2 geom_bar() + geom_text()

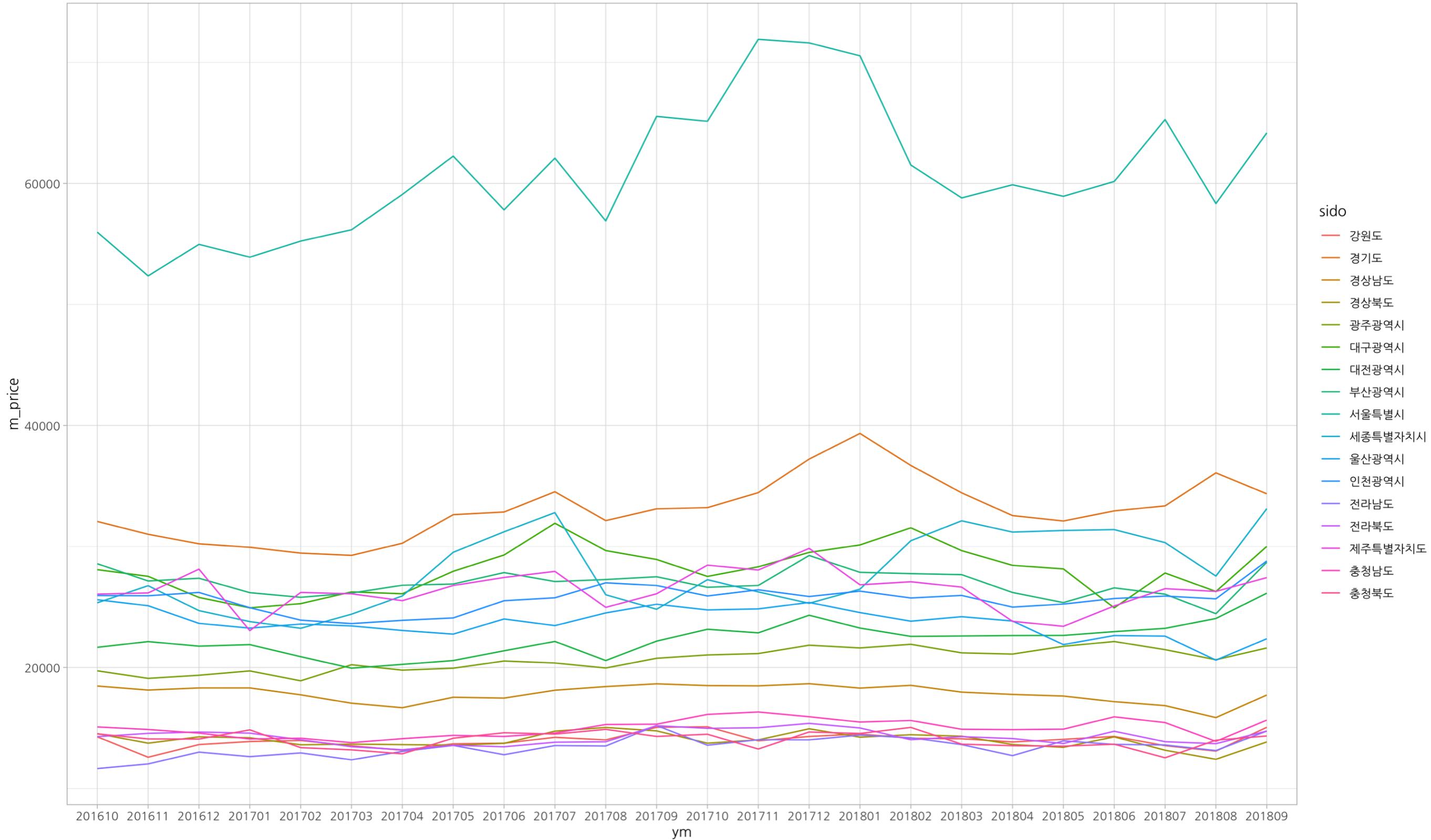
position_fill



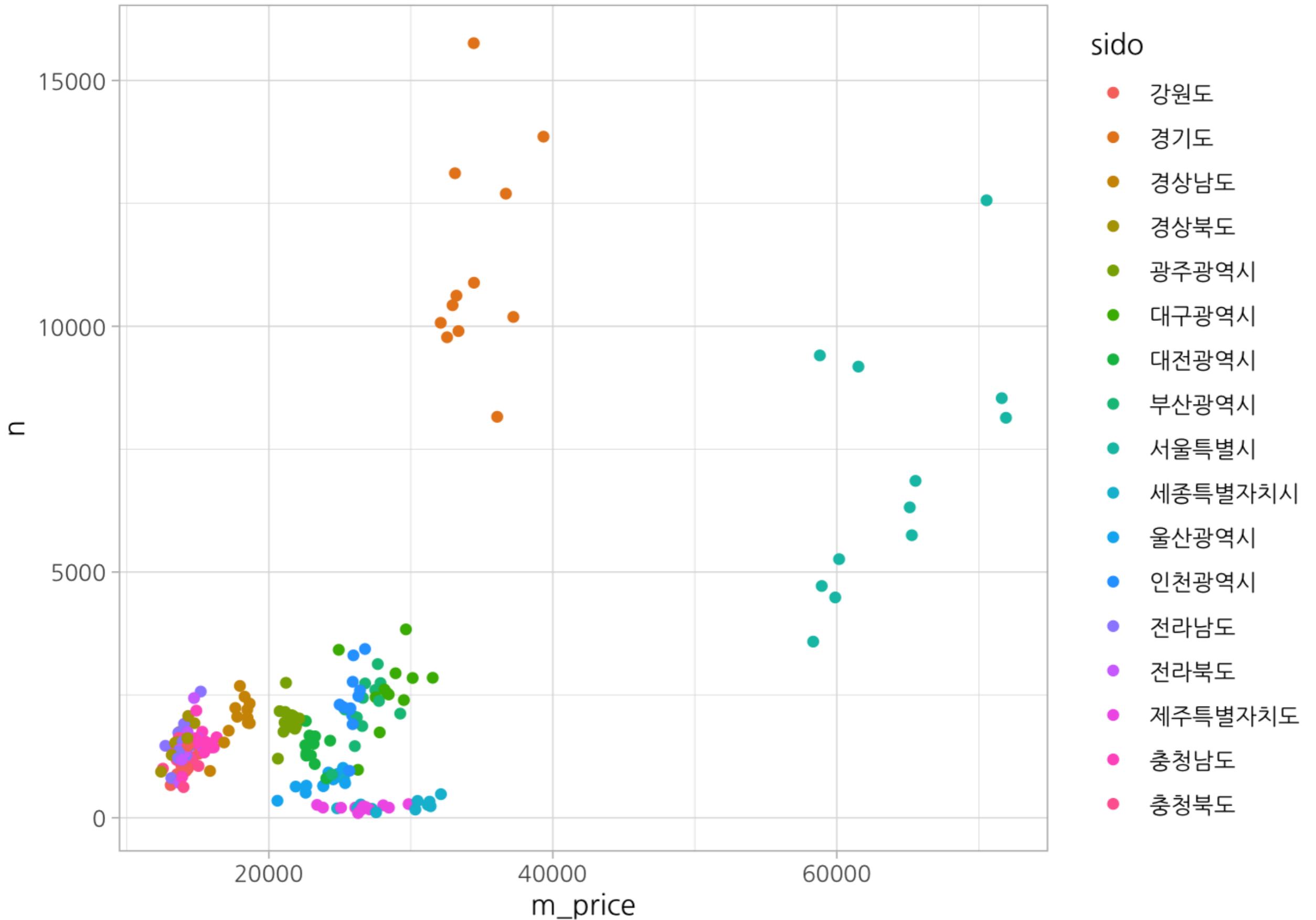
**ggplot(apt,aes(x=ym,fill=sido)) + geom_bar(position="fill")+
geom_text(data=percentData, aes(y=n, label=ratio), position=position_fill(vjust=0.5))**



ggplot2 geom_line()



```
ggplot(percentData, aes(x= ym, y= m_price, group=sido, color = sido)) + geom_line()
```

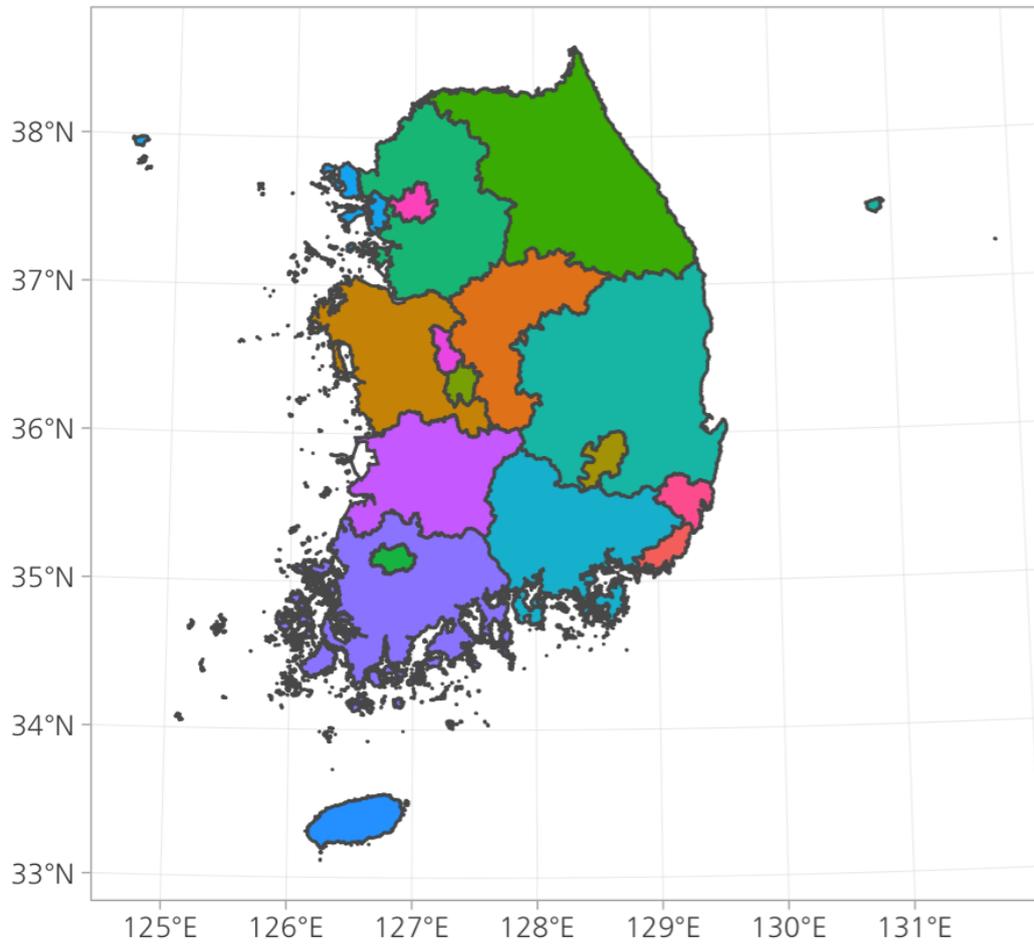


```
ggplot(percentData,aes(x=m_price,y = n, color=sido))+
  geom_point()
```

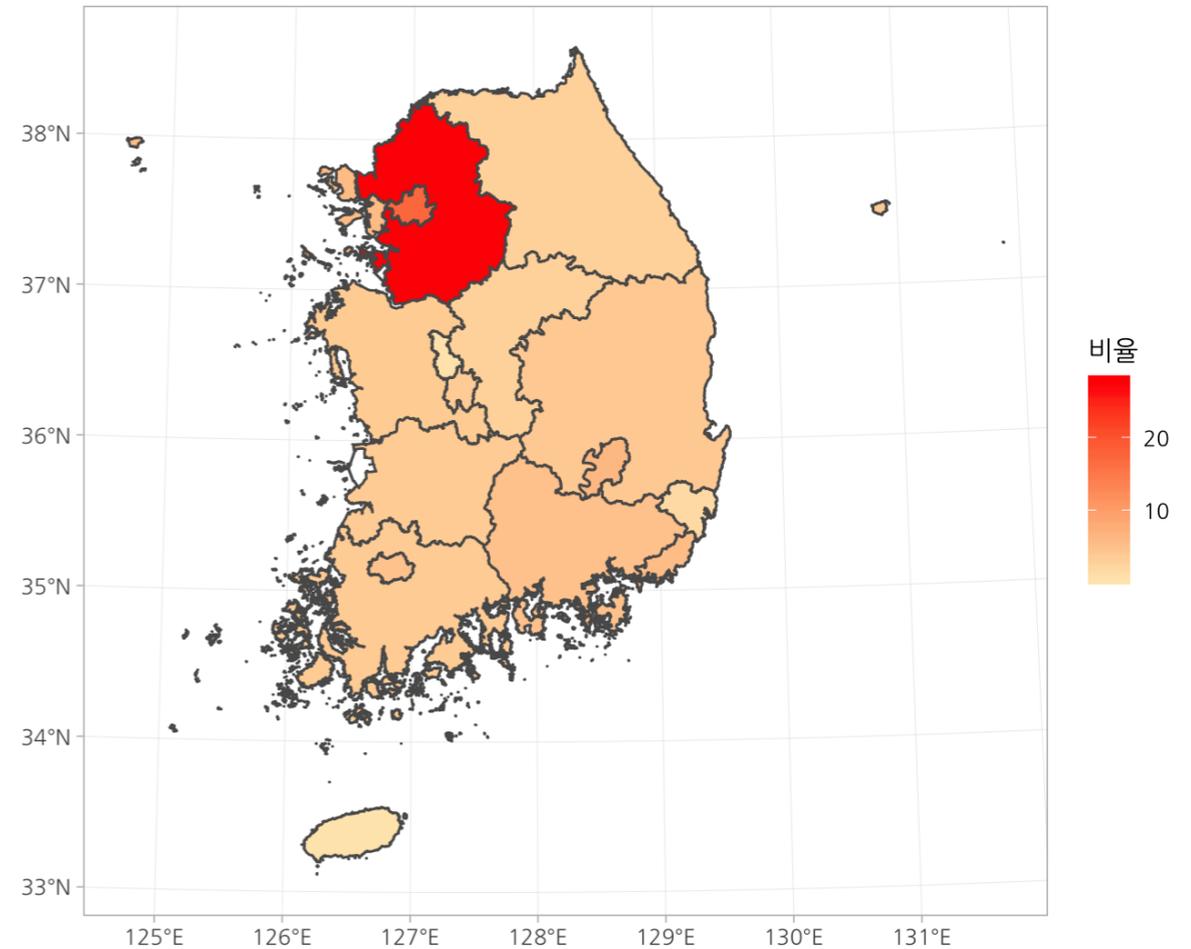



ggplot2 geom_sf()

광역시도 경계



아파트 거래량 비중



```
ggplot(data=sido_simp_shp, aes(fill=ratio)) +  
  geom_sf() +  
  labs(title="아파트 거래량 비중") +  
  theme(legend.position = "right") +  
  scale_fill_gradient(low = "wheat1", high = "red", name = "비율", labels = scales::comma)
```



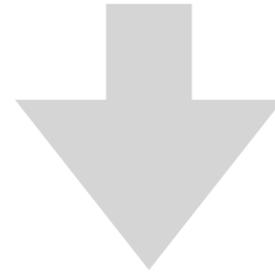
ggplot2 limitation

1. 구분 항목이 많아질 경우 2차원 그래프 상에서 색상 표현 및 연속적인 구분이 어렵다.
2. 그래프만 가지고 표현하려고 할 경우 누락되는 정보가 많다.



ggplot2 alternatives

1. 구분 항목이 많아질 경우 2차원 그래프 상에서 색상 표현 및 연속적인 구분이 어렵다.
2. 그래프만 가지고 표현하려고 할 경우 누락되는 정보가 많다.



1. 연속적인 흐름을 표현할 수 있도록 그래프를 그린다.
2. 데이터가 갖고 있는 정보들을 표현할 수 있도록 한다.

60,000

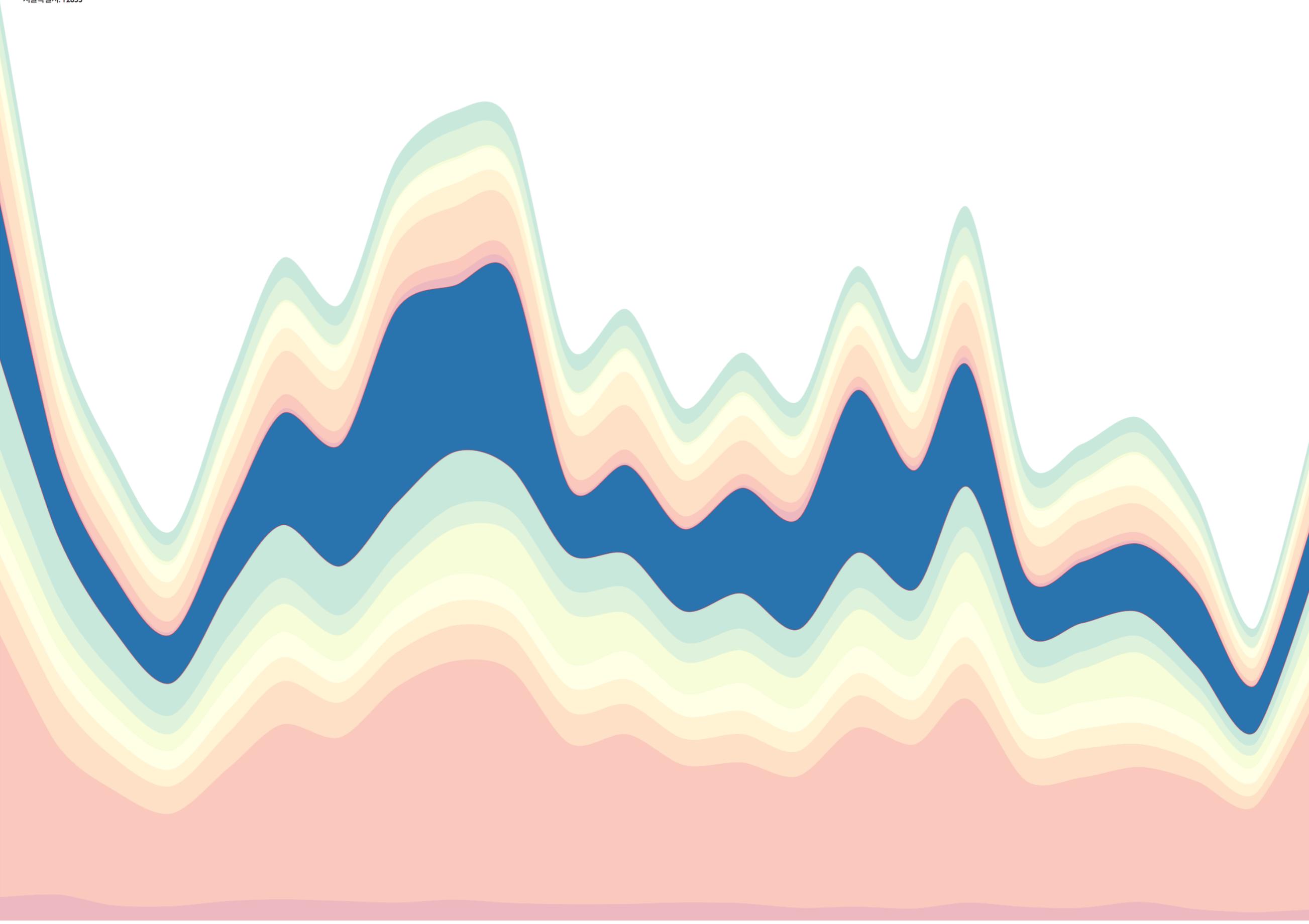
40,000

20,000

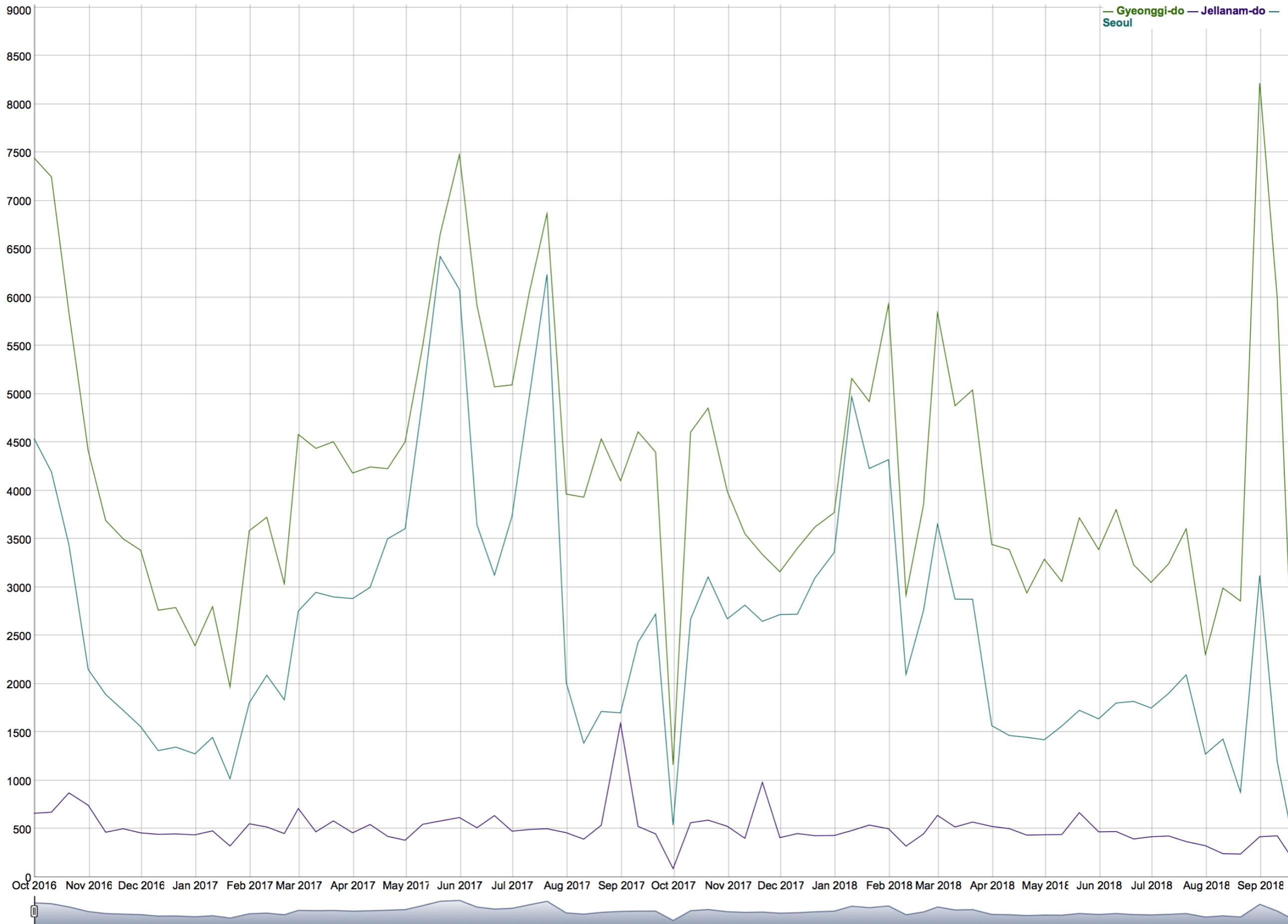
0

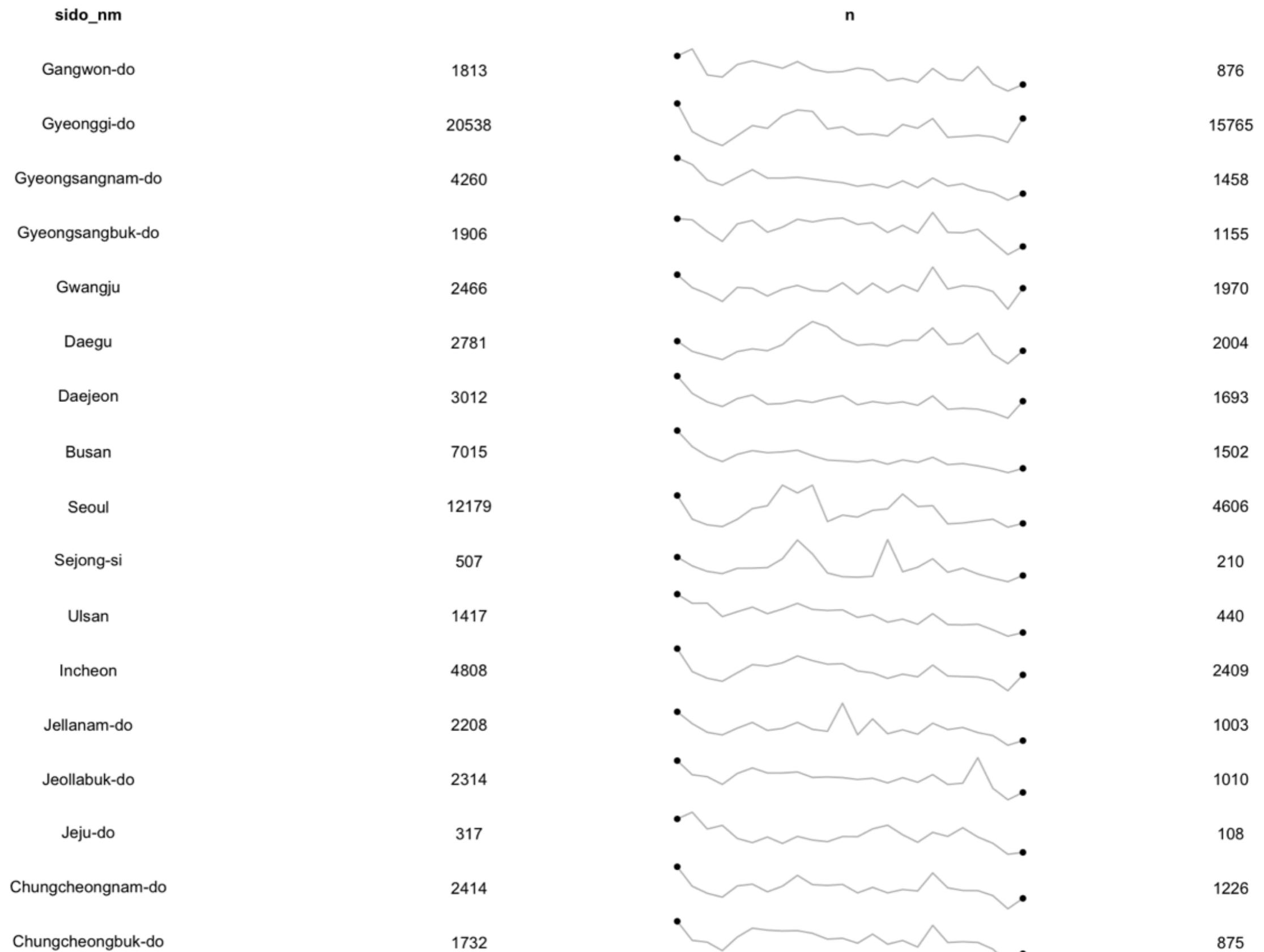
201610 201611 201612 201701 201702 201703 201704 201705 201706 201707 201708 201709 201710 201711 201712 201801 201802 201803 201804 201805 201806 201807 201808

시도:

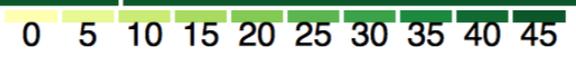
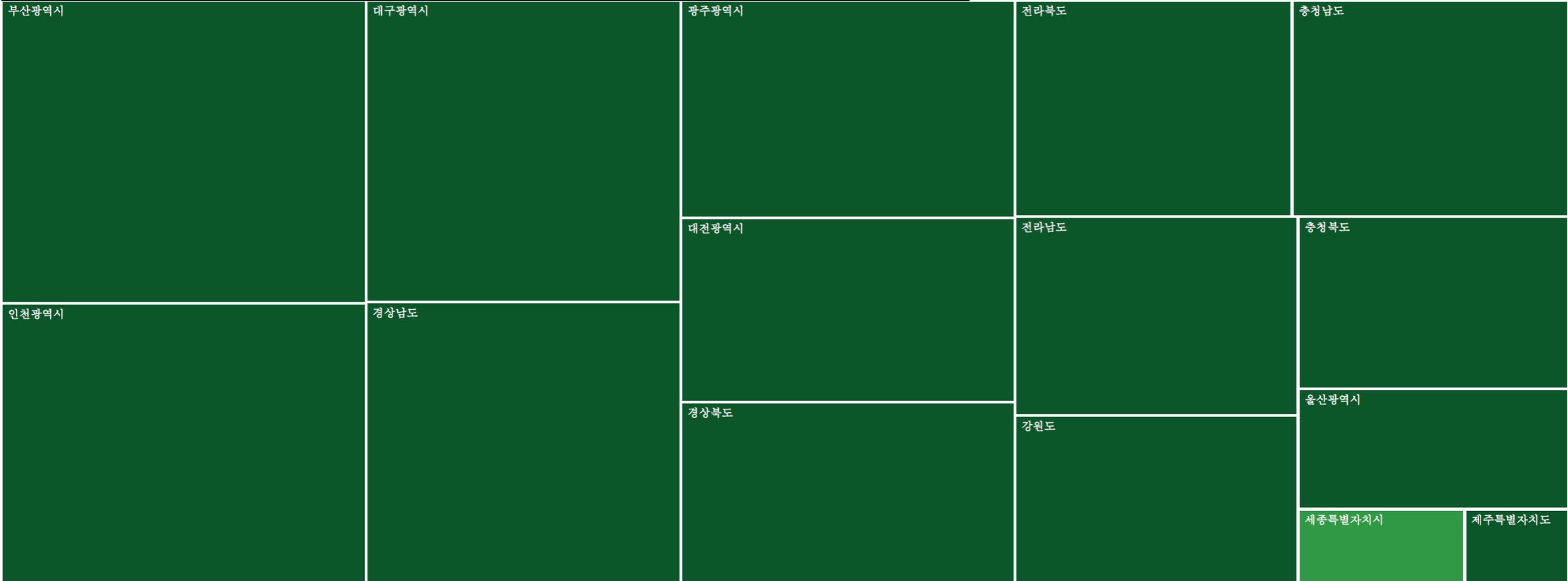
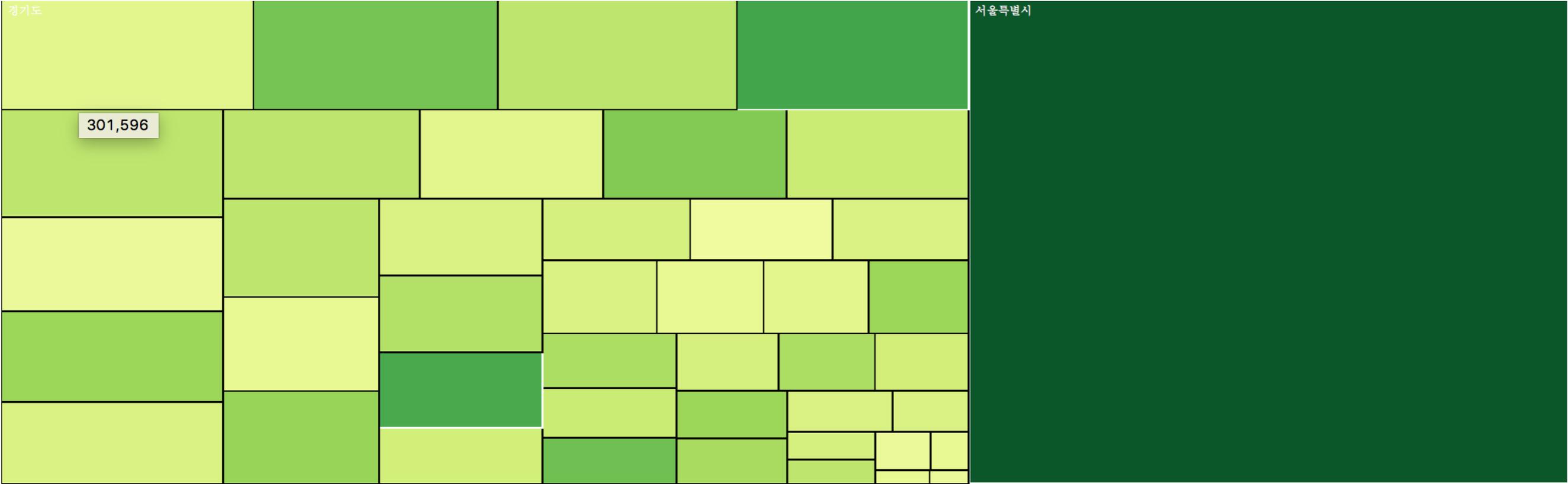


dygraphs for R

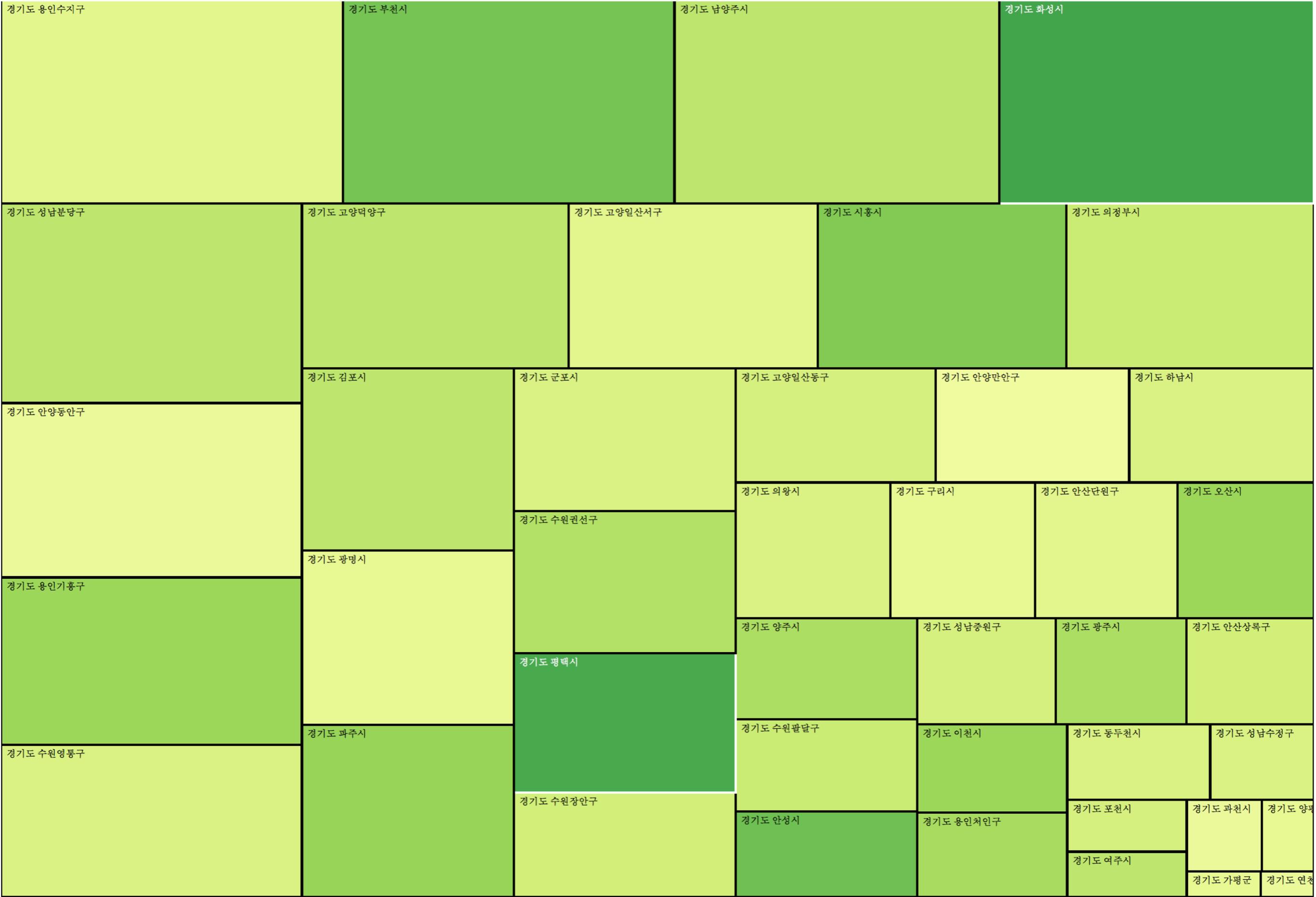




**plotSparklineTable(select(data.frame(percentData), sido_nm, n, ym),
row.var = 'sido_nm', col.vars = 'n')**



sidو.경기도



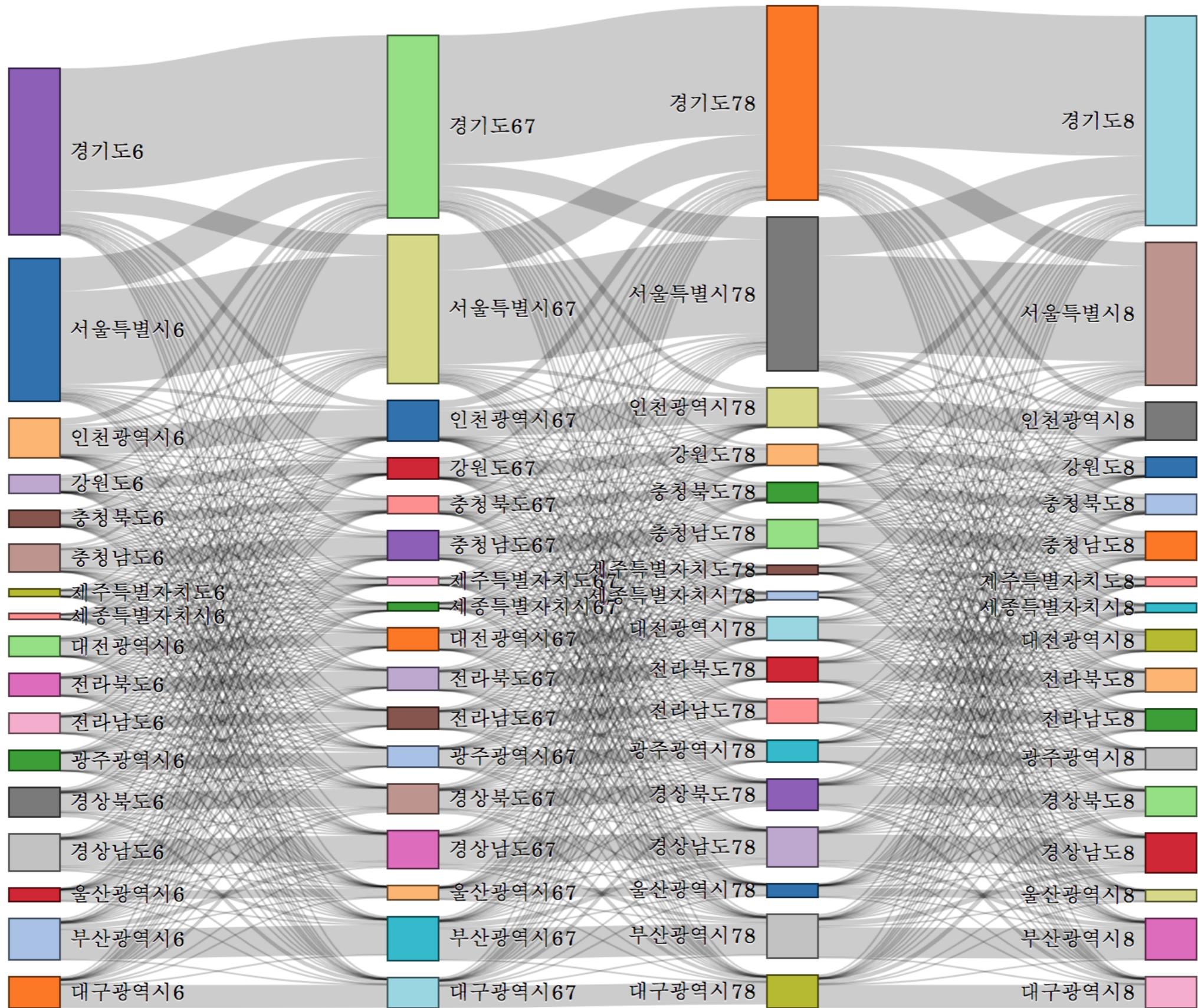
0 5 10 15 20 25 30 35 40 45

인구 이동

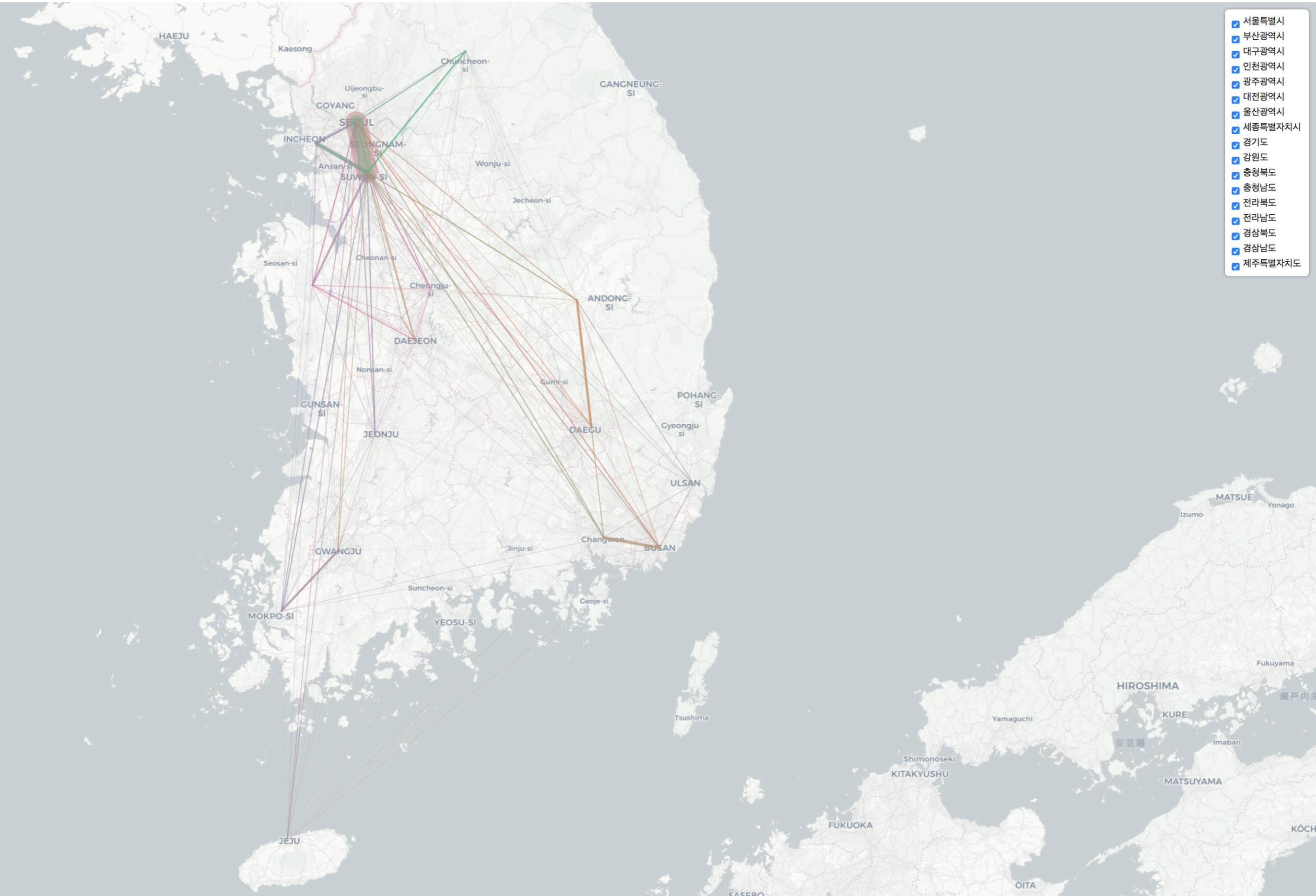
1) 시군구별 이동자수

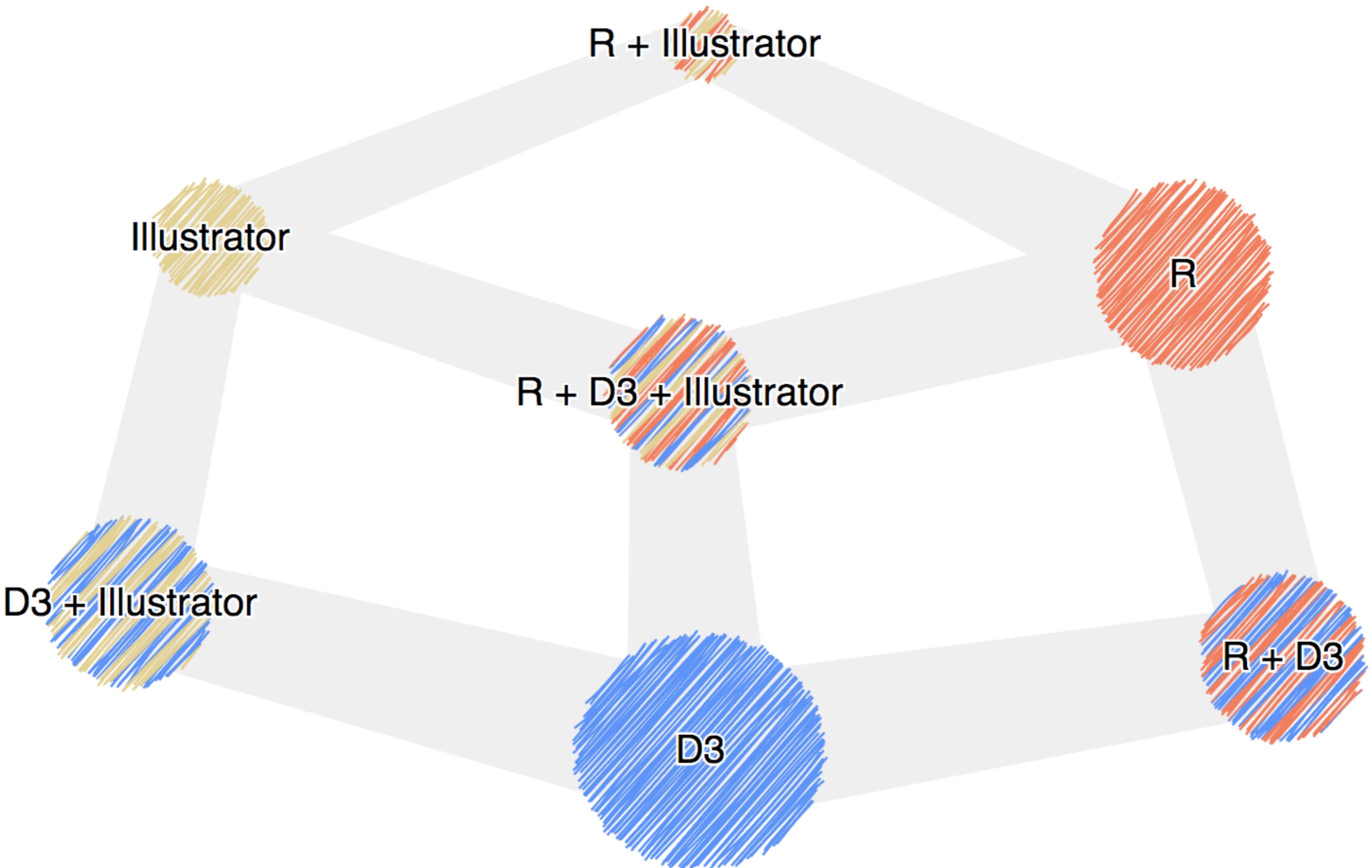
http://kosis.kr/statHtml/statHtml.do?orgId=101&tblId=DT_1B26001_A01

인구이동 Sankey



인구이동 leaflet





참고사이트

- <http://gallery.htmlwidgets.org/>
- <https://rmarkdown.rstudio.com/flexdashboard/>
- <http://personal.tcu.edu/kylewalker/interactive-flow-visualization-in-r.html>

