데이터 시각화

11. 그래프 디자인의 기본원칙(1)

최대영 교수





• 학습리뷰

- 1 변수 연관성의 시각화
 - ② 산점도(Scatter plot)
 - 데이터 하나를 x, y좌표 위에 점 하나로 표현
 - ত 버블 차트(Bubble chart)
 - 버블(점)의 크기를 이용해서 속성을 표시
 - ☑ 산점도 매트릭스(Scatterplot matrix)
 - 모든 변수 간의 관계를 일대일로 산점도로 표현

•• 학습리뷰

- 1 변수 연관성의 시각화
 - **Correlogram**
 - 상관계수를 색을 이용하여 시각화
 - ☑ 경사 차트(Slopegraph)
 - 두 변수의 값을 선으로 연결

3

◆ 학습리뷰

- 2 시계열 데이터의 시각화
 - ☑ 선 그래프(Line graph)
 - 이웃한 두 점(전후 시간의 데이터)을 연결
 - ☑ 이동 평균(moving average)을 이용한 평활화(smoothing)
 - 한 시점씩 이동하며 특정 구간의 데이터 평균값을 계산하여 표현
 - ☑ 국소 추정 산점도 평활(LOESS, locally estimated scatter plot smoothing)
 - 저차원의 다항식에 데이터의 부분집합을 적합(fitting)

Л

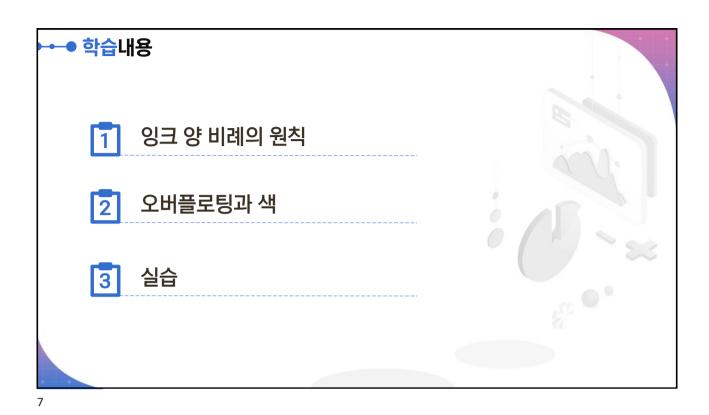
•• 학습리뷰

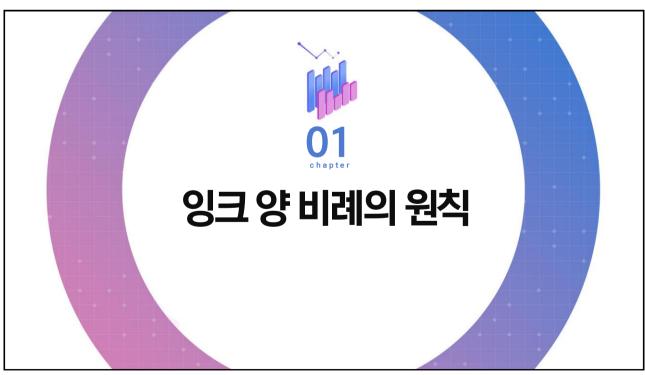
- 2 시계열 데이터의 시각화
 - 🖒 스플라인(Spline) 모형
 - 몇 개의 제어점을 기준으로 구간별로 적합
 - ☑ 정의된 함수 형식을 이용한 추세 시각화
 - 데이터에 잘 맞는 형식의 함수를 정하여 데이터를 적합
 - 🖒 추세 제거 (Detrending)
 - 장기적으로 지배하는 추세가 있을 때 추세를 제거하여 중요한 편차를 특정하여 강조

_

◆● 학습목표

- 🤷 잉크 양 비례의 원칙에 대해 설명할 수 있다.
- 오버플로팅과 효과적인 색 사용 방법에 대해 설명할 수 있다.
- 오버플로팅 관련 matplotlib 라이브러리를 이해하고 활용할 수 있다.





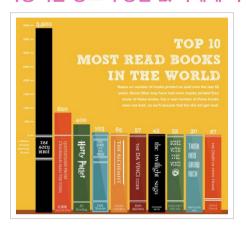
_

1. 축이 있는 도표

া 잉크 양 비례의 원칙(Principle of proportional ink)

☑ 시각화에서 음영 영역의 크기는 해당 영역이 나타내는 데이터 값에 비례해야 함

→ 값을 나타내는 데 사용되는 잉크의 양은 값 자체에 비례해야 함



[출처] Calling Bullshit

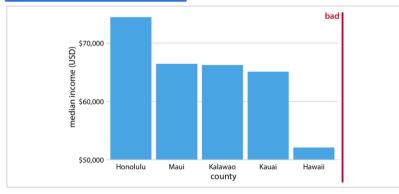
9

1. 축이 있는 도표

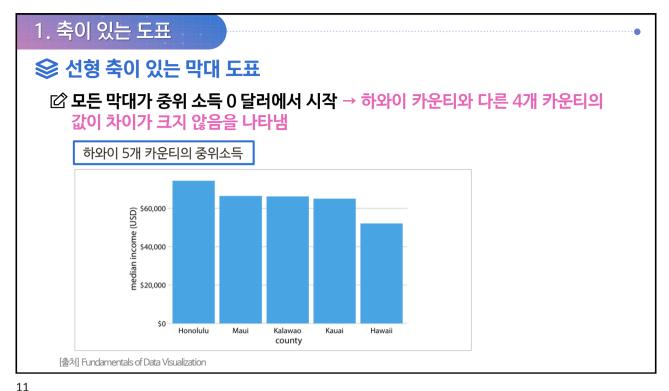
❤ 선형 축이 있는 막대 도표

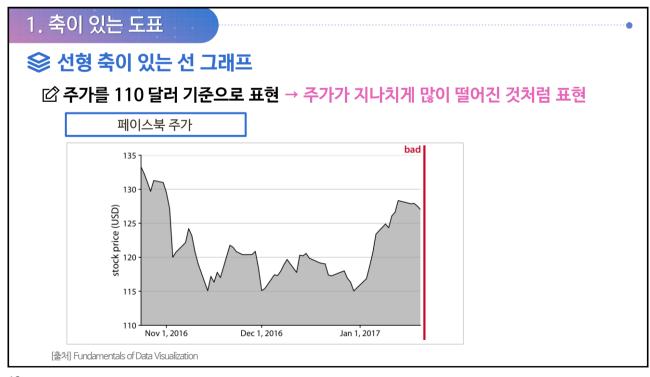
② 모든 막대가 중위 소득 5만 달러에서 시작 → 하와이 카운티와 다른 4개 카운티의 값의 차이가 실제보다 커 보임





[출처] Fundamentals of Data Visualization







(Age) 있는 도표

SPUTI 도표로 차이를 나타내는 방법

CP 0을 기준으로 중위소득의 변화량을 표현 → 0을 기준으로 하더라도 차이를 분명하게 표현할 수 있음

5년간 중위소득의 변화

SPUTI 도표로 차이를 나타내는 방법

(Age) 기준으로 하더라도 차이를 분명하게 표현할 수 있음

5년간 중위소득의 변화

SPUTI 도표로 차이를 나타내는 방법

(Age) 기준으로 하더라도 차이를 분명하게 표현할 수 있음

(Age) 기준으로 하더라도 차이를 보명하게 표현할 수 있음

(Age) 기준으로 하더라도 하는 기준으로 하는 기준으로 하는 기준으로 하더라도 하는 기준으로 하

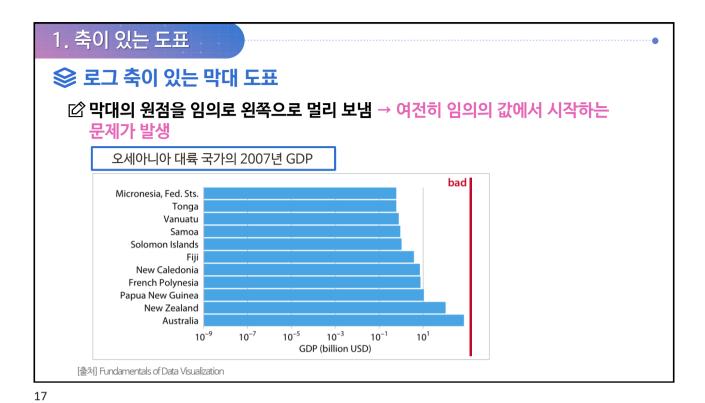
14

[출처] Fundamentals of Data Visualization

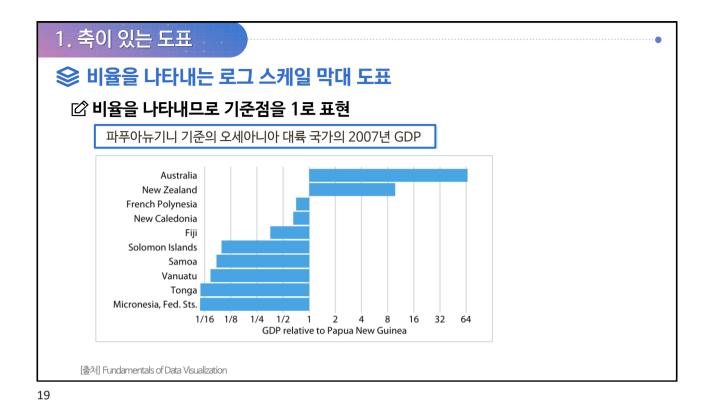
4 +01011 ---

15

1. 축이 있는 도표 ❤ 로그 축이 있는 막대 도표 [☑] 막대가 0.3에서 부터 시작 → 막대 길이가 데이터 값을 정확하게 보여주지 않고 0을 왼쪽으로 무한히 멀리에 표현해야 함 오세아니아 대륙 국가의 2007년 GDP bad Micronesia, Fed. Sts. Tonga Vanuatu Samoa Solomon Islands Fiji New Caledonia French Polynesia Papua New Guinea New Zealand Australia 0.3 1.0 30 100 300 10 GDP (billion USD) [출처] Fundamentals of Data Visualization



1. 축이 있는 도표 ★ 점으로 막대를 대체 [☑] 0을 기준으로 해야 하는 막대 도표의 길이 문제를 해결 → 길이의 개념이 없기 때문에 상대적인 크기를 적절히 표현 오세아니아 대륙 국가의 2007년 GDP Micronesia, Fed. Sts. Tonga • Vanuatu Samoa • Solomon Islands Fiji 🔵 New Caledonia French Polynesia Papua New Guinea New Zealand Australia 0.3 1.0 100 300 3.0 10 GDP (billion USD) [출처] Fundamentals of Data Visualization



2. 면적을 나타내는 도표

③ 파이 차트(Pie chart)

② 파이 조각의 면적은 각도에 비례, 각도는 데이터 값에 비례

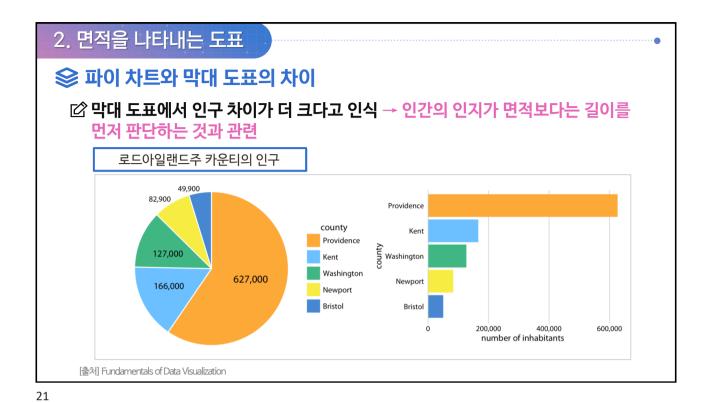
→ 잉크 양 비례의 원칙에 부합

로드아일랜드주 카운티의 인구

127,000
166,000
167,000
187,000
188,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000
189,000

20

[출처] Fundamentals of Data Visualization

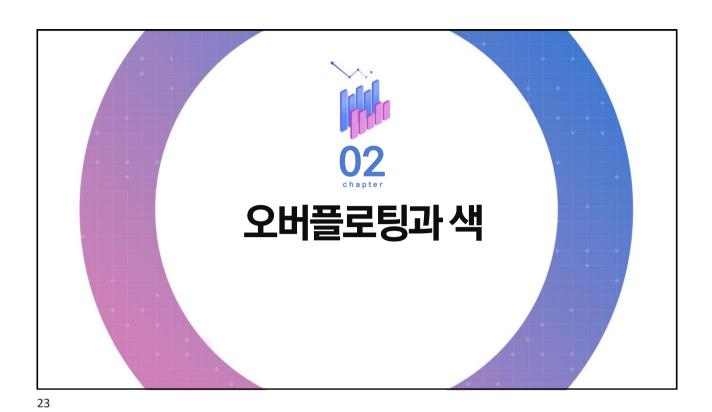


2. 면적을 나타내는 도표 ≫ 트리맵과 막대 도표의 차이 않 막대 도표에서 인구 차이가 더 크다고 인식 → 인간의 인지가 면적보다는 길이를 먼저 판단하는 것과 관련 로드아일랜드주 카운티의 인구 Providence Kent Washington Newport 166,000 127,000 82,900 Kent Washington Newport Providence 627,000 Bristol 200,000 400.000 600,000

22

[출처] Fundamentals of Data Visualization

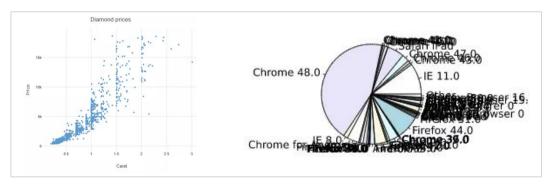
number of inhabitants



1. 오버플로팅

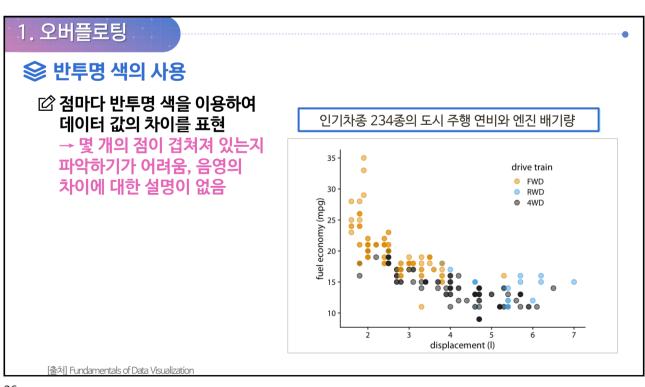
☑ 여러 개의 점이 같은 좌표에 겹쳐서 표현되어 있는 현상

- 데이터 포인트가 너무 많은 경우에 주로 발생
- 데이터 값이 부정확하거나 반올림 등으로 인해 같은 값이 많은 경우에도 발생

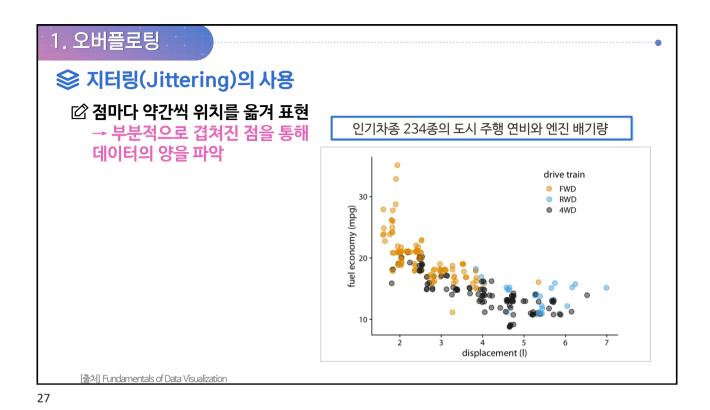


[출처] Fundamentals of Data Visualization

1. 오버플로팅 ≫ 데이터 값을 반올림하여 겹쳐진 데이터 예시 ☑ 연비와 엔진 배기량을 반올림 인기차종 234종의 도시 주행 연비와 엔진 배기량 → 겹쳐진 점으로 인해 데이터의 특성 파악이 어려움 35 -(차종별 데이터 포인트 수, drive train 노란점 뒤에 겹쳐진 drive train이 다른 차종이 FWD RWD 검은점이 존재 fuel economy (mpg) • 4WD 숨겨짐) 10 displacement (I) [출처] Fundamentals of Data Visualization



26



1. 오버플로팅

S 지터링(Jittering)의 사용

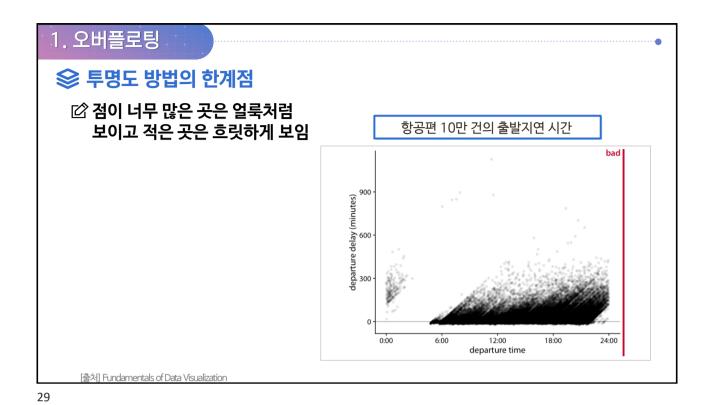
C 점마다 약간씩 위치를 옮겨 표현

→ 위치를 너무 많이 옮기면
데이터가 왜곡됨

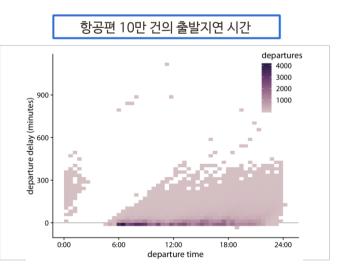
O기차종 234종의 도시 주행 연비와 엔진 배기량

OU기차종 234종의 도시 주행 연비와 엔진 배기량

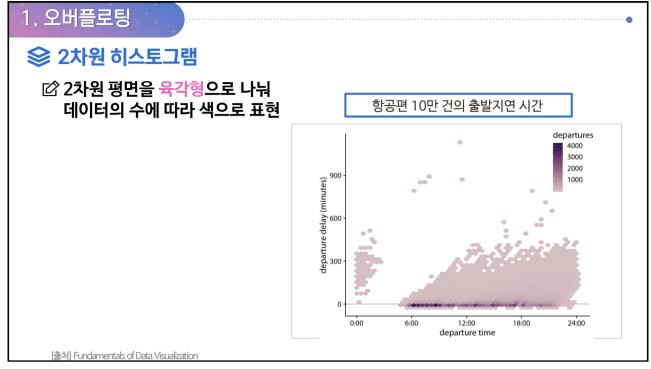
Output

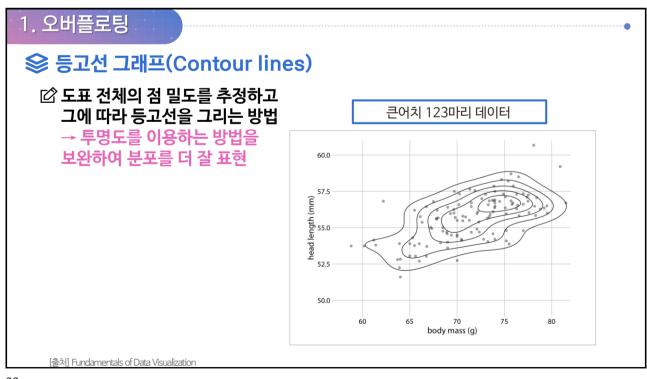


데이터의 수에 따라 색으로 표현 → 데이터의 특성을 요약해서 보여주는데 유리



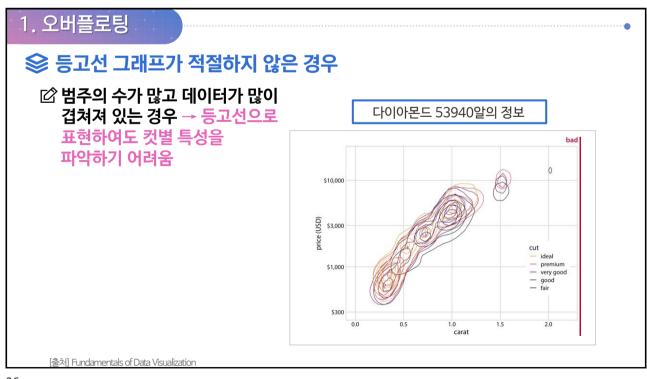
[출처] Fundamentals of Data Visualization





1. 오버플로팅 등고선 그래프(Contour lines) ☑ 등고선의 색을 다르게 하여 여러 큰어치 123마리 데이터 범주를 표현 → 범주의 수가 2~3개로 제한되고 분포가 겹쳐져 있지 않은 경우에 효과적 60.0 (mm) head length (r 52.5 50.0 female birds
 male birds 60 70 body mass (g) [출처] Fundamentals of Data Visualization

[출처] Fundamentals of Data Visualization



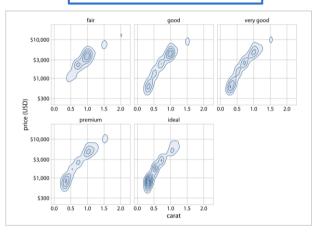
1. 오버플로팅

얼 범주별 등고선 그래프

☑ 범주별로 등고선을 그리고격자를 넣어 범주 간의 비교가가능하도록 표현

→ 여러 범주가 겹쳐 있을 때 유용

다이아몬드 53940알의 정보



[출처] Fundamentals of Data Visualization

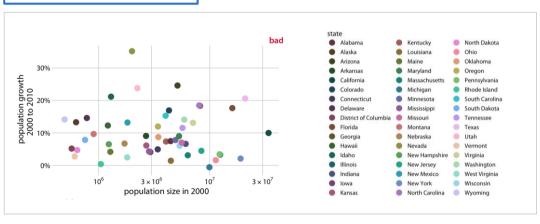
37

2. 효과적인 색 사용

❤️ 색상 스케일이 효과적인 범주의 수

☑ 3~5개의 범주를 색으로 나타낼 때 효과적, 8~10개가 넘어가면 효과가 떨어짐

미국의 주별 인구 성장율과 규모



[출처] Fundamentals of Data Visualization

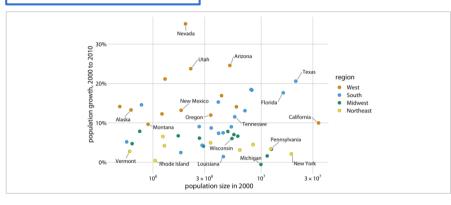
2. 효과적인 색 사용

❤ 색상 스케일이 효과적인 범주의 수

☑ 색상 스케일의 수를 한정하고 대표성을 띄는 중요한 데이터에만 레이블을 표시

→ 도표가 간결하고 메시지 전달력이 좋아짐

미국의 주별 인구 성장율과 규모



[출처] Fundamentals of Data Visualization

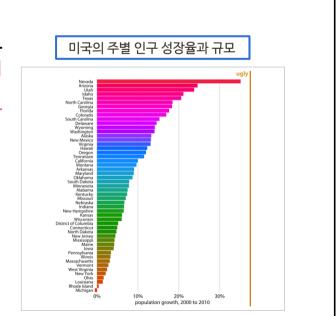
39

2. 효과적인 색 사용

목적이 명확하지 않은 색의 사용

☑ 색이 의미하는 바가 없고 채도가 높음

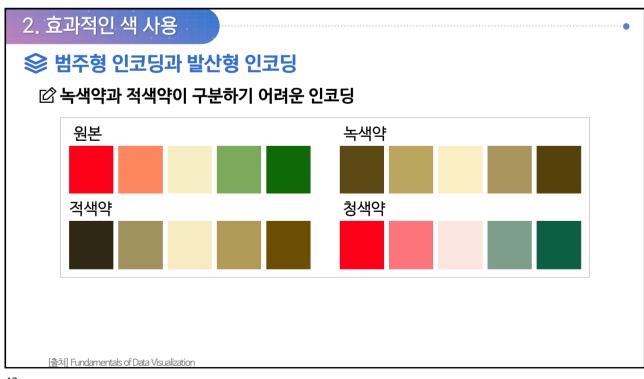
→ 메시지 전달에 도움이 되지 않으며 채도가 높으면 도표를 꼼꼼하게 살피기 어려움(넓은 면적에는 채도가 낮은 색 위주로 사용)



[출처] Fundamentals of Data Visualization

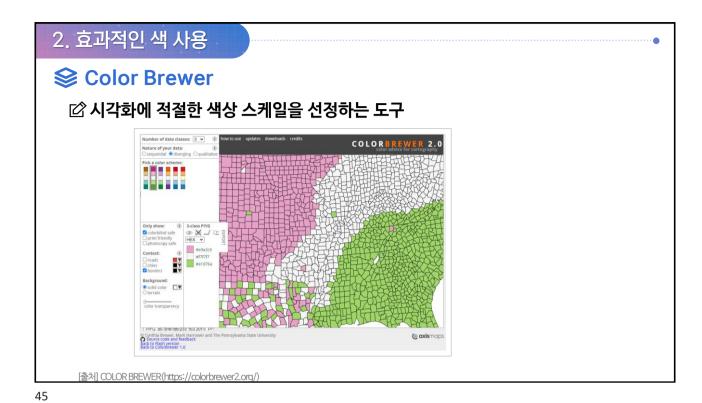


[출처] Fundamentals of Data Visualization





[출처] Fundamentals of Data Visualization





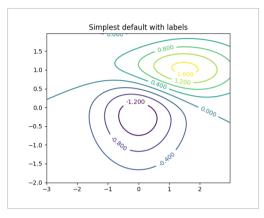
1, 오버플로팅 ◇ 산점도의 속성과 등고선 그래프 ⓒ 산점도의 속성 ■ scatter(x, y, s=area, c=color, alpha=transparency) 10 8 i xy - 4 2 1 2 3 4 X-Axis [출제 Matplotlib Tutorial - 파이썬으로데이터시각화하기

1. 오버플로팅

♦ 산점도의 속성과 등고선 그래프

🖒 등고선 그래프

contour(X, Y, Z)



[출처] Matplotlib Tutorial - 파이썬으로 데이터 시각화하기

49

● 학습정리

1 잉크 양 비례의 원칙

🖒 잉크 양 비례의 원칙(Principle of proportional ink)

■ 음영 영역의 크기는 해당 영역이 나타내는 데이터 값에 비례해야 함

🖒 축이 있는 도표

■ 잉크 양 비례의 원칙에 부합하기 위해 막대나 선은 축의 0부터 표현해야 함

[於 면적을 나타내는 도표

- 파이 조각의 면적은 각도에 비례, 각도는 데이터 값에 비례
- 인간의 인지가 면적보다는 길이를 먼저 판단하기 때문에 막대 도표가 파이 차트와 트리맵 보다 데이터 값의 차이를 더 두드러지게 나타냄

→ 학습정리

2 오버플로팅과 색

☑ 오버플로팅(Overplotting)

■ 여러 개의 점이 같은 좌표에 겹쳐서 표현되어 있는 현상

[♡] 반투명 색 사용

■ 점마다 반투명 색을 이용하여 데이터 값의 차이를 표현

② 지터링(Jittering)

■ 데이터를 왜곡하지 않는 한에서 점마다 약간씩 위치를 옮겨 표현

51

● 학습정리

2 오버플로팅과 색

☑ 2차원 히스토그램

■ 2차원 평면을 사각형 등으로 나눠 데이터의 수에 따라 색으로 표현

☑ 등고선 그래프(Contour lines)

■ 도표 전체의 점 밀도를 추정하고 그에 따라 등고선을 그리는 방법

● 학습정리

2 오버플로팅과 색

☑ 색상 스케일의 수와 의미

- 3~5개의 범주를 색으로 나타낼 때 효과적
- 색이 의미하는 바를 명확히 하고 채도를 조절하여 피로를 낮춰야 함

🏻 색각 이상자(color-vision deficiency)에 대한 고려

- 색각 이상자도 구별할 수 있는 색상 스케일로 인코딩
- Color Brewer 등 도구를 이용해 시각화에 쓰일 색상을 선정

53

▶ • • • 참고문헌

- 🔛 「데이터 시각화 교과서」, Claus O. Wilke, 책만, 2020.
- Fundamentals of Data Visualization, Claus O. Wilke, O'Reilly Media, 2019.

※ 서체 출처 | 넥슨Lv2고딕 - (넥슨코리아)www.levelup.nexon.com / 나눔바른고딕(네이버)

저작권 안내 이 강의록은 저작권법에 의해 보호받는 저작물로서 저작권자의 허락 없이 저작재산권 일체(복제권, 배포권, 대여권, 공연권, 공중전송권, 전시권, 2차적 저작물 작성권)를 침해 시 저작권법에 의거 처벌받을 수 있습니다. 계동캠퍼스(03051)서울특별시 종로구 북촌로 106 **안암캠퍼스**(02841)서울특별시 성북구 안암로 145 고려대학교