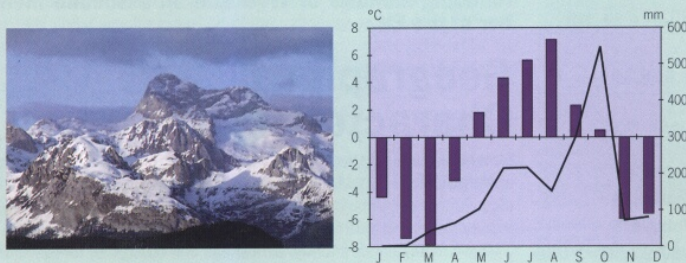
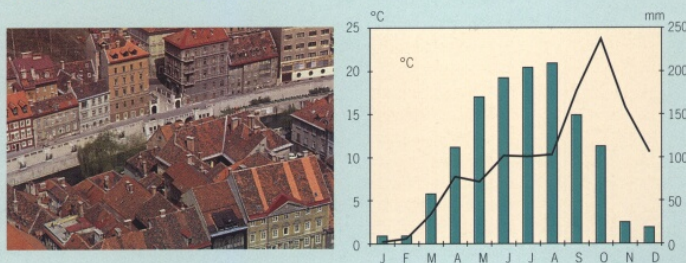


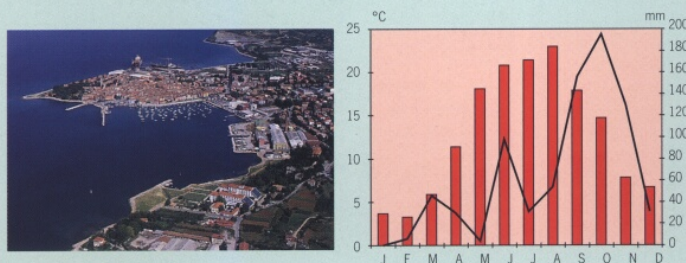
Kredarica 2514m above sea level



Ljubljana 300m above sea level



Izola 2m above sea level



Univerza v Ljubljani
Doktorski študijski program
Statistika

Sodobni statistični pristopi
Prikazi podatkov

Vladimir Batagelj

IMFM, Ljubljana

19. februar 2015

Kazalo

1	Merske lestvice	1
3	Preiskovanje in predstavitev podatkov	3
4	Osnove prikazov podatkov	4
7	Zgledi	7
14	Zgledi: ManyEyes	14
19	Zgledi: ggplot2	19
20	Knjižnice	20
21	Prikazni sestav	21
22	Izhodne naprave	22
23	Slikovne sestavine	23
24	Slikovne sestavine – globinska ostrina	24
25	Weber-Fechner-jev zakon	25
26	Stevens-ov zakon	26
27	Posledice	27
28	Posledice	28
29	Previdno	29
34	Podatkovja	34
35	Zunanji in notranji pogled	35

36	Krmilje	36
37	Očala in lupe	37
38	Prikazi večrazsežnih podatkovij	38
40	Ravni v podatkovju	40
41	Slikovne sestavine in vrste lestvic	41
42	Trirazsežne predstavitve	42
44	Hans Rosling – Gapminder	44
45	Viri	45

<http://zvonka.fmf.uni-lj.si/netbook/doku.php?id=pub:stat>

Merske lestvice

Merska lestvica je predpis $f : A \rightarrow \mathbb{R}$, ki posameznemu objektu x iz množice A priredi neko realno število $f(x)$. V analizi podatkov izmerjene vrednosti dane lastnosti na izbranih objektih določajo spremenljivko.

Merske lestvice, ki merijo isto lastnost so enakovredne – npr. vrednost izražena v USD oziroma v EUR. Za običajne lestvice velja, da lahko enakovredni lestvici predelamo eno v drugo z *dopustnimi transformacijami*. Te določajo *vrsto* (tip) lestvice (glej naslednjo prosojnico).

Vrste lestvic imajo v analizi podatkov pomembno vlogo – določajo, kaj lahko z danimi podatki počnemo.

VB: Teorija merjenja. F.S. Roberts: Discrete Mathematical Models.

Vrste merskih lestvic

dopustne transformacije	vrsta lestvice	primeri
$\varphi(x) = x$ (identiteta)	absolutna	štetje
$\varphi(x) = a.x, a > 0$ podobnost	razmernostna	masa temperatura (K)
$\varphi(x) = a.x + b, a > 0$	razmična	temperatura (C,F) čas (koledar)
$x \geq y \Leftrightarrow \varphi(x) \geq \varphi(y)$ strogo naraščajoča	urejenostna	šolske ocene, kakovost zraka, trdost kamnin
φ je povratno enolična	imenska	barva las, narodnost

Stevens (1946) *On the theory of scales of measurement.*

Preiskovanje in predstavitev podatkov

Z rastjo zmogljivosti namiznih računalnikov se je v 90. letih *prikazovanje* podatkov začelo uveljavljati kot pomembno orodje za *preiskovanje* podatkovij in za *predstavitve* dobljenih rezultatov.

Osnovna naloga *analize in prikazov podatkov* je razkrivanje lastnosti (obsežnih) podatkovij z uporabo zmogljivosti človekovega vida na prikazih ustvarjenih s sodobno računalniško grafiko. Končni cilj je dobiti vpogled v podatke in razumeti odnose med posameznimi deli.

Cilj (slikovnih) predstavitev je kar se da razumljivo posredovanje rezultatov. Pri (znanstvenih) prikazih uporabljamo slike, da bi razkrili posamezne značilnosti podatkov. Pogosto se oba vidika prekrivata.

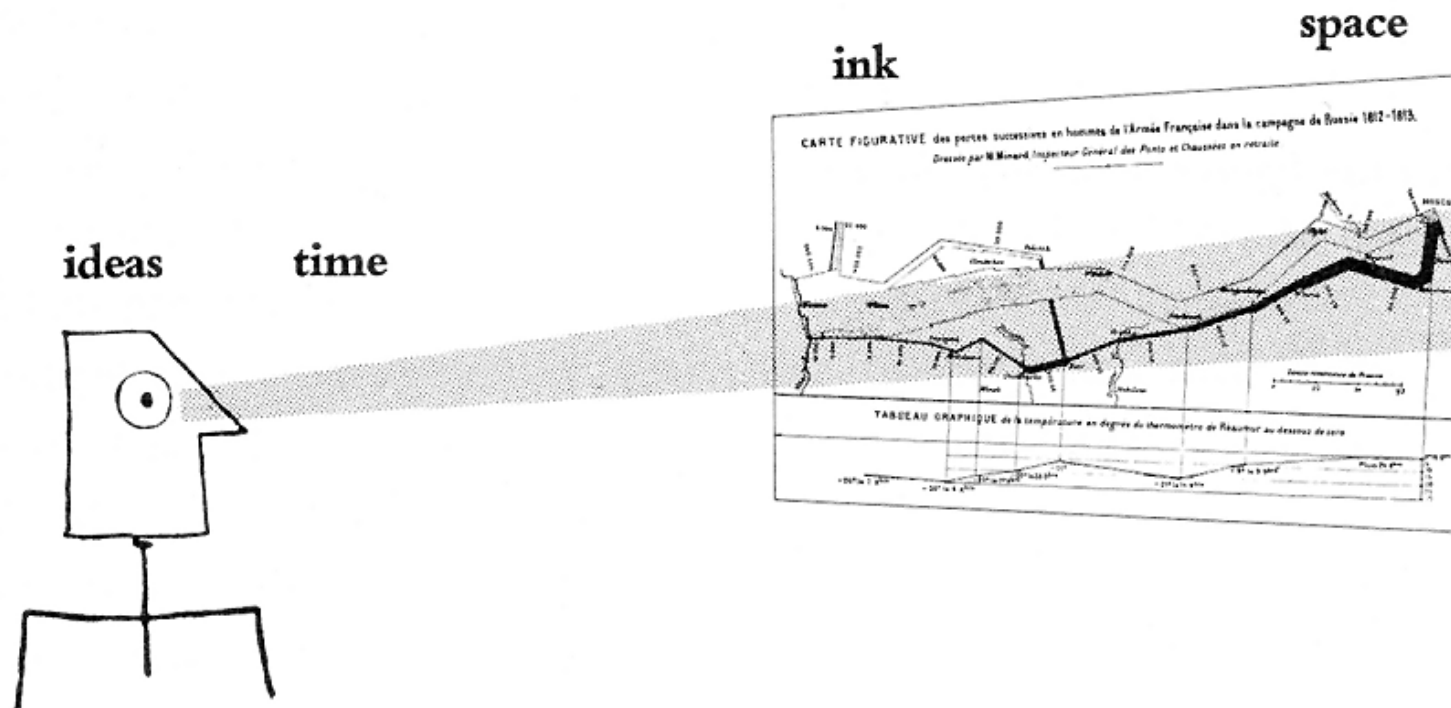
Milestones

Osnove prikazov podatkov

Teoretične osnove prikazov podatkov je postavil Jacques Bertin v svoji knjigi *Sémiologie graphique* (1967).

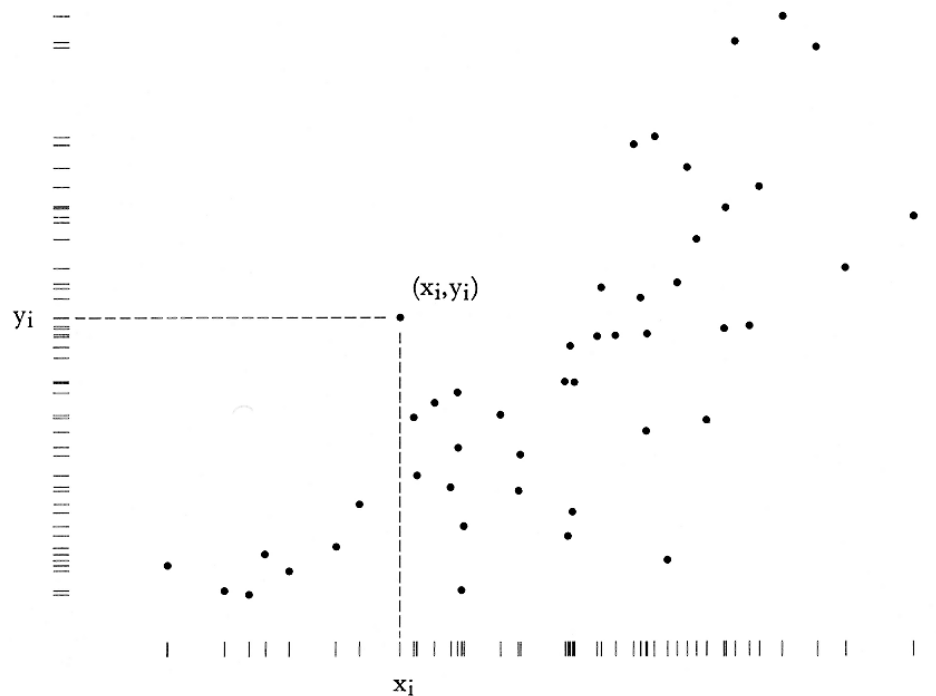
	<i>Points</i>	<i>Lines</i>	<i>Areas</i>	<i>Best to show</i>
<i>Shape</i>		<i>possible, but too weird to show</i>	<i>cartogram</i>	<i>qualitative differences</i>
<i>Size</i>			<i>cartogram</i>	<i>quantitative differences</i>
<i>Color Hue</i>				<i>qualitative differences</i>
<i>Color Value</i>				<i>quantitative differences</i>
<i>Color Intensity</i>				<i>qualitative differences</i>
<i>Texture</i>				<i>qualitative & quantitative differences</i>

Osnove prikazov podatkov (Edward R. Tufte)



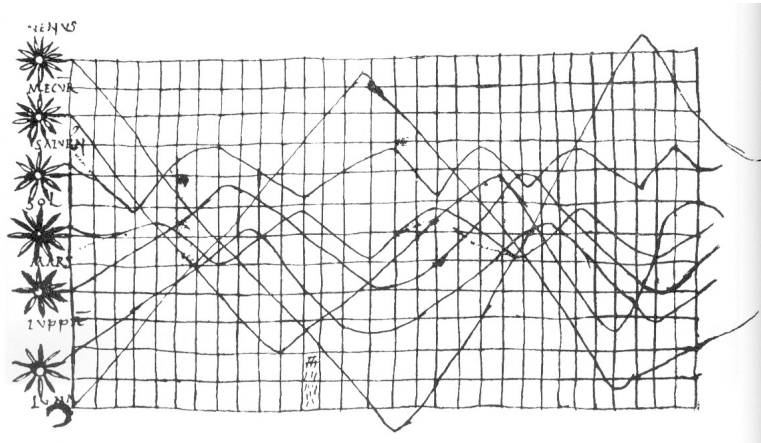
Graphical *excellence* consists of complex ideas communicated with *clarity*, *precision*, and *efficiency*. It gives to the viewer the greatest number of *ideas* in the shortest *time* with the least *ink* in the smallest *space*. It is telling the *truth* about the data.

Tufteova načela

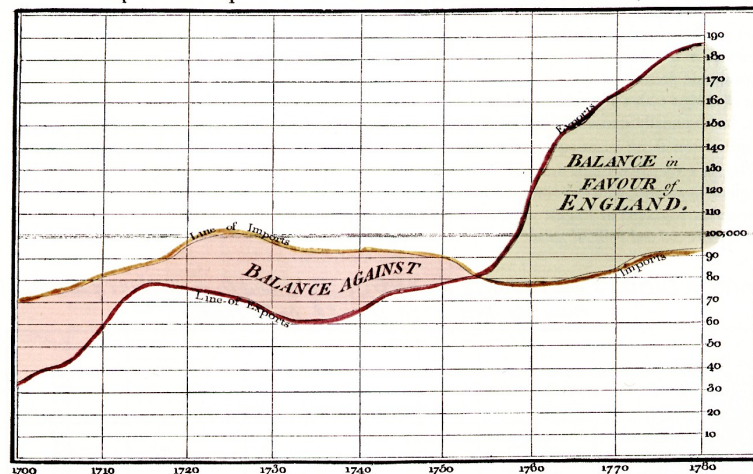


- Predvsem prikaži podatke.
- Povečaj razmerje med podatki in črnilom.
- Odstranuj nepodatkovno črnilo.
- Odstranuj odvečno podatkovno črnilo.
- Preglej in preuredi.

Zgledi



Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.



The Bottom line is divided into Years, the Right hand line into £10,000 each.
Published as the Act directs 11th May 1786. by W^m Playfair.
Sole vendor 532 Strand, London.

Najstarejši (10. stoletje) znani poskus prikaza zveznega spreminjanja neke količine je prikaz spreminjanja naklona planetarnih tirnic.

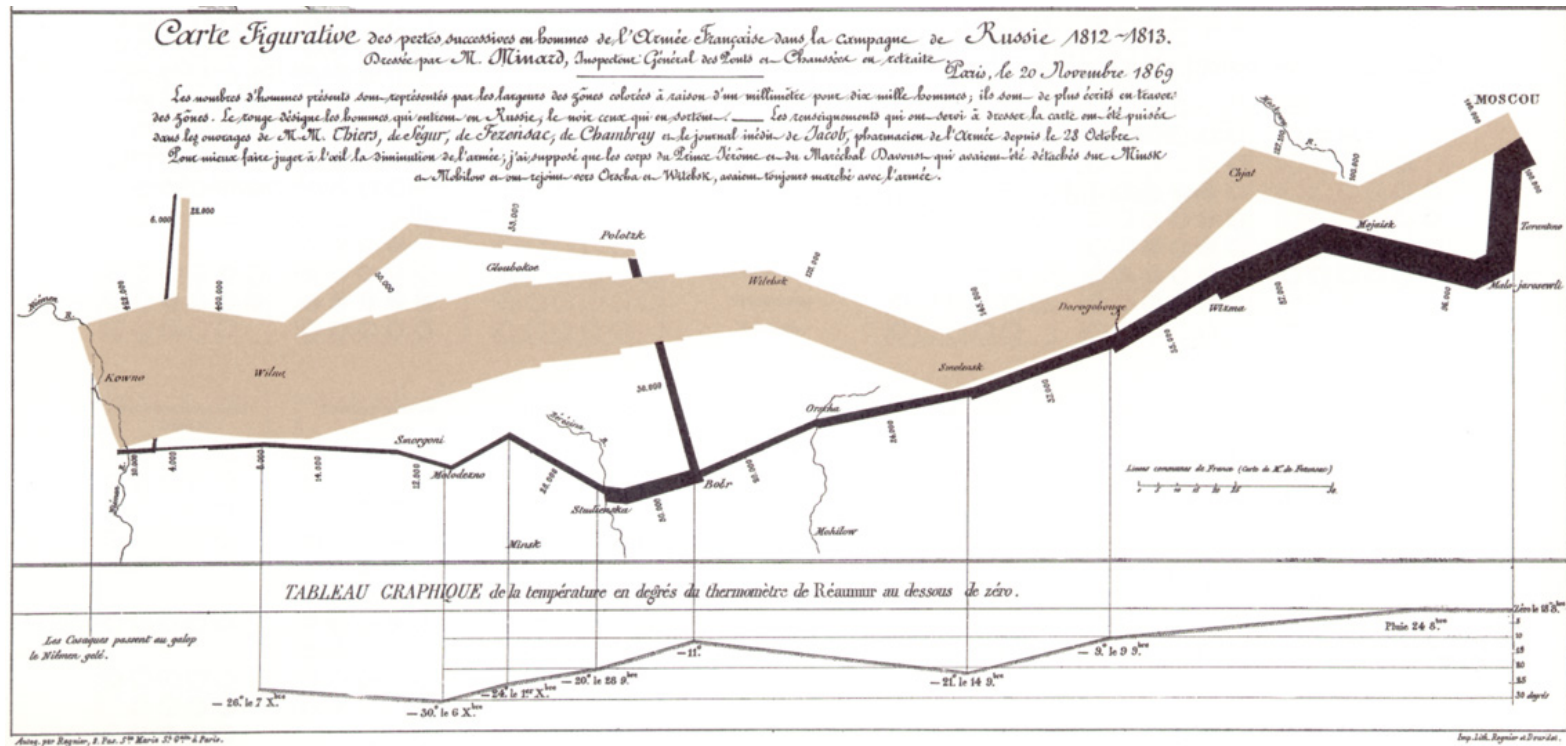
Eni izmed zgodnejših (1785) prikazov ekonomskih podatkov so Playfair-ovi prikazi izvoza in uvoza med državami.

Zgled: Kolera



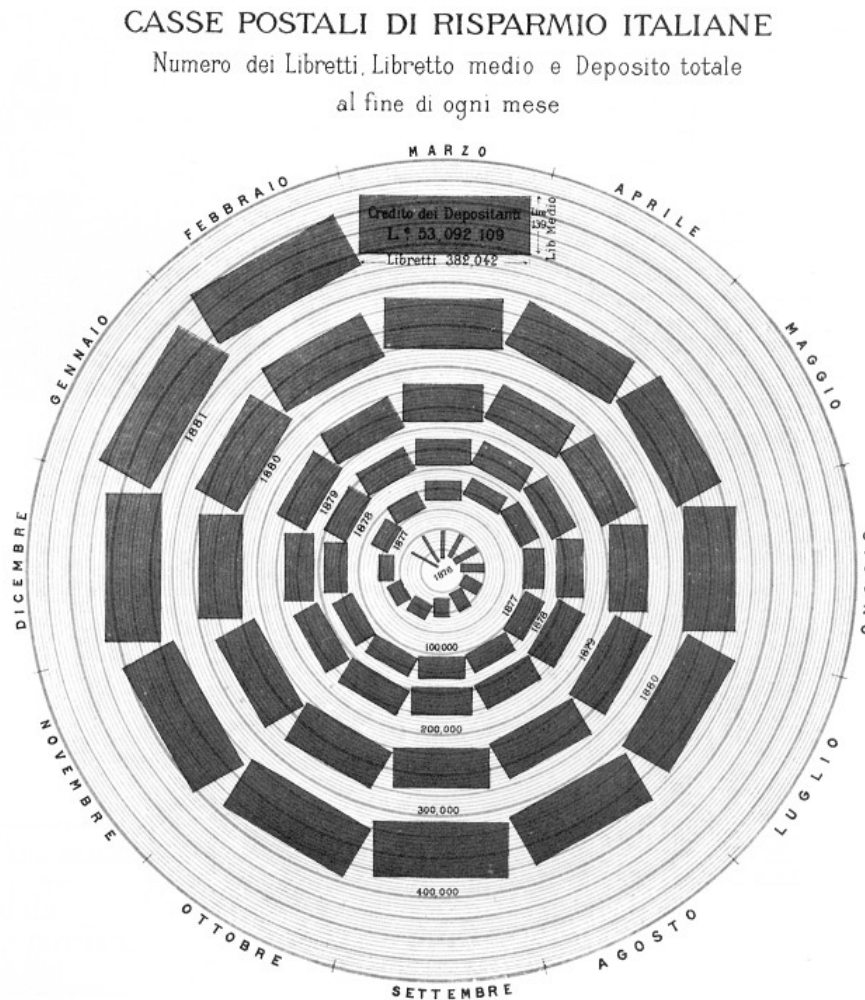
Leta 1854 je dr. John Snow na zemljevid mestnega predela, v katerem je izbruhnila kolera, vrisal s pikami mesta, kjer so posamezniki umrli. V središču pik se je nahajal vodnjak, iz katerega so se okoliški prebivalci oskrbovali z vodo. Dal ga je zapreti – bolezen je izginila.

Zgled: Napoleonov pohod v Rusijo



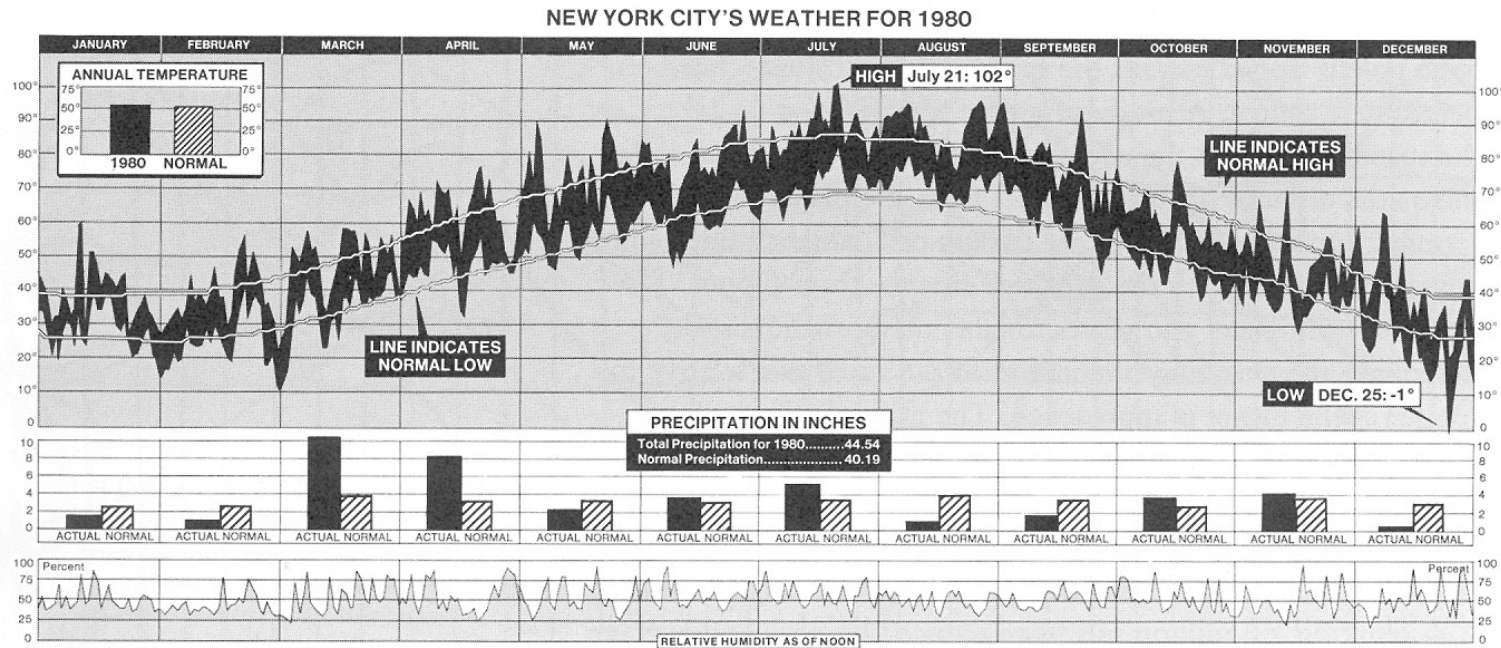
Za pravo mojstrovino velja Minardov (1861) prikaz Napoleonovega pohoda v Rusijo 1812-1813.

Zgled: Rast prihrankov



Zanimiv je prikaz mesečne rasti prihrankov v italijanski poštni hranilnici v letih 1876 do 1881 objavljen v knjigi Gabaglio A.: Teoria Generale della Statistica.

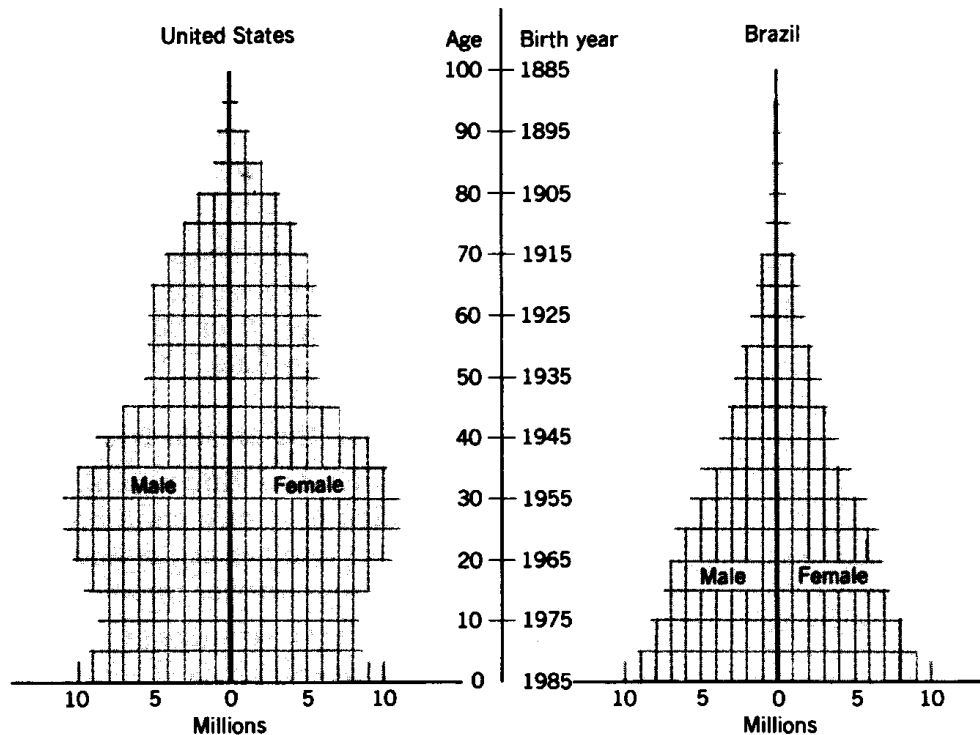
Zgled: Letni pregled vremena



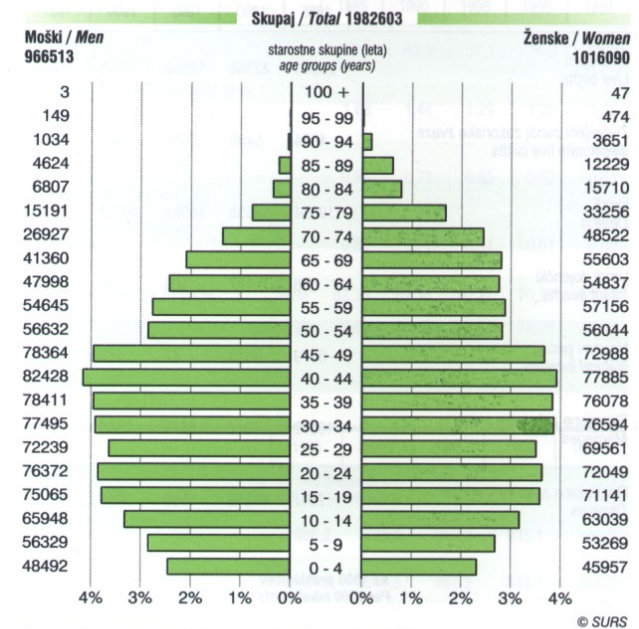
New York Times, January 11, 1981, p. 32.

Zelo veliko podatkov je vgrajenih tudi v prikaz vremena v New Yorku v letu 1980, objavljen v *New York Timesu*, 11. januarja 1981.

Zgled: Starostna porazdelitev prebivalstva



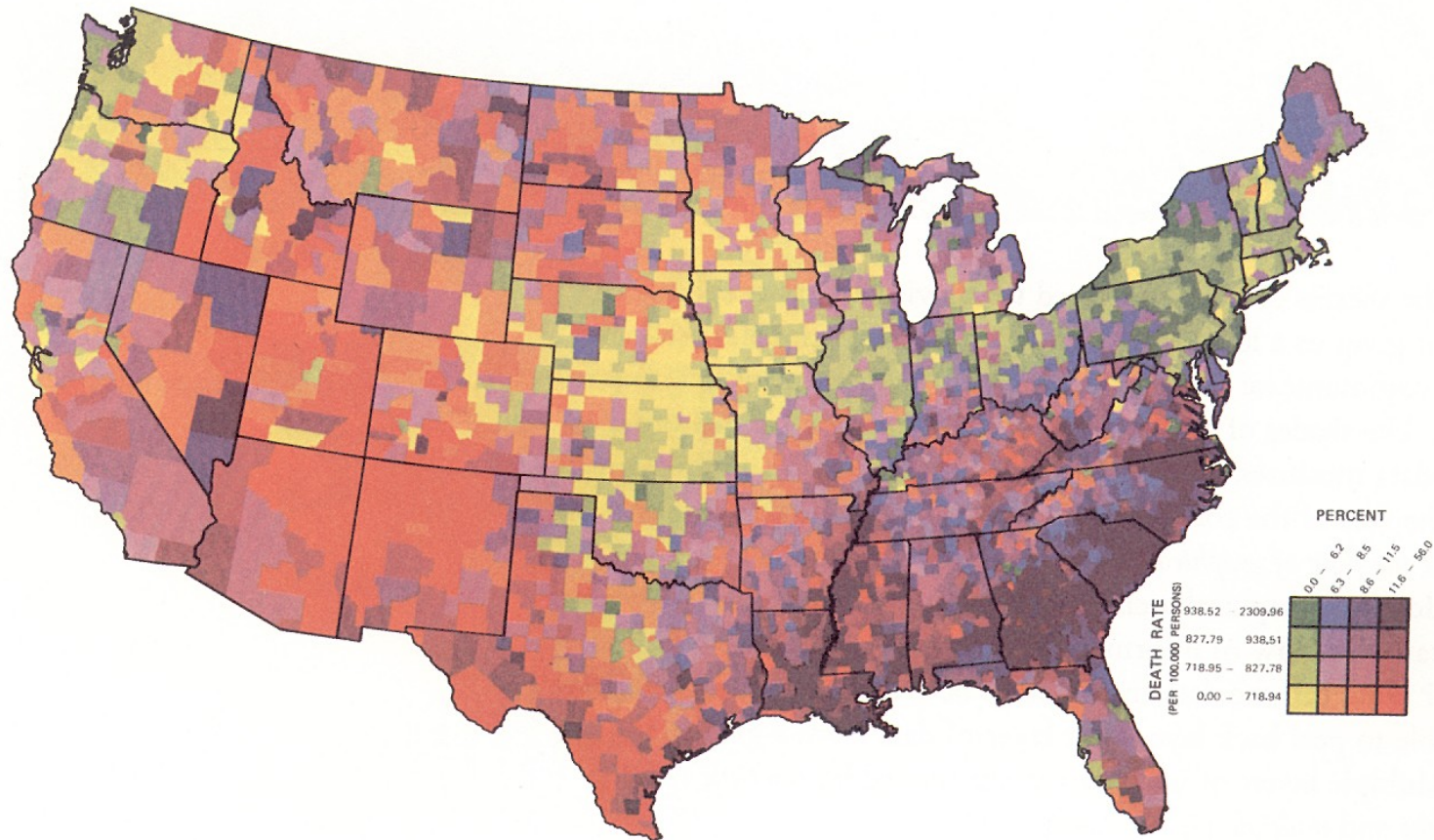
PREBIVALCI PO PETLETNIH STAROSTNIH SKUPINAH IN SPOLU, 30. JUNIJ 1998
POPULATION BY FIVE-YEAR AGE GROUPS AND SEX, 30 JUNE 1998



Vira: Ministrstvo za notranje zadeve - Centralni register prebivalstva (CRP), Uprava za upravne notranje zadeve.
Sources: Ministry of the Interior - Central Population Register (CPR), Administrative Internal Affairs Directorate.

Iz oblike porazdelitev lahko marsikaj izvemo o prebivalstvu in zgodovini neke dežele.

Zgled: Podatki na zemljevidu



Zelo učinkovito je tudi prepletanje zemljepisnih in statističnih podatkov.

Global Administrative Areas.

Zgledi: ManyEyes

ManyEyes; IBM Visual Communication Lab: Martin Wattenberg, Fernanda Viégas, Frank van Ham, ...

Visualizations : Distribution of US Foreign Aid over time, 1946-2005

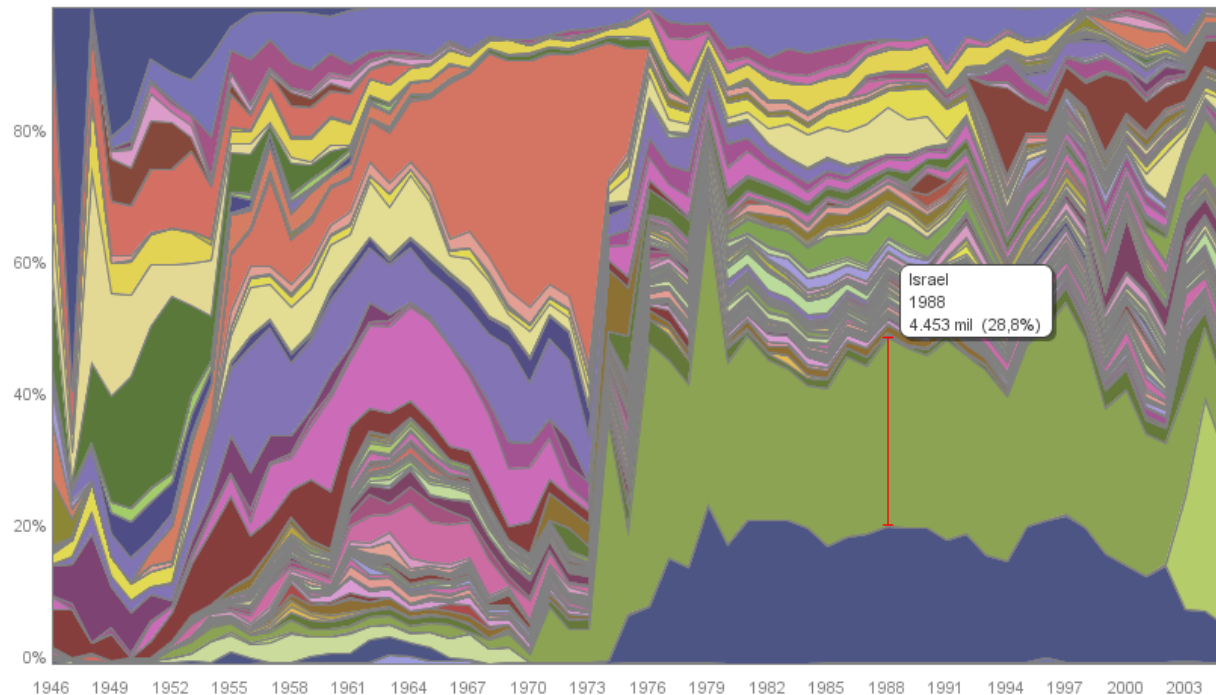
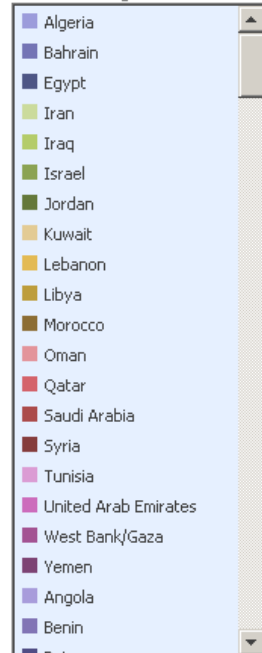
Uploaded by: **Emile Daigle**

Created at: Wednesday November 21 2007, 01:38 PM

Tags: **aid foreign**

Legend

Click to select,
Ctrl-Click: multiple
Shift-Click: range



Click or ctrl-click to highlight points on graph.

☒ % of items shown

Aggregate items with same label: Average

Sort: labels data order

Visualizations : Bubble Chart of Change in CO2 Emissions by State

Uploaded by: Anonymous

Created at: Monday June 25 2007, 03:06 PM

Tags: Emissions CO2 WITS,

Sector

Click to select,
Ctrl-Click: multiple
Shift-Click: range

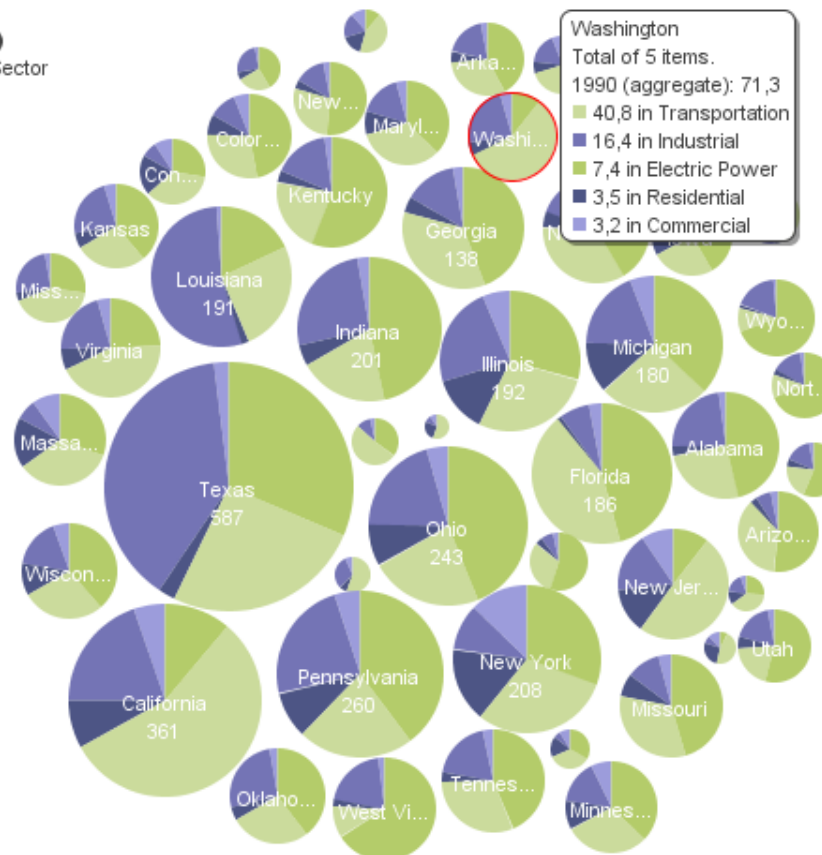
- Commercial
- Industrial
- Residential
- Transportation
- Electric Power

1990 (aggregate)

Disks colored by Sector

500
300
200
100
0

Search>>



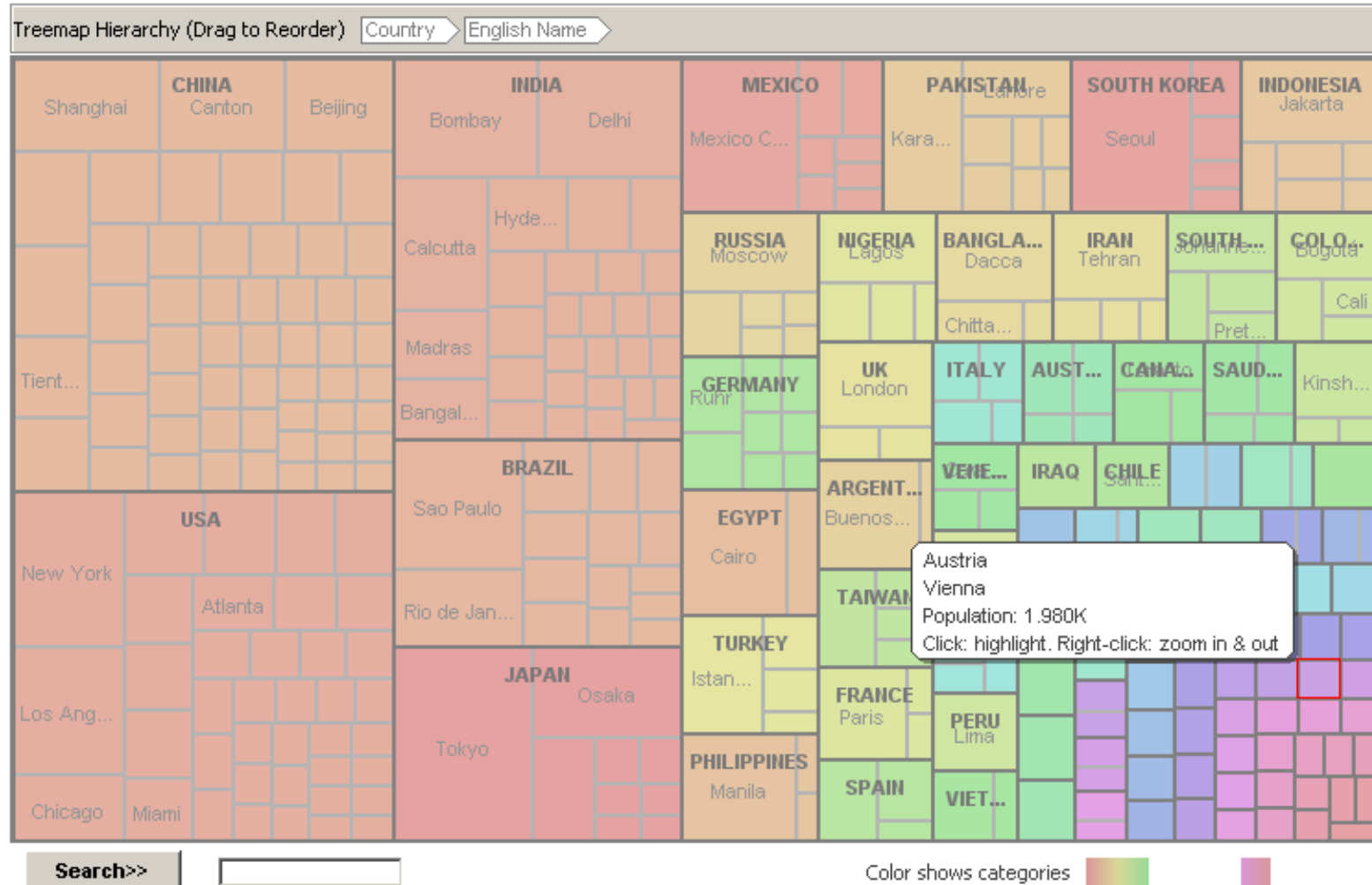
To highlight or find totals
click or ctrl-click.

Visualizations : World Cities with Populations over 1 million

Uploaded by: **crc stats**

Created at: Wednesday June 04 2008, 02:18 PM

Tags: **cities**

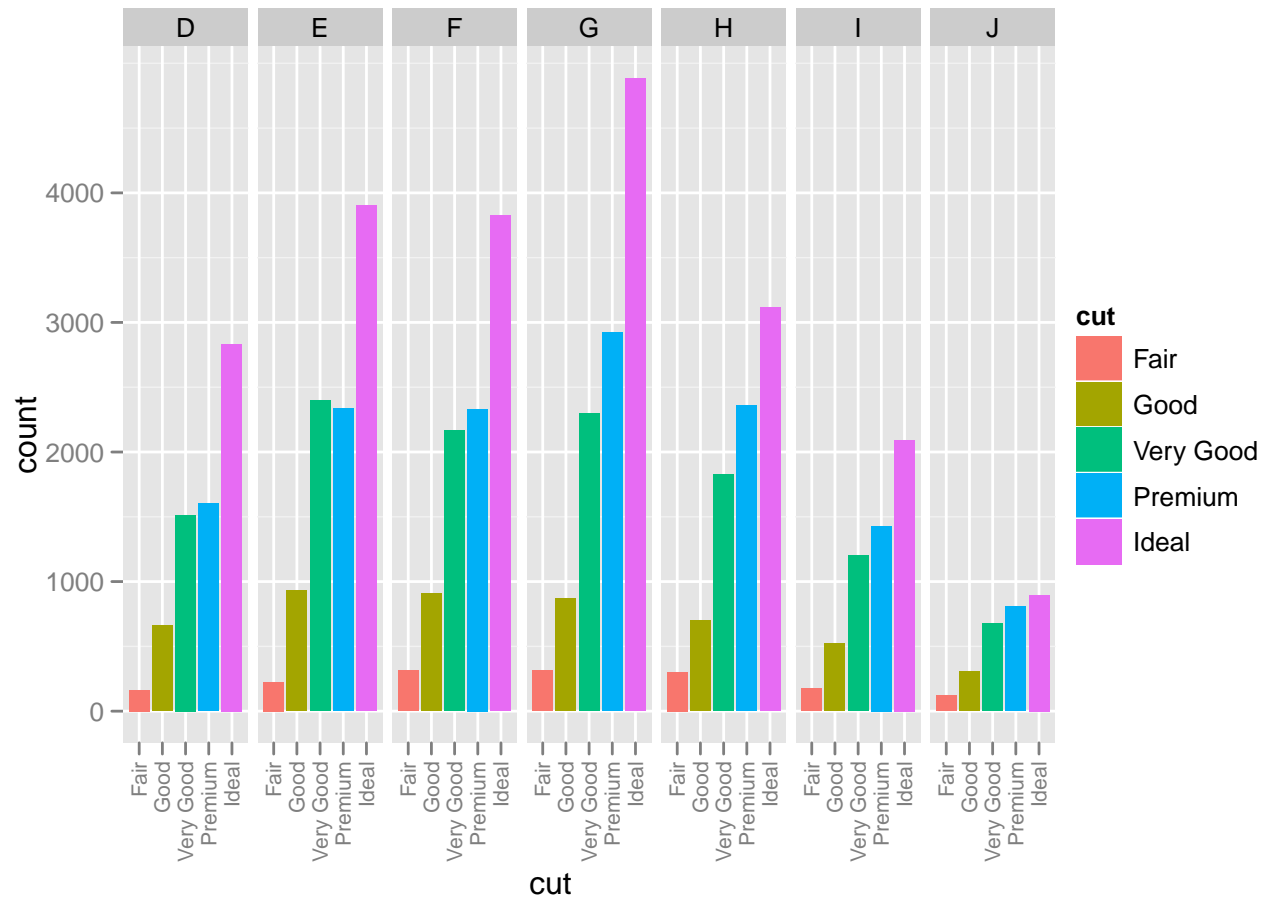


Tags: [music](#) [music,techno,electronic](#)



Zgledi: ggplot2

```
> library(ggplot2)
> qplot(cut, data = diamonds, geom = "bar", fill = cut) + facet_grid(. ~ color) +
+   opts(axis.text.x = theme_text(angle = 90, hjust = 1, size = 8, colour = "grey50"))
> ggsave("cut.pdf", width = 7, height = 5)
```



Knjižnice

Za prikaze podatkov imamo v R-ju na voljo nekaj knjižnic. Osnovna knjižnica je `graphics`. Pozna tudi S-ovo knjižnico `lattice`. V razvoju je nova knjižnica `grid`. Trirazsežne interaktivne prikaze omogočata knjižnici `rgl` in `rggobi` (vmesnik za program GGobi).

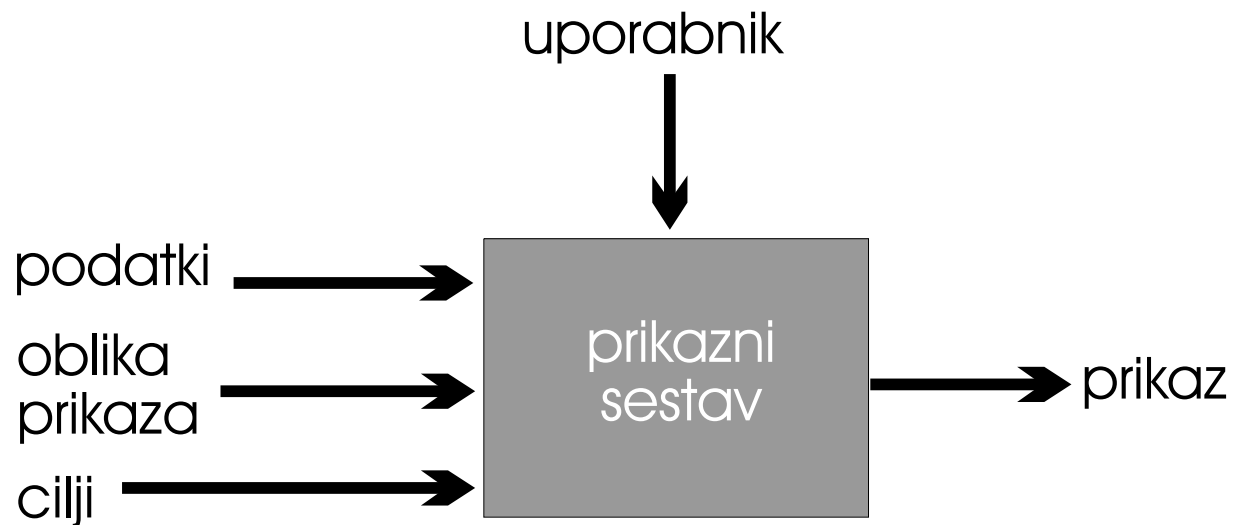
Različne posebne slikovne prikaze najdemo tudi v drugih knjižnicah. Na primer: drevesa razvrstitev in prikaze na zemljevidih.

Leland Wilkinson je napisal temeljno knjigo o prikazih podatkov *The grammar of graphics* in s sodelavci razvil okolje **nViZn** (za SPSS).

Na Wilkinsonovi knjigi temelji tudi R-jeva knjižnica `ggplot2`. Za prikaze podatkov na spletu so pred kratkim razvili v javascriptu knjižnico **Protovis** in njeno naslednico **d3**.

Prikazni sestav

Pri prikazovanju podatkov uporabnik sledi različnim ciljem: najde/določi, razlikuje, uvrsti, razvrsti, uredi, primerja, poveže, ...



Izhodne naprave

Večina starejših virov o prikazovanju podatkov je usmerjena na 'list papirja'.

Računalniški zaslon ponuja veliko novih možnosti

- vzporedni prikazi
- začasne sestavine (oznake, pojasnila, zastavice, ...)
- poudarjene izbire, povezave med prikazi
- sodejnost (interaktivnost)
- občutek trirazsežnosti (gibanje)
- programski nadzor (gibljive slike - animacije)

Te možnosti je potrebno kar se le da izkoristiti za podporo reševanja nalog analize podatkov. Prihodnost je v sodejnih, dinamičnih, trirazsežnih in večsličnih (lokalni/celotni pogled, izbire, primerjave, ...) prikazih.

GGobi

Slikovne sestavine

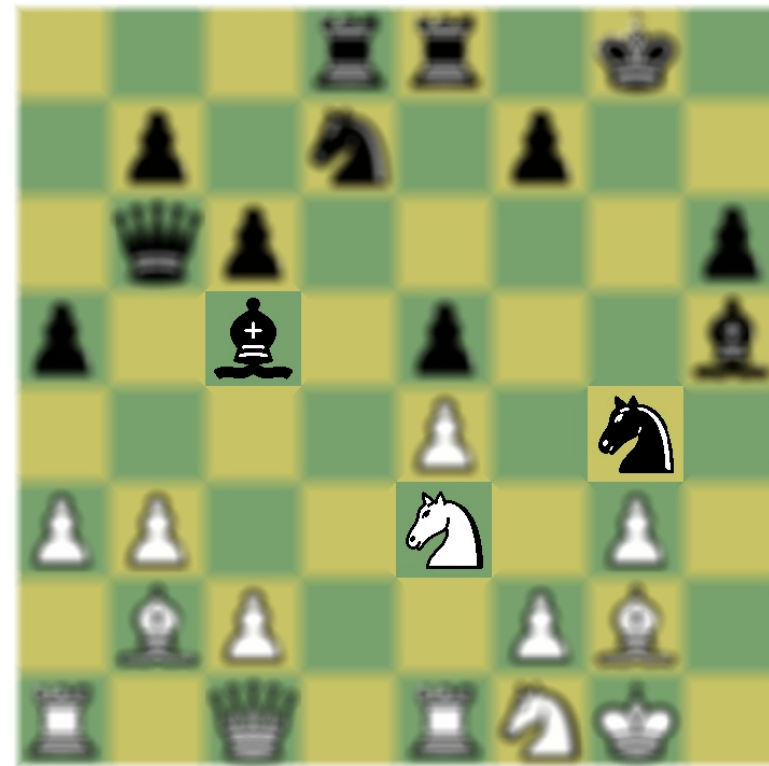
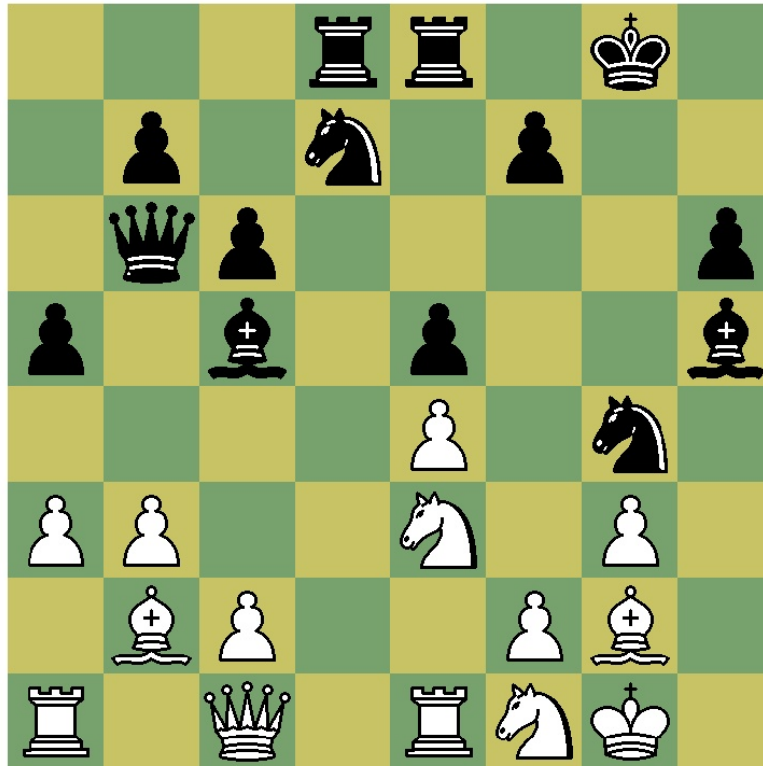
Slikovne sestavine so lastnosti s katerimi lahko predstavimo izbrano lastnost podatkov.

- polnila: barva (hue), nasičenost, svetlost, ostrina, vzorec
- prostorske sestavine:
 - splošne: dolžina, velikost, lik, mesto, smer, simetrija, položaj (spredaj / zadaj), zasuk, razdalja (med gradniki)
 - posebne 3D : zakrivanje, perspektiva, prozornost, globina, gibanje
- razsežnosti: pika, daljica, lik, telo

Slikovne sestavine so osnova različnih **slikovnih prikazov**.

Barve v R-ju

Slikovne sestavine – globinska ostrina



Semantic Depth of Field

Weber-Fechner-jev zakon

Ernst Weber (1795-1878) je opazil, da naj bi bila *najmanjša še zaznavna sprememba* občutka ψ sorazmerna razmerju spremembe dražljaja φ in njegove velikosti

$$d\psi = k \frac{d\varphi}{\varphi}$$

Zaznavamo torej relativne spremembe dražljajev.

Od tu je Gustav Theodor Fechner (1801-1887) izpeljal zvezo

$$\psi = k \ln \frac{\varphi}{\varphi_0}$$

kjer je φ_0 mejna velikost (ne)zaznavnosti dražljaja.

vid: $k = \frac{1}{60}$; bolečina: $k = \frac{1}{30}$; vonj: $k = \frac{1}{4}$; slanost: $k = \frac{1}{3}$.

Na primer: opaznost različnosti dveh dolžin je odvisna od razmerja teh dveh dolžin, ne pa od njune razlike. Opaznost razlik pri primerjavi dolžin lahko bistveno povečamo z vpeljavo oporne mreže.

Wikipedia

Stevens-ov zakon

Poskusi v 30. letih 20. stoletja so pokazali, da za večino dražljajev Fechnerjev zakon ne velja. Izkazalo se je, da velja potenčna zveza, ki jo je predlagal Stanley Smith Stevens (1906–1973):

$$\psi = k(\varphi - \varphi_0)^c$$

kjer so k , c in φ_0 konstante.

Natančnost naših zaznav je za lastnosti s c različnim od 1 odvisna od velikosti vrednosti – zaznave so pristranske.

Poskusi so pokazali, da je c za dolžino med .9 in 1.1; za ploščino med .6 in .9; in za prostornino med .5 in .8 .

[Wikipedia](#)

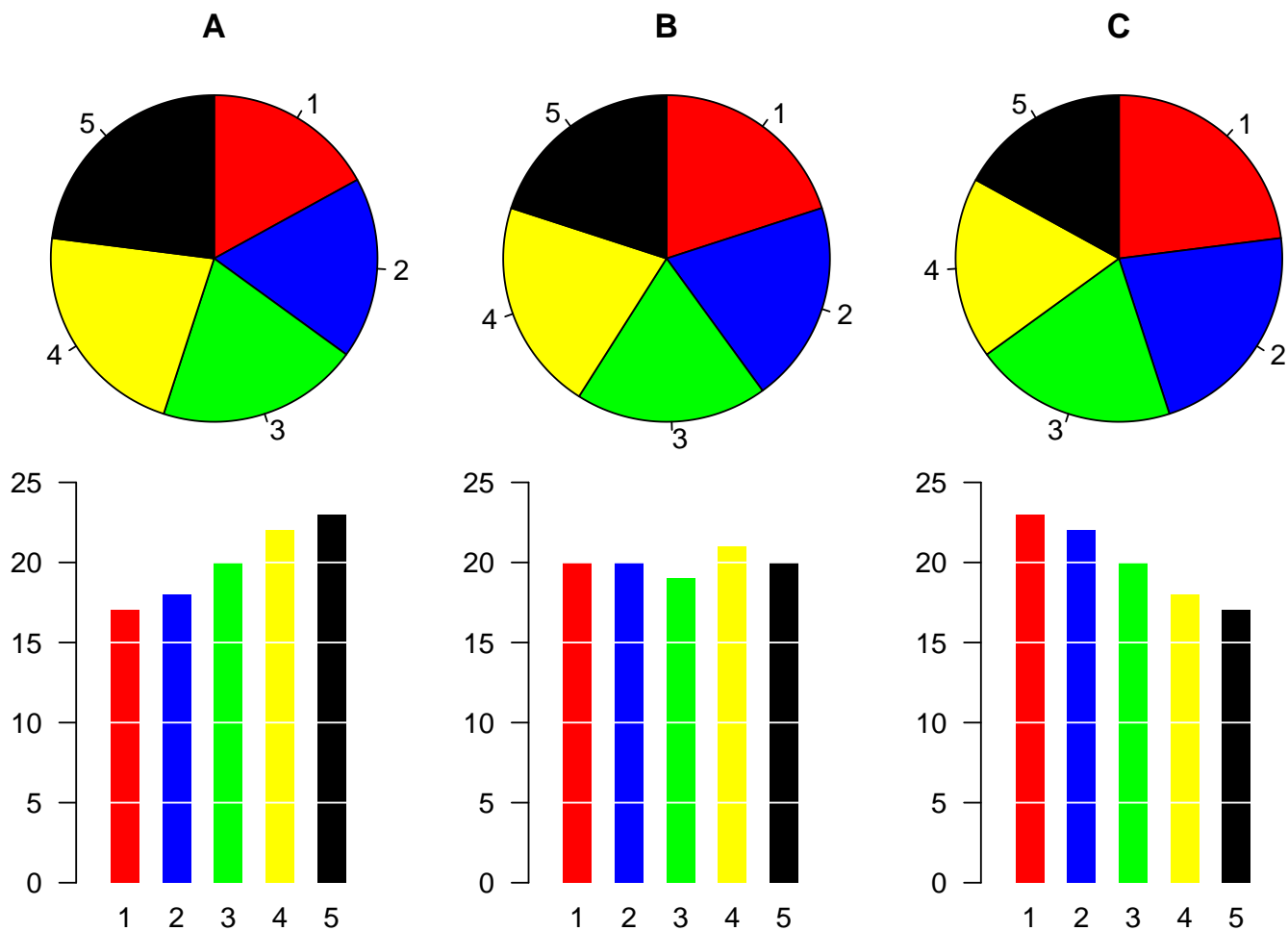
Posledice

Poskusi na osnovi Stevensovega zakona so dali naslednjo urejenost slikovnih sestavin glede na natančnost zaznavanja:

- mesto glede na skupno os (najbolj natančno)
- lega na dveh enakih, a neporavnanih oseh
- dolžina
- kot – naklon
- ploščina
- prostornina
- barva (hue) – barvna nasičenost – barvna svetlost (najmanj natančno)

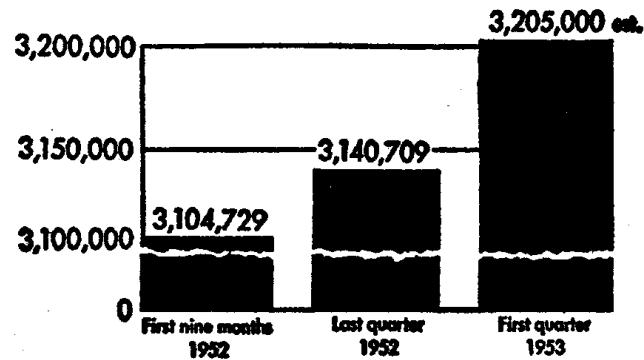
Ljudje smo različni: barvna slepota, navajenost, izurjenost, izobraženost, ...

Posledice

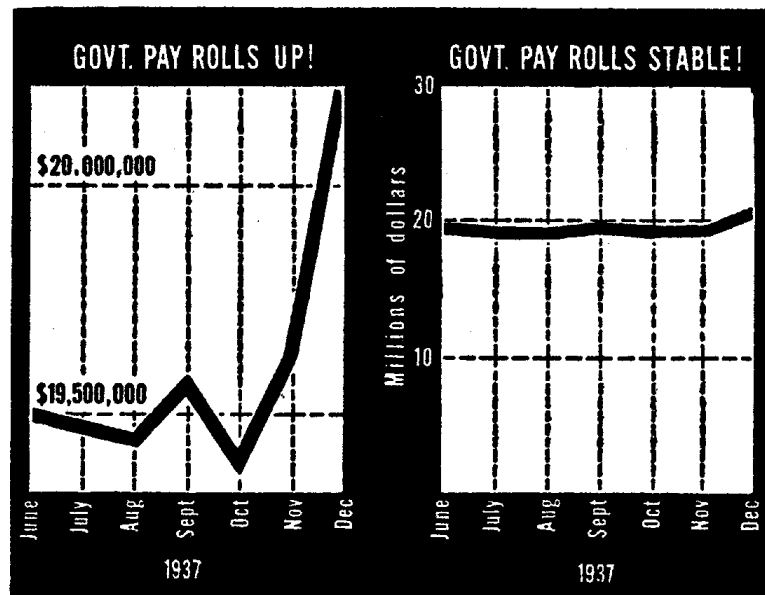


Wikipedia: Pie chart

Previdno



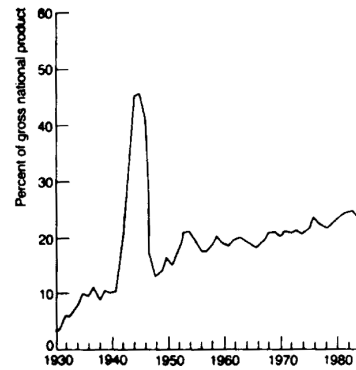
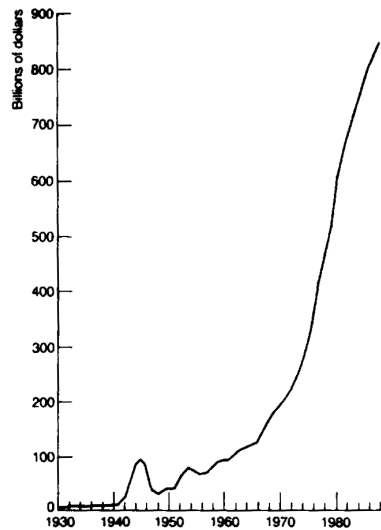
From an April 24, 1933, newspaper advertisement for COLLIER'S



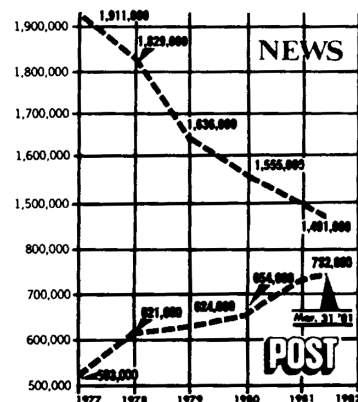
Pri tolmačenju prikazov moramo biti previdni in natančno pregledati prikaz, ne pa se kar prepustiti prvemu vtisu.

Pogosto je v prikazih podan samo zanimiv izrez iz podatkov.

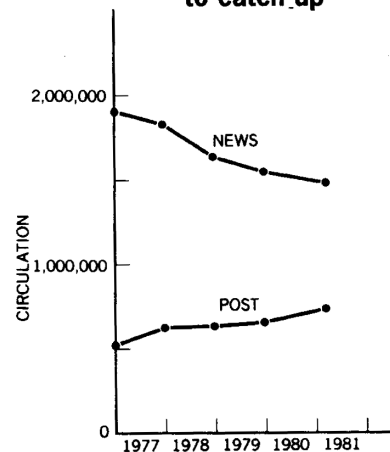
...Previdno



**The soaraway Post
— the daily paper
New Yorkers trust**



**The Post struggles
to catch up**

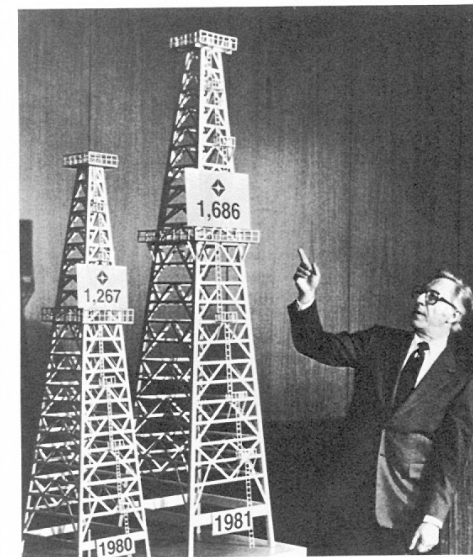
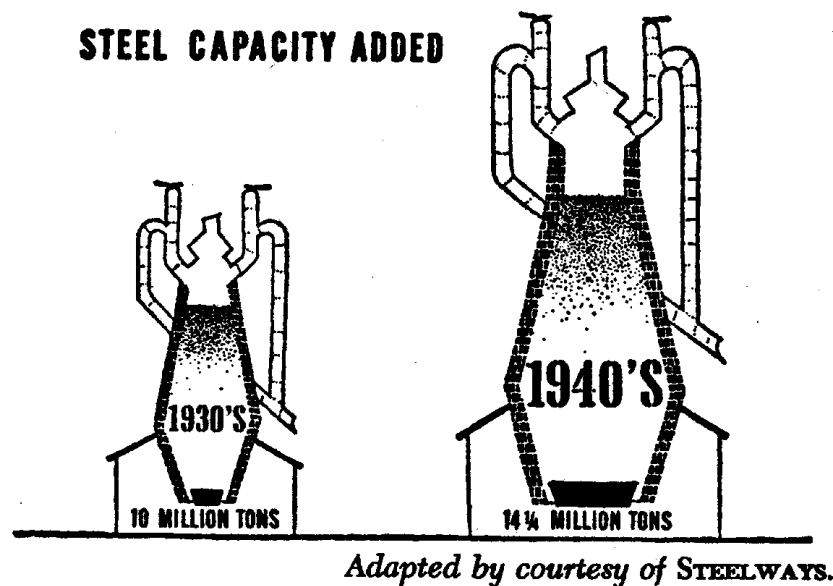


Kadar se pomen surovih podatkov spreminja skozi čas, jih moramo postaviti na 'skupni imenovalac'.

Neprekinjenost osi zavaja.

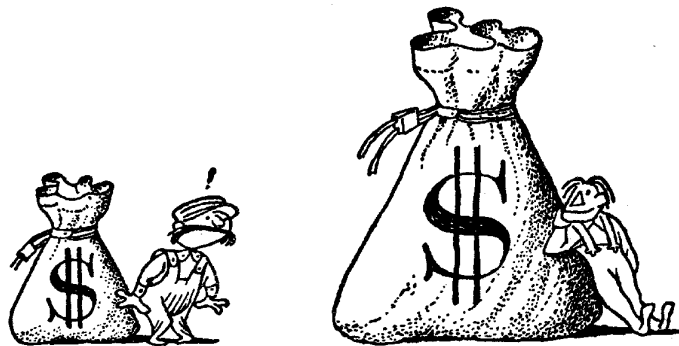
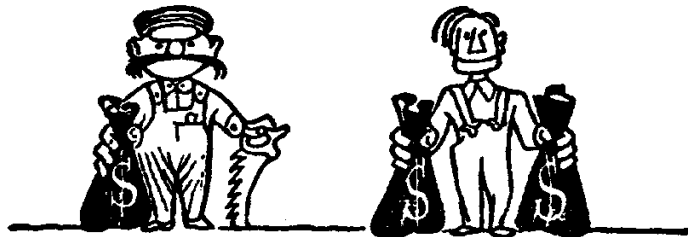
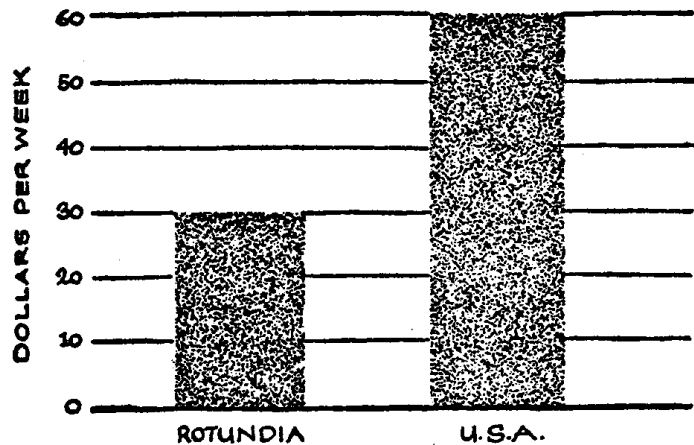
Slikovne prevare

Naše opažanje je odvisno tudi od **lastnosti** našega (človeškega) sistema videnja. Te lastnosti je potrebno upoštevati pri načrtovanju prikazov. **1, 2, 3, 4; iluzije**. Poglejmo za primer prikaz količine proizvodnje jekla (iz Darrell H.: *How to Lie with Statistics*).



Desna sličica je 1.4 krat višja od leve, ker se je proizvodnja povečala za 1.4 krat. Toda, ker pri primerjanju likov ljudje primerjamo njihove ploščine, zaznamo razliko kot dvakratno ($= 1.4^2$).

...Previdno: časopisi

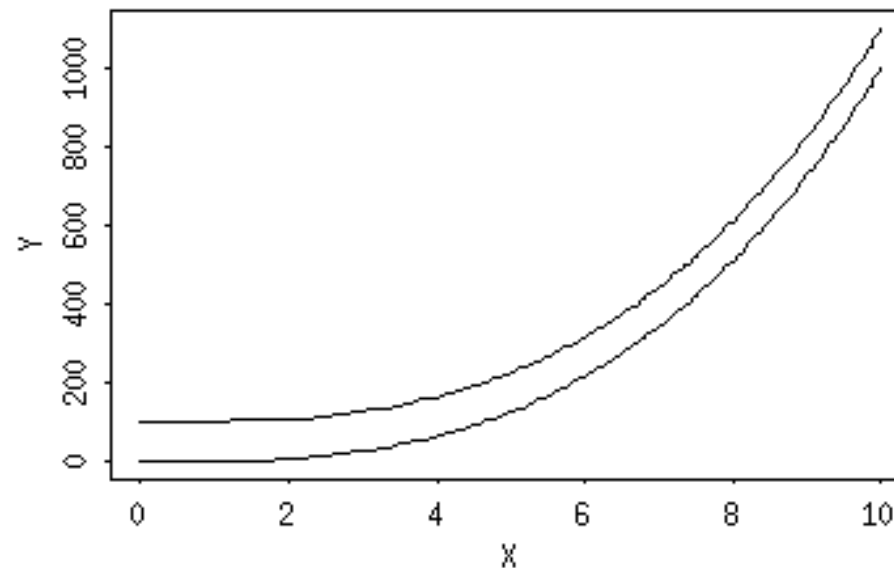


Pogosto, na primer v časopisih in predstavitev, je prikaz 'okinčan', da bi bil privlačnejši.

Toda pri tem moramo biti previdni – lik s podvojeno višino zaznamo kot štirikrat večjega.

...Slikovne prevare: krivulji

Pri primerjavi dveh krivulj nas običajno zanima **navpična razdalja** – razlika funkcijskih vrednosti; tisto, kar običajno opazimo pa je **najmanjša razdalja** med krivuljama.



Kaj je vaš *prvi vtis* o obeh funkcijah? Težita skupaj?

Podatkovja

Na izbiro oblik prikaza močno vplivajo:

- *velikost*: majhno, veliko, ogromno, 'neskončno';
- *gostota*: razpršeno, gosto, več skupin; in
- *dejavnost*: ustaljeno, spreminjajoče (deterministično, naključno).

Majhna podatkovja lahko prikažemo v celoti in v vseh podrobnostih na eni sliki. V celovitem prikazu obsežnega podatkovja se podrobnosti zgubijo; podrobni prikaz pa lahko zaobjame le del podatkovja.

Zunanji in notranji pogled

Običajni prikazi podatkov, kakršne najdemo v splošno namenskih programih (Excel, PowerPoint, ...) omogočajo praviloma pogled *od zunaj* – uporabnik je postavljen izven prikaza. V teh primerih je tretja razsežnost uporabljena predvsem za 'okraske', in ne za ustreznejši prikaz.

Ena od osnovnih značilnosti navidezne resničnosti (Virtual Reality) je podpora *notranjega* pogleda – uporabnik je vključen v prikaz kot njegov dejavni del in lahko potuje po prikazu. Slika je določena z njegovim mestom v prikazu.

Krmilje

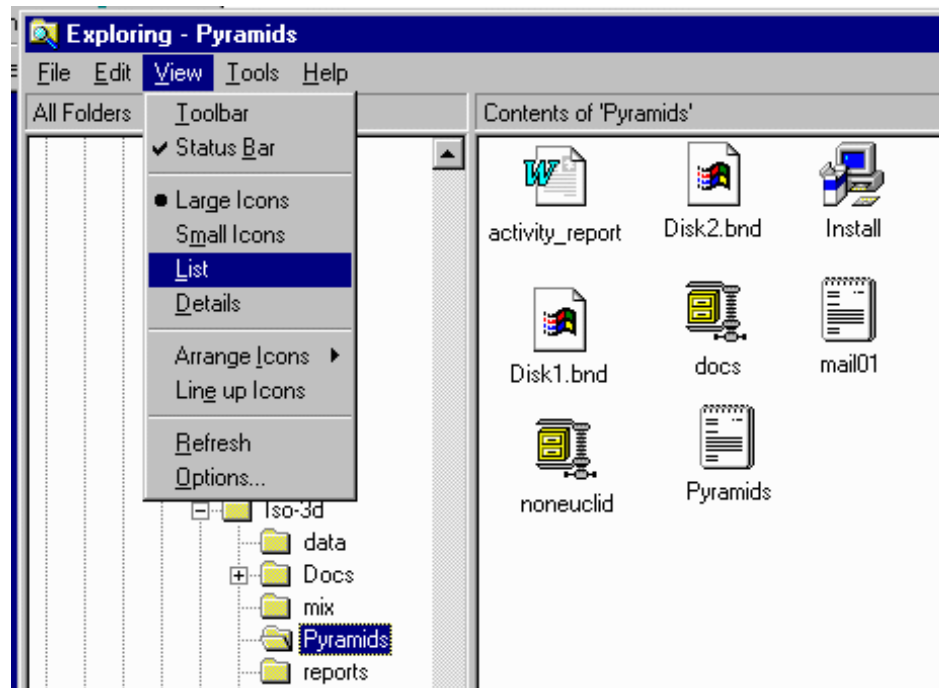
Pri pregledovanju obsežnih podatkovij se lahko uporabnik “*izgubi v gozdu*”. Pri usmerjanju uporabnika so na voljo različni prijemi:

- *ponoven začetek*: vrne uporabnika v izhodišče;
- vpeljava dodatnih *orientacijskih sestavin*: prikaz osi, mreže, sence, znamenja (stalna / uporabniška). Te sestavine lahko uporabnik vklaplja in izklaplja.
- *več slik*: prikaz na vsaj dveh slikah (oknih):
 - *zemljevid*: celovit (praviloma zunanji) pogled, ki prikazuje tudi trenutno mesto uporabnika in omogoča 'dolge' korake (skoke). Za zelo obsežna podatkovja lahko vključuje možnost povečave oziroma 'ribjega očesa'.
 - *trenutna okolica*: prikaz izbranega dela podatkovja.

Koristna je uvedba **sledi** / **nazaj** / **ponovi** in *vodenih sprehodov*.

Očala in lupe

Z večsličnimi prikazi so tesno povezani pojmi *očal*, *lup* in *povečevanja* (**Pad++**, **inXight**).



Z izbiro različnih očal dobimo različne poglede na iste podatke – ki podpirajo različne vidike/cilje prikaza. Očala delujejo na celo sliko, lupe pa samo na izbrano območje slike. Pristop poznamo iz slikovnih vmesnikov.

Prikaz lahko dopolnimo s *podpornimi sestavinami*: oznake, mreže, pojasnila, ...

Prikazi večrazsežnih podatkovij

Naj bo $E = \{X_i\}$ *množica enot*. *Enota* X je običajno opisana z naborom vrednosti izbranih lastnosti – *spremenljivk*

$$X_i = (V_1 = x_{i1}, V_2 = x_{i2}, \dots, V_m = x_{im})$$

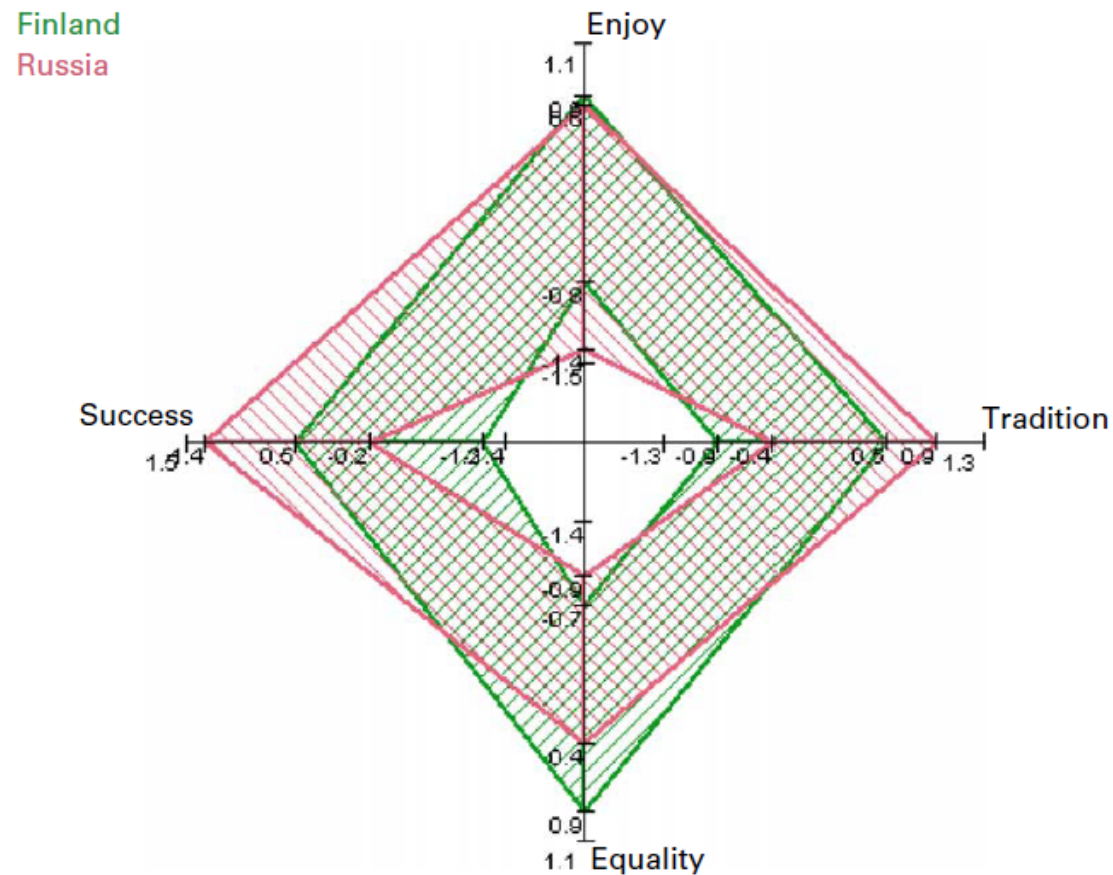
Posamezno enoto lahko prikažemo z uporabo *glifa*, ki v sebi združuje sestavine, ki predstavljajo vrednosti posameznih spremenljivk.

Znanih je več ravninskih (2D) glifov: točka v ravnini, krožni prikaz, stolpčni prikaz, zvezda, obraz Chernoffa, krivulja Andrews, tirnica, ... Večino teh rešitev je mogoče razširiti za prostorske prikaze. *

VisualComplexity

Zoom stars

Za prikaze simbolnih večrazsežnih podatkov so bile razvite **Zoom Stars** (Noirhomme-Fraiture M., 1997).



Ravni v podatkovju

Posamezne sestavine predstavitve enote podpirajo različna opravila (iskanje povezav, izbiranje, urejanje, vrednotenje, ...) nad enotami. Pri tem se pogosto prepletajo deli z različnih ravni

spremenljivka, enota, skupina, skupine, podatkovje

Večino postopkov analize podatkov je mogoče obravnavati kot transformacije ali odnose med temi ravnmi.

Slikovne sestavine in vrste lestvic

Posamezne vrste spremenljivk predstavljamo z ustreznimi slikovnimi sestavinami

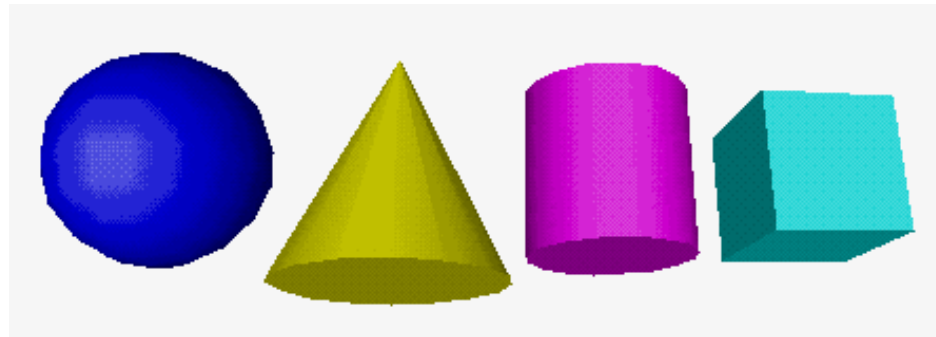
lestvica	slikovne sestavine
imenska	barva, oblika(lik)
urejenostna	stopnja, svetlost, gostota vzorca, razmestitev (mesto v)
številska	velikost, mesto, smer, kot

Ker so *številске* lestvice tudi *urejenostne*, in te tudi *imenske*, lahko ustrezne sestavine uporabimo tudi za predstavitev vsebovanih vrst lestvic – vendar pri tem zgubimo na natančnosti.

Uporaba neustreznih sestavin pa lahko privede do pojavitve neosnovanih povezav.

Trirazsežne predstavitve

Za razvoj trirazsežnih predstavitev se je uporabljal spletni jezik **VRML** (Virtual Reality Modeling Language), ki pa so ga v zadnjih letih nadomestili **X3D**, **Kinemage**, **WebGL**, ...



Pri tem lahko pri predstavitvah uporabimo naslednje sestavine:

- mesto v prostoru (x, y, z) ;
- obliko (krogla, kocka, stožec, valj, ravnina, ...);
- barva;
- velikost, kot, naklon, ploščina, prostornina;

- vzorec na ploskvah;
- smer (usmerjenost);
- besedilo;
- svetila (različni viri, sence, prozornost, odsevi, ...);
- zasuki;
- različni pogledi in načini sprehajanja po prostoru; lastnosti kamere (pravokotna, perspektiva, stereoskopska; zorni kot).
- sodejnost

Hans Rosling – Gapminder



- Gapminder
- Youtube: Hans Rosling
- The best stats you have ever seen
- New insights on poverty
- Let my dataset change your mindset
- Asia's rise – how and when
- The good news of the decade?

Viri

Jacques Bertin: **Graphics and Graphic Information-Processing**

Edward R. Tufte: **The Visual Display of Quantitative Information**

Leland Wilkinson: **The Grammar of Graphics**

Colin Ware: **Information Visualization: Perception for Design**

Paul Murrell: **R Graphics**

Hadley Wickham: **ggplot2: Elegant Graphics for Data Analysis**

Riccardo Mazza: **Introduction to Information Visualization**

Noah Iliinsky, Julie Steele: **Designing Data Visualizations**

Jock Mackinlay: **Information Visualization**

Marti Hearst: **Information Visualization and Presentation**

John Stasko: **Information Visualization**

Wojciech Basalaj: **Proximity Visualization of Abstract Data**

Chris North: **Information Visualization**

Gitta Domik: **Computer-generated Visualization**

Huff D.: **How to lie with statistics**, Norton, New York, 1993/1954.

Cleveland W.S.: **The Elements of Graphing Data**, AT&T, 1994

Friendly M., Denis D.J.: **Milestones in ... Data Visualization**

Toward a Perceptual Science of Multidimensional Data Visualization

A Periodic Table of Visualization Methods

Datavisualization: izbor orodij