

# Kennis maken met big data

In de wereldwijd explosief groeiende berg van digitale gegevens zitten nieuwe inzichten en kennis verstopt. Om die te vinden zijn slimme software en rekenkracht nodig, en heel veel data-experts. 'Over tien jaar is 80 procent van het onderzoek gebaseerd op het analyseren van databestanden.'

TEKST RIK NULAND ILLUSTRATIES KAY COENEN



**B**egin dit jaar werd Alphabet het meest waardevolle bedrijf ter wereld. Voor het eerst voerde niet een maak- of dielbedrijf de ranglijst aan, maar een nieuwkomer die handelt in informatie. Alphabet is het moederbedrijf van Google. De opmars van het bedrijf – tien jaar geleden kwam het bedrijf de top-100 binnen – is een duidelijk signaal hoeveel waarde er wordt toegekend aan informatie, informatie technologie en databasesystemen.

Op dat terrein groeit de hoeveelheid gegevens die digitaal worden opgeslagen en verwerkt met duizelingwekkende snelheid. Bijvoorbeeld op het gebied van DNA (zie kader), maar ook dichter bij ons alledaags leven.

Moderne auto's geven dag in dag uit informatie door aan de fabrikant over onder meer toerental en de lengte vanritten. Interessante informatie om het onderhoudspakket op te stemmen, maar ook handig voor studies naar bijvoorbeeld verschillen in rijgedrag. Leuk materiaal dus voor onderzoekers maar ook voor verzekeraars.

informatiestroom ook binnen het typisch Wageningse domein stiek zaaien. Dankzij moderne ICT worden nu sensoren ingezet die in de kas of op de trekker de gewasgroei in de gaten houden. Alle melkrobots samen weven bijna alles over honderdduizenden koelen. In die brit aan gegevens zitten nieuwe inzichten en kennis verstopt. Zeker als je databestanden kunt koppelen, bijvoorbeeld melkgift of vooropname met genetische informatie.

Het is noodzakelijk dat Wageningen die nieuwe kennis vorder aanhoort, vindt Karin Andeweg. Samen met Sander Janssen is zij aanjager en kwartiermaker voor big data, door de organisatie benoemd tot speerpunt. ‘Big data lijkt een hype, maar over een paar jaar is het gemanaged goed geworden. De verwachting is dat over tien jaar 80 procent van het onderzoek gebaseerd is op het analyseren en combineren van databestanden om zo nieuwe kennis te genereren.’

#### NIET MEER HET VELD IN

Die toekomst tekent zich nu al af. Wageningen UR doet daarmee een proef in Amsterdam om op drukke dagen aan de hand van het gebruik van honderdduizenden mobiele telefoons te bepalen waar de mensenmassa gevaarlijk aangroeit. Ook wordt zo gevoldigd of de mensenmassa niet alleen effect sorteert. Green warntemer hoeft de straat op.

Ook om de gewasgroei, zeg van tarwe te achterhalen in combinatie met wolkken. Niet wereldschockend, maar voor de groentesector een eyeopener over consumentengezag. Die sterk uitlijnende digitale gegevensberg om ons heen wordt aangeduid als big data. Big staat voor groot, maar de naamgeving is ook een verwijzing naar Big Brother, de alweten overheid die volgens schrijver George Orwell al lons doen en laten zou gaan bepalen. Wageningen UR zet immidels stappen om grote databestanden beter te exploreren. Verwacht wordt dat de

#### BROCCOLI WOKKEN

In de databestanden van Google, grote winkelketens, sms, twitter of autofabrikanten ligt een schat aan kennis opgeslagen over onze voorkeuren. Voor snel opbreken bijvoorbeeld waar we op internet naar zoeken; of wat we gezag staan. Bij een analyse van Nederlandse twetterberichten bleek het woord broccoli vaak voor te komen in combinatie met wolkken. Niet wereldschockend, maar voor de groentesector een eyeopener over consumentengezag.

Die sterk uitlijnende digitale gegevensberg om ons heen wordt aangeduid als big data. Big staat voor groot, maar de naamgeving is ook een verwijzing naar Big Brother, de alweten overheid die volgens schrijver George Orwell al lons doen en laten zou gaan bepalen. Wageningen UR zet immidels stappen om grote databestanden beter te exploreren. Verwacht wordt dat de

gegevens buitenkijken, niet van een drone of de satelliet in het veld in om metingen te doen. Nu al kan hij via sensoren in het veld maar ook van een drone of de satelliet in het veld. ‘Want dat is de belangrijkste verschillen tussen de verschillende disciplines’, zegt Dick de Ridder, hoogleraar bio-informatica aan Wageningen University. Miljarden gegevens over genomen, genen, eiwitten en andere moleculen worden in grote bestanden bij elkaar gebracht en systematisch onderzocht. ‘De verwachting is dat we dit jaar een miljoen miljard DNA-basen kunnen aflezen’, aldus De Ridder. ‘Het is de kunst van de bio-informaticus om op basis van de terabytes aan data nieuwe biologische hypotheses op te stellen. Waar nu biologen data-analyse vaak uittestden aan onderzoekers daarvoor jalang delaspelen onderzoeken. Uiteindelijk zal ook de boei daarvan profiteren als al die kennis kan worden vertaald naar handzame adviezen.’

#### BAKENS VERZETTEN

Om databestanden ook daadwerkelijk zo te kunnen gebruiken, dient de wetenschap de bakens te verzettten, denkt de EU. De Europese Commissie kondigde half april een miljardeninvestering aan in datamanagement.

## ‘Als je mee wilt doen, moet je investeren’



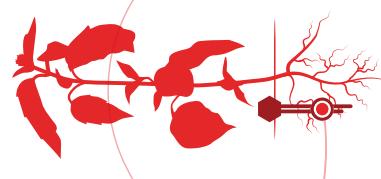
Foto: ANP

#### GROTE GETALLEN

Astronomen en fysici grossieren al heel lang in enorme hoeveelheden data. Tijdens de experimenten in de deeltjesversneller van CERN worden de resultaten van 600 miljoen botsingen per seconde geregistreerd. Deze seculaire biologie met een inhaarscène begonnen. In 2003 werd voor het eerst het DNA van een mens ontrafeld; in 2011 volgde de duizendste mens; volgend jaar komt waarschijnlijk 1 miljoen in de rek. Elk jaar verdrievoudigt de capaciteit om via sequencing DNA te ontleden. Een vergelijkbare daaexplozie vindt plaats rond eiwitten en stofwisselingsproducten in het lichaam. De biologie krijgt daardoor net karakter van een datawetenschap, zeft Dick de Ridder, hoogleraar bio-informatica aan Wageningen University. Miljarden gegevens over genomen, genen, eiwitten en andere moleculen worden in grote bestanden bij elkaar gebracht en systematisch onderzocht. ‘De verwachting is dat we dit jaar een miljoen miljard DNA-basen kunnen aflezen’, aldus De Ridder. ‘Het is de kunst van de bio-informaticus om op basis van de terabytes aan data nieuwe biologische hypotheses op te stellen. Waar nu biologen data-analyse vaak uittestden aan onderzoekers daarvoor jalang delaspelen onderzoeken. Uiteindelijk zal ook de boei daarvan profiteren als al die kennis kan worden vertaald naar handzame adviezen.’

#### PEUTER MET POTENTIE

Watson is een peuter die alles eerst moet leren, maar wel met een geweldige potentie’, stelt Richard Visser,



## De komende tien jaar zijn een half miljoen data-experts nodig

hogeplaatsen vanenveredeling bij Wageningen UR. Vorig jaar klopte hij aan bij IBM Research om de supercomputer in te zetten voor de veredeling van aardappel. Visser verwacht dat Watson Potato, zoals het systeem immiddels intiemeert, een belangrijk hulpmiddel wordt om efficiënter en niet een beter resultaat te verkrijgen. Maar eerst moet Watson leren zich te foussen. ‘Bij het woord ‘knol’ moet hij niet uitkomen bij ‘dahlia’, aldus Visser.

Na dat leerproces zal Watson Potato een waardevol hulpmiddel zijn om literatuur te scannen, bijvoorbeeld over de positie van een bepaald gen, verwacht Schaap. ‘Watson kan veel sneller zoeken dan wij en in veel meer literatuur. Wij kunnen er ook wel achter dat dat gen bijvoorbeeld bovenin chromosoom 3 ligt, maar dan heb je het over een gebied met wel duizend genen. Dat is onhandelbaar veel’, aldus Visser. ‘Watson kan verder graven; die duizend genen vergelijken met DNA bij andere planten, maar ook bij paddenstoelen, mossels of vogels, omdat wat bekend is over wat diegenen daar doen.’

Dan beperkt je de kans waarschijnlijk tot vijf of tien candidaatgenen, en naarmate de computer meer kennis vergaart, wordt het misschien wel neteens een schot in de roos.’

De computer leert van zijn fouten. Als hij van experts te horen krijgt dat de antwoorden in de goede richting wijzen, dan borduurt hij daarop voort. Bij negatieve feedback laat hij het gevuld spoor rusten. Hij leert dus, echter niet als een mens een bouw-expertise op’, aldus Visser. ‘Door kennis over uiterlijke kenmerken, DNA-samenstelling, ophengen en groeiomstandigheden te koppelen kunnen we veel beter ons uitgangsmateriaal kiezen.’

In op termijn zit er misschien meer in het vat, blijvend door Watson zijn tanden in complexe problemen te laten zitten. Welke genen zijn er verantwoordelijk voor dat de ene aardappel hardkokend is en de andere afkoekt? Daar hebben we nog eigenlijk geen idee van.’

#### AANSPRAAK MAKEN

Ondanks de mooie vooruitzichten met big data, zijn er ook beren op de weg. Wie mag aanspraak maken op de gegevens over een aardewaker? De eigenaar van de drone, degene die de gegevens kan interpreteren of de boer? Zijn de data die de melkrobotdag in dag uit verzamelt van de fabrikant of van de veehouder? Dat is nog hagenaeg ongongenom gebied. Van wie zijn de resultaten als een bedrijf en een universiteit hun datasets combineren?

Watson krijgt de antwoorden in de goede richting feedback door Watson zijn tanden in complexe problemen te laten zitten. Welke genen zijn er verantwoordelijk voor dat de ene aardappel hardkokend is en de andere afkoekt? Daar hebben we nog eigenlijk geen idee van.’

#### AANSPRAAK MAKEN

Ondanks de mooie vooruitzichten met big data, zijn er ook beren op de weg. Wie mag aanspraak maken op de gegevens over een aardewaker? De eigenaar van de drone, degene die de gegevens kan interpreteren of de boer? Zijn de data die de melkrobotdag in dag uit verzamelt van de fabrikant of van de veehouder? Dat is nog hagenaeg ongongenom gebied. Van wie zijn de resultaten als een bedrijf en een universiteit hun datasets combineren?

Volgens Ben Schaap is innovatie erme gedient als databestanden toegankelijk zijn voor iedereen. Schaap is door Wageningen UR gedetecteerd bij Global Open Data for Agriculture and Nutrition (GODAN), een internationale lobby-organisatie voor open data in de landbouw- en wedstrijdsector. Sponsors zijn onder meer de VS, Groot-Brittannië, Nederland en de FAO, en onder de 250 partners zijn bedrijven als IBM en Syngenta maar ook lokale Afrikaanse ngo's.

‘Openheid zorgt voor een gelijkwaardig playing field voor iedereen’, zegt Schaap. ‘Als gegevens niet toegankelijk zijn, hebben partijen niet meer macht en geld de grootste controle. Een multinational kan informatie kopen, een eenmansbedrijf kan daaraan niet tegenop. Als de informatie van aannemers of satellieten openbaar is, kan iedereen er mee aan de slag. Niet alleen een multinational, maar ook een slimme whizkid. Open data zorgt ervoor dat iedereen applicaties kan ontwikkelen en de boer niet afhankelijk is van één partij die bijvoorbeeld ook zaaizaad, meststoffen of gewasbeschermingsmiddelen levert.’

Daarnaast is het erg belangrijk dat publiek gefinancierde gegevens openbaar zijn zonder voorwaarden vooraf, vindt Schaap. ‘Daar horen ook de datasets bij die universiteiten en instituten verzamelen. De Nederlandse overheid en NWO stellen immiddels open science als voorwaarde voor subsidering, maar ook onderzoekers die geld willen van de Gates Foundation of van het Europees onderzoeksprogramma Horizon 2020 zijn verplicht hun data te publiceren.’

#### COMMERCIEEL BELANG

Dat bedrijven waarschijnlijk niet happy zijn, om eigen gegevens door anderen te laten gebruiken, begrijpt Schaap. Maar hij verwacht wel dat die bereid zijn uitzonderingen te maken als ze zelf geen competitiefelbelang hebben. ‘Syngenta heeft een dataset vrijgegeven over een pesticide tegen muggen, geen speerpunt voor het bedrijf en malariaonderzoekers waren er erg blij mee. Je moet dat zien als een bijdrage om de wereld te verbeteren, net zoets als aan CO<sub>2</sub>-reductie doen.’

Toch ziet hij openheid als moer dan liefdadigheid. ‘Er zijn ook bedrijven die alleen een open-data-landschap te creëren’, aldus Schaap. ‘Zij vinden dat data beter uitwisselbaar moeten zijn zodat meer partijen met elkaar gegevens uit de voeten kunnen voor de ontwikkeling van nuttige toepassingen, bijvoorbeeld in de precisielandbouw. Innewonen kun je regenwoordig niet meer alleen, is de gedachte. Open science zorgt voor een samenwerking met slimme uitvinders. Syngenta zegt tegen startups: maak maar gebruik van onze onderzoeksgegevens. Kom daar een interessante toepassing uit, dan willen we jullie misschien wel overnemen.’ Plantenverdeelaar Visser is ook voorstander van open access, maar onder voorwaarden. ‘Je wilt niet dat de eerste die bestewhizkid in Rusland ermee aan de haal gaat. Dat is op dit moment de frustratie van Amerikaanse onderzoekers die worden gefinancierd door de National Science Foundation. Sequenzen ze vandaag een genoom dan moet dat mogelijk op het web. Ze krijgen geen tijd er eerst zelf goed naar te kijken. Anderen zitten te wachten, zeggen dankjewel en publiceren een leuk resultaat.’

Ook bij open science horen spelregels, vindt Visser. ‘Wéllicht moet je als gebruiker eerst vertellen, wat je plan bent om de gegevens te gaan doen. Komt dat dicht op ons terrein dan is het logisch afspraken te maken over samenwerking of wij moeten eerst de tijd krijgen een paar artikelen te publiceren.’

#### HARDWARE NODIG

Om samenwerking in goede banen te leiden, zijn in maart door wetenschappelijke instellingen waar onder Wageningen UR in Nature de FAIR Guiding Principles gepubliceerd. Data moeten findable, accessible, interoperable en reusable zijn. En goed begin, vindt Visser. ‘Maar daarvan heb je ook hardware nodig. Het maakt niet uit waar je data huisvest, of dat bij IBM is, in de cloud, of bij SARA, maar we zullen in Wageningen ook eigen computerinfrastructuur moeten aanschaffen en mensen hebben die daar mee kunnen omgaan. Anders ben je voor alles afhankelijk van anderen. Als je mee wilt doen, moet je investeren.’ ■

[www.wageningenur.nl/bigdata](http://www.wageningenur.nl/bigdata)

## ‘Als gegevens niet toegankelijk zijn, hebben partijen met meer geld de controle’

