



Limperg Instituut

Data Science in de opleiding voor accountants

Ferry Geertman
Nyenrode



Agenda

- Statistieken
- Leerdoelen data science
- Verschillende raamwerken
- Technieken
- Bias Variance Trade-off
- Responsible AI



Statistieken

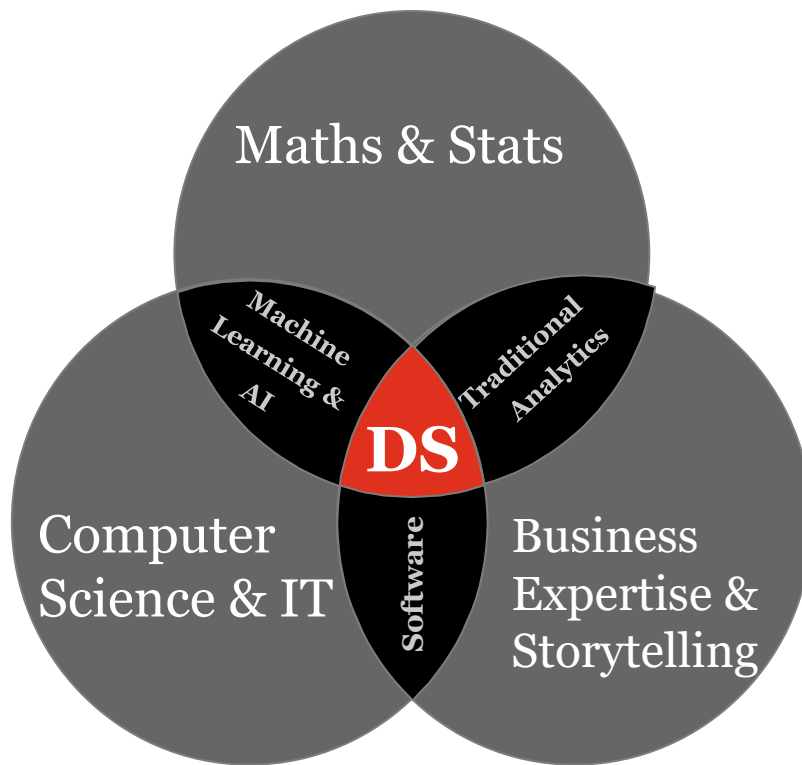
- Vak Data Science gegeven sinds 2020 op Nyenrode
- Ruim 1.100 studenten
- Studielast – 4 ECTS (112 uur)
- Gemiddeld cijfer (1ste poging) 7,4



Wat is Data Science?

Leerdoelen

- $y = f(x) + \varepsilon$
- Duiden van begrippen
- Kunnen beoordelen van kwaliteit van modellen/ algoritmen
- Zelf (eenvoudige) modellen in R bouwen





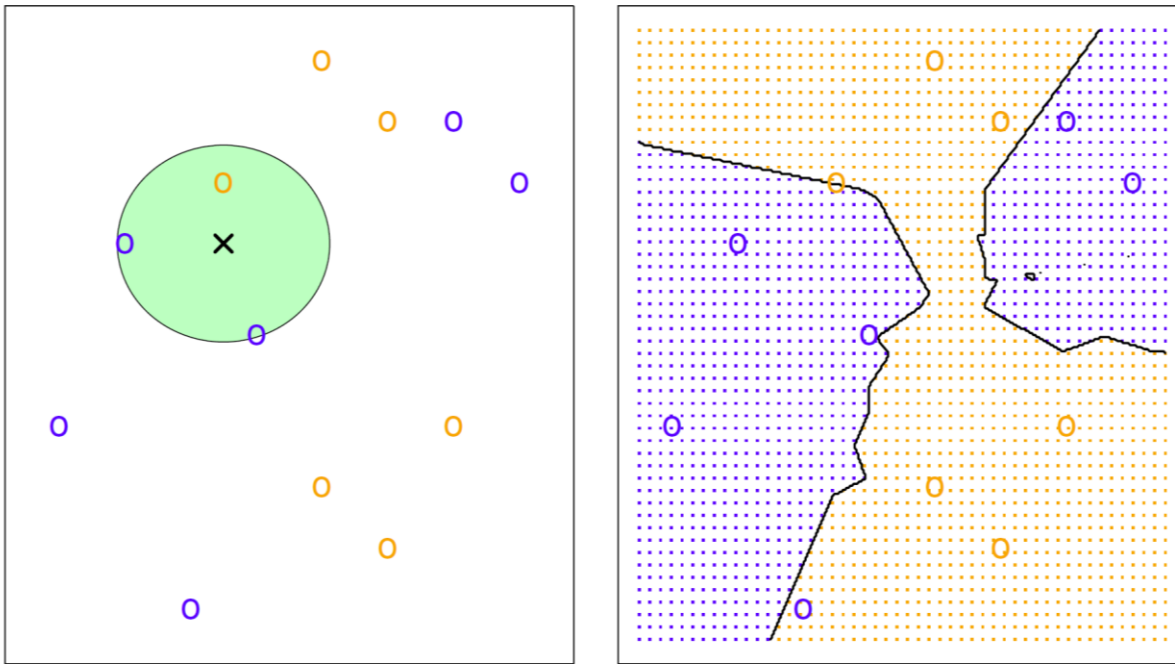
Data Science Lenzen

	Dimension	Low maturity	High maturity
1	Gartner	Descriptive	Prescriptive
2	Vs of big data	Small data	Big Data
3	Push-left	Reporting	Business partner
4	Wisdom pyramid	Deductive	Inductive



Technieken

- Regressie (parametrisch)
- KNN (non-parametrisch)





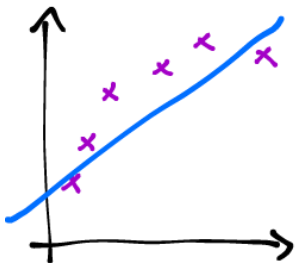
Beoordelen kwaliteit model

BIAS – VARIANCE TRADE-OFF



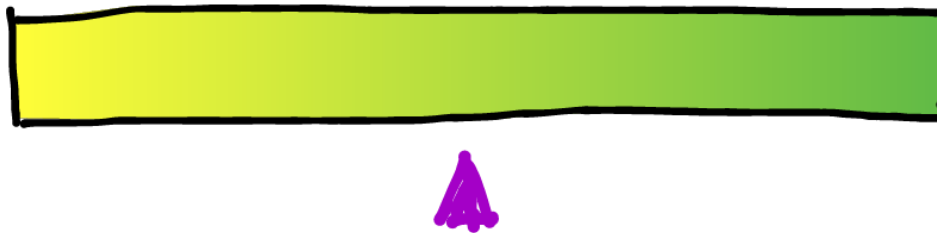
- Bias – fout geïntroduceerd door een complex probleem te benaderen door een simpel model
- Variance – Omvang waarmee de voorspelling gedaan met het model wijzigt als het model getraind zou zijn met andere data
- Trade-off – minder bias leidt tot meer variance en meer bias geeft minder variance
- Op zoek naar het optimum door de complexiteit (flexibiliteit) aan te passen

Underfitting

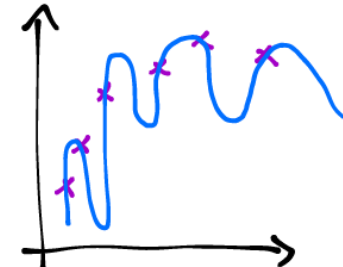


low complexity
high bias
low variance

Complexity



Overfitting



high complexity
low bias
high variance



Bias-Variance Trade-off

Bias - Variance Tradeoff

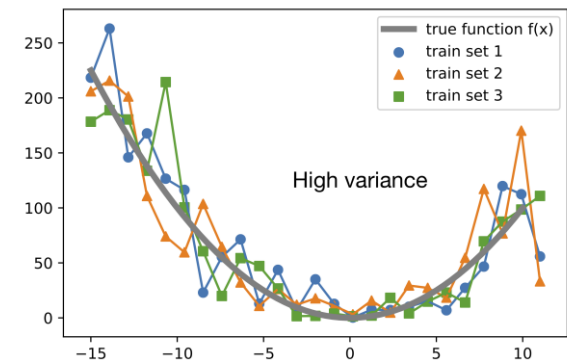
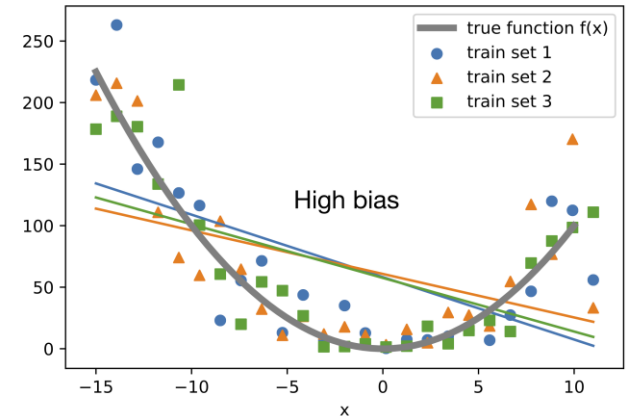
$$\text{Error}(x) = \left(E[\hat{f}(x)] - f(x) \right)^2 + E \left[\hat{f}(x) - E[\hat{f}(x)] \right]^2 + \sigma_e^2$$

predicted true predicted average predicted value irreducible error
 ↓ ↓ ↓ ↓
Bias² **Variance**

How much predicted values differ from true values.

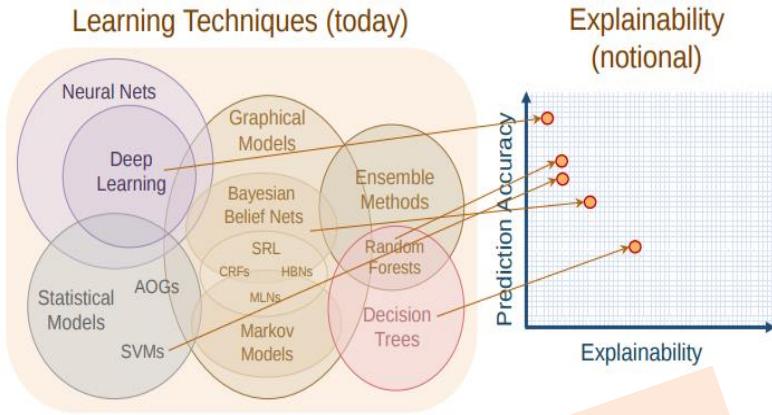
How predictions made on the same value vary on different realizations of the model

BY CHRIS ALBON





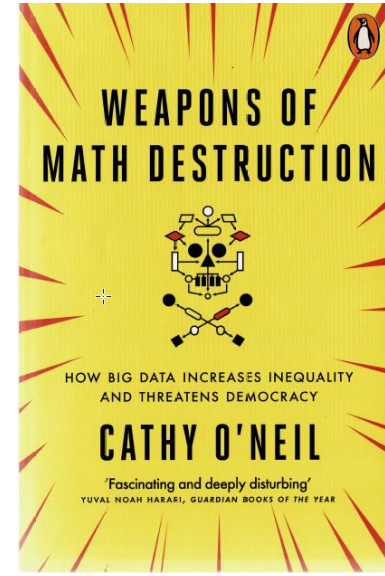
Responsible AI



OPINIE
Data beweegapps van zorgverzekeraars mogen niet gekoppeld worden aan polisvoorwaarden
 Opinie | Mark van Houdenhoven, bijzonder hoogleraar Economische bedrijfsvoering in de gezondheidszorg

11 apr 11:00 • Aangepast op 11 apr 16:02
AFM waarschuwt voor risico's van digitalisering verzekeringen
 Martijn Poels, Puck Sie

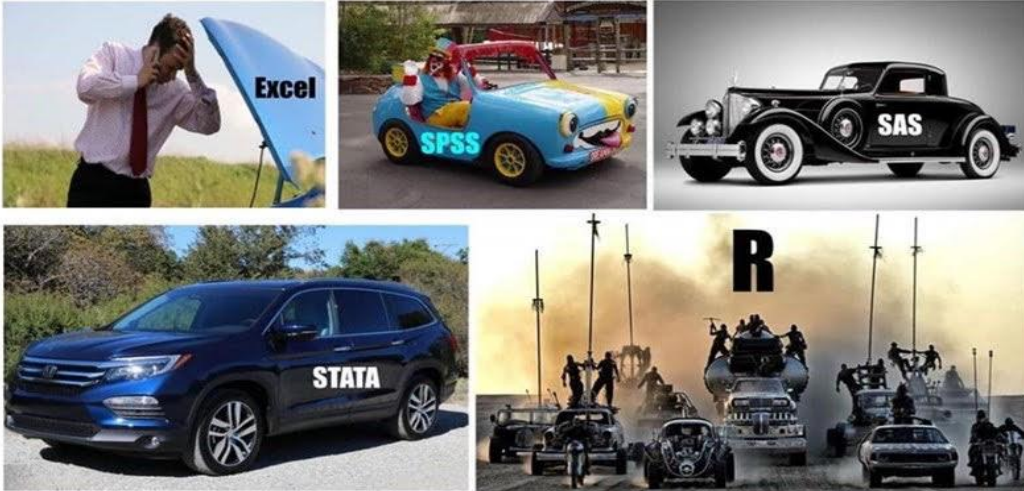
 Het kantoor van de AFM in Amsterdam. Foto: Harold Verstegen/ANP





Gebruik van tooling

If statistics programs/languages were cars...



```

> x=c(1,2,2,3)
> y=c(1,2,3,6)
> x
[1] 1 2 2 3
> y
[1] 1 2 3 6
> mean(x)
[1] 2
> mean(y)
[1] 3
> plot(x,y,col="red",pch=19,cex=2)
> linreg <- lm(y~x)
> abline(linreg,col="blue",lwd=5)
> summary(linreg)
Call: lm(formula = y ~ x)
Residuals:    1    2    3    4
             5.00e-01 -1.00e+00 1.11e-16 5.00e-01
Coefficients:
              Estimate      Std. Error  t value Pr(>|t|)
(Intercept) -2.00000      1.29900   -1.540  0.2635
x              2.50000      0.61240    4.082  0.0551 .

```

```

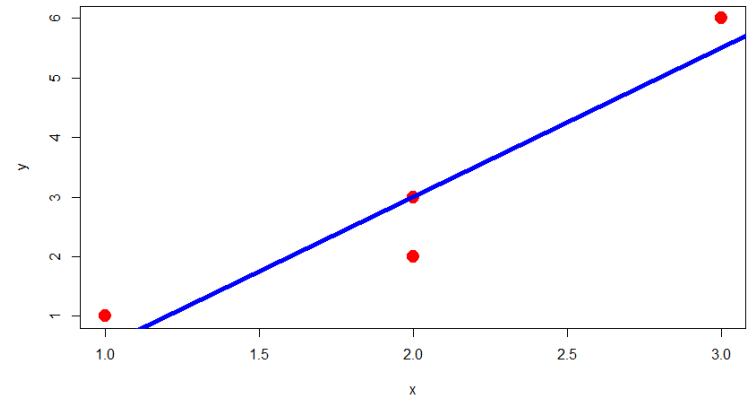
> fit[5,100,]
Error in fit[5, 100, ] : subscript out of bounds

```

```

> library(Hmisc)
Error: package or namespace load failed for 'Hmisc' in loadNamespace(j <- i[[1L]]@lib.loc,
.libPaths()), versionCheck = vI[[j]]):
there is no package called 'xfun'
> describe
Error: object 'describe' not found
> describe(COPD)
Error in describe(COPD) : could not find function "describe"

```





 Gemeente Amsterdam Amsterdam Algoritmeregister Beta

Algoritmeregister Meer informatie Reageer

Een algoritme kan echter zo goed zijn in het vinden van patronen, dat het direct uitsluiten van gevoelige gegevens niet voldoende is. We hebben daarom ook uitgezocht of de niet-gevoelige gegevens die het algoritme verwerkt indirect toch tot ongewenste verschillen in behandeling leiden. Het zou bijvoorbeeld kunnen dat in bepaalde wijken veel mensen wonen met een bepaalde nationaliteit, of dat een bepaalde groep gemiddeld genomen grotere gezinnen heeft. Als het algoritme dan gebruik kan maken van een gegeven als postcode of gezinsgrootte, kan het indirect alsnog onderscheid maken tussen bepaalde groepen, simpelweg door onderscheid te maken tussen wijken/gezinsgrootte. Zo kan een groep alsnog benadeeld worden door het algoritme, ook al is de groep niet expliciet bekend bij het algoritme.

We hebben ervoor gekozen om hier nader onderzoek naar te doen. **Het zogenaamd onderzoek naar onbewuste vooringenomenheid (bias) uitgevoerd en de resultaten hiervan zijn te vinden op:**

Blog post 1:

[Analyzing Bias in Machine Learning: a step-by-step approach \(amsterdamintelligence.com\)](https://amsterdamintelligence.com)

Er kan worden geconcludeerd dat het algoritme vrij is van een ongerechtvaardigde vooringenomenheid en daarmee geschikt is om een pilot mee te starten.

Gebruik van algoritmes binnen gemeente Amsterdam

Maak kennis met de diensten van gemeente Amsterdam waarbij algoritmes gebruikt worden.