

TDM	729.89	915.51	185.62	▲25.43%	FLR	660.27	745.28	85.01	▲12.88%
HUM	749.73	924.29	174.56	▲23.28%	UVD	155.59	181.57	25.98	▲16.70%
DMW	833.72	1004.01	170.29	▲20.43%	QUV	440.55	540.21	99.66	▲22.62%
YZJ	903.49	1127.46	223.97	▲24.79%	HZT	285.51	344.98	59.47	▲20.83%
GLY	982.07	1219.39	237.32	▲24.17%	PCW	811.44	1029.66	218.22	▲26.89%
VDA	113.74	143.41	29.67	▲26.09%	AIK	361.77	451.39	89.62	▲24.77%
UVV	468.08	535.41	67.33	▲14.38%	ZJJ	858.36	994.57	136.21	▲15.87%
HJS	545.49	659.05	113.56	▲20.82%	RHJ	894.79	1046.68	151.89	▲16.97%
ECC	586.36	654.68	68.32	▲11.65%	VDV	425.08	509.95	84.87	▲19.97%

Data & Data analyse

1038.36	125.73	▲13.78%	ZGK	991.59	491.48	99.89	▲25.51%		
1655.62	346.07	▲26.43%	BNY	959.21	1130.65	161.44	▲16.69%		
1641.68	346.49	▲26.75%	SOM	735.44	913.39	177.95	▲24.20%		
PNR	654.33	775.84	121.51	▲18.57%	TGQ	1429.91	1646.42	322.51	▲24.36%
171	101.11	115.11	14.00	▲13.84%	QIS	543.42	667.74	124.32	▲22.70%
171	101.11	115.11	14.00	▲13.84%	171	101.11	115.11	14.00	▲13.84%

Studiekring Gelderland & Overijssel

ERIC DANKAART

24-11-2022

Even voorstellen



Eric Dankaart

Fiscaal jurist

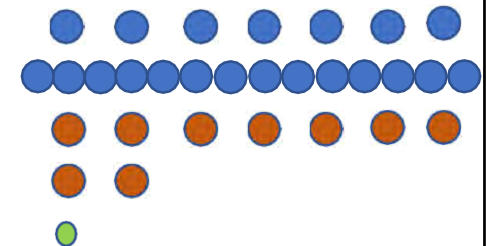
Loonheffingen inspecteur

Loonheffingen adviseur PwC

Tax Technology PwC

Tax Technology Dankaart BV

Loonheffingen Belastingdienst



De inhoud van deze presentatie is het persoonlijke gedachtegoed van de presentator, niet van organisatie(s) waarvoor hij werkt of gewerkt heeft.

Data & data-analyse

Data en de opkomst ervan

Data onderscheiden

'Data-denken'.

Data-analyse

Data-analyse, aandachtspunten

Aan de slag met data

Trends in data-analyse



Data en de opkomst ervan

Data kan verwijzen naar:

Een verzameling van tijdsaanduidingen, het meervoud van datum; zie [Datum \(dagtekening\)](#)

[Data \(Star Trek\)](#), een androïde uit *Star Trek*

[Data Records](#), een muzieklabel uit Groot-Brittannië

[Data \(geslacht\)](#), een vlindergeslacht

[Gegeven](#), een vastgelegde uitdrukking van een feit: *datum*, meervoud *data* gegevens

[Dataset](#), een gegevensverzameling

Data onderscheiden

Naar inhoud

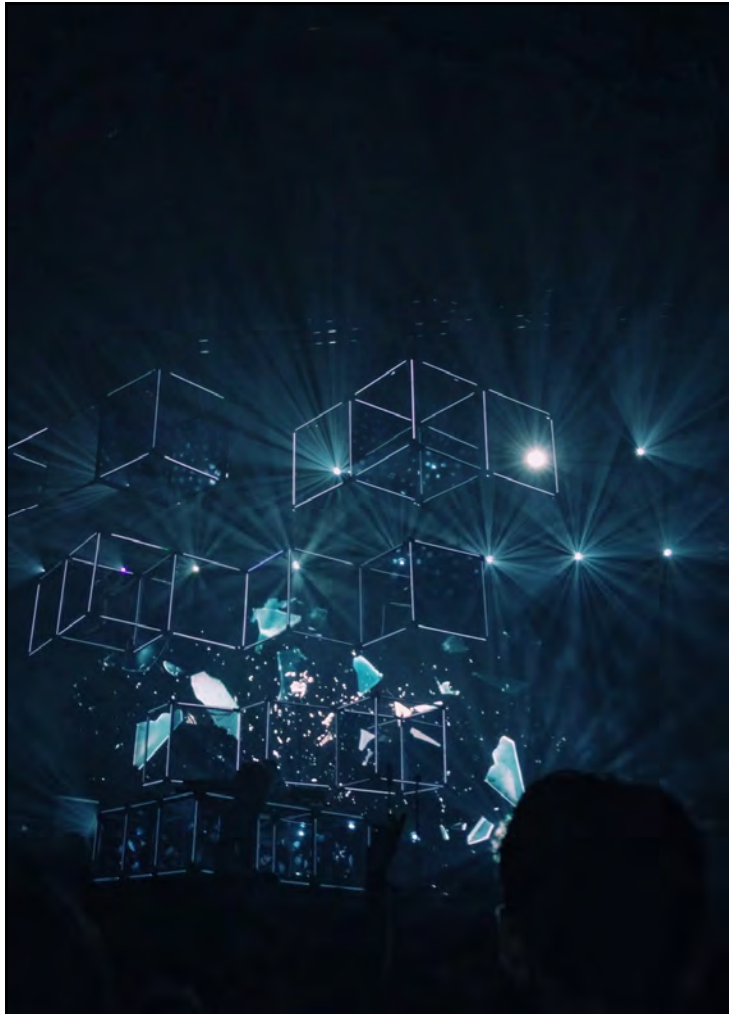
Gestructureerd - ongestructureerd

Naar omvang

Naar domein







Data naar inhoud

Een verder onderscheid is mogelijk naar bestandstype (fileformat). De meeste formats horen bij een bepaald type inhoud. Sommige formats worden door vele applicaties gebruikt, andere maar door 1 specifieke.

Kwisje! Bij welke applicatie hoort het format en/of welk soort inhoud bevat het?

.xls

MS Excel

.jpg

beeld/foto

.flac

geluid

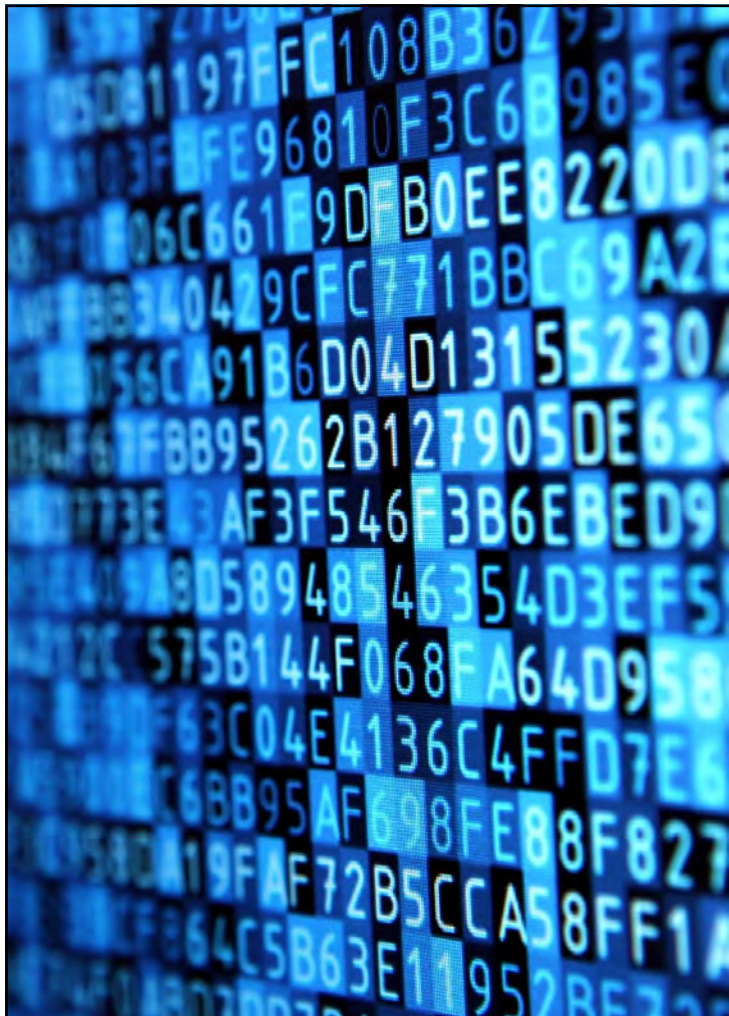
Data, gestructureerd vs. ongestructureerd

Gestructureerd

Noorwegen	56221	Rijst
Zweden	45027	Gerst
Oostenrijk	23992	Aardappel
Nederland	36892	Rijst
Italië	34568	Gerst
IJsland	23245	Rijst
Frankrijk	91102	-
Denemarken	16021	-
Duitsland	45289	Aardappel
België	49921	Rijst

Ongestructureerd

Zweden	12	Nee	28456
Aardappel	Nee	345289	
123	Misschien	Wenen	Gerst
6738	28972	18	Noordzee
Andalusië	1982		Zwart



Data naar omvang

1980: Kb

1990: Mb

2000: Gb

2010: Tb

2020 => Pb, Ex, Zb, Yb...

Data kunnen zeer verschillend in omvang zijn.

Voor de technische mogelijkheden is het van belang te weten hoe groot de databestanden zijn, maar inzicht in fileformat en opbouw van het bestand is voor de analyse belangrijker.



Databronnen

Open Data zoals CBS, data.overheid.nl

Commerciële data, company.info,
bureau van Dijk

Cliëntendata

Zelf verkregen data; sensordata,
internet



Logistieke data

HS-codes

Pos. No.	Naam	Bestelling	Schetsnr. / Ref.	Karakteristiek	Material
1	Stuk. Gehuize				433026
2	Stuk. Hoelwiel				140PMAL
3	Stuk. Schermerrol				15 - EUSIEN
4	Stuk. Schermerwiel				16PMAL15
5	Stuk. Spindel				16PMAL15
6	Stuk. Lagerdeksel groot				1225LH
7	Stuk. Lagerdeksel klein				1225LH
8	Stuk. Oorspronk				1225LH
9	Stuk. Wapelferster	BN 425 - 4209			1225LH
10	Stuk. Wapelferster	BN 720 - 9203			
11	Stuk. Propbeiler groot	BN 4405 - 8 12 x 8 x 22			
12	Stuk. Propbeiler klein	BN 4405 - 8 12 x 8 x 18			
13	Stuk. Versnellingsas	BN 908 - 90 x 15 - 12			
14	Stuk. Dichting	BN 7403 - A 8 x 18 x 14			
15	Stuk. Zwenkdeksel met versnellingsas	BN 4262 - 96 x 20 - 18			
16	Stuk. Zwenkdeksel met versnellingsas	BN 4262 - 96 x 18 - 18			
17	Stuk. Schermer	BN 9027 - 8 6.6			
18	Stuk. Wapelferster	BN 7760 - 41 x 13 x 13			
19	Stuk. Spindel	BN 1710-02015-4-400 10			
20	Stuk. Spindel	BN 1710-02015-4-400 10			
21	Stuk. Oorspronk	Kaafrol, geneesl. Zeebrugge			1225LH

Schneckengetriebe







Beroemd 'Big data' voorbeeld



SPORT



Politiek

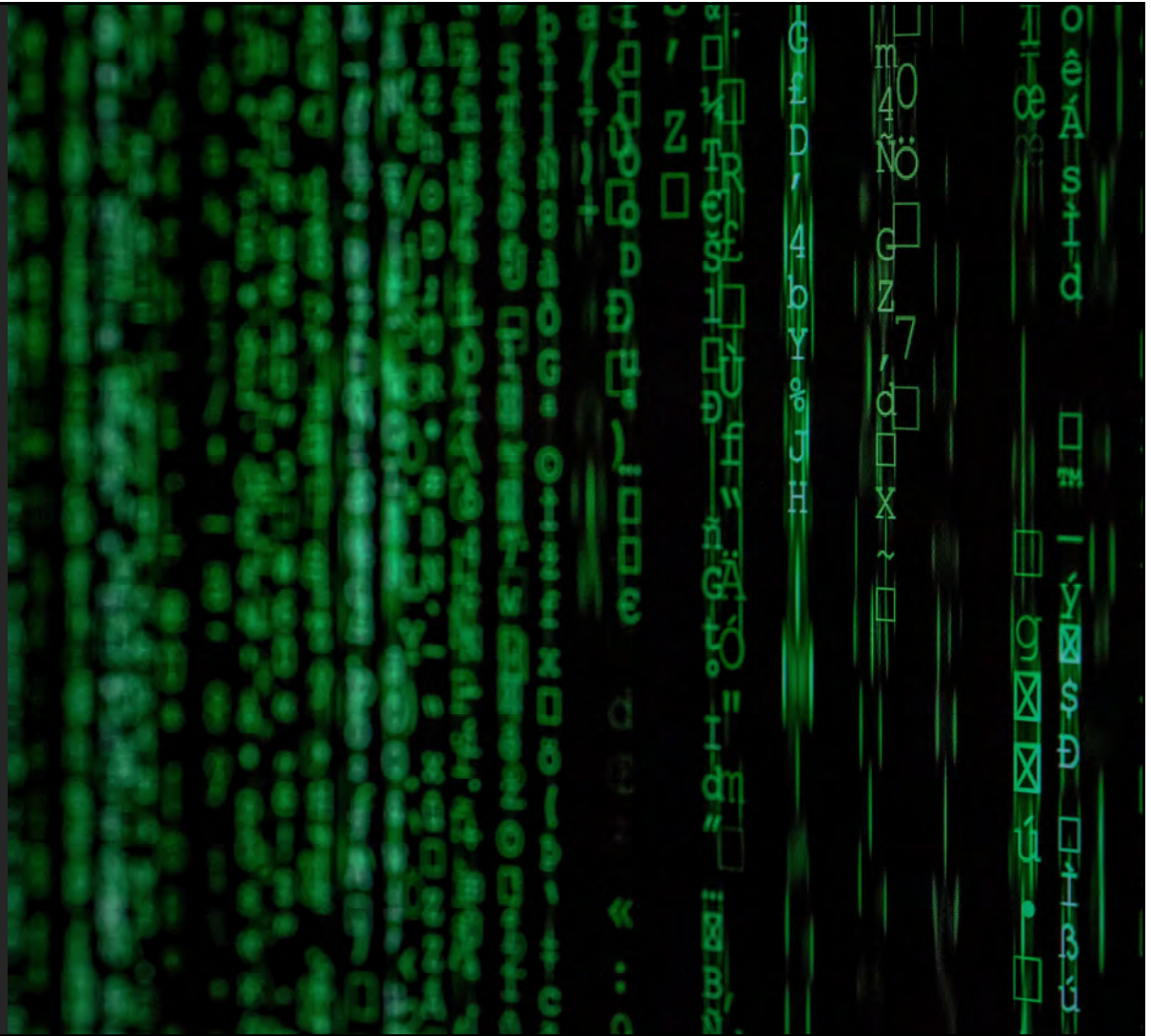




En
natuurlijk...

Datadenken

Om inzicht te krijgen in de mogelijkheden, is het goed om te leren denken in data; 'datadenken'.



Oefening 'datadenken'.



Bedenk welke data er gebruikt kunnen worden bij het houden van het toezicht op de belastingheffing.

Noem zo precies mogelijk welke data het zijn, bij wie ze opgevraagd kunnen worden en hoe de data geanalyseerd (beoordeeld) worden.




1. Auto van de zaak

Data?

Opvragen bij?

Hoe werkt de beoordeling?



NBA: 'Data-analyse is het ontdekken van patronen, afwijkingen, inconsistenties, en het onttrekken van andere nuttige informatie over het object van het onderzoek door middel van analyse, modellering en visualisatie met het oog op de planning of het uitvoeren van de opdracht.'

Data-analyse

Data-analyse is dus het proces van verkrijgen van informatie uit data door gebruik van methoden en technieken.

Voorafgaand aan de analyse zelf moeten de data verkregen worden. Dit proces wordt het ETL-proces genoemd:

Extraction, Transformation, Loading.

ETL-proces

Extraction => het verkrijgen van de data uit (diverse) systemen, bijvoorbeeld een ERP-systeem of een payrollpakket.

Transformation => data moeten zodanig aangepast worden dat het softwarepakket waarmee geanalyseerd gaat worden de data kan interpreteren.

Loading => de te analyseren gegevens worden ingeladen in een database of in de analysesoftware.

Data-analyse: 4 typen

- Descriptive analytics

- Diagnostic analytics

- Predictive analytics

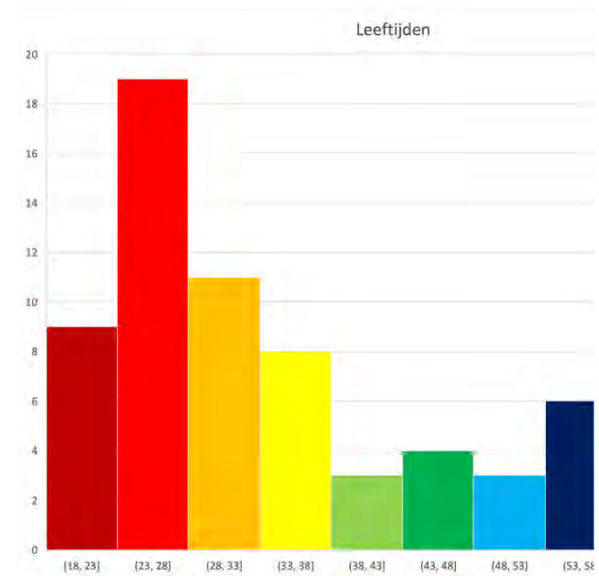
- Prescriptive analytics



Data-analyse in soorten

Descriptive analytics

‘Beschrijvende statistiek’, met kengetallen wordt een beschrijving gegeven van de data. Denk aan zaken als het gemiddelde, de mediaan, de modus en de standaarddeviatie. Aanvulling met grafische weergaven, bijvoorbeeld over de verdeling door middel van een histogram.



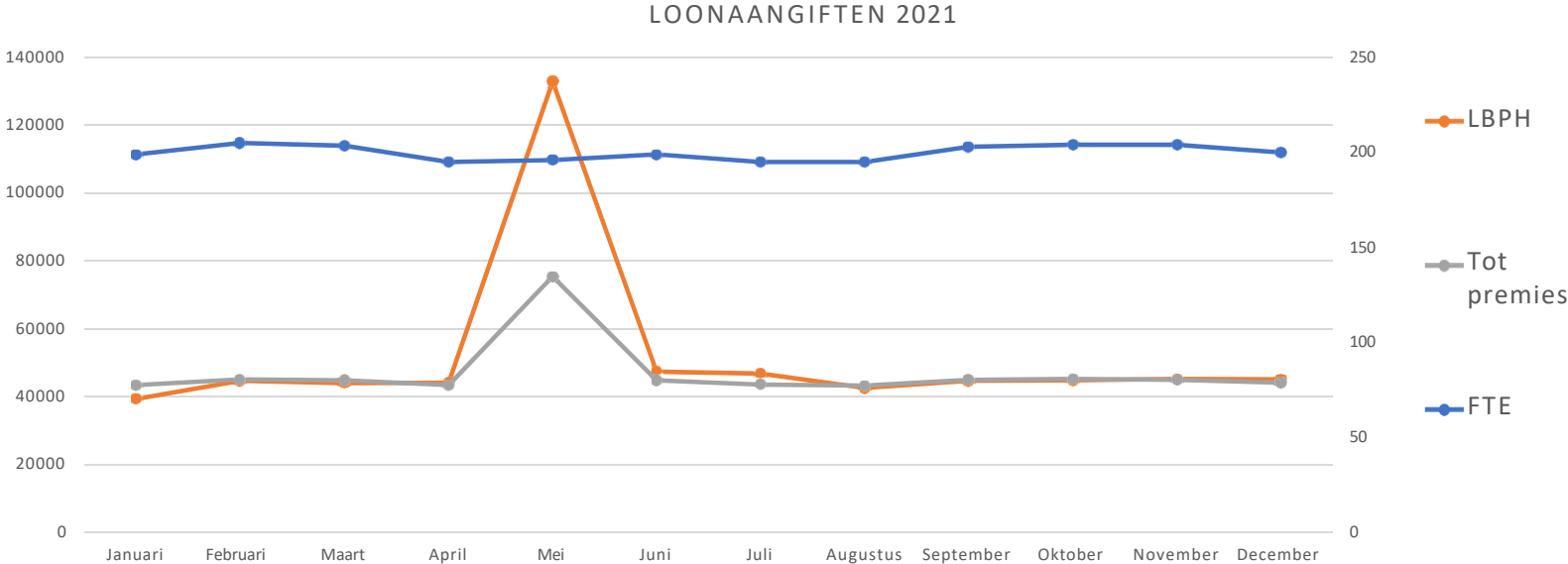
Descriptive analytics in Tax

Maand	FTE	LBPH	Tot premies	ZVW	WIA/WA	WW-Awf	Gem loon
Januari	199	39303	43311	2830	29477	11005	1975
Februari	205	44649	45022	3215	30443	11365	1980
Maart	203,5	44009	44777	3169	30006	11602	1966
April	195	44015	43324	3169	28958	11197	1980
Mei	196	132947	75412	9572	47481	18359	3230
Juni	199	47338	44780	3408	29835	11536	1999
Juli	195	46800	43540	3370	29250	10920	2000
Augustus	195	42510	43231	3061	29250	10920	2000
September	203	44615	44988	3212	30420	11357	1998
Oktober	204	44790	45165	3225	30539	11401	1996
November	204	45197	44991	3254	30539	11198	1996
December	200	45064	44122	3245	29910	10967	1994

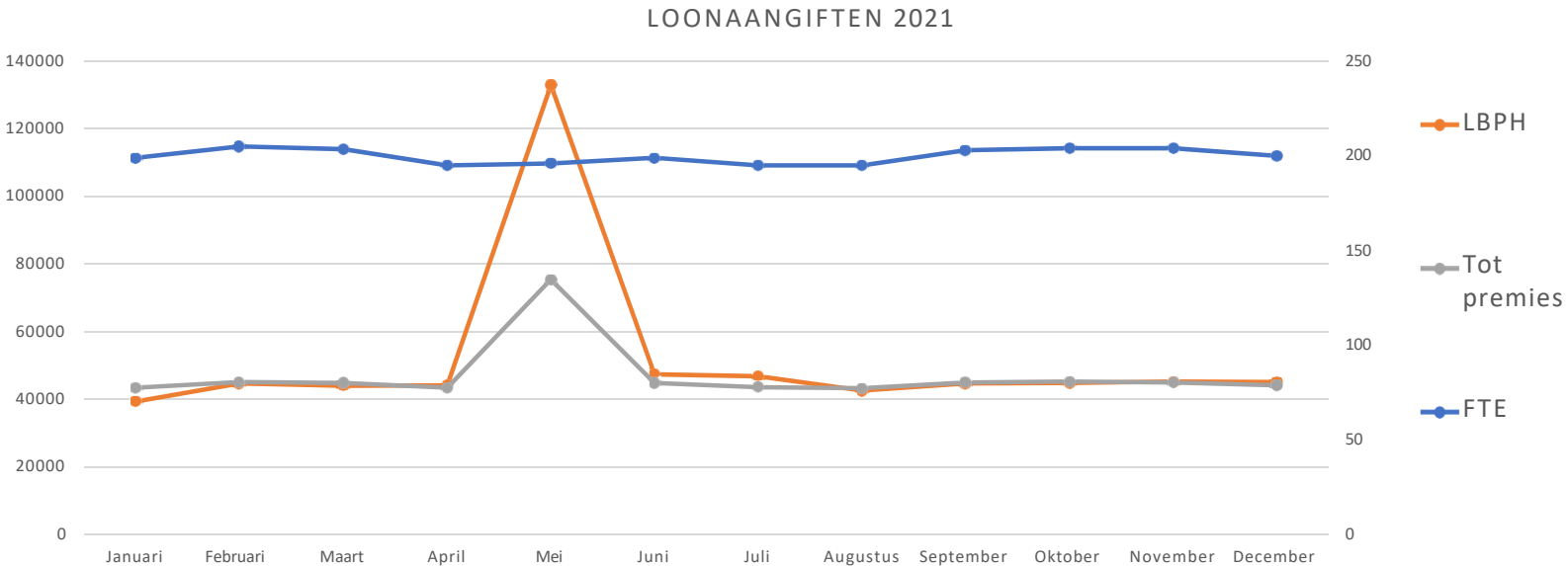
De beschrijving van de data kan worden verrijkt met kengetallen die meer inzicht geven.

Denk bijv. aan toegepaste gemiddelde tarieven, totalen etc.

Descriptive analytics in Tax



Descriptive analytics in Tax



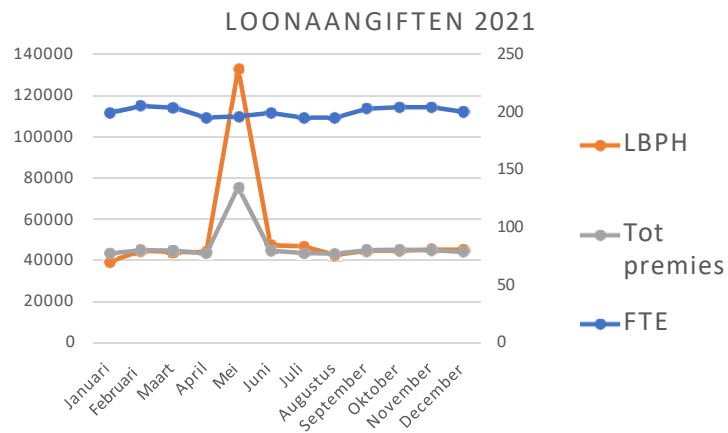
Diagnostic analytics

Aan de hand van data proberen te achterhalen wat de oorzaak of oorzaken van een bepaald fenomeen is of zijn. Diagnostic analytics zoekt naar het 'waarom' van een bepaalde gebeurtenis. Het beoordelen van de uitstoot van bepaalde gassen en de samenhang met klimaatverandering is een voorbeeld van diagnostic analytics.

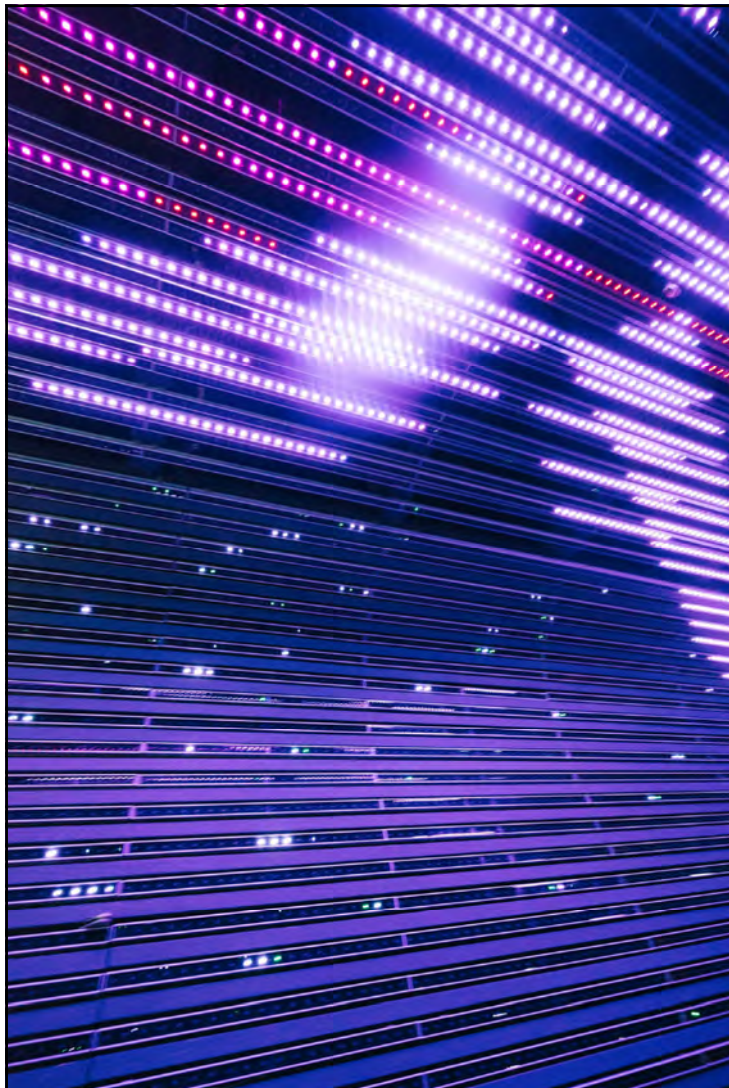


Diagnostic analytics in Tax

Wat is de oorzaak van de piek in LBPH en premies in mei?



Maand	FTE	LBPH	Tot premies	ZVW	WIA/WAO	WW-Awf	Gem loon
Januari	199	39303	43311	2830	29477	11005	1975
Februari	205	44649	45022	3215	30443	11365	1980
Maart	203,5	44009	44777	3169	30006	11602	1966
April	195	44015	43324	3169	28958	11197	1980
Mei	196	132947	75412	9572	47481	18359	3230
Juni	199	47338	44780	3408	29835	11536	1999
Juli	195	46800	43540	3370	29250	10920	2000
Augustus	195	42510	43231	3061	29250	10920	2000
September	203	44615	44988	3212	30420	11357	1998
Oktober	204	44790	45165	3225	30539	11401	1996
November	204	45197	44991	3254	30539	11198	1996
December	200	45064	44122	3245	29910	10967	1994



Predictive analytics

Het voorspellen op basis van data.

Vaak is voorafgaand aan het doen van de voorspelling het ontdekken van de wetmatigheid in de data noodzakelijk.

Als de wetmatigheid gevonden is, dan kan die vastgelegd worden in een formule (algoritme.)

Predictive analytics in Tax

Wat zal de volgende maanden de verschuldigde belasting worden? Een analyse van de patronen o.b.v. historische patronen en daarop vastgestelde wetmatigheden kan leiden tot nauwkeurige voorspellingen.

Maand	FTE	LBPH	Tot premies	ZVW	WIA/WAO	WW-Awf	Gem loon
Januari	199	39303	43311	2830	29477	11005	1975
Februari	205	44649	45022	3215	30443	11365	1980
Maart	203,5	44009	44777	3169	30006	11602	1966
April	195	44015	43324	3169	28958	11197	1980
Mei	196	132947	75412	9572	47481	18359	3230
Juni	199	47338	44780	3408	29835	11536	1999
Juli	195	46800	43540	3370	29250	10920	2000
Augustus	195	42510	43231	3061	29250	10920	2000
September	?	?	?	?	?	?	?
Oktober	?	?	?	?	?	?	?
November	?	?	?	?	?	?	?
December	?	?	?	?	?	?	?

Prescriptive analytics

Deze analyse borduurt voort op de predictive analytics en analyseert het 'recept' voor het bereiken van een bepaald resultaat.

Door gebruik te maken van meer data, dan in de voorspellende fase is gebeurd, kan met prescriptive analytics het resultaat verbeterd worden.





Descriptive analytics: beschrijving van de data



Diagnostic analytics: beantwoording van het waarom



Predictive analytics: levert een voorspelling



Prescriptive analytics: hoe kan het toekomstig resultaat bereikt (en verbeterd) worden?

Samenvatting data-analyse

Data-analyse, aandachtspunten

Is het doel voldoende duidelijk?

Juridische aspecten (AVG, opdracht)

Technische aspecten (ETL-traject)

Is er voldoende inzicht in de kwaliteit van de data?

Is duidelijk wat de databron is?





Data-analyse, nog meer aandachtspunten

Logging bewerkingen

Opslag 'deelproducten'

'Verborgen' persoonsgegevens

Bewaartermijnen

Overig?

Methoden & Techniken



"Lies, damned lies, and statistics"

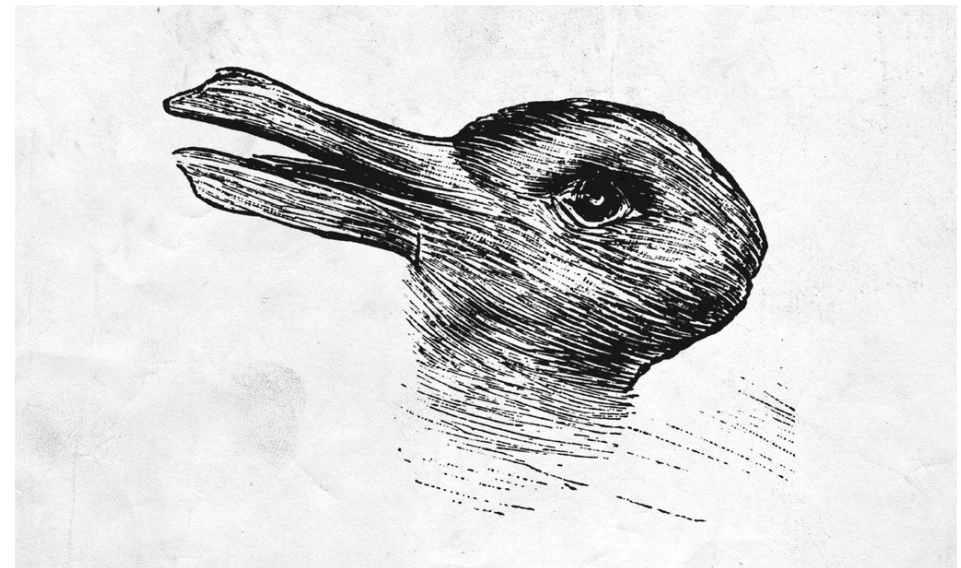
Mark Twain (Chapters from my Autobiography, 1907).

Ook voor het tijdperk van 'fake news' was de kracht van misleidende cijfers al bekend.

Grafische weergaven kunnen zeer misleidend zijn.

Hoe misinterpretatie voorkomen?

Methoden en technieken

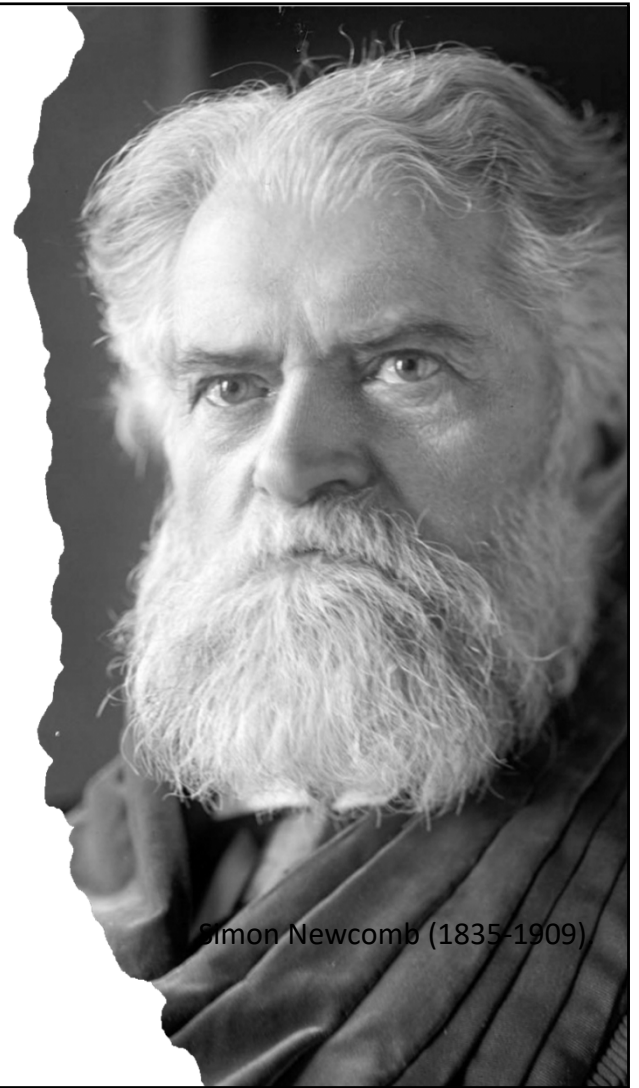
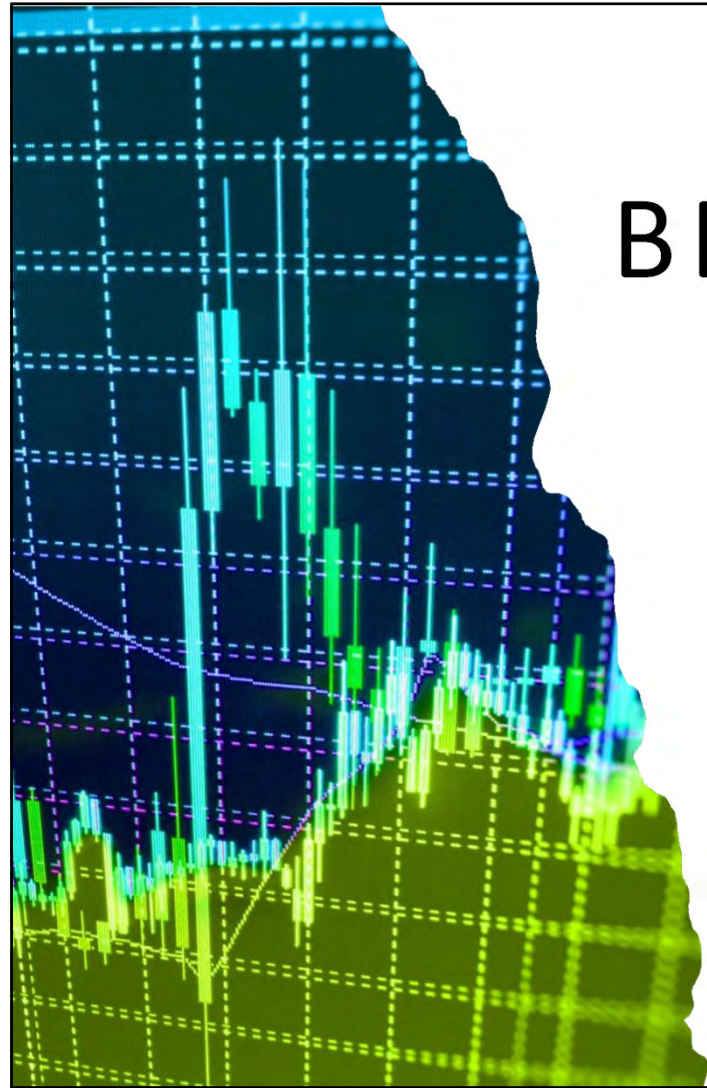




CHI² ?

	Cases	Deaths	New Deaths	Total Recovered	Active Cases	Serious, Critical
China	80,881	+21				
Italy	31,506	+3,526	3,226	68,715	8,940	3,226
Iran	16,169	+1,178	2,503	2,941	26,062	2,060
Spain	11,409	+1,467	988	5,389	9,792	
Germany	8,639	+1,367	510	1,028	9,871	563
S. Korea	8,320	+84	23	67	8,549	2
France	6,633		81	1,401	6,838	59
			148			
USA	5,704	+1,041	97	12	6,473	400
Switzerland	2,742	+389	27	74	5,533	12
UK	1,950	+407	71	15	2,700	
Netherlands	1,705	+292	43	52	1,827	20
Norway	1,452	+104	3	2	1,660	45
				1	1,448	27

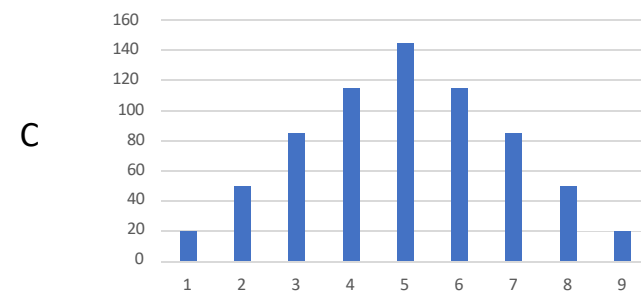
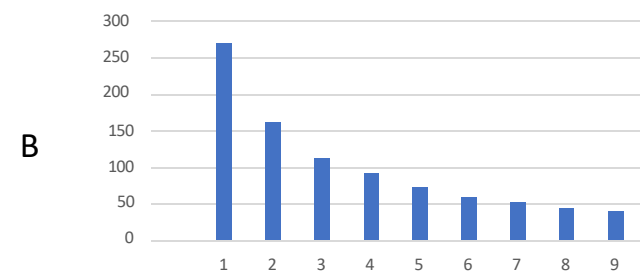
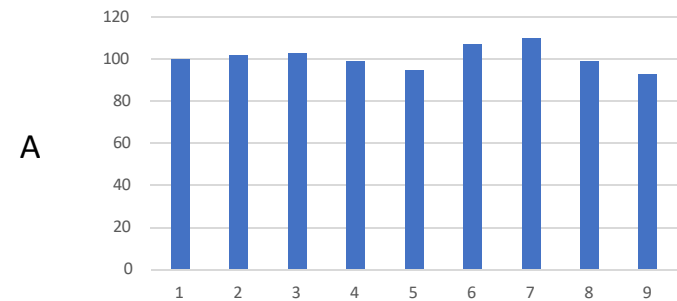
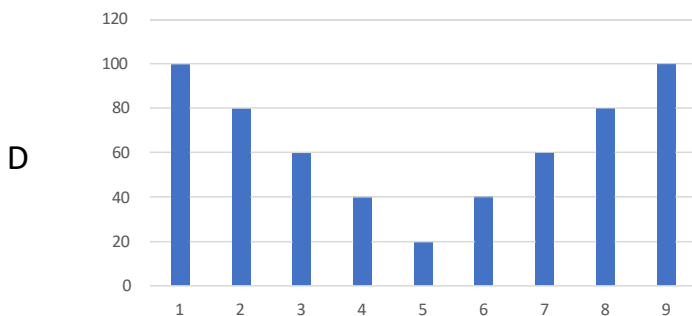
BENFORD'S LAW?



Simon Newcomb (1835-1909)

Kwisje!

Als je alle transacties van een onderneming sorteert naar omvang in geldbedrag, hoe ziet de verdeling van de *begincijfers* er dan waarschijnlijk uit?





CORRELATIE

Correlatie is de berekening van de
samenhang tussen (numerieke) gegevens.

Niet meer, niet minder.



Correlatie..

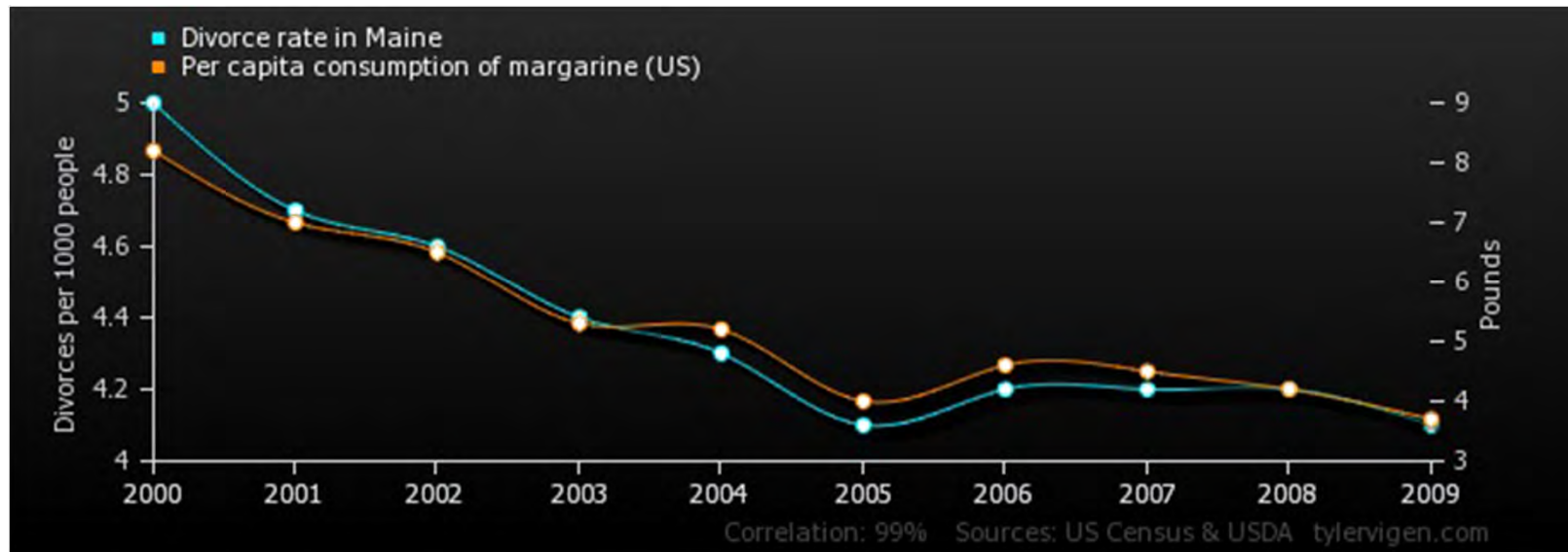


Causatie

Causatie is de wetenschap over oorzaak en gevolg.

1. Er is correlatie

2. Er is een toetsbare theorie die het verband verklaart



Correlatie vs causatie

Bron: www.tylervigen.com

Aan de slag met data

Om te voldoen aan de eisen van de fiscale wetten zijn juiste (en tijdige) aangiften van belang. Welke beoordelingen zouden er geautomatiseerd kunnen plaatsvinden? Welke informatie is daarover beschikbaar?

Voor de fiscalist is kennis van de mogelijkheden van belang.



Aan de slag met data, het landschap

Het 'landschap' is inzicht in de relevante software en de IT-omgeving waarin deze software functioneert. Denk aan het ERP-systeem, het payrollpakket maar bijvoorbeeld ook een geautomatiseerd toegangssysteem met pasjes.



Aan de slag met data.

1. Duidelijke scope; welke belasting(en), welk tijdvak, welke aangifteplichtige, welk doel?

2. Inzicht in het 'landschap', denk ook aan data voor andere toezichthouders!

3. Een overzicht van de 'queries' ofwel, de data-analyses die we willen uitvoeren

4. Welke data hebben we nodig? Aangiften en ?

5. 'Tooling'; software waarmede we kunnen analyseren en liefst ook visualiseren

Voorbeeld I

1. Scope: maandaangifte loonheffingen 2022. Doel: beoordeling op juistheid
2. Landschap: payrollpakket is AFAS, ERP is AFAS
3. Aangiftedata: XML-loonaangiften (kopie)
4. Queries, t.b.d.
5. Tooling: Excel

```
<?xml version="1.0" encoding="ISO-8859-1"?><Loonaangifte
xmlns="http://xml.belastingdienst.nl/schemas/Loonaangifte/2020/01"
version="2.0"><Bericht><IdBer>ADP012345R0004038071</IdBer><DatTdAa
nm>2020-04-23T13:52:00</DatTdAanm><ContPers>KC-ADP NL HR
DEMO </ContPers><TelNr>0201234567</TelNr><RelNr>swo00012</R
elNr><GebrSwPakket>Multipay
TWK</GebrSwPakket></Bericht><AdministratieveEenheid><LhNr>11223344
5L01</LhNr><NmIP>DEMO
BV </NmIP><TijdvakAangifte><DatAanvTv>2020-01-
01</DatAanvTv><DatEindTv>2020-01-
31</DatEindTv><VolledigeAangifte><CollectieveAangifte><TotLnLbPh>0000
596598</TotLnLbPh><TotLnSV>0000596598</TotLnSV><TotPrlnAwfAnwLg
>0000560577</TotPrlnAwfAnwLg><TotPrlnAwfAnwHg>0000018979</TotPrl
nAwfAnwHg><TotPrlnAwfAnwHz>0000000000</TotPrlnAwfAnwHz><PrlnUF
O>0000000000</PrlnUFO><IngLbPh>0000150106</IngLbPh><PrWAOAof>
0000042133</PrWAOAof><TotPrAwfLg>0000016480</TotPrAwfLg><TotPrA
wfHg>0000001507</TotPrAwfHg><TotPrAwfHz>0000000000</TotPrAwfHz>
<IngBijdrZvw>0000000000</IngBijdrZvw><TotWghZvw>0000038830</TotW
ghZvw><TotTeBet>0000249056</TotTeBet><TotGen>0000249056</TotGen
></CollectieveAangifte><InkomstenverhoudingInitieel><NumIV>0002</NumI
V><DatAanv>2011-08-
01</DatAanv><PersNr>0001</PersNr><NatuurlijkPersoon><SofiNr>0000000
01</SofiNr><Voort>A</Voort><SignNm>DEMO A</SignNm><Gebdat>1965-
06-
26</Gebdat><Nat>0001</Nat><Gesl>1</Gesl><AdresBinnenland><Str>DO
RPSSTRAAT</Str><HuisNr>1</HuisNr><Pc>6412WL</Pc><Woonpl>HEERL
EN</Woonpl></AdresBinnenland></NatuurlijkPersoon><Inkomstenperiode>
<DatAanv>2020-01-
01</DatAanv><SrtIV>15</SrtIV><CdAard>01</CdAard><CAO>0824</CAO>
<IndArbovOnbepTd>J</IndArbovOnbepTd><IndSchriftArbov>J</IndSchriftAr
bov><IndOprov>N</IndOprov><IndAvrLkvOudrWn>N</IndAvrLkvOudrWn><I
ndAvrLkvAgWn>N</IndAvrLkvAgWn><IndAvrLkvDgBafSb>N</IndAvrLkvDgB
afSb><IndAvrLkvHpAgWn>N</IndAvrLkvHpAgWn><IndLhKort>J</IndLhKort
><LbTab>012</LbTab></IndWAO>J</IndWAO>
```

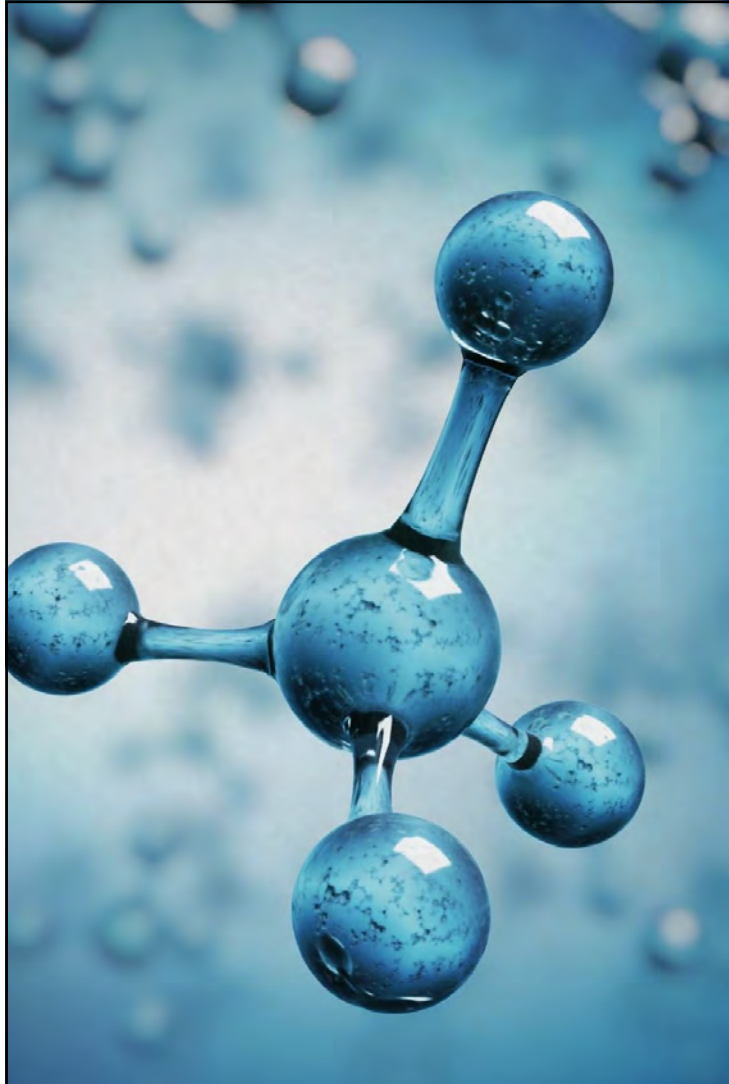
Data en queries

Het formuleren van de queries, ofwel het beschrijven van hetgeen onderzocht moet worden is natuurlijk essentieel. Dit is de uitwerking van het doel van het onderzoek: is het belang een algemeen beeld krijgen over de juistheid van de aangiften, of wordt er gezocht naar specifieke fouten?

In het voorbeeld is de scope nogal breed: beoordeling op juistheid. Hoe vertalen we dat in queries?

De mogelijke queries hangen af van de beschikbaarheid van de data, dus data en queries is enigszins een kip en ei vraagstuk..

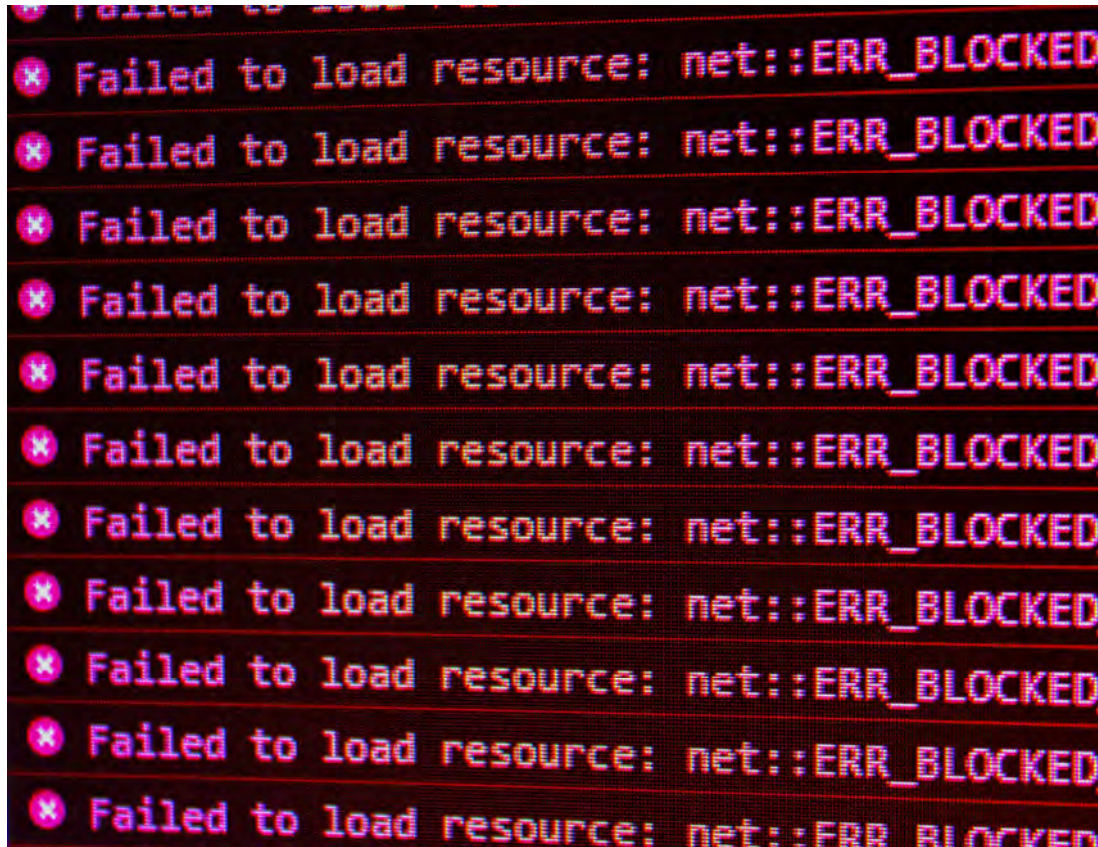




Queries, voorbeeld II.

Queries moeten zo scherp mogelijk geformuleerd worden. Stel u doet een beoordeling van de verkopen op juistheid van de in rekening gebrachte BTW. Dan heeft u deze gegevens nodig:

- Omschrijving geleverde dienst of geleverde goederen
- Datum factuur
- Datum prestatie
- Vestigingsland afnemer

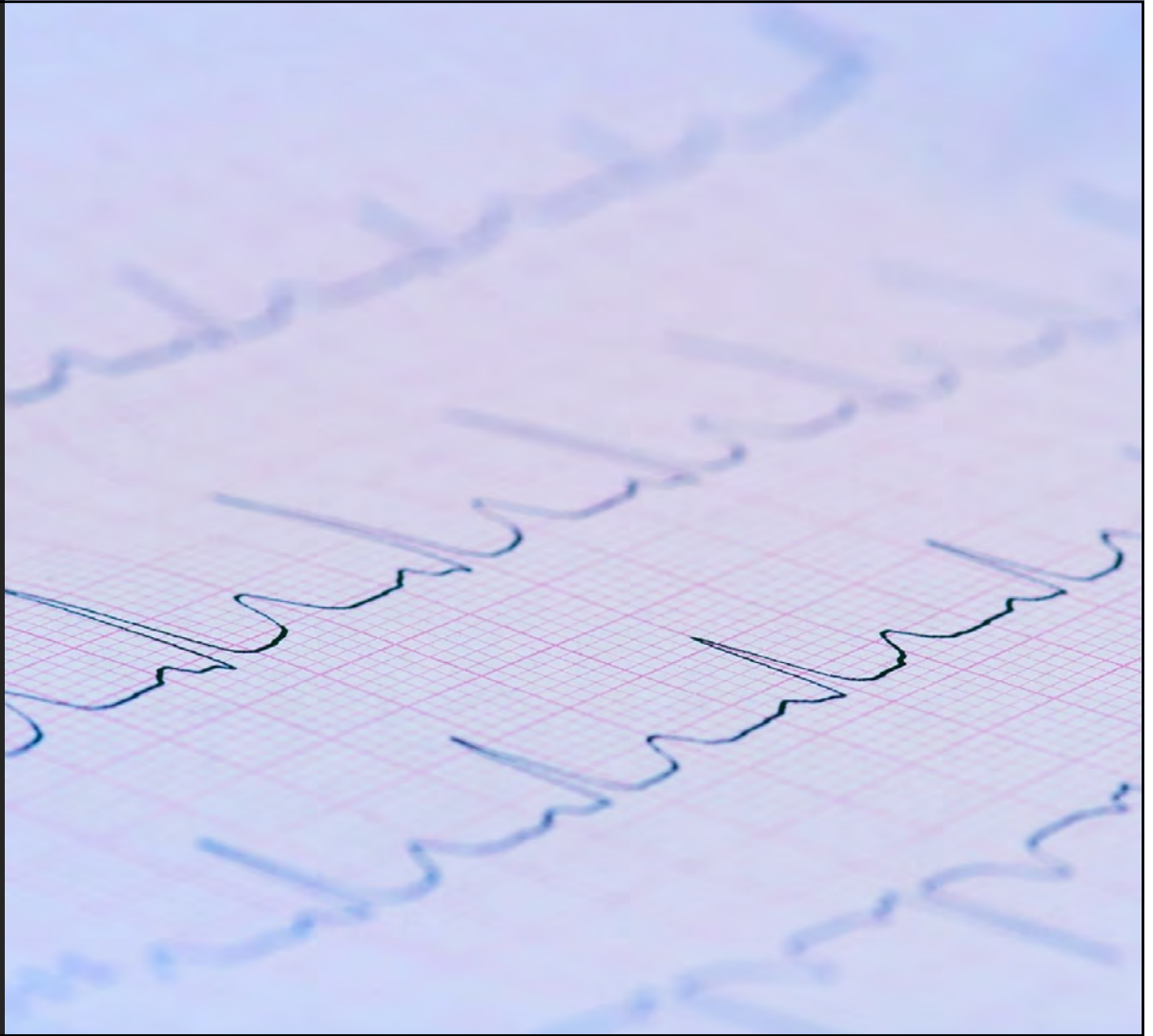


Queries, voorbeeld II.

Maar natuurlijk ook:

- het tijdvak waar de beoordeling op ziet
- hoe er gerapporteerd moet worden:
 - bij afwijking van 21% met omschrijving dienst/levering?
 - of bij 9% met omschrijving dienst/levering?
 - bij 0% maar geen afnemer in buitenland?

Trends in data-analyse



Continuous monitoring

Van incidenteel naar frequent naar doorlopend.

Veel organisaties hebben voor hun businessproces een grote informatiebehoefte en richten dan ook dit zelf al in (dashboards).

Voor (financiële) rapportagedoel-einden kan gebruik gemaakt worden van de bestaande informatiestromen.



Van geaggreerd naar transacties

Geaggreerde data zeggen minder dan de bouwstenen van deze data: de transacties. Veel meer informatie kan ontleend worden aan de transacties.

De beoordeling van de transactie-data is wel complexer; de samenhang tussen de data moet in kaart gebracht worden.



Van verzenden naar toegang systeem?

Verzenden van data is niet zonder risico. Een alternatief is het beoordelen van de originele data, dat vereist toegang tot deze data (systemen).

In sommige landen (Hongarije) is er al wetgeving die hiertoe verplicht.





Conclusie

Data is alomtegenwoordig en niet meer weg te denken uit de moderne tijd

De steeds toenemende hoeveelheid data biedt in combinatie met steeds snellere computers en daarmee samenhangende nieuwe technieken een enorm potentieel in alle domeinen.

Om bij te blijven is nodig:
inzicht in de benodigde data
inzicht in methoden en technieken (data-analyse)
inzicht in het juridisch/ethisch kader

Einde deel 1
