
Guia de lectures sobre visualització de dades

PID_00249216

Julià Minguillón Alfonso

Temps mínim de dedicació recomanat: 3 hores





Els textos i imatges publicats en aquesta obra estan subjectes –llevat que s'indiqui el contrari– a una llicència de Reconeixement-Compartir igual (BY-SA) v.3.0 Espanya de Creative Commons. Podeu modificar l'obra, reproduir-la, distribuir-la o comunicar-la públicament sempre que en citeu l'autor i la font (FUOC. Fundació per a la Universitat Oberta de Catalunya), i sempre que l'obra derivada quedi subjecta a la mateixa llicència que el material original. La llicència completa es pot consultar a <http://creativecommons.org/licenses/by-sa/3.0/es/legalcode.ca>

Índex

Introducció.....	5
1. Pioners de la visualització.....	9
2. Antecedents històrics.....	11
3. Representacions visuals.....	14
4. Què és una visualització?.....	20
5. Tipus de dades i operacions.....	26
6. Principis de disseny.....	29
7. Eines d'anàlisi visual.....	33
8. Introducció a D3.js.....	35

Introducció

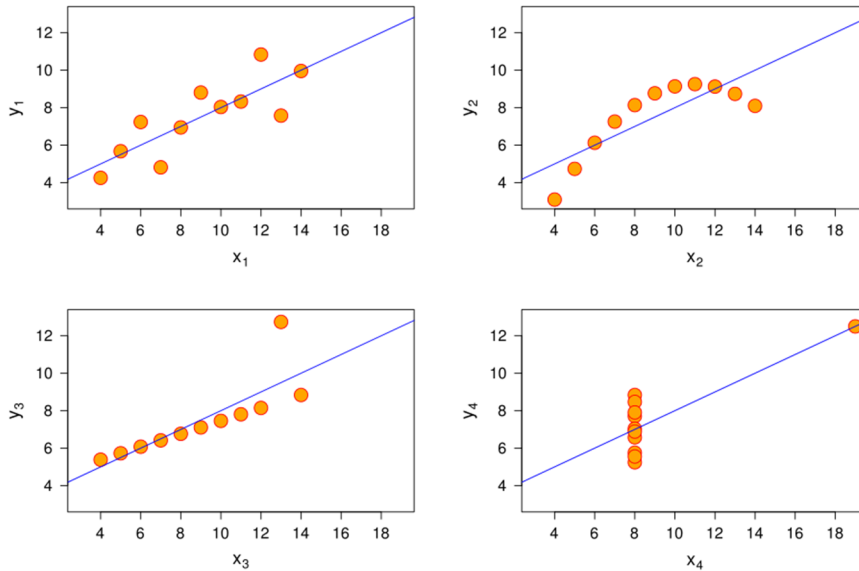
La visualització de dades és un àmbit de coneixement en evolució constant, que últimament s'ha vist impulsat a causa de la gran quantitat de dades disponibles, que esperen ser analitzades, interpretades i contextualitzades mitjançant una narrativa que combina text, imatges i altres recursos interactius, més enllà de presentacions estàtiques que fan servir certs elements gràfics com a suport. El *boom* de les xarxes socials, l'ús de dispositius mòbils i la pràctica digitalització de tots els serveis (consum, educació, lleure, etc.) permeten, a la pràctica, disposar de dades sobre qualsevol activitat humana que operi parcialment o totalment amb la tecnologia. Així, des de l'aparició d'internet a mitjan anys noranta, la capacitat de generar, processar i compartir dades ha anat creixent exponencialment, s'ha multiplicat per mil cada pocs anys i ha obligat a adoptar nous prefixos (*mega-*, *giga-*, *tera-*, etc.) que redueixin les xifres gestionades a quantitats raonables.

En aquest sentit, la visualització de dades és un dels mecanismes dels quals es disposa per a presentar tota aquesta informació raonablement per als usuaris finals, sense que aquests es vegin superats per una allau tan gran de dades. Una visualització de dades és un primer pas per a l'anàlisi i la projecció de les dades disponibles, utilitzant els recursos del sistema visual humà com a processador per detectar patrons, tendències o anomalies. En aquest sentit, una visualització permet, de manera eficient, mesurar i comparar dades, entre altres operacions. Fins i tot un simple resum d'un conjunt de dades pot ser més ben transmès i comprès mitjançant una visualització que mitjançant l'ús de text i taules numèriques.

Un exemple és el que es coneix amb el nom de «quartet d'Anscombe», creat per Francis Anscombe el 1973 per mostrar les limitacions de l'ús de descriptors estadístics per a resumir un conjunt de dades, tal com mostra la figura 1. Cadascun està compost per quatre conjunts d'onze punts en el pla de la forma (x, y) , de manera que la mitjana i la variància de cada variable, com també la correlació entre totes dues i el coeficient de la recta de regressió òptima són idèntics per als quatre conjunts, i són clarament diferenciables si s'utilitza una representació visual. Òbviament, cada conjunt representa el resultat de quatre processos diferents: una col·lecció de dades típica (figura superior esquerra), unes dades que segueixen una relació no lineal (figura superior dreta), unes dades que segueixen una relació lineal tret d'una, un possible *outlier* (figura inferior esquerra) i, finalment, unes dades que mostren una relació no lineal entre les dues variables, però en què un simple *outlier* genera un coeficient de correlació elevat. Sense la visualització d'aquestes dades fent servir un simple gràfic x - y és molt difícil fer-se a la idea de les quatre distribucions subjacents en cada cas. Encara que es tracta d'un exemple sintètic, mostra de manera

convincent les limitacions dels descriptors estadístics més habituals en els treballs de recerca i les possibilitats de la visualització com a eina d'anàlisi visual complementària.

Figura 1. El quartet d'Anscombe.



Font: Wikipedia

Per tant, acompanyar (o fins i tot substituir) les dades originals per una representació gràfica pot ser molt efectiu per a explicar-ne el perquè, especialment pel que fa a la causalitat. No obstant això, explicar històries per mitjà de les dades requereix combinar competències i habilitats de diferents àrees de coneixement, incloent-hi matemàtiques i estadística, informàtica, psicologia de la percepció i, per descomptat, disseny gràfic. Es tracta d'un àmbit clarament multidisciplinari on normalment es treballa en equip, encara que cal disposar de coneixements bàsics i d'un vocabulari comú en cadascuna de les àrees esmentades.

Per a ser un expert en visualització de dades és necessari, per tant, disposar d'un ampli ventall de coneixements, la qual cosa no és senzilla ni ràpida d'adquirir, sinó que exigeix un procés continu. Ja el 2010, Enrico Bertini (expert en visualització de dades i professor de la Universitat de Nova York) ho va descriure des de la seva experiència personal de la manera següent:

1) **Estudiar molt**, fent servir els recursos oberts disponibles a la xarxa i altres d'accessibles des d'una biblioteca (per exemple, articles en bases de dades). Bertini cita com a exemple un dels treballs clau en l'àmbit de la visualització escrit per Edward Tufte, anomenat *The Visual Display of Quantitative Information*, que, malgrat que és antic (data de 1983), continua sent una pedra angular de qualsevol treball relacionat amb la visualització de dades, tal com demostren les gairebé deu mil cites, segons Google Scholar, i la segona edició del llibre.

Enrico Bertini

Lloc web: <http://felinlovewithdata.com/>
Twitter: @FILWD

Edward Tufte

Lloc web: <https://www.edwardtufte.com>
Twitter: @EdwardTufte

2) **Robar** (o, millor dit, copiar) bones pràctiques, però no limitant-se a reproduir-les, sinó absorbint els detalls i els trucs que fan que una visualització de dades funcioni i sigui una bona pràctica. Bertini recomana estar al corrent de les principals revistes i congressos internacionals de l'àmbit, com també conèixer els treballs dels autors principals (com Nathan Yau, per exemple). Actualment, és possible seguir, mitjançant l'ús de xarxes socials (principalment Twitter), els millors especialistes de l'àmbit, atès que internet és més bon canal per a difondre aquest tipus de treballs, especialment quan demanen interactivitat.

3) **Criticar** i ser capaç de detectar els aspectes que diferencien el que és una bona pràctica del que no ho és, intentant, en aquest cas, detallar i corregir els aspectes erronis, la qual cosa permetrà veure les dificultats del problema que cal resoldre.

4) **Produir una bona visualització**, posant en pràctica els coneixements adquirits. Això pot exigir l'ús d'eines més o menys complexes (des d'Excel fins a Tableau, per exemple) o, fins i tot, llenguatges i biblioteques de programació (Processing o D3, entre altres), la qual cosa determinarà el grau de sofisticació assolible. Tal com diu Bertini, és possible crear visualitzacions excel·lents independentment de la tecnologia usada. També cal disposar de dades que es visualitzaran, encara que avui dia internet és un enorme repositori de dades sobre les quals es poden fer preguntes interessants i tractar de respondre-les mitjançant una visualització.

5) **Sortir de la zona de confort** exposant la visualització creada a l'escrutini públic, especialment mitjançant les xarxes socials o l'ús de repositoris i llocs web dedicats especialment a compartir treballs d'aquest tipus. Els tres beneficis esperats són haver de donar-li més importància i fer un bon treball abans d'exposar-lo, haver de pensar sobre l'ús que es donarà a la visualització i, especialment, obtenir retorn personalitzat d'altres usuaris interessats en la visualització de dades, de manera que sigui possible detectar i corregir els detalls que calgui.

Així doncs, una manera de començar el pla descrit per Bertini és conèixer alguns dels treballs d'autors clau en l'àmbit. Amb aquest objectiu, aquest material pretén guiar el lector per una sèrie de referències bàsiques en l'àmbit de la visualització de dades, amb l'objectiu de comprendre què és una visualització, els antecedents històrics anteriors a la situació actual d'abundància de dades i tecnologies, l'ús de visualitzacions interactives per a manipular i analitzar aquestes dades, com també els elements que componen una visualització i determinen, tant objectivament com subjectivament, la percepció de l'usuari final. Es tracta, per tant, de proporcionar un vocabulari bàsic i una introducció als principis que regeixen la visualització de dades com a mecanisme per a la creació i la transmissió de coneixement. És important recordar en aquest

Nathan Yau

Lloc web: <http://flowingdata.com/>
Twitter: @flowingdata

Processing

Lloc web: <https://processing.org/>

moment la seqüència «dades, informació, coneixement i saviesa» i evitar la discussió sobre si es visualitzen dades o informació, atès que sempre es tracta del segon, encara que s'utilitzin indistintament els dos conceptes.

Aquesta guia no és, de cap manera, exhaustiva, sinó que es tracta d'una sèrie de lectures escollides per la seva representativitat i per la importància dels seus autors dins de l'àmbit i que, d'alguna manera, estan relacionades entre si pel seu contingut. Al contrari, constantment es farà referència a altres recursos (principalment llocs web) que es poden fer servir com a punt de partida per aconseguir una visió més completa d'aquest àmbit en evolució constant. I, òbviament, aquesta guia no substitueix les lectures, sinó que serveix de punt d'entrada per a posar-les totes en context.

Els articles que constitueixen aquesta guia de lectura (per ordre d'aparició) són els següents:

- Edward Tufte (1983). «The Visual Display of Quantitative Information». *Graphics Press* (vol. 2, núm. 9). CT (EUA): Chesire.
- Michael Friendly (2006). «A brief history of data visualization». *Handbook of Data Visualization* (pàg. 15-56).
- Alan Blackwell (2011). «Visual Representation». *The Encyclopedia of Human-Computer Interaction* (2a. ed., cap. 5).
- Lev Manovich (2010). «What is Visualization?».
- Ben Shneiderman (1996). «The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations». *Proceedings of the IEEE Symposium on Visual Languages* (pàg. 336-343).
- Stephen Few (2011). «Data Visualization for Human Perception». *The Encyclopedia of Human-Computer Interaction* (2a. ed., cap. 35).
- Jeffrey Heer; Ben Shneiderman (2012). «Interactive Dynamics for Visual Analysis: a taxonomy of tools that support the fluent and flexible use of visualizations». *ACM Queue* (vol. 10, núm. 2).
- Michael Bostock; Vadim Ogievetsky; Jeffrey Heer (2011). «D³: Data-Driven Documents». *IEEE Transactions on Visualization and Computer Graphics* (vol. 17, núm. 12, pàg. 2301-2309).

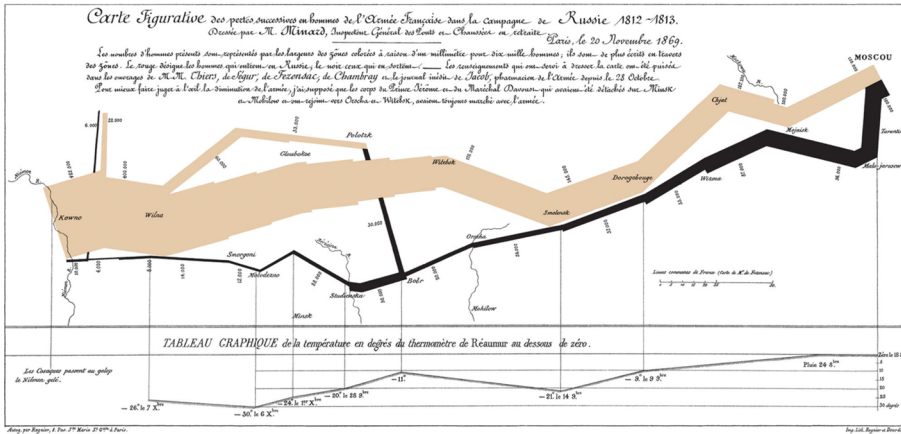
1. Pioners de la visualització

Un dels primers treballs en l'àmbit de la visualització de dades és *The Visual Display of Quantitative Information*, escrit per Edward Tufte el 1983. Tufte és professor de la Universitat de Yale i està considerat un dels pioners de la visualització de dades. Tufte va introduir l'ús de diagrames com a metodologia habitual per a la descripció de dades i la seva anàlisi preliminar com una eina més a banda de l'estadística. És l'inventor del concepte *chartjunk* (diagrama escombraries) per a criticar el mal ús de les visualitzacions quan no aporten res a les dades representades.

En aquest treball, Tufte desgrana una sèrie de visualitzacions de dades i les relaciona amb el tipus de dades que es visualitzaran i el context en el qual van ser creades, i també fa un repàs a algunes visualitzacions històriques, com el gràfic que descriu les pèrdues de l'exèrcit francès durant la campanya de Napoleó a Rússia (1812-1813), creat per Charles Minard i mostrat en la figura 2. Tufte descriu detalladament l'ús de mapes, sèries temporals i la seva combinació, cosa que ell anomena narratives espaciotemporals, com també l'ús de visualitzacions per a mostrar relacions entre elements.

Però el més interessant del treball de Tufte és el concepte d'*excel·lència*, que defineix com la comunicació d'idees complexes amb claredat, precisió i eficiència. L'excel·lència és el que proporciona a l'usuari de la visualització la quantitat d'idees més gran en l'espai de temps més curt mitjançant el mínim ús de tinta i en l'espai més petit possible. L'excel·lència és gairebé sempre multivariada, no depèn d'una sola variable. I, finalment, l'excel·lència demana explicar la veritat sobre les dades, cosa que desafortunadament sembla que s'ha perdut en moltes visualitzacions de caràcter polític, en les quals s'utilitzen les visualitzacions com a eina de manipulació. En aquest sentit, qualsevol persona interessada a ser un «periodista de dades» hauria d'adoptar les idees de Tufte com a principis bàsics.

Figura 2. Mapa figuratiu de les pèrdues successives d'homes de l'Armada Francesa en la campanya de Rússia 1812-1813, per Charles Minard (1869).



Font: Wikipedia

2. Antecedents històrics

Encara que hi ha introduccions excel·lents sobre què és una visualització de dades i el significat del concepte al llarg de la història, la lectura recomanada per a aquest apartat és l'article «A brief history of data visualization», de Michael Friendly, publicat l'any 2006 com un capítol d'un manual de visualització de dades que formava part d'una col·lecció de llibres d'estadística, la qual cosa mostra la importància de la visualització com a eina per a l'anàlisi de dades. Friendly és un autor molt prolífic en l'àmbit de la visualització, impulsor del lloc web DataVis, on es poden trobar molts altres recursos relacionats, incloent-hi articles, llibres i programari.

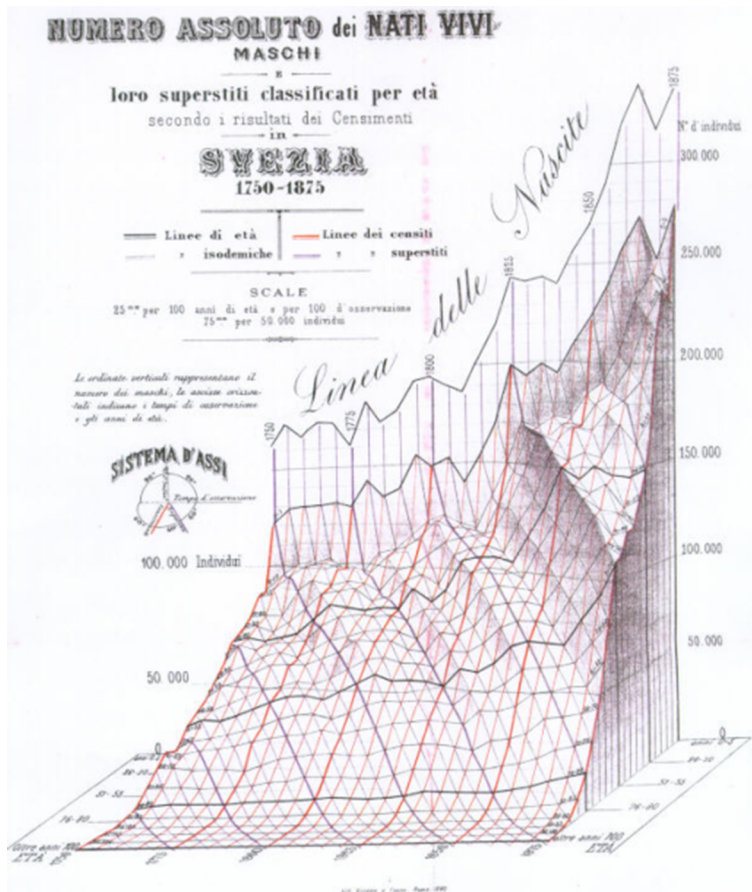
Michael Friendly

Lloc web: <http://www.datavis.ca/>

L'article de Friendly està estructurat segons una línia temporal, que inclou des de les primeres visualitzacions (de fet, mapes i diagrames) anteriors al segle XVII fins a l'actualitat (a partir de 1975), en què la tecnologia ha fet possible la creació massiva de visualitzacions. De fet, una revisió de l'article avui dia segurament afegiria una nova secció dedicada al *boom* que ha tingut la visualització recentment per la disponibilitat massiva de dades i l'aparició de tecnologies que permeten la creació de visualitzacions de dades dinàmiques, amb la utilització del web com a mitjà.

Una de les etapes més interessants destacades per Friendly és la que va tenir lloc en la segona meitat del segle XIX, quan es van desenvolupar moltes tècniques per a l'anàlisi estadística, que eren aplicades a tots els àmbits de la planificació social, la industrialització, el comerç i el transport. Això va provocar l'aparició de moltes innovacions en la visualització de dades, necessàries per a poder explicar les dades i els fenòmens tan complexos de la societat del moment. Un primer pas va ser la utilització d'elements 3D projectats com a via d'escapament del pla que fins aleshores limitava les possibilitats, com l'exemple de la figura 3. Un altre va ser la combinació de mapes amb dades de cada regió, de manera que en una mateixa representació es combinaven dades espacials amb altres de temporals. Finalment, l'ús de gràfics per a l'anàlisi estadística (la correlació era un concepte encara en desenvolupament) va permetre que Francis Galton i altres investigadors avancesin en la formalització de les observacions fetes i les convertissin en tècniques estadístiques, com mostra la figura 4. Un resultat de tota aquesta activitat va ser l'aparició d'atles estadístics, és a dir, informes de dades recopilades sobre gairebé tots els aspectes de la vida quotidiana acompanyats de gràfics detallats. Friendly destaca la col·lecció «Albums de Statistique Graphique», publicada anualment pel Govern francès entre 1879 i 1897, i que va ser discontinua pel seu alt cost de producció, com també la que va dur a terme el Govern dels Estats Units entre 1872 i 1874.

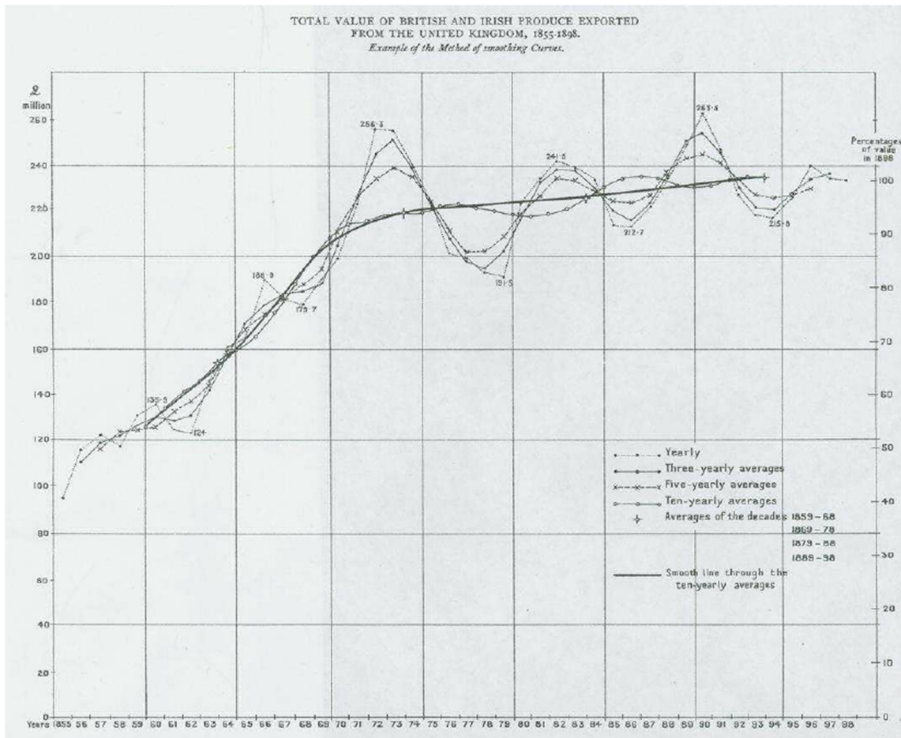
Figura 3. Població de Suècia entre 1750 i 1875, duta a terme per Luigi Perozzo el 1880.



Font: datavis.ca

Friendly destaca diversos elements que han estat clau en l'evolució del que actualment es coneix com visualització de dades. D'una banda, el desenvolupament i la formalització d'eines estadístiques per a l'anàlisi de dades i la necessitat consegüent de representar els resultats obtinguts mitjançant aquestes eines. De l'altra, l'aparició d'ordinadors i llenguatges de programació (com ara Fortran), que van permetre automatitzar càlculs i crear les primeres representacions gràfiques per a conjunts de dades amb centenars o milers d'elements.

Figura 4. Visualització de l'efecte del suavitzat en sèries temporals, per Arthur Bowley, 1901.



Finalment, Friendly resumeix l'evolució de la visualització de dades sobre la base de la necessitat de resoldre problemes concrets, relacionats amb el diseg de visualitzar fenòmens i relacions entre elements de manera diferent. Això és possible gràcies al desenvolupament de les metodologies (anàlisi estadística) i les tecnologies (ordinadors).

3. Representacions visuals

Una altra aproximació a l'evolució de l'ús de la imatge al llarg del temps es pot trobar en l'article «Visual Representation», d'Alan Blackwell, publicat el 2011 com un capítol d'una enciclopèdia de l'àmbit de la interacció persona-ordinador (en anglès *human-computer interaction*), la qual cosa no resulta estranya, ja que els éssers humans són, principalment, visuals i els ordinadors han evolucionat per maximitzar l'ús d'imatges com a interfície per a facilitar-ne l'ús.

Aquest treball té dues parts ben diferenciades. En la primera, Blackwell descriu, mitjançant l'ús d'exemples, els diferents tipus d'elements que s'han usat al llarg de la història per a la representació de dades i informació, des dels primers textos fins a l'ús de representacions basades en icones per a la creació d'interfícies d'usuari. En la segona, Blackwell presenta una aproximació holística que engloba tots els elements descrits en la primera, d'acord amb tres dimensions complementàries: el tipus d'element o recurs utilitzat, la correspondència amb la realitat representada i el seu ús en un context de disseny de la interacció.

Així doncs, Blackwell comença la primera part mitjançant la descripció dels elements que sempre han format part del conjunt d'eines bàsiques per a explicar històries, començant pel text, el qual es modula utilitzant estructures senzilles (paràgrafs, columnes i taules en pàgines) i les seves propietats (alineació, indentació, vores i ombrejos), juntament amb la tria d'una tipografia adequada (grandària, família i color). El concepte de text és molt ampli, i inclou, òbviament, l'ús de dígit i altres símbols propis d'altres sistemes de representació, com ara l'ús de lletres gregues en el cas d'equacions i fórmules matemàtiques o la notació musical, un llenguatge en ell mateix fortament visual.

L'element següent que destaca Blackwell són els mapes i diagrames, que són una evolució del concepte de símbol, ja que permeten concentrar més quantitat d'informació i establir relacions entre els elements que els componen, modificant paràmetres com la posició i la grandària, per exemple. En el cas concret dels mapes, usats des de temps ancestrals, el seu objectiu habitual és representar la realitat (una regió del nostre món 3D) mitjançant un conjunt de símbols i etiquetes, que representen i descriuen els elements que es volen destacar, com ara el mapa de Juan de la Cosa, creat just després del descobriment d'Amèrica, mostrat en la figura 5.

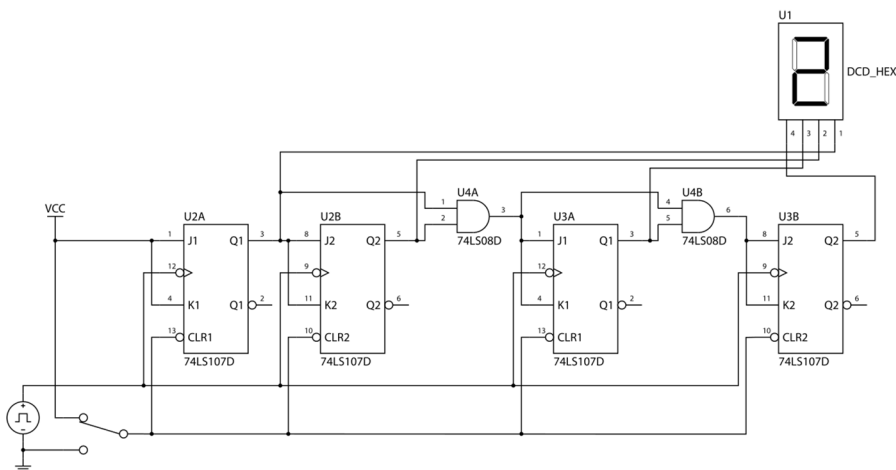
Figura 5. Mapamundi de Juan de la Cosa, cap a 1500, el primer mapa que representa el nou continent acabat de descobrir, Amèrica.



Font: Wikipedia

Per la seva banda, els diagrames són representacions esquemàtiques d'un element complex, amb l'objectiu de simplificar-ne la comprensió, que es focalitzen en els aspectes que determinen la seva naturalesa, habitualment amb un conjunt de símbols propi. Un exemple són els esquemes usats per a la representació de circuits electrònics, com el que apareix en la figura 6. En els diagrames la grandària i distància relatives dels elements que els componen no són tan importants com en un mapa, atès que s'està fent una abstracció de la realitat i se n'eliminen els detalls innecessaris.

Figura 6. Diagrama de circuit d'un comptador TTL de 4 bits.



Font: Wikipedia

Un tipus especial de diagrames són els que fan èmfasi en les relacions entre els elements que els componen, de manera que la posició exacta perd força pel que fa a l'ordre o proximitat entre els elements. Un exemple característic són els diagrames usats per a la representació de les xarxes de transport (especialment el metro, com el mostrat en la figura 7), on la posició absoluta dels elements no importa, sinó que és la posició relativa (entre aquests elements) la que aporta informació. No són un mapa pel que fa a la precisió, però serveixen igualment d'orientació. Els diagrames solen representar estructures tipus

arbre o graf, que permeten descriure relacions (en alguns casos jeràrquiques) entre elements. Blackwell esmenta la dificultat que pot comportar per a alguns usuaris entendre aquest tipus de visualització, especialment si no es relativitza el concepte de distància entre elements.

Figura 7. Diagrama de la xarxa de transport subterrani de Nova York.

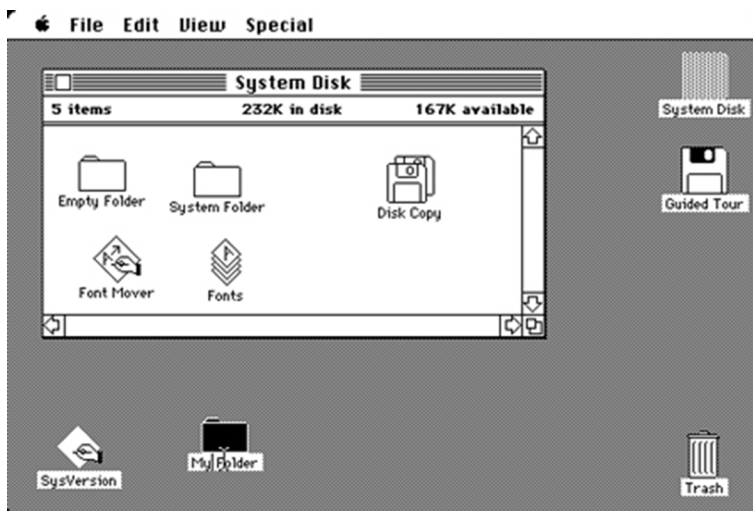


Font: Wikipedia

El pas següent que descriu Blackwell és no representar la realitat mitjançant un esquema, sinó fer-ho directament mitjançant una imatge, sigui natural o sintètica, de manera que sigui possible fer-se una idea fidedigna de l'element representat. Per *natural* s'entén una imatge que reproduïx un fragment de la realitat («analògica»), mentre que per *sintètica* s'entén una imatge que ha estat generada mitjançant un algorisme o procés («digital»). En el primer cas les imatges es capturen mitjançant la tecnologia existent, sigui la pintura o la fotografia. En el segon es tracta d'imatges generades per ordinador, com les que hi ha en un videojoc o el resultat d'un algorisme. Ara bé, la sofisticació dels gràfics generats per ordinador i la seva popularització en el món audiovisual han causat la difuminació de les barreres entre els dos tipus, de manera que en alguns casos és difícil determinar què és real i què no ho és en una imatge.

L'ús d'imatges per a representar la realitat ha evolucionat també en el que es coneix com a *icones*, és a dir, símbols que representen esquemàticament un element de la realitat quotidiana, suplantant-la i simplificant-la, fent-la fàcilment reconeixible. Un exemple són els logotips de les marques comercials, que substitueixen el nom mateix de la marca, com en el cas de Nike, el logo del qual (anomenat Swoosh) és reconegut sense necessitat de cap text. En un altre sentit, les interfícies d'usuari actuals són un compendi d'icones que representen accions (serveis) o continguts (recursos), d'acord amb certs criteris. Això sorgeix del concepte de metàfora visual, que va ser explotat inicialment pels investigadors de Xerox PARC, a Palo Alto (Califòrnia), per al desenvolupament d'interfícies d'usuari visuals, en un context en el qual els ordinadors personals es feien servir i es programaven mitjançant una línia d'ordres. Aquestes interfícies, basades en l'ús d'icones que representen elements típics d'un entorn de treball o d'una oficina, permetien un ús més senzill en reduir la quantitat de text que s'havia de llegir i escriure (és a dir, les ordres usades per a les operacions bàsiques). La metàfora d'escriptori, desenvolupada per Alan Kay el 1970, converteix la pantalla de l'ordinador en un espai virtual que s'assembla o reproduïx un entorn de treball físic habitual. Apple va ser qui ho va popularitzar el 1984 amb la interfície d'usuari de l'Apple Macintosh, mostrada en la figura 8.

Figura 8. Escriptori de l'Apple Macintosh el 1984.



Font: Wikipedia

Finalment, Blackwell es basa en la tesi doctoral de Yuri Engelhardt per a classificar tots aquests elements segons tres dimensions: el tipus de recurs gràfic utilitzat, la seva correspondència amb algun concepte del món real que vol ser representat i la seva utilització com a part del disseny que determina la visualització. Els elements s'agrupen en quatre categories:

- **Marques:** són els atributs de nivell més baix, incloent-hi forma, orientació, grandària, textura, saturació, color i tipus de línia. En aquest cas, la correspondència amb el concepte representat pot ser literal (una imitació d'alguna de les característiques físiques), un mapatge a una escala relativa o convencional (arbitrària). Els usos d'aquest tipus de recursos són mar-

Yuri Engelhardt

Lloc web: <http://datagood.org/>
Twitter: @YuriEngelhardt

car la posició, identificar categories (mitjançant formes, textures i colors) i indicar direcció i magnitud, mitjançant l'ús de símbols i codis de color senzills.

- **Símbols:** inclouen elements geomètrics, text, logotips i icones, elements pictòrics i elements de connectivitat. La correspondència pot ser topològica (enllaçant elements), descriptiva (usant convencionalismes pictòrics), figurativa (mitjançant metonímia o bromes visuals), connotativa (associada amb aspectes culturals o professionals) o adquirida (en el cas d'alfabets especialitzats). Es fan servir per a la representació de textos i el càlcul simbòlic, diagrames, marques o logotips, retòrica visual i la definició de regions.
- **Regions:** aquest concepte inclou reixetes per a alinear elements, vores, marcs, farciments, l'ús de l'espai en blanc i la idea d'integració de la Gestalt per a donar més valor al tot que a les parts que el componen. Així, es poden crear contenidors per a elements de nivell més baix, separar-los, enquadrar-los (en un sentit fotogràfic, integrant-los en una composició) i crear capes. Els usos són identificar una pertinença compartida en un cert conjunt o categoria, la separació d'elements diferents en panells o bé la col·locació d'etiquetes, subtítols o llegendes.
- **Superfícies:** és la generalització de nivell més alt, i inclou els objectes físics (3D) on es mapa (superposa) l'element gràfic. En aquest cas, la correspondència pot ser literal (com en un mapa), euclidiana (respectant escales i angles), mètrica (segons uns eixos quantitius), juxtaposada o ordenada (en regions o categories), esquemàtica o ben situada en un context. En aquest nivell, els usos són els habituals en la visualització de dades: la creació de dissenys tipogràfics, gràfics i diagrames, incloent-hi els relacionals, interfícies visuals, etc.

Amb aquesta classificació, Blackwell proporciona uns criteris bàsics per a l'anàlisi de qualsevol representació visual, indicant quins elements s'han d'identificar i quin és el seu significat, seguint l'ordre proposat. Un exemple indicat per l'autor és una partitura (mostrada en la figura 9), que parteix d'unes formes bàsiques (línies, cercles...) i una disposició (el pentagrama) per a representar una cosa tan complexa com una simfonia. Hi ha un ordre determinat per a la lectura del pentagrama, d'esquerra a dreta i de dalt a baix, incloent-hi múltiples capes (per exemple, les anotacions que fa el director o autor sobre l'obra).

Figura 9. Pentagrama d'una composició de Mozart en la seva infància.



Font: AP.

4. Què és una visualització?

Aquest apartat pren el nom d'un famós article del no menys famós Lev Manovich, teòric conegut pels seus treballs sobre els mitjans de comunicació i les transformacions motivades per l'adopció de les noves tecnologies. L'article va ser publicat el 2010, amb el títol «What is Visualization?». En aquest article l'autor presenta una anàlisi dels principis que han estat clau en el desenvolupament de l'àmbit, especialment pel que fa als nous mitjans que han impulsat i popularitzat l'ús de visualitzacions per a narrar històries. El lloc web de Manovich és també una referència per a tots els interessats en la visualització de dades, però des d'una perspectiva més àmplia, i inclou projectes relacionats amb l'ús social de la imatge (per exemple, la moda de les *selfies* o autofotos) així com altres de caràcter més teòric.

Lev Manovich

Lloc web: <http://manovich.net/>

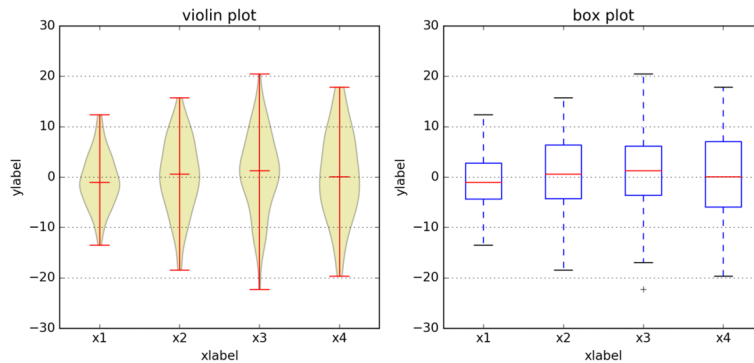
Twitter: @manovich

En aquest article, Manovich revisa la seva definició inicial de 2002, en la qual defineix una visualització com «una transformació de dades quantificades no visuals en una representació visual [d'aquestes dades]». L'autor es preocupa per definir *visualització d'informació* (*infovis*, abreujadament) de la manera més inclusiva possible, tenint en compte la diversitat de treballs que es podrien acollir en un terme com aquest. Manovich comença mostrant la diferència entre visualització d'informació i visualització científica, i indica que aquesta última (segons altres autors) es limita a dades numèriques, mentre que la primera engloba altres conceptes de semàntica més complexa, com el text o les xarxes i els grafs. Manovich no fa cap distinció segons aquest criteri, ja que, segons ell, la majoria de les visualitzacions combinen dades numèriques i no numèriques. Per a Manovich, la diferència principal entre visualització d'informació i visualització científica és l'ús de tecnologies diferents i el fet que provinguin de cultures diferents, del disseny en el primer cas i de l'àmbit científic en el segon. Igualment, Manovich es pregunta si la visualització d'informació és diferent del disseny d'informació; en aquest cas (abusant del llenguatge), és una qüestió de visualitzar dades contra visualitzar informació, respectivament. No obstant això, Manovich rebutja una distinció taxativa i considera que tots els termes s'encavalquen.

Segons Manovich, des de la segona meitat del segle XVIII fins avui hi ha hagut dos principis clau que han donat forma a la visualització d'informació. El primer és el principi de reducció, que consisteix en l'ús de gràfiques primitives (punts, línies, formes geomètriques simples...) per a la representació d'elements i les seves relacions, que revelen patrons i estructures subjacents, sense necessitat de visualitzar les dades originals. Això ha comportat una pèrdua d'importància de les dades pel que fa a les seves representacions, massa esquemàtiques en alguns casos. Un exemple és el resum d'un conjunt de dades mitjançant descriptors estadístics, com el quartet d'Anscombe ja esmentat. El més senzill és la mitjana, acompanyat habitualment de la variància, que

indica fins a quin punt les dades estan centrades al voltant de la mitjana. El pas següent és fer servir diagrames de caixa per a descriure els quartils, mostrant la distribució de les dades i l'existència de possibles *outliers*. Actualment s'utilitzen els diagrames de violí, que integren l'histograma com a part de la visualització i afegeixen informació sobre la distribució real de les dades, tal com mostra la figura 10. Segons el nivell de detall desitjat i de la naturalesa de les dades, es pot optar per una representació o l'altra.

Figura 10. Exemple de diagrama de violí com a extensió del diagrama de caixa equivalent.



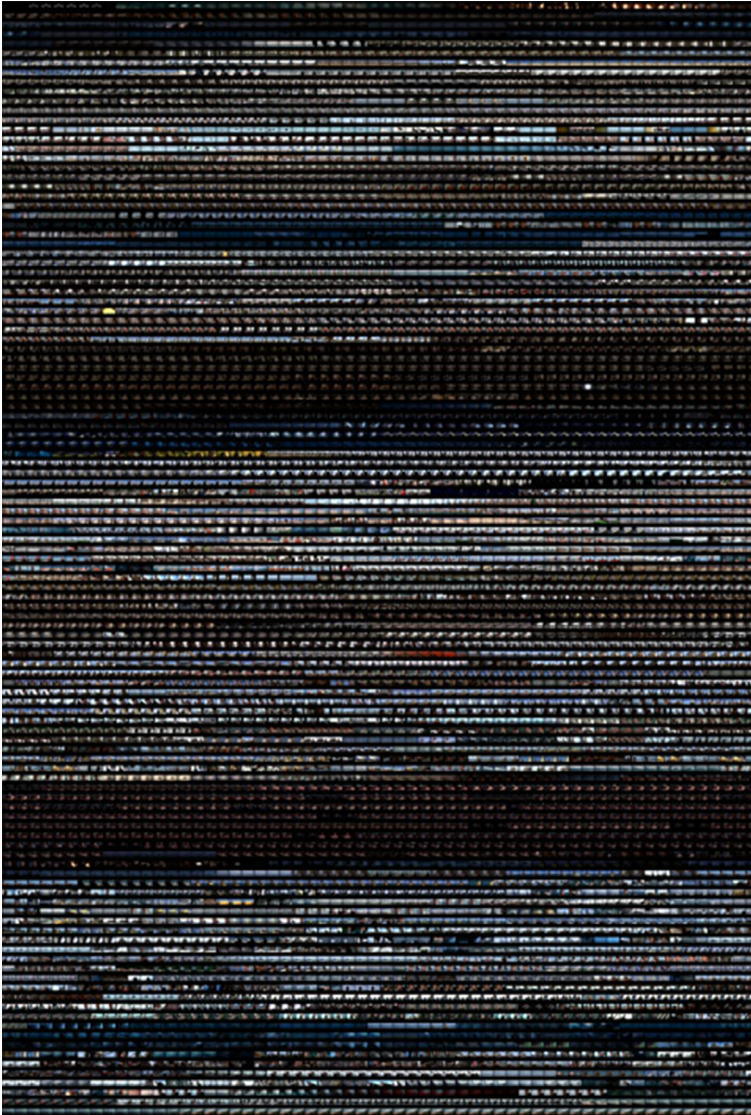
Aquest reduccionisme, present en tots els àmbits de les ciències, proposa que el món pot ser analitzat sobre la base dels elements simples que el componen i les regles que regeixen les seves interaccions, de manera que es pugui comprendre la totalitat mitjançant una descripció simplificada o reduïda. Així, durant el segle XIX es van desenvolupar tots els gràfics típics per a representar aquestes dades «reduïdes» que permeten explicar aspectes socials, demogràfics, etc. Va ser en aquesta època quan van aparèixer els gràfics de barres i de pastís, els histogrames, etc., tots conceptualitzats des d'aquesta visió reduccionista i fent servir els mateixos elements.

Així, el segon principi és l'ús de variables espacials (posició, grandària, forma, etc.) per a representar diferències en les dades i revelar, així, els patrons i les relacions existents més importants. En l'exemple (fictici) de la figura 10 es poden observar les diferències entre quatre classes diferents quant a la distribució d'una certa variable pel que fa a cada classe. Manovich fa notar que la visualització d'informació privilegia les dimensions espacials sobre altres i dona més importància a la topologia i a la geometria i menys a altres aspectes com el color, la saturació o la transparència. Així, per a representar un conjunt de dades, les dimensions més importants són assignades a la disposició espacial (anomenada *layout*), mentre que la resta de les dimensions es mapen habitualment en la resta de les variables visuals (color, etc.). En aquest cas, el color o la forma es fan servir per a dividir els elements d'un conjunt de dades en diferents classes.

Manovich fa una reflexió teòrica sobre el perquè d'aquesta distinció, és a dir, per què l'organització geomètrica dels elements en una representació ha de ser més important (per a la percepció humana) que les altres dimensions visuals (color, etc.)? Manovich esmenta el fet que tot objecte ocupa una única part de

l'espai com a possible raó, atès que el cervell utilitza aquesta informació per a segmentar el món tridimensional en una col·lecció d'objectes diferents que probablement conformen identitats diferents (per exemple, la gent, el cel, el terra, els edificis, etc.). El mateix concepte d'art modern i contemporani en corrents com l'abstracció s'ha basat a trencar amb la tradició d'identificar i representar entitats en el seu context, i donar més importància al color que a la forma, per exemple. Manovich esmenta també la dificultat de reproduir representacions gràfiques mitjançant la tecnologia existent, la qual cosa limita l'ús del color, la transparència, etc. Han estat els ordinadors els que han permès crear i manipular representacions més complexes, potenciant l'ús d'altres dimensions visuals.

Manovich prossegueix llavors amb el concepte de visualització sense reducció (o visualització directa), en la qual les dades adopten més importància que en el cas anterior. Un exemple són els núvols d'etiquetes, popularitzats per diferents eines i xarxes socials aparegudes amb el web 2.0, com Flickr o els blogs, entre altres. Un núvol d'etiquetes mostra la freqüència d'aparició de cada paraula en un text, de manera que ràpidament és possible fer-se una idea dels conceptes que hi apareixen. La figura 11 mostra un núvol d'etiquetes amb les dues-centes paraules més freqüents de l'article de Manovich. Encara que seria possible fer servir un gràfic de barres per a representar la mateixa informació, la visualització directa proporcionada pel núvol d'etiquetes és molt més rica, i alhora més agradable estèticament. No és important saber el percentatge de vegades que apareix la paraula *visualization*, sinó que es pot veure ràpidament que és una de les més usades. De fet, una visualització directa interactiva podria proporcionar aquesta informació si l'usuari se situa damunt d'una paraula en concret, per exemple; és a dir, aportant detalls només quan són requerits, a diferència d'una visualització basada en la reducció, en la qual aquests detalls són, precisament, els que governen la visualització.

Figura 12. Resum visual de «Jaws», per Brendan Dawes, 2004.

Font: Brendan Dawes.

Aquest tipus de visualització és impossible de fer sense una tecnologia que permeti accedir al contingut i manipular-lo segons unes regles. Dawes va utilitzar el llenguatge de programació Processing per a capturar els fotogrames de cada pel·lícula, reduir-los de grandària i reordenar-los en forma de matriu. Aquestes noves visualitzacions apareixen de les possibilitats que ofereix la tecnologia per a manipular dades íntegrament, no mitjançant l'ús de les reduccions habituals. Manovich destaca que Dawes fa servir fotogrames reals, no el color mitjà. No cal reduir les dades (a un sol nombre, usant descriptors estadístics) per a destacar-ne patrons, sinó que es pot utilitzar mostreig per mitjà de les dades originals. A més, no cal seguir cap configuració espacial concreta, sinó que es fa servir l'ordre natural de les dades.

Manovich acaba el seu article amb una reflexió al voltant del concepte de visualització directa, i es planteja si és, o no, un mètode diferent de la visualització d'informació, que actualment es continua basant en l'ús de gràfiques primitives. Atès que les visualitzacions directes també permeten extreure patrons i revelar estructures en les dades, Manovich planteja que es poden conside-

rar visualitzacions tradicionals, tot i que en la majoria dels casos vinguin d'autors amb referències molt diferents del tradicional de l'àmbit estadístic o científic. Ha estat la disponibilitat de tecnologies com Processing el que ha permès que altres autors creessin visualitzacions en les quals els paràmetres que en determinen l'aspecte no són els tradicionals, aprofitant les possibilitats de manipulació de la mateixa manera en tots. No obstant això, encara hi ha una limitació en forma de l'amplada de banda necessària per a transmetre imatges (mapes de bits, dades completes), així que els gràfics tradicionals (vectorials, basats en reduccions) continuaran sent omnipresents, encara que cada vegada més complexos. En aquest sentit, llenguatges com ara Processing o, millor encara, D3 són la prova que Manovich té raó quan diu que la visualització directa o sense reducció és un nou paradigma per a ser explorat, especialment des d'àmbits com les ciències socials i les humanitats.

5. Tipus de dades i operacions

Encara que és relativament antic (data de 1996), l'article «The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations», de Ben Shneiderman, continua sent una referència (ha estat citat més de tres mil set-cents vegades segons Google Scholar) pel que fa al disseny d'interfícies visuals per a la manipulació de dades. Citant el crític d'art Ernst Hans Josef Gombrich, Shneiderman va presentar aquest treball en un simposi sobre llenguatges visuals i va mostrar la necessitat de categoritzar els elements que componen una representació visual i les operacions que se'n deriven.

En aquest treball, Shneiderman introdueix el seu famós mantra sobre com s'ha de plantejar la cerca visual d'informació: «Overview first, zoom and filter, then details-on-demand»; és a dir, primer una visió general, després seleccionar i ampliar la zona d'interès, i finalment afegir el detall necessari segons les necessitats de l'usuari. A partir d'aquesta idea, Shneiderman defineix set tasques que es poden aplicar a conjunts de dades, que també són (coincidentment) de set tipus diferents:

- Dades unidimensionals: dades organitzades seqüencialment, que són consumides segons l'ordre en el qual s'han generat. Per exemple, documents de text, llistes d'elements, etc.
- Dades bidimensionals: dades que representen una certa estructura espacial, incloent-hi, per tant, el concepte de posició. Un exemple obvi són els mapes o diagrames.
- Dades tridimensionals: en alguns casos, es vol representar la posició d'un conjunt d'elements en el món real, i en aquest cas cal utilitzar tres dimensions. Molts dels primers esforços en l'àmbit de la visualització de dades pertanyien a aquesta categoria, com els sistemes d'informació geogràfica, els entorns de disseny assistits per ordinador (CAD) per al modelatge de peces o en l'arquitectura, o les imatges resultants de processos d'escombratge del cos humà.
- Dades temporals: són les relacionades amb esdeveniments que tenen un inici i un final i que es poden encavalcar. És important no confondre-les amb les dades unidimensionals, que són generades per un procés amb un cert ritme, com ara un conjunt de mesuraments de la temperatura diària en un lloc.
- Dades multidimensionals o n -dimensionals: en la majoria dels casos els ítems emmagatzemats en una base de dades es descriuen mitjançant n

atributs, la qual cosa en dificulta la visualització, ateses les limitacions del nostre món tridimensional.

- **Arbres (dades jeràrquiques):** en aquest cas es tracta de dades que tenen una relació entre elles, en la qual cada element o node té un enllaç al seu (únic) *ancestor*, excepte en el cas del node arrel.
- **Xarxes:** quan les relacions entre nodes són més generals i no hi ha restriccions, els arbres es generalitzen en grafs, que poden ser de molts tipus diferents (dirigits, bipartits, acíclics, etc.). Actualment l'anàlisi i la visualització de dades de xarxes és una de les àrees més interessants (per exemple, l'anàlisi del flux de contingut a Twitter).

Sobre aquests tipus de dades, Shneiderman defineix set operacions bàsiques, segons la idea repetida en el mantra («Overview first, zoom and filter, then details-on-demand»). Òbviament per a cada tipus de dades el significat de l'operació pot ser lleugerament diferent:

- **Overview:** es tracta de crear una vista general de totes les dades disponibles, com a punt d'entrada a la seva exploració. La vista general pot no incloure totes les dades alhora, però llavors ha de proporcionar un mecanisme perquè l'usuari pugui seleccionar altres dades al mateix nivell de detall.
- **Zoom:** es tracta d'un ajustament de la vista anterior, seleccionant un subconjunt de dades, però mantenint la sensació de posició i el context, apropant-s'hi o allunyant-se'n. El *zoom* complementa l'operació anterior per a poder fer-se una idea de l'estructura subjacent en les dades.
- **Filter:** l'usuari pot seleccionar un subconjunt de les dades, de manera que només visualitzi les que compleixin uns criteris determinats i elimini la resta. Els criteris es poden especificar de moltes maneres, segons el tipus de dades: en el cas d'una variable numèrica unidimensional es pot especificar un rang, mentre que en una de bidimensional s'utilitza el concepte de *bounding box*, o caixa que tanca un conjunt de dades.
- **Details-on-demand:** quan, mitjançant les operacions anteriors, ja s'ha seleccionat un subconjunt reduït de dades, l'usuari pot sol·licitar informació de cadascuna, normalment mitjançant l'ús del ratolí, passant-hi per damunt o fent clic a les dades, de manera que aparegui la informació addicional continguda en cadascuna.
- **Relate:** en el cas de dades n -dimensionals és possible establir criteris de distància entre els elements del conjunt, segons la seva tipologia i l'ús que es vulgui donar a la visualització. Mitjançant aquesta distància és possible, llavors, establir relacions entre elements (propers) que comparteixen una o més característiques.

- *History*: es tracta de mantenir una llista de les operacions dutes a terme, de manera que sigui possible desfer algun dels passos fets durant el procés de visualització de les dades.
- *Extract*: finalment, una vegada s'han seleccionat les dades i la informació addicional requerida, es tracta de poder bolcar-les per a poder reutilitzar-les en un altre context, emmagatzemant-les com un nou conjunt de dades.

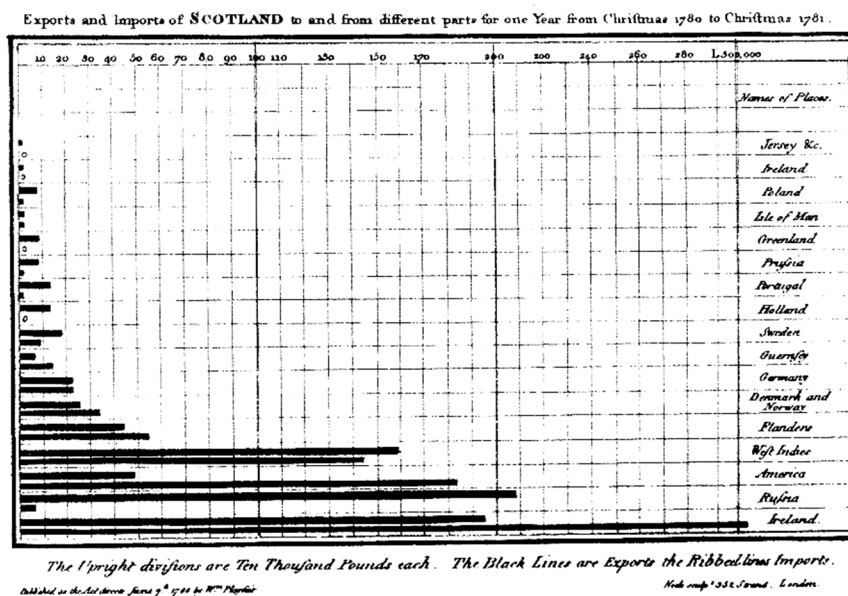
Finalment, Shneiderman proposa una manera senzilla perquè usuaris no experts puguin crear seleccions complexes (és a dir, filtres) a partir de les dades disponibles, combinant operadors booleans senzills dinàmicament per obtenir visualitzacions de dades adaptades a les seves necessitats. Si la visualització no és usable i no permet a l'usuari respondre les seves preguntes, llavors no és útil. Shneiderman, de manera gairebé profètica, acaba el seu article destacant que l'ús dels ordinadors ha generat una explosió de dades difícilment gestionable, però que, alhora, és possible utilitzar-los per a visualitzar aquestes dades mitjançant diferents eines conegudes (que combinen les operacions descrites anteriorment) com també altres (usant paraules textuais de l'autor) que encara han de ser «domades i validades».

6. Principis de disseny

També forma part de l'enciclopèdia de la interacció persona-ordinador l'article «Data Visualization for Human Perception», publicat el 2011, que presenta els aspectes relacionats amb la percepció que són més importants a l'hora de prendre decisions pel que fa al disseny d'una visualització. En aquest article, Stephen Few descriu uns principis bàsics de percepció visual i cognició que s'han de respectar, de manera que sigui possible traduir una idea abstracta a un conjunt d'atributs físics, entre altres la forma, la posició, la grandària i el color, dels elements que compondran la visualització.

És una dita popular que «una imatge val més que mil paraules», però això no-més és cert si la història es pot narrar millor visualment i està ben dissenyada. Per exemple, analitzar taules amb dades numèriques no dona una idea de la seva estructura; és molt millor visualitzar-les i detectar-hi ràpidament tendències, patrons, màxims i mínims, etc. Aquest és el poder de la visualització de dades, poder captar ràpidament els aspectes més característics d'un conjunt de dades mitjançant les capacitats del sistema visual humà. El mateix serveix per a altres tipus de dades, com les xarxes o les estructures d'elements relacionats entre si. Few descriu l'evolució de l'ús de gràfics per a la representació de dades des del treball pioner de William Playfair (1759-1823), que va fer servir per primera vegada gràfics de diferents tipus (de línies, de barres i de pastís), tal com mostra l'exemple de la figura 13.

Figura 13. Importacions i exportacions d'Escòcia a altres regions del món, per William Playfair, 1786.



Fent servir diferents exemples, Few introdueix les idees que haurien de regir qualsevol visualització de dades:

- Una visualització hauria d'indicar clarament com els valors mostrats es relacionen (comparen) els uns amb els altres, o cadascun amb la totalitat.
- Representa les quantitats de manera precisa.
- Facilita la comparació entre quantitats.
- Facilita establir l'ordre dels elements segons la quantitat que representen, és a dir, detecta màxims i mínims.
- Deixa clar com s'hauria de fer servir la visualització i quins són els seus objectius, i per a què s'hauria d'usar.

Així, es pot observar que un gràfic tan estès com el gràfic de pastís no és eficient, ja que no facilita la comparació entre quantitats ni permet reordenar fàcilment els elements una vegada representats. És molt més efectiu un gràfic més senzill com el de barres, principalment perquè es percep millor una dimensió lineal (una barra) que una d'angular (un segment de pastís).

Segons Few, la visualització de dades funciona perquè canvia l'equilibri entre la percepció i la cognició, aprofitant les capacitats visuals del cervell humà. Few es basa en la psicologia de la Gestalt, un corrent de la psicologia moderna que tracta (entre altres) el concepte de percepció per descriure com els éssers humans perceben patrons, les seves formes i la seva organització quan visualitzen alguna cosa. Alguns dels principis que determinen aquests aspectes són els següents:

- Principi de proximitat: els objectes que es troben a prop els uns dels altres es perceben com un grup.
- Principi de similitud: els objectes que comparteixen atributs similars (per exemple, la forma o el color) es perceben com un grup.
- Principi d'adjunció: els objectes que semblen tenir una frontera o vora al seu voltant (una línia o una àrea de color) es perceben com un grup.
- Principi de clausura: una estructura oberta (parcialment) es percep com tancada, completa i regular sempre que hi hagi la possibilitat raonable d'interpretar-la d'aquesta manera.
- Principi de continuïtat: els objectes que estan alineats o que apareixen un a continuació de l'altre es perceben com un grup.

- Principi de connectivitat: els objectes que estan connectats (per exemple, mitjançant una línia) es perceben com un grup.

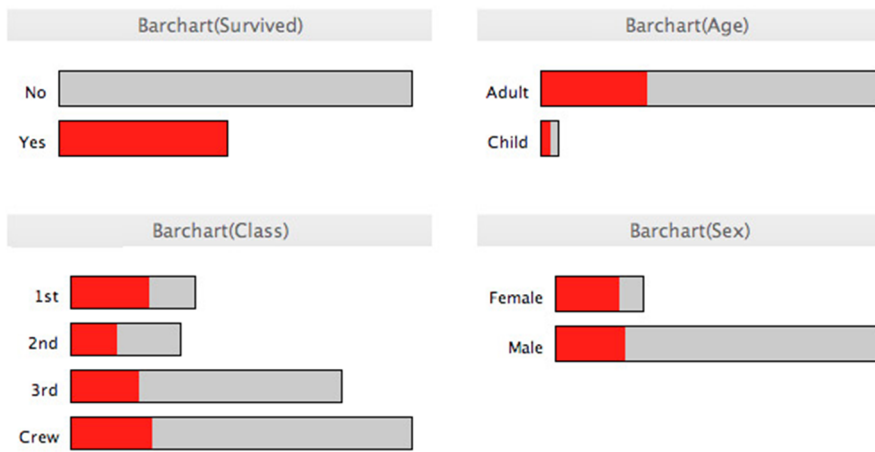
Aquests principis i les idees esmentades anteriorment es poden entendre millor mitjançant els avenços en dos àmbits d'estudi: el processament visual preventiu (en anglès, *preemptive*) i els mecanismes i les limitacions de l'atenció i la memòria. És sabut que el processament visual és més ràpid que el verbal. En part, perquè el sistema visual humà fa tasques de nivell baix molt ràpidament, atès que estan codificades mitjançant circuits neuronals específics, per detectar atributs bàsics com ara longitud, grandària, color (to i intensitat), angle, textura i forma, entre altres. A més, les limitacions per a processar i recordar múltiples elements simultàniament fan que sigui més eficient utilitzar una visualització. Segons Few, la visualització de dades pot estendre la capacitat d'anàlisi, sigui mitjançant l'ús de visualitzacions simples però efectives o, en un futur, mitjançant nous usos, incloent-hi l'ús d'interfícies complexes per a interactuar amb visualitzacions de manera senzilla, integrant l'anàlisi estadística i l'ús de mineria de dades per a l'extracció de coneixement.

Finalment, es pot destacar un comentari de Robert Kosara en relació amb el concepte de metàfora visual, molt estès actualment i esmentat amb anterioritat en l'article d'Alan Blackwell. Kosara, un dels responsables actuals de l'eina de visualització de dades Tableau i del blog *eagereyes*, presenta breument la idea de metàfora visual com una manera de forçar la visualització, d'acord amb altres criteris (potser purament estètics), especialment pel que fa a la interacció. Kosara descriu diferents operacions que es poden fer servir per a manipular dades en una visualització interactiva, des de les més senzilles associades al moviment del ratolí sobre els diferents elements que la componen fins al *linking* (quan una visualització mostra el mateix subconjunt de dades en diferents vistes mitjançant els mateixos atributs i estableix una relació o enllaç entre ells) i el *brushing* (quan l'usuari selecciona diferents elements per crear un subconjunt de dades d'interès), tal com mostra la figura 14. Així, una vegada s'han seleccionat (*brushing*) els passatgers que van sobreviure a l'accident del *Titanic* en la figura superior esquerra, la visualització mostra els mateixos passatgers repartits segons les altres categories (dimensions) usades per a classificar-los, la qual cosa permet veure, per exemple, que el percentatge de dones que van sobreviure és molt més gran que el d'homes.

Robert Kosara

Lloc web: <http://kosara.net/> i
<https://eagereyes.org/>
Twitter: @eagereyes

Figura 14. Exemple de *brushing* i *linking* en una visualització de dades sobre els passatgers i la tripulació del *Titanic*.



7. Eines d'anàlisi visual

Juntament amb Ben Shneiderman, Jeffrey Heer revisa en l'article «Interactive Dynamics for Visual Analysis: a taxonomy of tools that support the fluent and flexible use of visualizations» (publicat el 2012) les operacions bàsiques descrites anteriorment, però des d'una perspectiva de les eines i tecnologies amb les quals s'executen. De fet, aquest article desenvolupa el treball previ de Shneiderman esmentat amb anterioritat i classifica dotze accions o dinàmiques en tres grans blocs, afegint-hi algunes de noves i revisant-ne d'altres respecte a la llista inicial de set tasques de l'article de Shneiderman:

- **Especificació de les dades i la vista:** inclou visualitzar dades mitjançant una representació pictòrica, filtrar les dades no rellevants per a centrar-se en les que sí que ho són, ordenar les dades per a exposar els possibles patrons i, finalment, derivar valors o models a partir de les dades.
- **Manipulació de la vista:** inclou seleccionar elements per a poder-los destacar, aplicar-hi filtres o manipular-los, navegar per les dades per a poder detectar-hi patrons de nivell alt i també el detall de nivell baix, coordinar vistes per a explorar dades multidimensionals i organitzar espais de treball i finestres múltiples.
- **Processament i procedència de les dades:** inclou emmagatzemar l'històric d'anàlisis fetes per a poder-lo revisar, visitar i compartir, crear anotacions sobre la base dels patrons descoberts, compartir vistes i anotacions per a promoure la col·laboració i guiar els usuaris per tasques i històries (narratives).

Els autors descriuen cadascuna de les operacions i actualitzen el treball anterior de Shneiderman, tenint en compte els avenços tecnològics que hi ha hagut des de llavors. Així, per a cadascuna de les operacions esmentades, els autors proporcionen també exemples usant diferents eines per a la visualització de dades, incloent-hi R i ggplot2, Protovis (el precursor de D3), Tableau o IBM Many Eyes, entre altres. L'objectiu dels autors és mostrar la visualització de dades com una eina més per a l'anàlisi i la presa de decisions, com també la narració d'històries basades en les dades.

En aquest sentit, un dels objectius de qualsevol visualització de dades és permetre'n una anàlisi preliminar, de manera que sigui possible detectar patrons, tendències, etc. Tanmateix, això pot ser difícil amb les dades originals, a causa de la seva naturalesa, per exemple. La visualització hauria d'integrar la transformació de les dades originals, la creació de noves variables, el càlcul i la integració (en la mateixa visualització) de descriptors estadístics, i també la creació de models (estadístics o de mineria de dades) per a l'extracció de

Jeffrey Heer

Lloc web: <http://homes.cs.washington.edu/jheer/>

Twitter: @jeffrey_heer

coneixement, facilitant la detecció de grups o de les variables més rellevants, per exemple. Els autors destaquen que aquesta és una de les àrees (coneguda com *visual analytics*) en les quals encara cal avançar, combinant (superposant) dos àmbits que fins ara s'havien plantejat seqüencialment, l'anàlisi i la visualització.

La navegació per les dades que formen part de la visualització és una altra de les operacions que pren més importància, especialment en escenaris de dades massives, on hi ha un gran volum de dades molt diverses i canviant en el temps. El mantra de Shneiderman («Overview first, zoom and filter, then details-on-demand») pot no ser adequat si la primera operació ha de bregar amb un hipercub (el producte cartesià de les tres dimensions esmentades, volum, varietat i velocitat de les dades, les tres «ves» de les dades massives) enorme, per la qual cosa cal establir mecanismes per a facilitar una aproximació inversa: «search, show context, expand-on-demand», és a dir, anar de l'element seleccionat (trobat en una cerca) fins a la totalitat. El canvi de paradigma provocat pel que es coneix com dades massives ha causat també la necessitat de repensar el que s'entén per una visualització, atès que no és possible (ni té sentit) visualitzar-ho tot. Encara que la capacitat de càlcul dels ordinadors és brutal, ho és encara més el volum de dades que es generen o es capturen, de manera que cal replantejar-se totes les operacions descrites pels autors per a fer-les eficaces i eficients.

Finalment, els autors plantegen aquesta taxonomia com un punt de partida per a les persones interessades en l'àmbit de la visualització de dades, i proporcionen una gran quantitat de referències bibliogràfiques i exemples, identificant els aspectes clau per al desenvolupament de visualitzacions i les àrees de recerca més interessants, com ara nous mètodes per a l'especificació de vistes interactives, la integració d'anàlisi i visualització de dades o l'ús d'anotacions per a afegir semàntica a les visualitzacions, amb la qual cosa en faciliten també la utilització.

8. Introducció a D3.js

Finalment, una pregunta òbvia una vegada s'arriba a aquest punt és «com?», és a dir, de quina manera es poden construir aquestes visualitzacions més enllà de l'ús d'una eina concreta, des de la perspectiva d'un desenvolupador d'aplicacions o, més ben dit, de visualitzacions de dades. Encara que ja hi ha diferents biblioteques de programari per a crear gràfics i visualitzacions més o menys complexes (per exemple, *prefuse* o *Protovis*; aquest últim és l'antecessor de D3, o el llenguatge de programació *Processing*), és a finals de 2011 quan Michael (Mike) Bostock, Vadim Ogievetsky i Jeffrey Heer presenten D3 en l'article «D³: Data-Driven Documents», publicat en la revista *IEEE Transactions on Visualization and Computer Graphics*, una de les més prestigioses de l'àmbit.

El nom donat a aquesta nova biblioteca, D3 (o també D3.js, ja que es tracta d'una biblioteca escrita en JavaScript), no és gratuït, ja que porta implícit un canvi de paradigma en la manera com es construeix una visualització. D3 permet crear gràfics interactius combinant elements de les diferents capes que componen una pàgina web: llenguatge HTML, manipulació dels objectes que constitueixen el document mitjançant nodes DOM (*document object model*), aplicació d'estils mitjançant CSS (*cascading style sheets*) i, òbviament, JavaScript, amb el qual s'obtenen gràfics en format SVG (*scalable vector graphics*) que poden ser visualitzats directament per qualsevol navegador web. Així, *data-driven documents* (D3) fa referència al fet d'enllaçar directament dades amb elements del DOM, que permet la seva actualització immediata segons els canvis que es produeixin en les dades i genera, d'aquesta manera, visualitzacions interactives i dinàmiques. D3 no és un nou llenguatge de programació per a la creació de gràfics, l'aproximació habitual d'altres biblioteques predecessores, sinó que permet pensar en el document (és a dir, la pàgina web resultant) com una representació visual d'un conjunt de dades, segons una configuració (*layout*) i un conjunt de paràmetres que la determinen.

En aquest article, de caràcter predominantment tècnic, els autors descriuen els aspectes més importants de D3 pel que fa al funcionament intern:

- Selecció: permet accedir a un conjunt d'elements (un, diversos o tots) del document, sigui per nom, classe, identificador, atribut, etc. És possible combinar seleccions fent servir interseccions o unions.
- Operadors: sobre els elements seleccionats es poden aplicar operadors, com ara establir o canviar atributs, estils, propietats i continguts, tant textuals com HTML.
- Dades: l'operador *data* permet establir vincles entre les dades d'entrada i els elements que conformaran la visualització. Les dades es poden ordenar

i filtrar segons diferents criteris. La manera com D3 processa les dades i les vincula amb els elements és una de les qüestions clau que cal entendre per a dominar-ne el funcionament. Les dades poden «entrar» (*enter*) en el document (és a dir, ser mapades als nodes), ser actualitzades (*update*) o bé ser eliminades (*exit*).

- Esdeveniments: es proporcionen mecanismes per a supervisar la interacció mitjançant els elements habituals (teclat i ratolí).
- Configuracions (*layouts*): un dels aspectes més interessants de D3 és l'existència de diferents *layouts* per a crear visualitzacions, com si fossin plantilles, incloent-hi la majoria dels més habituals; entre altres, diagrames de Sankey, mapes d'arbre, grafs de força dirigits i mapes, tal com mostra la figura 15.
- Altres: finalment, D3 també inclou una col·lecció de gràfiques primitives, l'ús de transicions animades o l'ús d'interpolació per a crear nous elements a partir d'uns de bàsics (colors, fonts, etc.).

Els autors també analitzen un aspecte clau en el desenvolupament de visualitzacions de dades interactives (la raó per la qual fer servir D3): el rendiment o cost computacional necessari per a visualitzar un conjunt de dades. És obvi que una visualització que necessiti molt de temps per a renderitzar i mostrar un conjunt de dades no serà ben rebuda pels usuaris, així que D3 ha estat pensat per a ser molt eficient en aquest aspecte, tant en la càrrega inicial com en l'actualització de les dades. Ara bé, D3 no està pensat per a visualitzar conjunts de dades amb centenars de milers o milions d'elements; cal no oblidar que s'executa en el client (és a dir, el navegador web), per la qual cosa cal transmetre i carregar les dades a la memòria, i també processar-les, la qual cosa pot ser costós per a certs *layouts* com els grafs de força dirigits, per exemple.

Figura 15. Exemples de *layouts* en D3.js.

Font: D3js.org

D3 s'ha convertit en l'estàndard *de facto* per a la creació de visualitzacions interactives en línia, en part a causa de la gran quantitat de bones pràctiques que s'han desenvolupat a partir d'aquest programa i que l'han popularitzat, juntament amb un encertat disseny visual, molt lleuger i molt modern. El mateix Mike Bostock va liderar diversos projectes de visualització de dades mentre va ser a *The New York Times*, com ara l'anàlisi de la disputa electoral l'any 2012 per a la presidència dels Estats Units entre el demòcrata Barack Obama i el republicà Mitt Romney.

D3 continua sent un projecte viu que ha estat actualitzat fa poc i el nombre de tutorials i exemples d'ús ha anat creixent des del seu llançament el 2011, per la qual cosa és una molt bona aposta per als desenvolupadors que vulguin crear visualitzacions de dades interactives. Encara que la corba d'aprenentatge de D3 és empinada (en el sentit que és complicat començar a fer-lo servir des de zero), sí que és possible reutilitzar els exemples existents, adaptant les dades disponibles als requeriments de cada *layout*. El fet d'executar-se en un navegador web també facilita la creació de visualitzacions sense necessitat de disposar d'un entorn de desenvolupament complex.

Mike BostockLloc web: <https://bost.ocks.org>

Twitter: @mbostock

D3.js: <https://d3js.org/>

