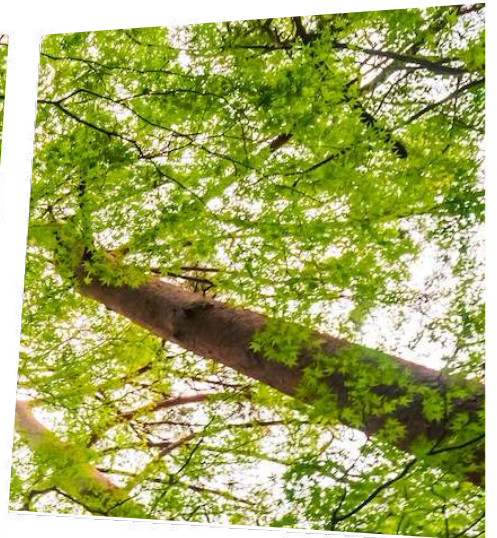


No Fairness without Awareness

Een statistisch onderzoek naar kansengelijkheid aan De HHs



let's change
YOU. US. THE WORLD.

Inhoudsopgave

1	Inleiding	3
1.1	Wat hebben opleidingen en studenten aan dit onderzoeksproject?	3
1.2	Wat heeft De HHs aan dit onderzoeksproject?	3
1.3	Positie van het project in de onderzoekslijnen van het lectoraat, het kenniscentrum en De HHs	5
1.4	Doel en opbouw van dit document	5
2	Onderbouwing van het belang	6
2.1	Theoretische onderbouwing	6
2.2	Ethische overwegingen	8
3	Operationalisering van de benaderingen van kansengelijkheid	13
3.1	Gelijke kansen op de verwezenlijking van het leerpotentieel c.q. gelijke kansen op instroom	13
3.2	Gelijke leer- en diplomakansen bij gelijk potentieel c.q. gelijke kansen op succesvolle doorstroom en uitstroom	14
3.3	Gelijke kansen op een goede plek in de samenleving voor verschillende talenten c.q. gelijke kansen op een vervolgstudie of positie op de arbeidsmarkt	16
3.4	Afbakening	17
3.5	Aansluiting bij bestaand onderzoek binnen De HHs	17
4	Databronnen	18
4.1	Toelichting op bewerkingen per bron	18
4.1.1	Basisgegevens van De HHs	18
4.1.2	Verrijking op basis van aanvullende bronnen	19
4.1.3	Redenen voor de verwerking per bron	21
4.1.4	Schematische weergave bronnen en koppelingen	23
5	Proces van levering en bewerking	24
5.1	Aanvraag aan levering of download	24
5.2	Datamanagement	24
5.3	Bewerking en verrijking	24
6	Methoden van analyse	25
6.1	Deelonderzoek I – Analyse van (mogelijke) studiekeuzes	26
6.1.1	Analyse van verwachte en werkelijke marktaandeelen	26
6.1.2	Analyse van gewogen proportionaliteit in de studentenpopulatie	27
6.2	Deelonderzoek II - Analyse van bottlenecks en gelijke kansen tijdens de studie	28
6.2.1	Modelontwikkeling	28

6.2.2	Detectie, visualisatie en mitigatie van bias in modellen	31
6.3	Deelonderzoek III - Analyse van succes in een vervolgstudie of op de arbeidsmarkt . .	37
7	Verwachte resultaten	38
8	Reproduceerbaarheid	39
	Referenties	41
	Versiegeschiedenis	43
	Repository	43
	Bijlage 1 - Data Science Ethics Checklist	44
8.1	A. Dataverzameling	44
8.2	B. Data-opslag	45
8.3	C. Analyse	46
8.4	D. Modelleren	47
8.5	E. Inzet	48

1 Inleiding

Het project No Fairness without Awareness van het lectoraat Learning Technology & Analytics (LTA) heeft tot doel gelijke kansen voor studenten binnen De Haagse Hogeschool (De HHs) in kaart te brengen en – waar nodig – adviezen te ontwikkelen om deze te verbeteren. Analyse van instroom, doorstroom en uitstroom (*student journeys*) is belangrijk om zicht te krijgen op de positie en betekenis van onze hogeschool voor de studenten en de regio.

1.1 Wat hebben opleidingen en studenten aan dit onderzoeksproject?

Studenten krijgen door deze analyses antwoord op de vraag of De HHs hen gelijke kansen biedt op toelating en het behalen van een diploma binnen afzienbare tijd. Ook krijgen zij zicht op de mogelijke bottlenecks die er zijn en kunnen daarop hun opleiding of de ondersteuning van De HHs bevragen. Ook wordt voor hen duidelijk wat hun kansen zijn op een baan of vervolgstudie op niveau na hun studie aan De HHs. Deze inzichten kunnen vanuit een analyse van de opleidingen komen, maar ook uit onderzoek naar de student journey van groepen studenten. Een voorbeeld van dit laatste is een onderzoek naar de overgang van studenten met een vooropleiding op de Caribische eilanden naar Den Haag en hun student journey.

Opleidingen kunnen deze analyses gebruiken om hun onderwijsbeleid of onderwijs te verbeteren. Zo kan een opleiding met zicht op de ontwikkeling van de instroom van de afgelopen 10 jaar daar mogelijke verbeterpunten uit halen om de opleiding relevant te houden. Of we tonen de impact aan van een wijziging in het curriculum op de samenstelling van nieuwe instroom en de effectiviteit van deze wijziging. En met partners als het Centraal Bureau voor de Statistiek (CBS) kunnen we bijvoorbeeld een analyse maken van het succes van afgestudeerden op de arbeidsmarkt.

Eventuele bottlenecks in de student journey worden zo geïdentificeerd, waarmee opleidingen deze kunnen wegnemen.

1.2 Wat heeft De HHs aan dit onderzoeksproject?

Dit onderzoeksproject draagt bij aan het instellingsplan van De HHs '[Onderzoekend leren met impact](#)'. In het bijzonder aan de strategische thema's I. Kwaliteit van onderwijs en onderzoek (ambitie 1) en IV. Een inclusieve community (ambities 11 en 12).

AMBITIE 1

Continue verbetering van de kwaliteit van de opleidingen.

We streven naar een zo hoog mogelijke kwaliteit van ons onderwijs om de impact en waarde voor studenten en het werkveld te maximaliseren. Kwaliteit is daarnaast ook essentieel voor het

voorkomen van onnodige studieuitval. We kiezen voor een kortcyclische en resultaatgerichte aanpak waarbij we de potentie van studiedata benutten voor de verdere ontwikkeling van studeerbare en doceerbare curricula.

AMBITIE 11

Welzijn van studenten en medewerkers staat voorop.

Wij bieden een omgeving waarin het welzijn van studenten en medewerkers voorop staat, zodat iedereen met plezier en trots aan De Haagse Hogeschool kan studeren en werken. Dat doen we met oog voor individuele achtergronden, behoeftes en ambities van de leden van onze community. Onze community biedt saamhorigheid, vertrouwen en verantwoordelijkheid – sterk motiverende factoren die onze studenten en medewerkers in staat stellen om bij te dragen aan onze community en de samenleving als geheel.

AMBITIE 12

Inclusieve cultuur.

Wij waarderen een verscheidenheid aan perspectieven en we doorbreken uitsluitingsmechanismen. Dit vraagt om openheid en moed om elkaar aan te spreken als er sprake is van discriminatie, ook als dit onbedoeld is, en we erkennen dat het ieders verantwoordelijkheid is om dit te doen. Wij beseffen dat studenten en medewerkers verschillende startposities hebben en daarom ook verschillende bottlenecks tegen kunnen komen. Het is onze verantwoordelijkheid om deze bottlenecks te identificeren en ze samen met studenten en medewerkers weg te nemen of alternatieve routes aan te bieden.

Daarnaast draagt het bij aan de doelstellingen van De HHs ten aanzien van inclusiviteit zoals verwoord in het visiedocument ‘Een inclusieve hogeschool: samenwerking voor inclusief onderwijs en onderzoek’ d.d. 26 november 2021, versie 1.0.

Visie op diversiteit en inclusie

(...) Een inclusieve hogeschool is een omgeving die gelijke kansen biedt voor eenieder, gelijkwaardigheid hoog heeft staan en daarmee fysieke, sociale, culturele, onderwijs- en werkinhoudelijke toegankelijkheid waarborgt voor studenten en medewerkers van diverse achtergronden. In een omgeving als deze wordt scherp toegezien op het voorkomen van belemmeringen of vormen van uitsluiting en discriminatie. Op dit vlak is sprake van een gezamenlijke verantwoordelijkheid voor iedereen in de hogeschool.

Principe van antidiscriminatie

(...) Toegang tot onze hogeschool (en daarmee toegang tot het onderwijs) wordt maximaal gegarandeerd. Faculteiten en diensten nemen de nodige maatregelen om drempels weg te nemen en te allen tijde kritisch te zijn op impliciete en expliciete criteria, regels en procedures

die (onbedoeld) leiden tot uitsluiting en discriminatie. De hogeschool monitort en treedt snel corrigerend op wanneer beleidsmatige uitgangspunten, procedures, werkinstructies en handelwijzen (onbedoeld) uitsluiting en discriminatie in de hand werken.

1.3 Positie van het project in de onderzoekslijnen van het lectoraat, het kenniscentrum en De HHs

Het project draagt bij aan de onderzoekslijnen Student Analytics, Institutional Analytics en Inclusion Analytics van het lectoraat. Niet alleen met nieuwe inzichten, maar ook met de methode die we ontwikkelen om tot die inzichten te komen. Aanvullende onderzoeksvragen die we hierbij stellen zijn: welke vormen van rapportages geven direct inzicht en welke kosten meer moeite? Wat zet aan tot handelen of juist niet? Wat leidt tot impact en wat niet? Hoe kunnen we hierin verschillende doelgroepen het beste bedienen? Deze vragen worden in deze projectbeschrijving verder buiten beschouwing gelaten.

Het onderzoek valt verder in de onderzoekslijn Transformative Technology van het Kenniscentrum Global & Inclusive Learning en sluit vanwege het snijvlak Machine Learning en gelijke kansen aan bij de onderzoeksthema's Digitale Toekomst en Rechtvaardige Samenleving van De HHs.

1.4 Doel en opbouw van dit document

Dit document beschrijft de verdere onderbouwing van het belang van het project (hoofdstuk 2), de operationalisering en onderzoeksvragen van het project (hoofdstuk 3), de bronnen die we daarvoor gebruiken (hoofdstuk 4), het proces voor levering en bewerking (hoofdstuk 5), de methoden van analyse (hoofdstuk 6), de verwachte resultaten (hoofdstuk 7) en de reproduceerbaarheid van het onderzoek (hoofdstuk 8).

2 Onderbouwing van het belang

Dit hoofdstuk bevat een theoretische en ethische onderbouwing van het belang van het onderzoeksproject.

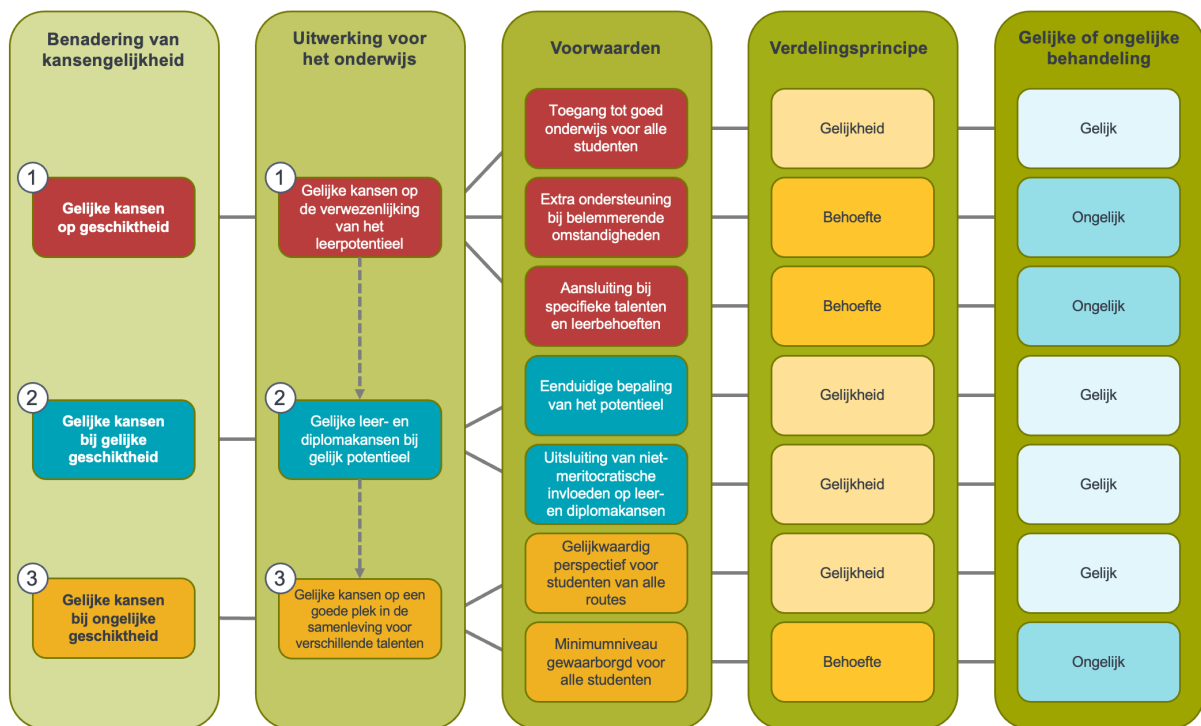
2.1 Theoretische onderbouwing

Het onderzoeksprogramma No Fairness without Awareness heeft tot doel gelijke kansen voor studenten binnen De Haagse Hogeschool in kaart te brengen en – waar nodig – adviezen te geven hoe deze te verbeteren.

Het is een gegeven dat er verschillen bestaan tussen studenten en dat deze een effect hebben op de instroom, doorstroom, uitstroom van studenten in het hoger onderwijs en hun latere succes op de arbeidsmarkt. Daar waar dit de verschillen in capaciteiten van de student weerspiegelt is dat niet altijd direct een probleem. Maar waar dit niet het geval is, kan er sprake zijn van een gebrek aan gelijke kansen. Dit is het onderwerp van studie van het onderzoeksproject ‘No Fairness without Awareness’.

Het Nederlandse onderwijs kent – vergeleken met andere landen in Europa – een bijzondere structuur in het voortgezet onderwijs. Differentiatie in de studiepaden van leerlingen wordt al op 12 jarige leeftijd gemaakt, waar in andere Europese landen dit pas is bij 15 jaar. Dit werkt ongelijkheid in de hand ([Copier, 2022](#)).

De term ‘gelijke kansen’ wordt veel gebruikt, maar het is belangrijk de inhoud daarvan nader te definiëren. Voor gelijke kansen in het hoger onderwijs zijn er drie benaderingen die in de tijd en reikwijdte opschalen ([Elffers, 2022](#)): 1) gelijke kansen op de verwezenlijking van het leerpotentieel, 2) gelijke leer- en diplomakansen bij gelijk potentieel en 3) gelijke kansen op het vinden van een goede plek in de samenleving voor verschillende talenten. Zie voor de nadere uitwerking [Figuur 1](#).



Figuur 1: Framework Gelijke Kansen naar Elffers (2022)

In het Engels is er een verschil tussen de termen *equality* en *equity* (Espinoza, 2007). Terwijl ‘equality’ is gericht op groepsgelijkheid als verdelingsprincipe, is ‘equity’ gericht op de individuele behoeften. In dit document en onderzoek kiezen we als lens het eerste principe, de groepsgelijkheid. In de tweede benadering - gelijke kansen bij gelijke geschiktheid - is een nader onderscheid te maken (Espinoza, 2007), p. 347:

1. **Overleving** (*survival: educational attainment*) - De kans dat studenten uit verschillende sociale groepen in het schoolsysteem blijven tot een bepaald niveau. Hier valt te denken aan de propedeuse of het diploma.
2. **Uitkomst** (*output: educational achievement based on test performance*) - De kans dat studenten uit verschillende sociale groepen op een bepaald moment in het schoolsysteem dezelfde dingen leren op hetzelfde niveau. Hier valt te denken aan de uitkomst op vakniveau.

Wij richten ons in dit onderzoek op ‘overleving’, waarbij we verschillende momenten van de student journey zullen nemen als referentiepunten c.q. bottlenecks. Zie voor een uitwerking hoofdstuk 3.2: [Gelijke leer- en diplomakansen bij gelijk potentieel c.q. gelijke kansen op succesvolle doorstroom en uitstroom](#).

2.2 Ethische overwegingen

Vanuit het principe van academische integriteit is dit document opgesteld om transparant te zijn in het onderzoek en verantwoording af te leggen over de verzamelde data en toegepaste methoden. In deze paragraaf lichten we op hoofdlijnen toe welke ethische overwegingen zijn gemaakt. In bijlage 1 is aan de hand van de [deon checklist voor data scientists](#) een breed scala aan ethische vragen behandeld.

Het onderzoek wil bijdragen aan de oplossing van een maatschappelijk en ethisch vraagstuk: ongelijke kansen in het hoger onderwijs

Het onderzoek is gericht op het beantwoorden van een ethisch vraagstuk: in welke mate hebben studenten van De Haagse Hogeschool gelijke kansen in hun studie en – na afloop – in een vervolgstudie of op de arbeidsmarkt. Het doel is dat De Haagse een veilige en stimulerende omgeving is voor alle studenten in de regio. Daarbij sluiten we ons aan bij de [agenda](#) van de Gelijke Kansen Alliantie van de gemeente Den Haag (2019-2021):

Alle kinderen en jongeren hebben recht op gelijke kansen in het onderwijs. Het is belangrijk dat de talenten van alle kinderen en jongeren optimaal benut worden. Dat is niet alleen van groot belang voor de toekomst van kinderen en jongeren zelf, maar ook voor de toekomst van onze samenleving.

Om gelijke kansen te onderzoeken, bestuderen we achtergrondkenmerken van studenten die mogelijk van invloed zijn op gelijke kansen

Om de onderzoeksvraag te kunnen beantwoorden is een inschatting nodig van de startsituatie van studenten. Dit zijn zowel meritocratische, als eerder behaalde resultaten en vooropleidingen, en niet-meritocratische kenmerken, zoals sociaal-economische status, geslacht en leeftijd. Voor een deel zijn deze direct gerelateerd aan de student – zoals geslacht, leeftijd en vooropleiding. Voor een ander deel zijn ze indirect afgeleid – zoals de sociaal-economische status op basis van de geografische herkomst van de student op buurt- en wijkniveau; zonder deze gegevens kunnen we geen inschatting maken van verschillen in sociaal-economische status. Om analyses te kunnen maken wordt in de modellering gebruik gemaakt van de SES-WOA-scores op buurniveau van het CBS.

Het betrekken van deze factoren is niet onomstreden. Ook al zijn het geen bijzondere persoonsgegevens en is de data die onderzoeken geanonimiseerd, volgens sommigen is het toch beter deze gegevens niet te betrekken in analyses, om bias te voorkomen ([Pedreshi et al., 2008](#)). Dit is echter naïef ([Hardt et al., 2016](#)): het verwijderen van deze variabelen sluit niet uit dat deze kenmerken alsnog via andere variabelen in de dataset als proxy aanwezig zijn en alsnog invloed hebben op de uitkomsten van een analyse of afgeleid besluit. We willen mogelijke bias in bestaande data te ontdekken om te voorkomen dat deze bij de ontwikkelingen van statistische in het hoger onderwijs tot regel worden verheven en daarmee institutionele discriminatie tot gevolg kunnen hebben ([Pedreshi et al., 2008](#)). Hier geldt: “Fairness through awareness” ([Dwork et al., 2011](#)). Hiervoor is het nodig om de data uit

processen van werving, toelating, onderwijs en toetsing van De Haagse Hogeschool te onderzoeken om te zien of er mogelijk sprake is van ongelijke kansen die terug te voeren zijn op deze kenmerken.

We vergroten het inzicht in het onderzoek tot nu toe naar kansengelijkheid in het onderwijs

In een verkennende studie naar kansengelijkheid in het onderwijs zijn 30 factoren in 5 niveaus onderscheiden: leerling (9), familie (6), school (7), wijk (5) en samenleving (3) (Badou & Day, 2021). Niet alle factoren laten zich vertalen in beschikbare data, daarom nemen we in dit onderzoek een select aantal van deze factoren mee, met name op leerling-, school- en wijkniveau. Een van de aanbevelingen uit de studie was vervolgonderzoek naar de onderlinge interactie en rangorde van deze factoren. Ter illustratie kunnen geslacht en leeftijd met elkaar interacteren bij studiesucces na 1 jaar, maar blijkt wellicht dat daarin de leeftijd belangrijker is dan het geslacht. Het statistische onderzoek dat we voorstellen met behulp Machine Learning draagt hieraan bij, omdat het zowel interactie-effecten kan onderzoeken als de relatieve bijdrage van elke factor aan een prognosemodel.

Gelijke kansen zijn niet gelijk aan discriminatie

Het onderscheid tussen gelijke kansen en discriminatie is van belang. Het onderzoek heeft tot doel om gelijke kansen te onderzoeken in de context van het hoger onderwijs en De Haagse Hogeschool in het bijzonder; dit kan los staan van discriminatie. Discriminatie is het moedwillig maken van onderscheid op basis van een of meer niet-meritocratische achtergrondkenmerken van een persoon met als doel hen nadelig en ongelijk te behandelen. Gelijke kansen is in dit onderzoek een neutrale term die stelt of een student in transitie in de student journey gelijke uitkomsten kan verwachten op basis van gelijke merites. Hiermee doelen we op persoonlijke capaciteiten, zoals intelligentie, inzet, executieve vaardigheden en motivatie.

In de presentatie van de onderzoeksresultaten houden we rekening met het voorkomen van onthullingsgevaar en stigmatisering

Het onderzoek mag zelf discriminatie of stigmatisering niet in de hand werken. Daarom wordt rekening gehouden met mogelijk onthullingsgevaar en betrouwbaarheid afgeleid van de [richtlijnen van het CBS](#):

- Bij tabellen en grafieken wordt een minimaal aantal van 10 observaties per cel gehanteerd.
- Percentages worden alleen gepubliceerd wanneer de noemer tenminste 100 waarnemingen bevat.
- Visualisaties worden geaggregeerd naar hogere niveaus (bijv. naar cohort, naar opleiding, naar wijk).

In de toepassing van Machine Learning in dit onderzoek hanteren we de ethische principes van de EU

De EU heeft [Ethische richtlijnen voor betrouwbare AI](#) opgesteld (8 april 2019). Deze bevatten vier ethische richtlijnen waaraan we ons verbinden in dit onderzoek:

- **Respect voor menselijke autonomie** - De Machine Learning algoritmen uit dit onderzoek worden niet toegepast voor praktijktoepassingen, zoals een voorspelmodel voor individueel advies aan studenten, maar juist om de 'menselijke cognitieve, sociale en culturele vaardigheden te vergroten, aan te vullen en te versterken'.

De grondrechten waarop de EU is gefundeerd, zijn erop gericht respect voor de vrijheid en autonomie van mensen te waarborgen. Mensen die met KI-systemen werken, moeten hun volledige en effectieve zelfbeschikking kunnen behouden en kunnen deelnemen aan het democratische proces. KI-systemen mogen mensen niet onterecht onderwerpen, dwingen, misleiden, manipuleren, conditioneren of drijven. Ze moeten veeleer worden ontworpen om de menselijke cognitieve, sociale en culturele vaardigheden te vergroten, aan te vullen en te versterken. De verdeling van functies tussen mensen en KI-systemen moet gebeuren volgens ontwerpbeginselen waarbij de mens centraal staat, en moet zinvolle mogelijkheden voor menselijke keuze openlaten. Er moet dus worden gezorgd voor menselijk toezicht en menselijke controle op de werkprocessen in KI-systemen.

- **Preventie van schade** - De Machine Learning uit dit onderzoek wordt toegepast juist om vast te stellen waar mogelijk het risico voor negatieve gevolgen kunnen hebben voor minderheden. In de loop van het onderzoek (Deelonderzoek II) zullen we studenten betrekken bij de ontwikkeling van de analyse op bottlenecks, waar ML voor gebruikt zal worden. We verkennen met het Inclusion Office en het Partner Up! programma wat de mogelijkheden zijn.

KI-systemen mogen geen schade veroorzaken of vergroten of anderszins negatieve gevolgen hebben voor mensen. Dat betekent bescherming van de waardigheid, alsook de geestelijke en fysieke integriteit van mensen. KI-systemen en de omgeving waarin zij werken, moeten veilig en zeker zijn. Ze moeten technisch robuust zijn en er moet voor worden gezorgd dat ze geen ruimte bieden voor kwaadwillig gebruik. Kwetsbare personen moeten meer aandacht krijgen en moeten worden betrokken bij de ontwikkeling en installatie van KI-systemen. Er moet specifiek aandacht worden besteed aan situaties waarin KI-systemen negatieve gevolgen kunnen veroorzaken of vergroten vanwege ongelijkheid wat betreft macht of beschikking over informatie, bijvoorbeeld tussen werkgevers en werknemers, tussen bedrijven en consumenten of tussen overheden en burgers. Preventie van schade betekent ook dat er rekening moet worden gehouden met de natuurlijke omgeving en alle levende wezens.

- **Rechtvaardigheid** - Het onderzoek is ingesteld om rechtvaardigheid en gelijke kansen te bevorderen. Machine Learning wordt ingezet om biases te ontdekken in historische gegevens van De HHs om van daaruit af te leiden welke bias in toekomstige regelgeving of toepassingen van Machine Learning voorkomen dient te worden.

De ontwikkeling, de installatie en het gebruik van KI-systemen moeten rechtvaardig zijn. Wij erkennen dat er veel verschillende interpretaties van rechtvaardigheid bestaan, maar zijn ervan overtuigd dat rechtvaardigheid zowel een inhoudelijke als een procedurele dimensie heeft. De inhoudelijke dimensie impliceert een toezegging om de gelijke en rechtvaardige verdeling van zowel voordelen als kosten te waarborgen en ervoor te zorgen dat personen en groepen vrij zijn van onrechtvaardige vertekening, discriminatie en stigmatisering. Indien onrechtvaardige vertekening kan worden voorkomen, zouden KI-systemen zelfs de maatschappelijke rechtvaardigheid kunnen vergroten. Gelijke kansen wat betreft toegang tot onderwijs, goederen, diensten en technologie moeten ook worden bevorderd. Daarnaast mag het gebruik van KI-systemen nooit tot gevolg hebben dat de (eind)gebruikers worden misleid of worden beperkt in hun keuzevrijheid. Verder impliceert rechtvaardigheid dat beroepsbeoefenaars op het gebied van KI het beginsel van evenredigheid tussen middelen en doelen moeten eerbiedigen en zorgvuldig moeten afwegen hoe ze tegengestelde belangen en doelstellingen in evenwicht kunnen brengen. De procedurele dimensie van rechtvaardigheid omvat het vermogen om beslissingen die worden genomen door KI-systemen en door de mensen die deze systemen beheren, aan te vechten en er effectief beroep tegen in te stellen. Om dat te kunnen doen moet de entiteit die verantwoordelijk is voor de beslissing, identificeerbaar zijn en moet het besluitvormingsproces verklaarbaar zijn.

- **Verantwoording** - Naast de verantwoording die dit document zelf vormt voor het onderzoeksproject, zullen de toegepaste modellen en inputfactoren ter verantwoording beschreven worden en komt de broncode van deze modellen publiek beschikbaar.

Verantwoording is cruciaal voor het scheppen en behouden van het vertrouwen van gebruikers in KI-systemen. Dat betekent dat processen transparant moeten zijn, dat de capaciteiten en het doel van KI-systemen openlijk kenbaar moeten worden gemaakt en dat beslissingen – voor zover mogelijk – verklaarbaar moeten zijn aan degenen die er direct of indirect de gevolgen van ondervinden. Zonder die informatie kan een beslissing niet naar behoren worden aangevochten. Het is niet altijd mogelijk om te verklaren waarom een model een bepaald resultaat of een bepaalde beslissing heeft opgeleverd (en welke combinatie van inputfactoren daaraan heeft bijgedragen). Deze gevallen worden “blackbox”-algoritmen genoemd en vereisen speciale aandacht. In die situaties kunnen andere verantwoordingsmaatregelen (zoals traceerbaarheid, controleerbaarheid en transparante communicatie over de capaciteiten van het systeem) nodig zijn, mits het systeem als geheel de grondrechten eerbiedigt. De mate waarin verantwoording nodig is, hangt sterk af van de context en de ernst van de gevolgen, mocht het resultaat onjuist of anderszins onnauwkeurig zijn.

In de presentatie van de uitkomsten gebruiken we verschillende invalshoeken met uitwerking van voor- en nadelen

Er is niet een definitie van eerlijkheid; je kunt en moet er meervoudig naar kijken. Om te voorkomen

dat er één dominante, te eenvoudige kijk komt op eerlijkheid in het onderwijs, beleid of begeleiding ontstaat binnen De HHs, zullen we in de presentatie van de resultaten en besprekingen met stakeholders die nuances inbrengen. We zullen de verschillende invalshoeken naast elkaar presenteren met voor- en nadelen. Daarnaast zullen we de beperkingen over het onderzoek meenemen, zoals bijvoorbeeld relationele eerlijkheid (Fish & Stark, 2022) die niet met de gehanteerde methode beantwoord kunnen worden.

We geven studenten een stem in ons begrip van eerlijkheid en de mogelijke vertalingen van de inzichten naar toepassingen in De HHs

We willen voorkomen dat de onderzoeksresultaten in de vertaling naar de praktijk kunnen leiden tot een versmalling van beelden over groepen studenten.

Volgens de Haagse inclusiviteitsprincipes ‘inclusieve omgangsvormen’ en ‘inclusieve governance’ willen we studenten niet reduceren tot een categorie en voor de mogelijke verbetering van onderwijsbeleid en praktijk actieve participatie stimuleren. De onderzoeksresultaten bespreken we daarom via diverse gremia en kanalen met studenten. Mogelijke studentnetwerken zijn faculteitsraden, het netwerken van het Inclusion Office en van het Kenniscentrum Global & Inclusive Learning, en studenten van de opleiding CMD.

We toetsen met deze studenten vooraf de beelden over eerlijkheid en de selectie van de data uit het onderzoeksvoorstel, en gaandeweg het onderzoek de uitkomsten en de communicatie daarover. Studenten van de opleiding CMD zullen we vragen de visualisaties van de software die we gaan gebruiken te onderzoeken op gebruiksvriendelijkheid en bruikbaarheid voor studenten. Daarnaast zullen we aan faculteiten een handreiking bieden hoe zij het eventuele gebruik van de inzichten kunnen communiceren naar hun studenten en medewerkers die we bij studenten toetsen.

3 Operationalisering van de benaderingen van kansengelijkheid

Om gelijke kansen en eventuele bottlenecks in kaart te brengen onderzoeken we de drie benaderingen van Elffers op gelijke kansen op basis van 1) herkomst en instroom, 2) doorstroom en uitstroom met of zonder diploma, en 3) succes in een vervolgopleiding of de arbeidsmarkt.

3.1 Gelijke kansen op de verwezenlijking van het leerpotentieel c.q. gelijke kansen op instroom

Dit eerste perspectief houdt zowel een gelijke als een ongelijke behandeling in:

“Daarvoor is in de eerste plaats nodig dat alle leerlingen gelijke toegang hebben tot onderwijs van voldoende kwaliteit om zicht te kunnen ontwikkelen, wat neerkomt op een gelijke behandeling. Maar het vereist evenzeer dat er aanpassingen worden gedaan wanneer leerlingen te maken hebben met omstandigheden die de verwezenlijking van hun potentieel mogelijk begrenzen, en dat er maatwerk wordt geboden aan leerlingen met uiteenlopende talenten en leerbehoeftes. Dat vereist juist ongelijke behandeling.” (Elffers, 2022)

Hiervoor onderzoeken we de **herkomst** en **instroom** van onze studenten en hun **kans op instroom** in De HHs. De aanname is dat alle studenten uit het toeleveringsgebied van De HHs evenveel kans hebben om een opleiding te volgen aan De HHs. De verwachting is dat dit waarschijnlijk niet zo zijn. Eventuele bottlenecks kunnen om verschillende redenen ontstaan (Elffers, 2022): a) geen toegang tot goed onderwijs voor alle studenten, b) te weinig ondersteuning bij belemmerende omstandigheden, en c) geen aansluiting bij de specifieke talenten en leerbehoeftes.

Tabel 1: Operationalisering gelijke kansen op de verwezenlijking van het leerpotentieel

Voorwaarde	Operationalisering	Mogelijke data
a) Geen toegang tot goed onderwijs voor alle studenten	1. Kans op instroom op basis van demografische en sociaal-economische achtergrond	Geslacht, leeftijd, sociaal-economische achtergrond, datum van aanmelding
	2. Selectieve werving door De HHs	Wervingscriteria van De HHs
	3. Selectie van studenten bij selectieve opleidingen	Selectiecriteria

Voorwaarde	Operationalisering	Mogelijke data
b) Te weinig ondersteuning bij belemmerende omstandigheden	1. Ondersteuning van studenten met een functiebeperking	Ontwikkeling omvang en aard van aanvragen van studenten met een functiebeperking
	2. Ondersteuning voor studenten met een niet-reguliere vooropleiding (staatsexamen of 21+)	Ontwikkeling omvang en aard van instroom van studenten met een niet-reguliere vooropleiding
	3. Ondersteuning voor studenten met deficiënties	Criteria voor toekenning van voorzieningen
c) Geen aansluiting bij de specifieke talenten en leerbehoeften	1. Selectieve oriëntatie en studiekeuze door studenten	Kans op instroom vanuit de regio afgezet ten opzichte van de daadwerkelijke instroom
	2. Kwaliteit van aansluiting afhankelijk van vooropleiding (havo, mbo, buitenlands diploma, etc.)	Vooropleiding, land van de hoogste vooropleiding, soort aansluiting
	3. Kwaliteit van aansluiting afhankelijk van soort aansluiting (direct na vooropleiding, tussenjaar, switch)	Aansluiting en instroom van Caribische studenten

De **onderzoeksvragen** die we op basis van de studiedata van De HHs in dit perspectief willen onderzoeken zijn:

1. In welke mate hebben aankomende studenten gelijke of ongelijk kansen tot de toelating van het onderwijs van De HHs?
2. Worden studenten in het toelatingsproces door De HHs gelijk of juist ongelijk behandeld waar dit – gezien hun kansen – nodig is?

3.2 Gelijke leer- en diplomakansen bij gelijk potentieel c.q. gelijke kansen op succesvolle doorstroom en uitstroom

Het tweede perspectief vraagt om een gelijke behandeling van studenten bij gelijk potentieel.

Dat vergt enerzijds een eenduidige bepaling van het potentieel van leerlingen en anderzijds het zoveel mogelijk uitsluiten van de invloed van niet-meritocratische factoren op de kansen van leerlingen om te leren en een diploma te halen in een bepaalde vorm van onderwijs. (Elffers,

2022)

Hiervoor onderzoeken we de **doorstroom** van studenten door hun studie en **uitstroom met of zonder diploma**. De aanname is dat studenten bij gelijke geschiktheid evenveel kans hebben om per vak, fase of voor de hele studie door te stromen en evenveel kans hebben een diploma te behalen of uit te vallen na ieder studiejaar.

Tabel 2: Operationalisering gelijke leer- en diplomakansen bij gelijk potentieel

Voorwaarde	Operationalisering	Mogelijke data
a) Eenduidige bepaling van het potentieel	1. Kwaliteit van onderwijs en toetsing	Achtergrond- en vooropleidingskenmerken als bij perspectief 1
b) Uitsluiting van niet-meritocratische leer- en diplomakansen bij gelijk potentieel	1. Kans op doorstroom en diplomering	Doorstroom bij vakken, transitiepunten (zoals propedeuse c.q. BSA, volgende fasen in de studie) en afstuderen binnen een afzienbare tijd (nominale studietijd of nominale studietijd + 1 jaar)
	2. Kans op uitval	Uitval na ieder studiejaar
	3. Kans op het volgen van keuzeonderdelen	Toegankelijkheid van het gewenste onderwijs
	4. Kans op voorzieningen voor studenten met een functiebeperking	Beschikbaarheid van voorzieningen voor het volgen van onderwijs of afleggen van toetsen
	5. Kans op het volgen van alternatieve onderwijs en toetsvormen voor het behalen van leeruitkomsten.	Beschikbaarheid van alternatieve vormen voor het behalen van leeruitkomsten

De **onderzoeksvragen** die we op basis van de studiedata van De HHs in dit perspectief willen onderzoeken zijn:

1. In welke mate hebben studenten gelijke of ongelijk kansen op doorstroom in het onderwijs van De HHs?
2. Worden studenten in de doorstroom in het onderwijs door De HHs gelijk of juist ongelijk behandeld waar dit – gezien hun kansen – nodig is?

3.3 Gelijke kansen op een goede plek in de samenleving voor verschillende talenten c.q. gelijke kansen op een vervolgstudie of positie op de arbeidsmarkt

Dit derde perspectief houdt zowel een gelijke als een ongelijke behandeling in:

Dat stelt in de eerste plaats eisen aan de perspectieven die verschillende routes in het onderwijs bieden. De routes moeten een gelijkwaardig perspectief bieden op het bereiken van een goede plek in de samenleving. (...) Ongeacht de mate van ongelijkheid in beloning, vereist het waarborgen van de kansen op een goed leven voor iedereen in elk geval dat het onderwijs alle leerlingen in alle routes tot op een niveau weet te brengen dat minimaal nodig is om volwaardig te kunnen participeren in het onderwijs, op de arbeidsmarkt en in de samenleving in brede zin. Gegeven de uiteenlopende niveaus en leerbehoeften van leerlingen vergt het bereiken van dat minimumniveau een ongelijke behandeling. (Elffers, 2022)

Hiervoor onderzoeken we het **succes** van studenten in een **vervolgopleiding** of de **arbeidsmarkt**. De aanname is dat studenten bij ongelijke geschiktheid evenveel kans hebben om een vervolgopleiding te volgen of een baan te vinden op niveau, binnen het afstudeerdomein en binnen afzienbare tijd.

Tabel 3: Operationalisering gelijke kansen op een goede plek in de samenleving voor verschillende talenten

Voorwaarde	Operationalisering	Mogelijke data
a) Gelijkwaardig perspectief voor studenten van alle routes	1. Kans op het volgen van een vervolgstudie na de verschillende onderwijsvormen en opleidingsvormen van De HHs 2. Kans op een baan na de verschillende onderwijsvormen en opleidingsvormen van De HHs	Inschrijfvormen (voltijd, deeltijd, duaal) en opleidingsvormen (bachelor, ad, master) ICHO data over doorstroom naar een vervolgopleiding Inschrijfvormen (voltijd, deeltijd, duaal) en opleidingsvormen (bachelor, ad, master) CBS-microdata over arbeidsmarktsucces
b) Minimumniveau gewaarborgd voor alle studenten	[nader te bepalen]	[nader te bepalen]

De **onderzoeksvragen** die we op basis van de studiedata van De HHs in dit perspectief willen onderzoeken zijn:

1. In welke mate hebben studenten gelijke of ongelijk kansen op doorstroom naar een vervolgopleiding of een baan op niveau, binnen het afstudeerdomein en binnen afzienbare tijd na het onderwijs van De HHs?
2. Worden studenten in de doorstroom na het onderwijs door De HHs gelijk of juist ongelijk behandeld waar dit – gezien hun kansen – nodig is?

3.4 Afbakening

Voor alle formele KPI's en rapportages van De HHs geldt dat deze niet uit het onderzoek van het lectoraat af te leiden zijn, maar uitsluitend gebaseerd kunnen zijn op de managementdashboards en -rapportages van de diensten Onderwijs, Kennis & Communicatie (OKC) en Business & Control (B&C).

Omdat vanuit de operationalisering blijkt dat we voor deze gegevens gebruik maken van Nederlandse bronnen (DUO, CBS) beperkt dit onderzoek zich in beginsel tot instroom van studenten met een Nederlandse vooropleiding in het voortgezet onderwijs¹ die instromen in een bachelor of associate degree opleiding aan De HHs, voltijd, deeltijd of dual.

Om te kunnen bepalen of er (significante) verschillen zijn met studenten met een internationale vooropleiding zullen die inschrijvingen wel worden meegenomen in exploratieve analyses.

Het onderzoek wordt voor nu een eerste keer uitgevoerd; of het vaker herhaald zal worden en met welke frequentie is nu nog onbekend en zal afhankelijk zijn van de uitkomsten, behoefte van De HHs en de ontwikkeling van het onderzoeksprogramma van het lectoraat.

3.5 Aansluiting bij bestaand onderzoek binnen De HHs

Definities zullen geëxpliciteerd worden en vergeleken worden met de formele definities die gehanteerd worden binnen rapportages van OKC en B&C. Het is mogelijk dat in het onderzoek afgeweken wordt van definities om meer dynamiek in de studentenpopulatie te kunnen onderzoeken, zoals uitval voor 1 oktober, of omdat data (nog) niet voorhanden is, zoals nationaliteit. In al deze gevallen worden afwijkingen onderbouwd en gerapporteerd.

Waar mogelijk of gewenst vergelijken we uitkomsten met bestaande onderzoeken van OKC of de VH, zoals de [HBO-monitor](#) (uitgevoerd door ROA van de Universiteit Maastricht en DESAN Research Solutions).

¹ Dit is dus inclusief Caribische studenten.

4 Databronnen

Afhankelijk van het perspectief op gelijke kansen onderzoeken we verschillende bronnen:

Tabel 4: Bronnen per perspectief

Perspectief	Fase	Bronnen
1) Gelijke kansen op de verwezenlijking van het leerpotentieel	Oriëntatie en instroom	<p>De HHs: OKC - studiedata over aanmeldingen en inschrijvingen uit het Studiedataproject</p> <p>De HHs: OKC - CRM-data over werving en oriëntatie</p> <p>CBS: Open data over demografische ontwikkeling en sociaal economische status van wijken en buurten in Nederland</p> <p>OCW: Open data over leerlingen en afgestudeerden in het voortgezet onderwijs</p>
2) Gelijke leer- en diplomakansen bij gelijk potentieel	Doorstroom en behalen diploma	<p>De HHs: OKC - studiedata over studievoortgang, uitval en diplomering uit het Studiedataproject</p>
3) Gelijke kansen op een goede plek in de samenleving voor verschillende talenten	Uitstroom naar vervolgstudie of arbeidsmarkt	<p>CBS: Microdata over arbeidsmarktsucces</p> <p>1CHO: Studiedata over doorstroom naar een vervolgopleiding</p>

4.1 Toelichting op bewerkingen per bron

De basis voor het onderzoek vormt de dataset van het studiedataproject van OKC. Deze dataset noemen we 'de studiedata analyseset'.

4.1.1 Basisgegevens van De HHs

De **HHs studiedata analyseset**:

- bevat per student per jaar per jaar inschrijving een serie variabelen voor de cohorten 2012 tot en met 2022 (zie Bijlage 2 - Variabelen_Analyseset_202306068 voor een beschrijving van deze variabelen) ;

- wordt op verzoek geanonimiseerd geleverd door het team IR & Analytics van OKC is daarmee een secundaire bron:
 - gegevens zoals studentnummer, naam, geboorteplaats, geboortedatum of andere direct identificeerbare kenmerken zijn uitgesloten;
 - bevat geen bijzondere persoonsgegevens.

4.1.2 Verrijking op basis van aanvullende bronnen

De studiedata analyseset wordt door het lectoraat LTA tot een tertiaire bron vanuit aanvullende, secundaire bronnen verrijkt voor een beter begrip van de context van de student en van de opleiding die de student volgt: het Customer Relationship Management (CRM) van De HHs, en data van de Vereniging Hogescholen (VH), de Dienst Uitvoering Onderwijs (DUO), het Centraal Bureau voor de Statistiek (CBS), Studiekeuze123, OpenStreetmaps (OSM) / OpenTripPlanner (OTP), en 1 cijfer HO (1CHO; eveneens van DUO).

- De **HHs CRM data**:
 - bevatten per opleiding van De HHs gegevens over aantallen bezoekers aan oriëntatieactiviteiten en conversiepercentages van oriëntatie naar aanmelding;
 - zijn per collegejaar geaggregeerd naar opleidingsniveau;
 - worden op aanvraag geleverd door het team Marketing & CRM van de unit Marketing & Communicatie van de dienst OKC van De HHs;
 - zijn beschikbaar vanaf collegejaar 2022-2023;
 - worden gekoppeld op opleidingsniveau op basis van de naam van de opleiding.
- De **VH opleidingsdata**:
 - bevatten landelijke, publieke gegevens over eerstejaars instroom, doorstroom en uitstroom van studenten in opleidingen aan hogescholen;
 - zijn per collegejaar geaggregeerd naar opleidingsniveau;
 - zijn publiek beschikbaar via de website van de VH via de dashboards [instroom](#), [inschrijvingen en diploma's](#) (cohorten 2017 tot en met 2022) en [studiesucces, uitval en studiewissel](#) (cohorten 2017 tot en met 2020);
 - worden gekoppeld op opleidingsniveau op basis van de naam van de opleiding.
- De **DUO open onderwijs data**:
 - bevatten in de vorm van open onderwijs data publieke, landelijke gegevens over a) (verwachte) instroom, doorstroom en uitstroom van scholen uit het voortgezet en hoger onderwijs (cohorten 2017 tot en met 2022; prognoses: 2022 tot en met 2041) en b) vestigingsadressen;
 - zijn per collegejaar geaggregeerd naar schooltype en onderwijsinstelling;

- zijn publiek beschikbaar via de website van [DUO](#);
 - worden gekoppeld op schoolniveau op basis van de BRIN6 code: school en vestiging.
- De **CBS open data**:
 - bevatten publieke statistieken van het CBS uit Statline naar buurt- en wijkniveau;
 - [kerncijfers wijken en buurten 2004-2022](#) worden gebruikt per buurt en wijk van kenmerken betreffende bevolking, wonen, energie, onderwijs, arbeid, inkomen, sociale zekerheid en voorzieningen als proxies van deze kenmerken van de student bij aanvang van de studie;
 - [sociaal-economische status; scores per wijk en buurt, regio-indeling 2021](#) worden per buurt en wijk gebruikt als proxy voor de sociaal-economische status² van de student bij aanvang van de studie;
 - voor een koppeling op buurt en -wijkniveau³ wordt de 4-cijferige postcode van het woonadres op het moment van afstuderen van de vooropleiding van de student geaggregeerd naar een buurt- en wijkcode binnen een gemeentecode⁴.
 - De **Studiekeuze123 NSE data**:
 - zijn het [NSE benchmarkbestand](#): een meerjarig bestand en bevat studentoordelen en zogenaamde duidingsgegevens (zoals contacttijd);
 - worden gebruikt voor het benchmarken van de verschillen in studentoordelen uit de Nationale Studenten Enquête per opleiding;
 - worden op verzoek ter beschikking gesteld aan belangstellenden, waaronder onderzoekers;
 - worden gekoppeld op opleidingsniveau op basis van de naam van de opleiding.
 - De **OSM en OTP data**:
 - worden gebruikt voor het berekenen van de publieke reistijden tussen de 4-cijferige postcode van het woonadres op het moment van het behalen van het diploma van de vooropleiding en een mogelijke vervolgopleiding bij alternatieve onderwijsinstellingen met eenzelfde opleidingsaanbod en daarmee de kans op instroom naar De HHs;
 - worden gekoppeld per student op basis van de 4-cijferige postcode van het woonadres op het moment van het behalen van het diploma.
 - De **DUO ICHO data**:
 - worden door DUO geleverd aan De HHs;

² CBS: "De sociaal-economische status (SES-WOA) van gemeenten, wijken en buurten in Nederland. Deze status wordt beschreven in termen van de financiële welvaart, het opleidingsniveau en het recente arbeidsverleden van particuliere huishoudens op 1 januari van het verslagjaar."

³ CBS: "Onderdeel van een gemeente, bestaande uit één of meerdere buurten. Vaak komt een wijk overeen met een woonplaats of een deel van een grotere woonplaats."

⁴ Daar waar een postcode in meerdere buurten, wijken of gemeenten ligt, wordt de postcode gekozen met het hoogste aantal woningen, conform de methode die het CBS hier zelf voor hanteert.

- worden gebruikt voor het analyseren van doorstroom naar een vervolgstudie per opleiding;
- bevatten alle deelnames en resultaten van studenten in het hoger onderwijs, vanaf 1991 tot en met 1 oktober van het afgelopen jaar;
- worden gekoppeld op opleidingsniveau op basis van de naam van de opleiding.

Naast de verrijking met bestaande bronnen zal het lectoraat – in samenwerking met het studiedatateam van IR&A van OKC en medewerkers van de faculteiten – een primaire bron ontwikkelen, de **opleidingendataset**. Deze dataset:

- bevat specifieke kenmerken per opleiding per cohort zoals formele naam, populaire naam, BSA-grens, aanmelddeadlines, ingangseisen voor delen van het curriculum, cijferschalen, tracks, etc.
- wordt in de basis vanuit de **studiedata analyseset** opgebouwd - zoveel mogelijk kenmerken worden vanuit beschikbare data overgenomen of afgeleid;
- wordt gekoppeld op opleidingsniveau op basis van de naam van de opleiding.

4.1.3 Redenen voor de verwerking per bron

Tabel 5: Onderbouwing voor de verwerking per bron

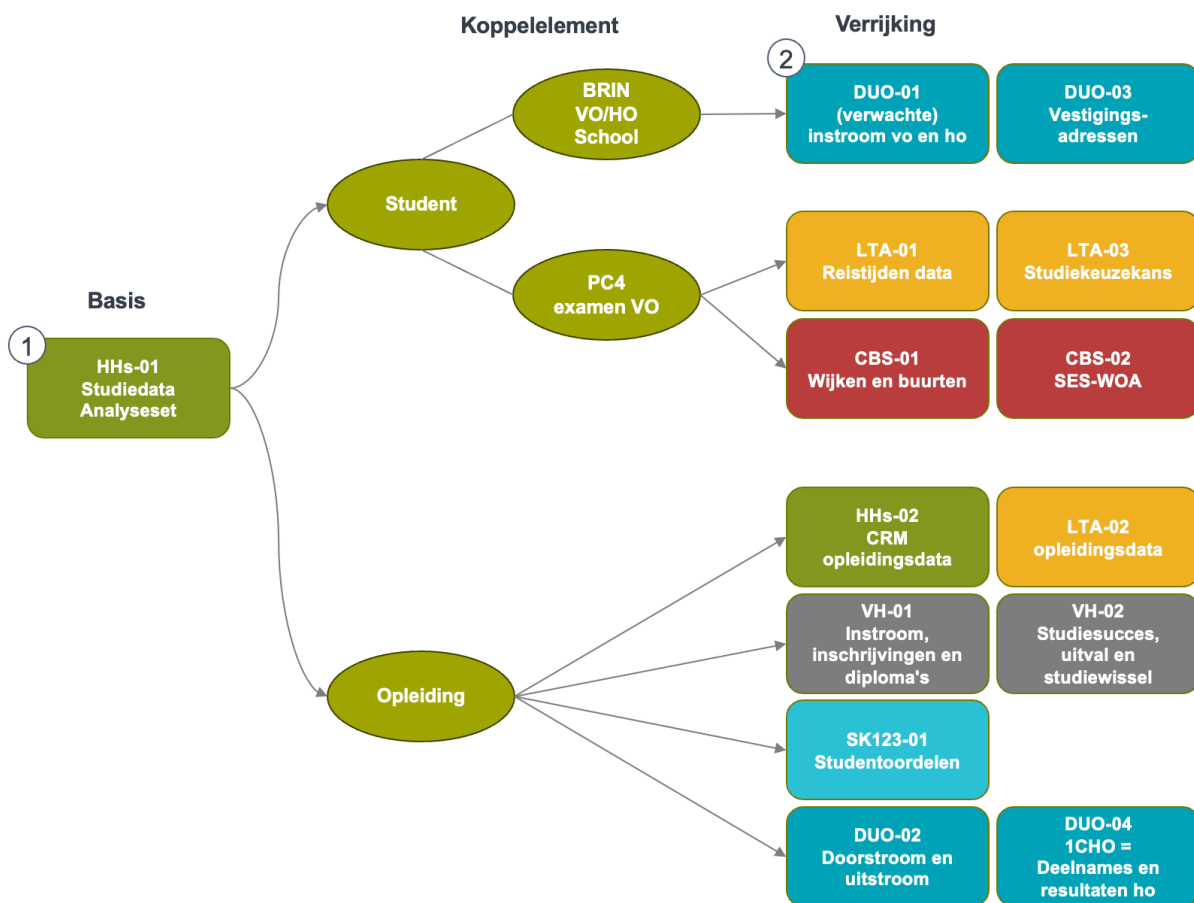
Bron	Data	Reden
Studiedata analyseset (De HHs) = HHs-01	Basisgegevens per student per studie per jaar over persoonlijke demografische kenmerken, oriëntatie, vooropleiding, aanmelding, toelating, inschrijving, studievoortgang, studieuitval of studiesucces	Onderzoek naar student journeys: gelijke kansen bij mogelijke bottlenecks (instroom, doorstroom, uitstroom)
CRM (De HHs) = HHs-02	Deelname aan oriëntatie van De HHs, conversie van oriëntatie naar aanmelding per opleiding	Onderzoek naar student journeys: gelijke kansen bij mogelijke bottlenecks (instroom)
VH = VH-01 + VH-02	Landelijk overzicht met equivalente opleidingen ten aanzien van 1) instroom, inschrijvingen en diploma's, en 2) studiesucces, uitval en studiewissel per opleiding	Benchmark onderzoek naar landelijke ontwikkelingen per opleiding om onderscheid te kunnen maken tussen trends en uitzonderingen (instroom, doorstroom, uitstroom)

Bron	Data	Reden
DUO = DUO-01 + DUO-02	Basisgegevens per school over 1) (verwachte) instroom en 2) doorstroom en uitstroom van scholen uit het voortgezet en hoger onderwijs	Onderzoek naar de mate waarin de studentpopulatie van De HHs een afspiegeling is van de regio op basis van marktaandeelen en proportionaliteit van instroom bij opleidingen van De HHs (instroom)
DUO = DUO-03	Vestigingsadressen per school	Onderzoek naar de afstand van deze scholen naar de vestigingen van De HHs (instroom)
DUO = DUO-04	1CHO - Deelnames en resultaten van studenten in het hoger onderwijs	Onderzoek naar student journeys: onderzoek naar doorstroom naar een vervolgstudie (uitstroom)
CBS = CBS-01	Kenmerken betreffende bevolking, wonen, energie, onderwijs, arbeid, inkomen, sociale zekerheid en voorzieningen	Onderzoek naar student journeys: gelijke kansen bij mogelijke bottlenecks (instroom, doorstroom, uitstroom)
CBS = CSB-02	Sociaal-economische status	Onderzoek naar student journeys: gelijke kansen bij mogelijke bottlenecks (instroom, doorstroom, uitstroom)
Studiekeuze123 = SK123-01	Studentoordelen en zogenaamde duidingsgegevens (zoals contacttijd)	Onderzoek naar student journeys: berekening van de kans op studiekeuze van een opleiding aan De HHs (instroom)
Reistijden-dataset (LTA) = LTA-01	Reistijden vanaf PC4 codes	Onderzoek naar student journeys: berekening van de kans op studiekeuze van een opleiding aan De HHs (instroom)
Opleidingen-dataset (LTA) = LTA-02	Kenmerken per opleiding per cohort	Onderzoek naar student journeys: gelijke kansen bij mogelijke bottlenecks (instroom, doorstroom, uitstroom)

Bron	Data	Reden
Studiekeuze-kans (LTA) = LTA-03	Kans voor een keuze voor een opleiding van De HHs op basis van PC4	Onderzoek naar student journeys: gelijke kansen bij mogelijke bottlenecks (instroom, doorstroom, uitstroom)

4.1.4 Schematische weergave bronnen en koppelingen

Samenvattend zijn de bronnen en koppelingen als volgt geordend:



Figuur 2: Bronnen en koppelingen

5 Proces van levering en bewerking

5.1 Aanvraag aan levering of download

- Data van De HHs wordt geleverd vanuit de betreffende teams op basis van een formele aanvraag. Dit document is bedoeld ter onderbouwing en documentatie van dit proces.
- Niet publieke data van publieke instellingen wordt separaat aangevraagd volgens de daartoe ingerichte procedures; het betreft hier vooralsnog alleen Stichting Studiekeuze123.
- Alle overige data is publiek beschikbaar en wordt via publieke toegang gedownload.

5.2 Datamanagement

- De data wordt opgeslagen op de Research Drive van de lector, conform de procedure voor [wetenschappelijk datamanagement](#) die is ingericht door de hogeschoolbibliotheek van De HHs.

5.3 Bewerking en verrijking

- De data wordt in meerdere stappen omgevormd tot onderzoeksresultaten:
 - **documentatie** van de bron conform de *Leidraad mappenstructuur en datadocumentatie* van De HHs (versie 1.0, februari 2021, Hogeschoolbibliotheek).
 - **inlezen** van de studiedata analyseset;
 - **opschonen** van de dataset: uniformeren van namen en velden waar nodig en opvullen van missende waarden waar nodig (met bijv. het gemiddelde of de mediaan van een veld⁵);
 - **opslag** naar een .rds en .fst bestandsformaat;
 - **koppelen** met aanvullende bronnen op basis van PC4, opleidingsnaam en cohort;
 - **verrijken** van de dataset: aanvullende velden ontwikkelen die nog niet beschikbaar zijn maar af te leiden, zoals ‘aantal dagen aanmelding na de aanmelddeadline’, ‘reistijd tot alternatieve studie bij een andere onderwijsinstelling’.
- Iedere stap in de bewerking wordt aanvullend gedocumenteerd en elk tussenresultaat wordt als separaat bestand opgeslagen, zodat bewerkingen naderhand goed te traceren zijn en hierover gerapporteerd kan worden in een verantwoording of publicatie.

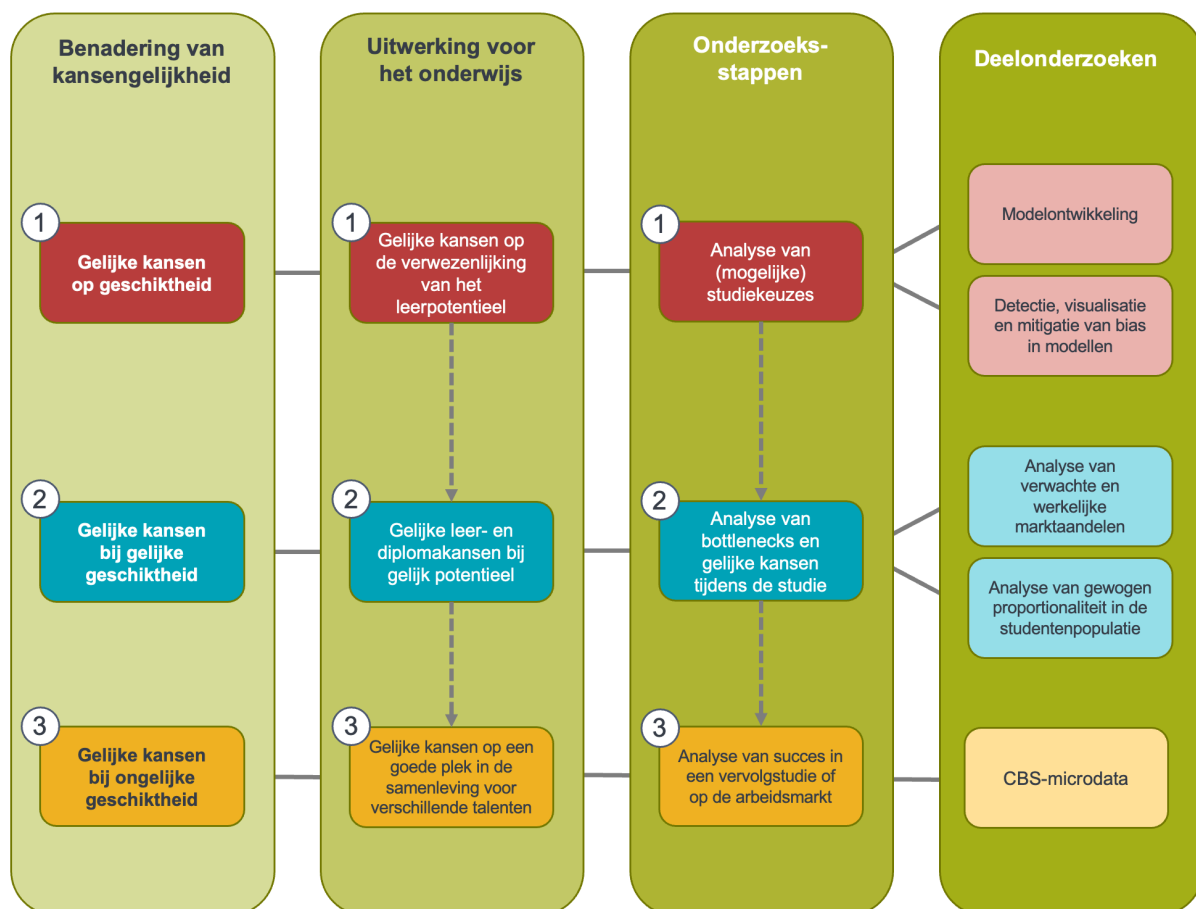
⁵ Het opvullen van missende waarden, imputatie, is nodig om statistische berekeningen te kunnen uitvoeren. De beste methode voor imputatie is op voorhand niet te zeggen; dit is afhankelijk van de variabele en de context ervan (Cox et al., 2014). Verschillende methoden zullen hiervoor onderzocht worden.

6 Methoden van analyse

In dit hoofdstuk beschrijven we per benadering van kansengelijkheid de methode van analyse.

1. Een **analyse van (mogelijke) studiekeuzes** om te bepalen welke (groepen) studenten naar verhouding meer of minder vertegenwoordigd zijn aan De HHs dan dat we op basis van de geografische ligging van de vestigingen van De HHs zouden verwachten. Hiermee krijgen we antwoord op de vraag welke bias er is in de instroom van De HHs.
2. Een **analyse van bottlenecks en gelijke kansen in de student journeys** van opleidingen aan De HHs ten opzichte van verschillende groepen studenten. Hiermee krijgen we antwoord op de vraag welke bias er is in student journeys van De HHs.
3. Een **analyse van succes in vervolgopleidingen of de arbeidsmarkt** om te bepalen welke groepen studenten meer of minder kansen hebben na hun studie. Hiermee krijgen we antwoord op de vraag welke bias er is in het succes na de studie van studenten van De HHs.

Zie voor een visuele uitwerking [Figuur 3](#).



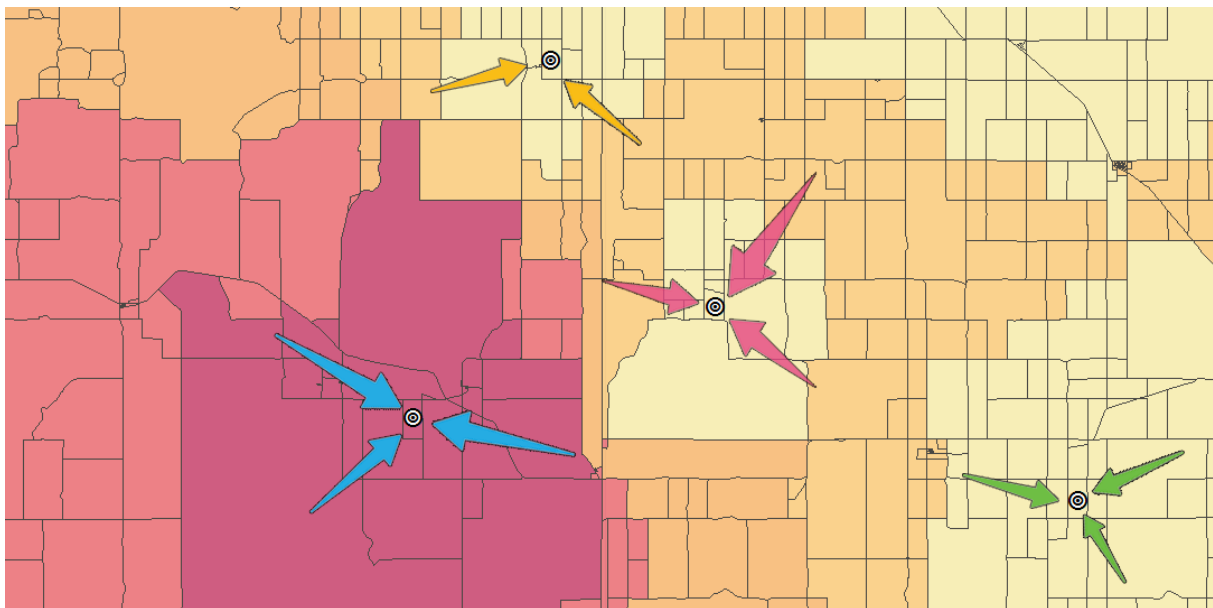
Figuur 3: Deelonderzoeken

6.1 Deelonderzoek I – Analyse van (mogelijke) studiekeuzes

Dit deelonderzoek geeft antwoord op de onderzoeksvragen uit het eerste deel van de operationalisering van [Gelijke kansen op de verwezenlijking van het leerpotentieel c.q. gelijke kansen op instroom](#).

6.1.1 Analyse van verwachte en werkelijke marktaandeelen

De eerste stap van de deelonderzoek bestaat uit het berekenen van studiekeuzekansen van de instroom van De HHs ten opzichte van het spreidingsgebied c.q. de regio. Hiervoor gebruiken we een [Huff Gravity model](#) (Huff, 1963) met behulp van het [REAT](#) package en QGis. Dit type model wordt gebruikt in de retail bij het berekenen van de kans dat klanten kiezen voor een winkel. Naarmate een winkel verder weg ligt, neemt de bereidheid van de klant om ernaartoe te reizen af (afstandsverval - *distance decay*). Er zijn echter factoren die dit verval dempen: goede bereikbaarheid met het OV of eigen vervoer (*time-space compression*) en de aantrekkelijkheid van de winkel, bijvoorbeeld de omvang of aanbod van de winkel. Zie ter illustratie [Figuur 4](#).



Figuur 4: Voorbeeld van een Huff Gravity Model (bron: GISGeography)

Dit concept en model passen we toe op de keuze van studenten voor een van de onderwijsinstelling in de regio Den Haag.

- De aanname is dat ook studenten voor een onderwijsinstelling kiezen waar de minste belemmeringen zijn om hun gewenste studie te volgen.
- Aspecten die bevorderend zijn voor het wegnemen van belemmeringen zijn een beperkte reistijd, met als dempende factoren een goede reputatie en beperkte toegangseisen.

- Omdat het OV voor studenten gratis is, worden de kosten van het OV niet meegenomen.
- Verder nemen we aan dat de meeste studenten in hun eerste studiejaar nog thuis wonen of dat het in alle steden even moeilijk is een kamer te vinden, waardoor beschikbare huisvesting als factor buiten beschouwing gelaten kan worden.

Voor deze analyse worden vier datasets bewerkt of ontwikkeld en geanalyseerd:

- De historische **uitstroomdata** van scholen in het voortgezet onderwijs (havo en mbo) binnen een straal van 90 minuten reistijd in de spits op een maandag van elke vestiging van De HHs (bronnen: DUO-02 + DUO-03 + LTA-01). De aanname is dat studenten maximaal bereid zijn om 90 minuten te reizen naar De HHs. Vervolgens wordt berekend welk marktaandeel De HHs heeft in deze uitstroom.
- Voor elke opleiding wordt vervolgens bepaald welke **alternatieve, gelijke opleidingen** er aangeboden worden bij andere HBO-instellingen in een straal van 2x 90 minuten. Dit zijn in theorie de equivalente opleidingen die een student kan kiezen in plaats van een opleiding aan De HHs; deze opleidingen bevinden zich geografisch in Noord-Holland, Zuid-Holland, Zeeland, Utrecht en Brabant. Van elk van deze opleidingen worden de NSE scores verzameld over de periode 2010 t/m 2023 (bron: SK123-01).
- Voor elke inschrijving worden de **reistijden** berekend in het eerste jaar aan de opleiding aan De HHs en equivalente opleidingen vanaf de PC4 ten tijde van het behalen van de vooropleiding (bronnen: LTA-01 en LTA-02).
- Vervolgens wordt een **studiekeuzekans** (een gewogen kans op aanmelding) berekend op basis van de reistijden, de reputatie en de ingangseisen per opleiding en opgeslagen (bron: LTA-03). De weging van de onderdelen hierin moet nog nader bepaald worden.

6.1.2 Analyse van gewogen proportionaliteit in de studentenpopulatie

Nadat we de verwachte en werkelijke marktaandelen hebben berekend per buurt, kunnen we de mogelijke **instroom-bias** berekenen op basis van een weging per buurt gebaseerd op de kansberekening vanuit de eerdere Huff-analyse. De bias werkt naar twee kanten toe: studenten die uit gebieden komen met een lagere kans tellen zwaarder mee, terwijl studenten die uit gebieden komen met een hogere kans minder zwaar meetellen. Door de verdeling van de kansen te normaliseren, kunnen we de weging projecteren rondom het getal 1, waarmee we de totale populatie gelijk houden. Op deze manier hebben we een **propensity score**⁶ ontwikkeld voor instroom aan De HHs en **balanceren** we de totale studentenpopulatie.

Vervolgens **vergelijken** we **de gewogen populatie met de ongewogen populatie** op achtergrondkenmerken, zoals vooropleiding, leeftijd, geslacht en sociaal-economische status, waarmee we de **achtergrond-**

⁶ Een *propensity score* (neigingsscore) is een score die uitdrukt welke kans een item heeft om in een bepaalde categorie te vallen. In deze analyse is het de neiging om de studeren bij een van de vestigingen van De HHs.

bias kennen voor de totale populatie en beter zicht hebben op minderheden en meerderheden.

We voeren deze analyse uit per vestiging en per opleiding op basis van 1 nader te bepalen peiljaar.

6.2 Deelonderzoek II - Analyse van bottlenecks en gelijke kansen tijdens de studie

Dit deelonderzoek geeft antwoord op de onderzoeksvragen uit het tweede deel van de operationalisering van [Gelijke leer- en diplomakansen bij gelijk potentieel c.q. gelijke kansen op succesvolle doorstroom en uitstroom](#) in andere woorden: gelijke kansen in de student journey. De data worden volgens de 'Cross-industry Standard Process for Data Mining' (CRISP-DM) methodiek geanalyseerd ([Chapman et al., 2000](#)), verbijzonderd volgens de methode 'Exploratory Model Analysis' ([Biecek & Burzykowski, 2021](#)) met het bijbehorende DALEX package⁷ en het `fairmodel` package. Deze packages maken deel uit van [DrWhy.AI](#), een ecosysteem van packages c.q. analyses voor explainable en fair AI.

Eerst worden een of meerdere algemene modellen ontwikkeld, waarna deze op rechtvaardigheid (fairness) gecontroleerd en gecorrigeerd worden. Deze aanpak geeft zowel zicht op de mate waarin verschillende groepen gelijke kansen of een gebrek daaraan hebben in De Haagse Hogeschool, als aan welke factoren 'gedraaid' moet worden om ongelijke kansen tegen te gaan.

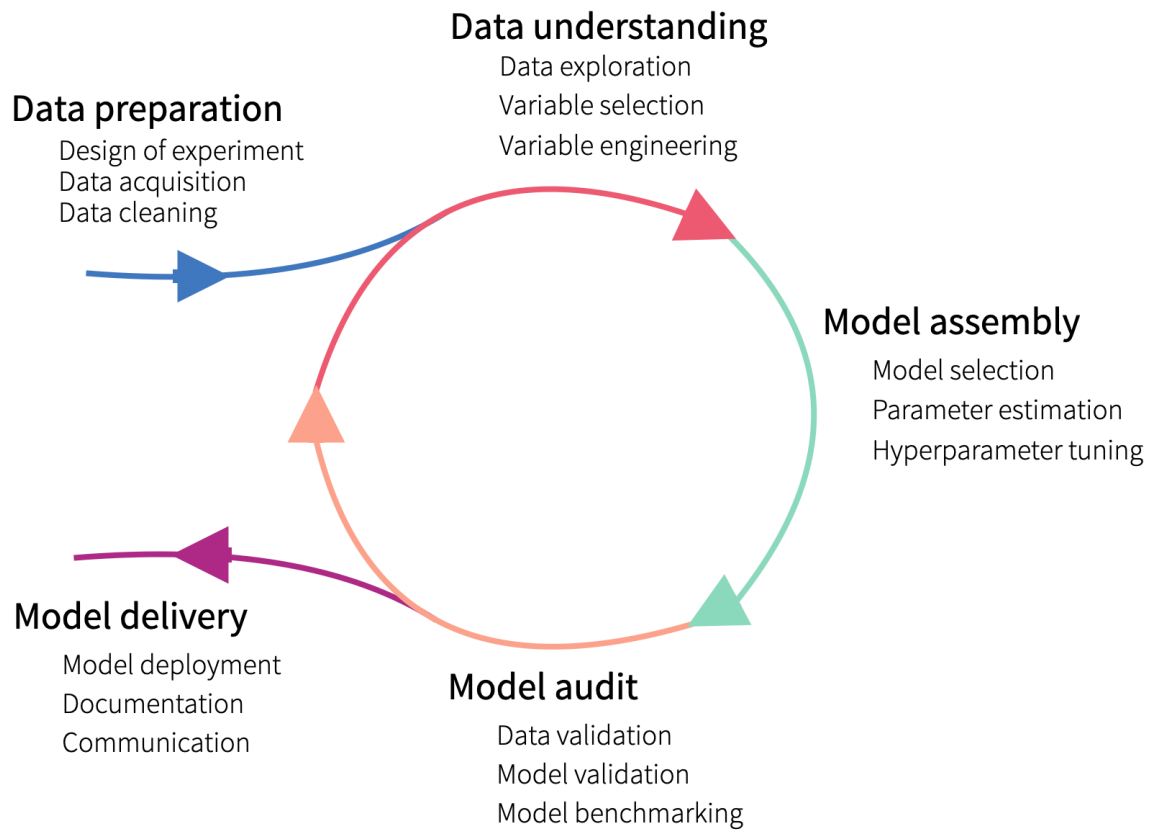
6.2.1 Modelontwikkeling

Bij de ontwikkeling van een voorspelmodel zijn drie voorwaarden van belang:

1. **Validatie van de voorspelling.** Voor elke voorspelling van een model moet we na kunnen gaan worden hoe sterk het bewijsmateriaal is dat de voorspelling ondersteunt.
2. **Verantwoording van de voorspelling.** Voor elke voorspelling van een model moeten we kunnen uitleggen welke variabelen de voorspelling beïnvloeden en in welke mate.
3. **Speculatie van de voorspelling.** Voor elke voorspelling van een model moeten we kunnen uitleggen hoe de voorspelling zou veranderen als de waarden van de variabelen in het model zouden veranderen.

De ontwikkeling is een iteratief proces dat op basis van de uitkomsten in de verschillende stappen meerdere keren wordt herhaald naarmate het begrip over de data groeit (zie [Figuur 5](#)). Het DALEX package helpt om modellen transparant te ontwikkelen, waardoor de interpreteerbaarheid van de uitkomsten zo optimaal mogelijk is.

⁷ Zie voor een populaire beschrijving van deze methode: *The Hitchhiker's Guide to Responsible Machine Learning The R version* ([Kozak et al., 2022](#)).



Figuur 5: CRISP-DM methodiek; visualisatie van Biecek en Burzykowski (2021)

Tabel 6: Methoden van analyse

Stappen	Onderdelen	Toelichting
1. Data preparatie	a. Onderzoeksontwerp	Dit document en in het bijzonder dit hoofdstuk.
	b. Data verzameling	De verzameling en bewerking van de data zoals beschreven in Databronnen .
	c. Data cleaning	
2. Data begrip	a. Data exploratie	Verkenning van basiskenmerken op basis van frequentietellingen in tabelvorm, grafieken of kaarten.

Stappen	Onderdelen	Toelichting
		<p>Berekening en visualisatie van statistische verbanden, rekening houdend met de statistische samenstelling van de populatie (verdelingen) en bijbehorende assumpties, en correcties voor meervoudige testen.</p> <p>Onderzoek van missende waarden, uitbijters (<i>outliers</i>) en disbalansen in verdelingen (scheefheid of disproportionaliteit).</p>
	b. Variabele selectie	Keuze voor variabelen waarin de verschillen tussen groepen significant is of van algemeen belang zijn vanuit de literatuur.
	c. Variabele engineering	<p>De volgende statische bewerkingen vinden plaats:</p> <ul style="list-style-type: none"> • Opvulling van missende waarden • Verwijdering van uitbijters • Correctie van scheefheid of disproportionaliteit met behulp van transformaties
3. Model assembly	<p>a. Model selectie</p> <p>b. Parameterestimatie</p> <p>c. Hyper-parameter tuning</p>	<p>Ontwikkeling van explainable Machine Learning modellen voor het voorspellen van transitie in de studie: instroom c.q. selectie, doorstroom in de studie langs verschillende punten (zoals het BSA, de propedeuse, aanvullende onderdelen van curricula waarvoor ingangseisen gelden zoals een stage, exchange met het buitenland, afstudeerproject) en uitval of diploma.</p> <p>Toetsing van een scala aan modellen afhankelijk van de uitkomstvariabele (binair – bijv. wel of niet geslaagd – of continu – bijv. het aantal EC of de hoogte van cijfers).</p> <p>Mogelijke modellen zijn onder meer lineaire regressie (GAM/GLM), classification and regression trees (CART), random forest (RF), stochastic gradient boosting, bagged CART.</p> <p>Afhankelijk van het soort model is extra tuning van de instellingen (hyperparameters) nodig (bijv. bij RF). Deze kunnen automatisch gezocht en geselecteerd worden.</p>

Stappen	Onderdelen	Toelichting
4. Model audit	a. Data validatie b. Model validatie c. Model benchmarking	Onderzoek van het model op: <ul style="list-style-type: none"> • De verschillende voorspellende onderdelen • Een sensitiviteitsanalyse voor de voorspellingen • De voorspelkracht: de uiteindelijke selectie wordt gebaseerd op vergelijking van de performance van deze modellen, gebaseerd op standaard performance maten zoals contingency tables voor classificaties (bijv. man / vrouw versus diploma ja / nee) of de Receiver Operating Characteristic (ROC) en de Area Under the Curve (AUC) voor analyses op rangordes (kansen op basis van meerdere variabelen). • De variable importance per verklarende variabele; dit geeft inzicht in verschillen tussen opleidingen en groepen studenten, die inzicht geven in de achterliggende mechanismen (Shmueli, 2010). • Het effect van de verklarende variabelen op de voorspellingen; hiervoor worden onder meer what-if scenario's via Ceteris Paribus analyses gemaakt, welke kunnen aangeven hoe en voorspelling verandert als een kenmerk net wat anders zou zijn, bijv. de leeftijd van een student • Eventuele residuen (niet verklaarde effecten)
5. Model oplevering	a. Model uitrol b. Documentatie c. Communicatie	De uitkomsten kunnen gedeeld worden via de meegeleverde publicatietools voor deze modellen (modelstudio en arena). Het gebruik van deze tools zal bepaald worden als het onderzoek is aangekomen bij deze fase en is afhankelijk van de mate van aggregatie op de data.

6.2.2 Detectie, visualisatie en mitigatie van bias in modellen

Nadat een of meer modellen volgens de bovenstaande methode in de eerste fase zijn ontwikkeld, onderzoeken we in de tweede fase de mogelijke bias in deze modellen met behulp van het [fairmodels](#) package ([Wiśniewski & Biecek, 2022](#)). We richten ons op binaire classificaties, bijv.

‘een student heeft de BSA wel of niet behaald’. Hier zijn vier hoofdvragen ([Wiśniewski & Biecek, 2022](#)): Hoe meten we bias? Hoe ontdekken we bias? Hoe visualiseren we bias? Hoe gaan we bias tegen?

Het vaststellen van criteria voor bias

Er zijn drie criteria voor het meten van bias c.q. eerlijkheid: onafhankelijkheid, separatie en toereikendheid ([Barocas et al., 2019](#)). Ter illustratie lichten we deze begrippen toe aan de hand van het behalen van een BSA van 50 EC voor mannen en vrouwen.

1. **Onafhankelijkheid** – Het onafhankelijkheids criterium houdt in dat de kans op het behalen van 50 EC in jaar 1 (de BSA grens) even groot moet zijn tussen mannen en vrouwen. Dit is eerlijk vanuit het maatschappelijke perspectief.
2. **Separatie** – Bij de toewijzing van een negatief BSA kan er sprake zijn van een foutieve toewijzing: het kan voorkomen dat een vrouwelijke student een negatief BSA haalt vanwege een nog niet nagekeken tentamen. Het separatiecriterium houdt in dat de verhouding terecht versus onterecht een negatief BSA gelijk moet zijn in de subgroepen (mannen en vrouwen). Dit is eerlijk vanuit het student perspectief.
3. **Toereikendheid** – Het toereikendheids criterium houdt in dat de verhouding terecht of onterecht een negatief of positief BSA even groot moet zijn in subgroepen (mannen en vrouwen). Dit is eerlijk vanuit het perspectief van de organisatie (in dit voorbeeld de onderwijsinstelling).

Een probleem en uitdaging is dat geen van deze criteria tegelijkertijd volledig op kan gaan ([Barocas et al., 2019](#)). Dit vraagt om het balanceren van de verschillende criteria of aan een van de criteria de voorkeur te geven.

Het ontdekken en visualiseren van bias met een cut-off criterium en fairness checks

Het is gangbaar in Machine Learning om bij het balanceren uit te gaan van de 4/5 regel⁸: er is sprake van discriminatie als de selectie van leden van een minderheidsgroep lager is dan 80% van leden van de meest gekozen groep ([Code of Federal Regulations. Section 4d, uniform guidelines on employee selection procedures \(1978\), 1978](#)). Het `fairmodels` package hanteert daarom een marge van 0,8. Dit is een analyse van ongelijk effect (*disparate impact analysis*).

Ter verduidelijking:

Stel bij een selectieve opleiding worden 96 van 120 havo kandidaten geselecteerd en 30 van 50 mbo kandidaten, dan is er volgens deze regel sprake van een ongewenst negatief effect.

De selectie van havo studenten is 80% (96/120); de selectie van mbo studenten is 60% (60/100); het minimum voor een eerlijke ratio ten opzicht van mbo studenten is 64% (80% van 80%); 60%

⁸ Voor zover bekend is er geen Europees of Nederlands equivalent van dit criterium. Vandaar dat we in dit onderzoek aansluiten bij deze regel.

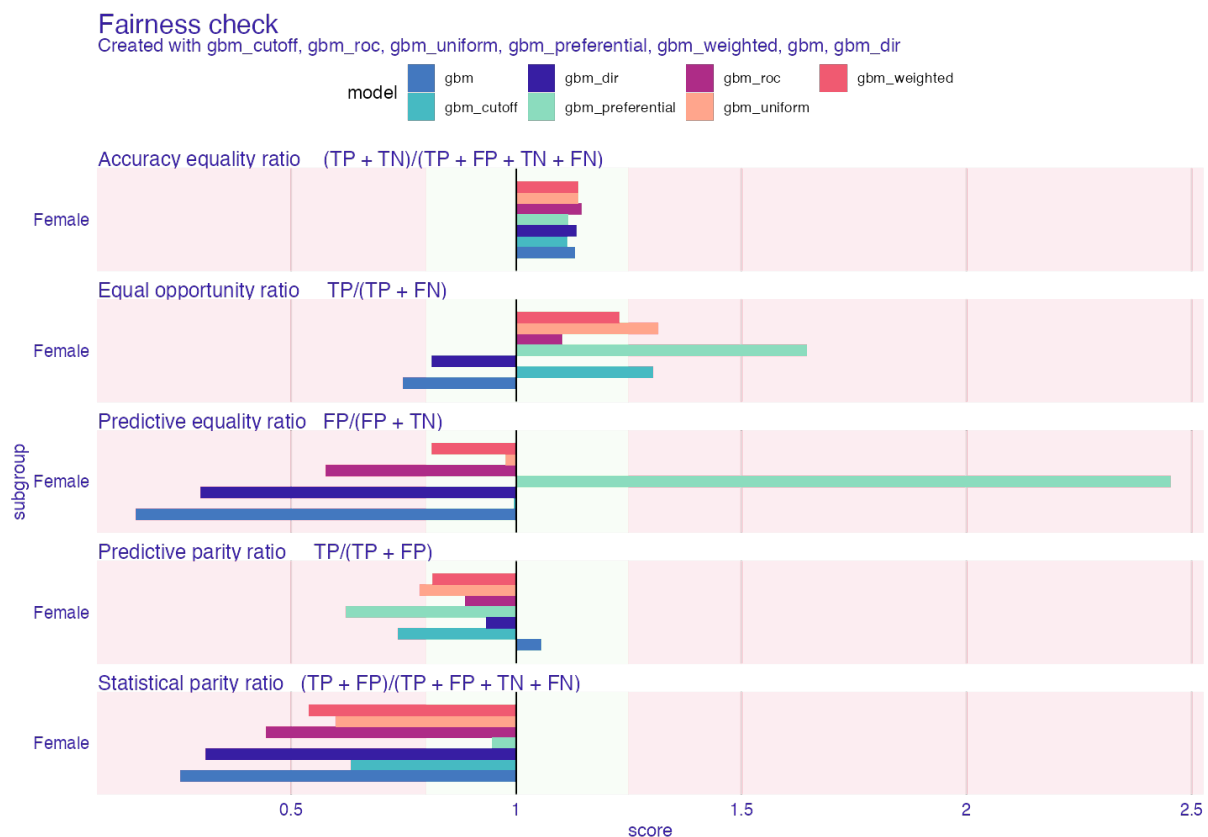
is lager dan 64% en daarom is er sprake van een negatief effect volgens de 4/5 regel.

Om de drie criteria te toetsen worden 5 *fairness checks* uitgevoerd (zie [Figuur 6](#)), die worden afgebeeld in een *fairness plot* (zie [Figuur 7](#)):

1. **Accuracy equality ratio (ACC)** = $(TP + TN)/(TP + TN + FP + FN)$ = accuraatheid: de verhouding van de correcte voorspellingen (zowel positief als negatief) ten opzichte van alle voorspellingen - het aantal studenten dat correct een positief of negatief BSA heeft ontvangen ten opzichte van het totaal aantal studenten.
2. **Equal opportunity ratio (EO)** = $TP/(TP + FN)$ = sensitiviteit: de verhouding van de correct positieve voorspellingen ten opzichte van de vals negatieve voorspellingen - het aantal studenten dat terecht een positief BSA heeft ontvangen ten opzichte van het aantal studenten dat ten onrechte een negatief BSA heeft ontvangen.
3. **Predictive opportunity ratio (PO)** = $FN/(FN + TN)$: de verhouding van de fals negatieve voorspellingen ten opzichte van alle negatieve voorspellingen - het aantal studenten dat ten onrecht een negatief BSA heeft ontvangen ten opzichte van het totaal aantal studenten dat een negatief BSA heeft ontvangen.
4. **Predictive parity ratio (PPV)** = $TP/(TP + FP)$ = precisie: de verhouding van de correct positieve voorspellingen ten opzichte van alle positieve voorspellingen - het aantal studenten dat terecht een positief BSA heeft ontvangen ten opzichte van het totaal aantal studenten dat een positief BSA heeft ontvangen.
5. **Statistical parity ratio (STP)** = $(TP + FN)/(TP + TN + FP + FN)$: de verhouding van de correct positieve voorspellingen en de vals negatieve voorspellingen ten opzichte van alle voorspellingen - het aantal studenten dat terecht een positief BSA heeft ontvangen of ten onrechte een negatief BSA ten opzichte van het totaal aantal studenten.

		Predicted Class							
		Positief	Negatief						
Actual Class	Positief	True Positive (TP)	False Negative (FN) Type II Error	2	Sensitivity Equal opportunity ratio (EO) $\frac{TP}{(TP + FN)}$				
	Negatief	False Positive (FP) Type I Error	True Negative (TN)		Specificity $\frac{TN}{(TN + FP)}$				
		4	Precision Predictive parity ratio (PPV) $\frac{TP}{(TP + FP)}$		Negative Predicted Value $\frac{TN}{(TN + FN)}$	1	Accuracy Accuracy equality ratio (ACC) $\frac{TP + TN}{(TP + TN + FP + FN)}$	5	Statistical parity ratio (STP) $\frac{TP + FN}{(TP + TN + FP + FN)}$
				3	Predictive opportunity ratio (PO) $\frac{FN}{(FN + TN)}$				

Figuur 6: Confusion matrix met fairness checks



Figuur 7: Visualisatie van fairness checks voor het kenmerk geslacht (Bron: fairmodels - dataset: German credit data)

Het tegengaan van bias

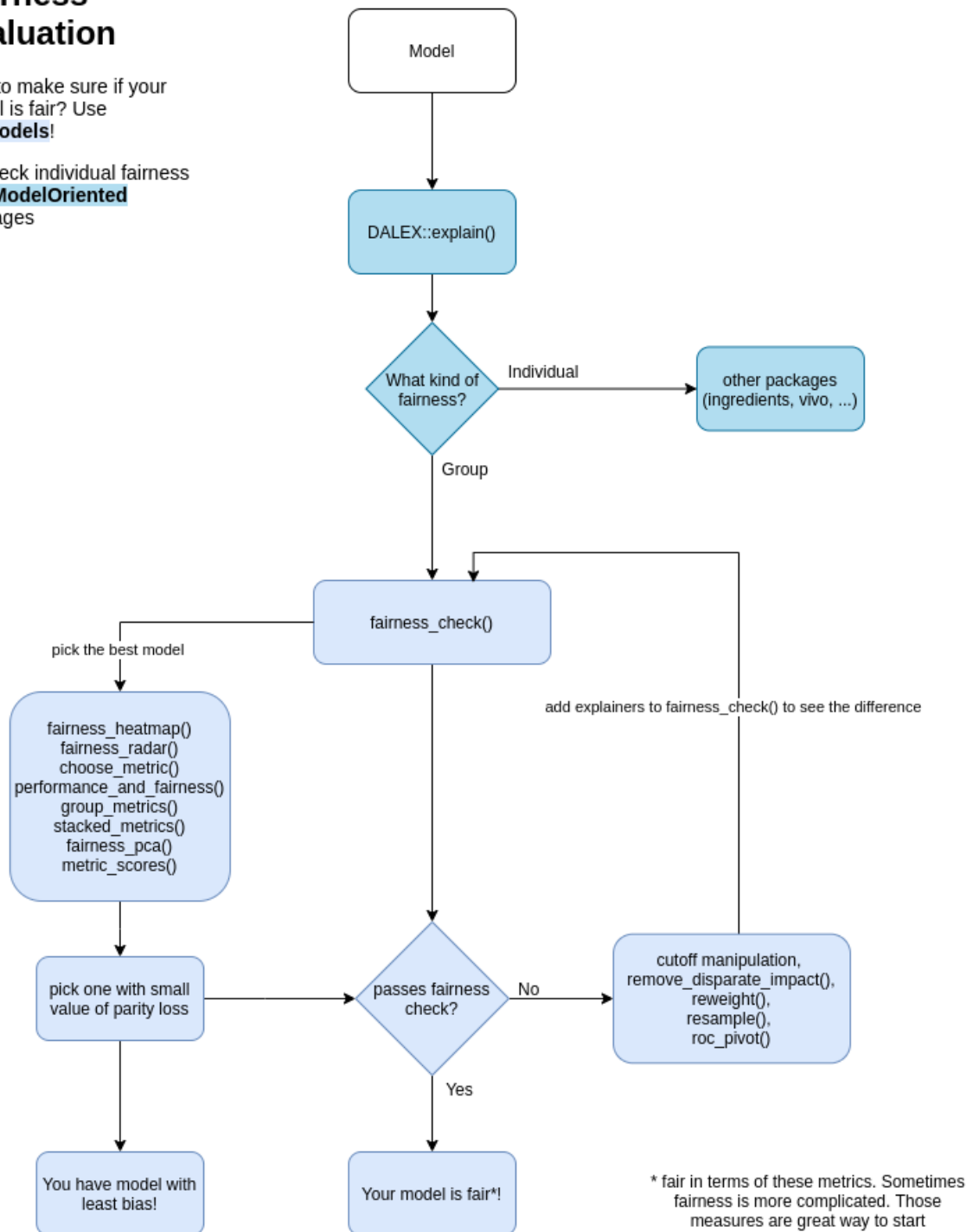
De modellen die ontwikkeld zijn in de eerste fase met het **DALEX** package, kunnen in de tweede fase getoetst worden op fairness voor statistische minderheden en – waar nodig – gecorrigeerd kunnen worden aan de hand van de bovenstaande ratio's. Hiervoor zijn methoden a) voorafgaand aan het modelleren (*pre-processing*): het toepassen van transformatie van een distributie (disparate impact remover), gewichten (reweighting) of hersamplen van de data, of b) na afloop van het modelleren (*post-processing*): reject option based classification pivot (roc_pivot) of het aanpassen van cutoffs per groep (cutoff manipulation). De visualisatie van de uitkomsten vindt plaats met behulp van het **arenar** package ([demo](#)) en het **modelstudio** package ([demo](#)).

Zie [Figuur 8](#) voor de complete onderzoeksmethode.

Fairness evaluation

How to make sure if your model is fair? Use **fairmodels!**

Or check individual fairness with **ModelOriented** packages



Figuur 8: Groepsgewijze evaluatie van fairness in het modelleren

6.3 Deelonderzoek III - Analyse van succes in een vervolgstudie of op de arbeidsmarkt

Dit deelonderzoek geeft antwoord op de onderzoeksvragen uit het derde deel van de operationalisering van [Gelijke kansen op een goede plek in de samenleving voor verschillende talenten c.q. gelijke kansen op een vervolgstudie of positie op de arbeidsmarkt](#).

Hiervoor zullen we een onderzoek uitvoeren met behulp van CBS-microdata. De invulling van dit deel van het onderzoek dient nog nader bepaald te worden. Op hoofdlijnen houdt het onderzoek in dat we via het CBS onderzoeken wat het succes is studenten die aan De HHs hebben gestudeerd en zijn uitgestroomd met of zonder diploma:

1. In een **vervolgstudie**: Naar welke vervolgstudie zijn zij na het verlaten van De HHs doorgestroomd? Op welk niveau was deze studie? Hebben zij in deze vervolgstudie een diploma behaald? In hoeveel tijd hebben zij dat diploma behaald?
2. Op de **arbeidsmarkt**: Hebben zij na het verlaten van De HHs binnen drie, acht en twaalf jaar betaald werk gevonden? In welk werkveld hebben zij dit gevonden? Op welk opleidingsniveau en met welk inkomen?

In dit onderzoek sluiten we aan bij de gangbare onderzoeksmethoden van nog nader te verkennen arbeidsmarktonderzoek, zoals het CBS en het Researchcentrum voor Onderwijs en Arbeidsmarkt (ROA).

Voor dit onderzoek zal een separaat onderzoeksplan uitgewerkt worden.

7 Verwachte resultaten

De resultaten van het onderzoek zullen op de volgende manieren beschikbaar worden gesteld:

- **Beroepsproducten:** Gerichte adviezen aan het management van opleidingen, onderzoeksrapporten op deelgroepen, methodes voor het onderzoeken van gelijke kansen in studiedata. Borging van de toepassing van de inzichten via het Inclusion Office van De HHs, de uitvoeringsagenda van het Instellingsplan en aansluiting bij de Gelijke Kansen Alliantie in de regio Den Haag. Aansluiting bij landelijke agenda's via de SURF SIG Learning Analytics, de Informatiehub Studiedata & AI van Npuls, SURF Studiedata en de Nederlandse AI Coalitie voor Onderwijs.
- **Intreerede:** Het onderzoek zal de basis vormen voor de intreerede van de lector (verwacht april 2024).
- **Publicaties:** “Gelijke kansen in het Haags hoger onderwijs” (mogelijk uitgesplitst naar instroom, doorstroom en uitstroom); mogelijk afgeleide artikelen over subthema's, zoals “Invloed van reistijden op gelijke kansen in het hoger onderwijs”, “Studievoortgang en -succes van studenten met een Caribische vooropleiding” zowel in journals zoals Higher Education als populaire wetenschappelijke bladen of magazines voor doelgroepen zoals Thema voor managers in het hoger onderwijs.
- **Presentaties:** Op conferenties zoals de SURF Onderwijsdagen, de DAIR, de HO-link, De HHs AI Fest, etc.
- **Blogposts:** Op de site van De HHs en van SURF.
- **Broncode:** De broncode van het onderzoek wordt via GitHub onder CC-licentie ter beschikking gesteld om soortgelijk onderzoek te kunnen uitvoeren: [Naamsvermelding-NietCommercieel-GelijkDelen 4.0 Internationaal \(CC BY-NC-SA 4.0\)](#).

8 Reproduceerbaarheid

Het onderzoek volgt de [FAIR principes](#) voor reproduceerbaarheid van onderzoeksresultaten: **F**indability (vindbaarheid), **A**ccessability (toegankelijkheid), **I**nteroperability (uitwisselbaarheid), en **R**euse (herbruikbaarheid) van digitale assets. We betrekken dit op broncode, metadata, data en beroepsproducten.

Tabel 7: Toepassing van de FAIR principes

	Broncode	Metadata	Data	Beroepsproducten
Findability	Wordt publiek gepubliceerd op github	Wordt publiek gepubliceerd op github	De toegang wordt beschreven en gepubliceerd op github ; bronnen die niet te herleiden zijn tot De HHs of studenten van De HHs worden gepubliceerd	Wordt publiek in de vorm van artikelen; interne adviesrapporten zullen afhankelijke van de mate van herleidbaarheid tot De HHs wel of niet publiek gemaakt worden.
Accessibility	Is vrij toegankelijk	Is vrij toegankelijk	Is op verzoek in te zien	Is afhankelijk van de mate van publieke toegankelijkheid
Interoperability	Is te lezen met behulp van open software (R en RStudio)	Is te lezen met behulp van open software (R en RStudio)	Is te lezen met behulp van open software (R en RStudio)	Is te lezen met een PDF lezer of via vakjournals, presentaties, blogposts

	Broncode	Metadata	Data	Beroepsproducten
Reuse	Kan naar eigen inzicht onder CC BY-NC-SA 4.0 licentie hergebruikt worden	Kan naar eigen inzicht onder CC BY-NC-SA 4.0 licentie hergebruikt worden	Kan voor zover publiek naar eigen inzicht onder CC BY-NC-SA 4.0 licentie hergebruikt worden	Kan voor zover publiek naar eigen inzicht onder CC BY-NC-SA 4.0 licentie hergebruikt worden

Referenties

- Badou, M., & Day, M. (2021). *Kansengelijkheid in het onderwijs. Verkennend onderzoek naar factoren die samenhangen met onderwijs(on)gelijkheid* (p. 34). Gelijke kansen alliantie. <https://www.gelijke-kansen.nl/documenten/publicaties/2021/10/05/verkennend-onderzoek-naar-factoren-die-samenhangen-met-kansengelijkheid>
- Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness in Machine Learning Limitations and Opportunities*. fairmlbook.org. <http://www.fairmlbook.org>
- Biecek, P., & Burzykowski, T. (2021). *Explanatory Model Analysis*. Chapman; Hall/CRC, New York. <https://pbiecek.github.io/ema/>
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T. P., Shearer, C., & Wirth, R. (2000). *CRISP-DM 1.0: Step-by-step data mining guide*.
- Code of Federal Regulations. Section 4d, uniform guidelines on employee selection procedures (1978)*. (1978). <https://www.govinfo.gov/content/pkg/CFR-2011-title29-vol4/xml/CFR-2011-title29-vol4-part1607.xml>
- Copier, J. (2022). *Tussen idealen en dwalingen. Verhalen over onderwijs*. Garant.
- Cox, B. E., McIntosh, K., Reason, R. D., & Terenzini, P. T. (2014). Working with Missing Data in Higher Education Research: A Primer and Real-World Example. *The Review of Higher Education*, 37(3), 377–402. <https://doi.org/10.1353/rhe.2014.0026>
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2011). Fairness Through Awareness. *arXiv*. <https://doi.org/10.48550/arxiv.1104.3913>
- Elffers, L. (2022). *Onderwijs maakt het verschil - kansengelijkheid in het Nederlandse onderwijs*. Walburg Pers B.V.
- Espinoza, O. (2007). Solving the equity–equality conceptual dilemma: a new model for analysis of the educational process. *Educational Research*, 49(4). <https://doi.org/10.1080/00131880701717198>
- Fish, B., & Stark, L. (2022). It's Not Fairness, and It's Not Fair: The Failure of Distributional Equality and the Promise of Relational Equality in Complete-Information Hiring Games. *Equity and Access in Algorithms, Mechanisms, and Optimization*. <https://doi.org/10.1145/3551624.3555296>
- Hardt, M., Price, E., & Srebro, N. (2016). Equality of Opportunity in Supervised Learning. *arXiv*. <https://doi.org/10.48550/arxiv.1610.02413>
- Huff, D. L. (1963). A Probabilistic Analysis of Shopping Center Trade Areas. *Land Economics*, 39(1), 81. <https://doi.org/10.2307/3144521>
- Kozak, A., Biecek, P., & Zawada, A. (2022). *The Hitchhiker's Guide to Responsible Machine Learning The R version*. Scientific Foundation SmarterPoland.pl.
- Pedreshi, D., Ruggieri, S., & Turini, F. (2008). Discrimination-aware data mining. *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. <https://doi.org/10.1145/1401890.1401959>
- Shmueli, G. (2010). To Explain or to Predict? *Statistical Science*, 25(3), 289–310. <https://doi.org/10.1214/10-sts330>

Wiśniewski, J., & Biecek, P. (2022). fairmodels: a Flexible Tool for Bias Detection, Visualization, and Mitigation in Binary Classification Models. *The R Journal*, 14(1), 227–243. <https://doi.org/10.32614/rj-2022-019>

Versiegeschiedenis

- 10-05-2023: versie 0.9.0 - eerste conceptversie
- 18-05-2023: versie 0.9.1 - tweede conceptversie; aanvulling met bronnen, levering en bewerking, methoden van onderzoek, verwachte resultaten
- 24-05-2023: versie 0.9.1 - derde conceptversie; uitbreiding methoden van onderzoek, ethische principes, reproduceerbaarheid
- 25-05-2023: versie 0.9.2 - aanscherping Inleiding
- 27-05-2023: versie 0.9.3 - aanvulling met CBS-microdata onderzoek + extra afbeelding ter verduidelijking van de deelonderzoeken
- 30-05-2023: versie 0.9.4 - aanscherping terminologie
- 31-05-2023: versie 0.9.5 - toevoeging tools voor visualisatie in deelonderzoek II
- 01-06-2023: versie 0.9.6 - aanscherping duiding gelijke kansen op basis van Espinoza (2007)
- 02-06-2023: versie 0.9.7 - aanvulling met bijlage 1 - Data Science Ethics Checklist - versie voor de Ethische Adviescommissie van De HHs
- 03-06-2023: versie 0.9.8 - aanvulling van achtergrondinformatie bij het DALEX package
- 06-06-2023: versie 0.9.9 - aanpassing op foutieve versiegeschiedenis
- 07-06-2023: versie 0.9.9.1 - verbetering typo's
- 07-06-2023: versie 0.9.9.2 - aanvullingen op de implementatie van de resultaten
- 08-06-2023: versie 0.9.9.3 - toevoeging lijst van variabelen
- 30-06-2023: versie 1.0 - verwerking van de feedback van de Ethische Adviescommissie van De HHs
- 05-07-2023: versie 1.0.1 - verwerking van feedback van team IR&A, OKC
- 18-03-2023: versie 1.0.2 - wijziging naamgeving project
- 04-07-2024: versie 1.0.3 - verbetering operationalisatie

Repository

De broncode voor dit document kan bewerkt worden via [GitHub](#). Leden van de kenniskring van het lectoraat kunnen op verzoek toegang krijgen en meewerken aan de inhoud.

Bijlage 1 - Data Science Ethics Checklist

ethics checklist deon

8.1 A. Dataverzameling

Aandachtspunt	Overwegingen
A.1 Informed consent	<p>Als er menselijke proefpersonen zijn, hebben deze dan geïnformeerde toestemming gegeven, waarbij de proefpersonen er zelf voor kiezen en een duidelijk begrip hebben van het gegevensgebruik waarvoor ze toestemming geven?</p> <p>Er is in dit onderzoek geen sprake van menselijke proefpersonen. De studiedata dataset wordt geanonimiseerd geleverd door OKC aan het lectoraat. Door OKC verwijderde elementen zijn: studentnummer, geboortedatum, voornamen en achternaam.</p> <p>Formeel valt de dataset vanwege de anonimisering niet onder de AVG omdat anonieme data geen persoonsgegevens zijn. Op deze dataset is daarom geen consent van toepassing.</p> <p>Als grondslag voor het onderzoek geldt het gerechtvaardigd belang van de onderzoeker en De HHs.</p>
A.2 Bias in data verzameling	<p>Hebben we bronnen van vooringenomenheid overwogen die geïntroduceerd zouden kunnen worden tijdens het verzamelen van gegevens en het ontwerpen van enquêtes en hebben we stappen ondernomen om deze te beperken?</p>

Aandachtspunt	Overwegingen
<p>A.3 Beperkte onthulling van persoonlijk identificeerbare informatie</p>	<p>Zie Databronnen. De gekozen bronnen komen voort uit de administraties van De HHs en openbare, administratieve bronnen en de nationale studenten-enquête. Er zijn door het lectoraat geen aanvullende bronnen verzameld op basis van eventuele enquêtes onder de studenten van De HHs.</p> <p>Hiaten in de bronnenverzameling zijn die gegevens die informeel van aard zijn: sociale contacten, motivatie, verwachtingen, etc. Deze gegevens zijn binnen De HHs niet beschikbaar of niet ontsloten (bijv. sociale interacties in een online leeromgeving).</p> <p>Alle administraties en enquêtes maken een selectie aan informatie en zijn daarmee vooringenomen. De mate van vooringenomenheid is onderwerp van het onderzoek vanuit de weging van modellen en de onverklaarde variantie op mogelijke uitkomstvariabelen.</p>
<p>A.4 Mitigatie van downstream bias</p>	<p>Hebben we manieren overwogen om de blootstelling van persoonlijk identificeerbare informatie (PII) te minimaliseren, bijvoorbeeld door anonimisering of door geen informatie te verzamelen die niet relevant is voor analyse?</p> <p>De data is geanonimiseerd (zie boven). Informatie die per postcode beschikbaar komt wordt geaggregeerd naar buurten en wijken. Ter illustratie van mogelijke voorspelkracht zullen fictieve studentprofielen gemaakt worden.</p> <p>De verzameling van specifieke gegevens wordt beperkt tot die gegevens die bekend zijn vanuit literatuur of eerder onderzoek van de lector.</p> <p>Hebben we manieren overwogen om downstreamresultaten te testen op bevooroordeelde uitkomsten (bijv. het verzamelen van gegevens over beschermde groepsstatus zoals etniciteit of geslacht)?</p> <p>Deze testen zijn het cruciale onderwerp van dit onderzoek om vast te stellen welke bias er mogelijk is op deze kenmerken. Vooralsnog wordt etniciteit niet meegenomen ten gunste van de SES-WOA kenmerken.</p>

8.2 B. Data-opslag

Aandachtspunt	Overwegingen
B.1 Data beveiliging	<p>Hebben we een plan om gegevens te beschermen en te beveiligen (bijv. encryptie in rust en in transit, toegangscontroles op interne gebruikers en derde partijen, toegangslogs en up-to-date software)?</p> <p>De gegevens worden opgeslagen conform de richtlijnen van de hogeschoolbibliotheek van De HHs via SURF Research Drive. Eventuele overdracht van data vindt plaats met behulp van SURF Filesender, waarbij de bewijzen van transacties separaat worden opgeslagen. De code wordt separaat opgeslagen in een private omgeving van github voor leden van de kenniskring van het lectoraat Learning Technology & Analytics.</p> <p>Het versiebeheer op de software (R) vindt plaats via het renv package in R.</p>
B.2 Het recht om vergeten te worden	<p>Hebben we een mechanisme waarmee een persoon kan vragen dat diens persoonlijke gegevens worden verwijderd?</p> <p>Omdat de data is geanonimiseerd is het niet mogelijk om een student naderhand te verwijderen omdat de student niet geïdentificeerd kan worden.</p>
B.3 Data bewaar plan	<p>Is er een schema of plan om de gegevens te verwijderen als ze niet langer nodig zijn?</p> <p>De data wordt na afloop van het onderzoek en publicatie van de resultaten na 1 jaar gearchiveerd in DARK-store. De HHs heeft deze voorziening nog niet, maar hierover zal het lectoraat in gesprek gaan met de data stewards van de hogeschoolbibliotheek.</p>

8.3 C. Analyse

Aandachtspunt	Overwegingen
C.1 Missende perspectieven	<p>Hebben we geprobeerd om blinde vlekken in de analyse aan te pakken door samenwerking met relevante belanghebbenden (bijv. het controleren van aannames en het bespreken van implicaties met betrokken gemeenschappen en materiedeskundigen)?</p>

Aandachtspunt	Overwegingen
	De uitkomsten zullen met studenten en opleidingen besproken worden (docenten, opleidingsdirecteuren, onderwijsadviseurs, etc.). Voor het contact met studenten gaan we de samenwerking met het Partner Up! programma van het Kenniscentrum Global & Inclusive Learning en het Inclusion Office van De HHs.
C.2 Dataset bias	<p>Hebben we de gegevens onderzocht op mogelijke bronnen van vooroordelen en stappen ondernomen om deze vooroordelen te verminderen of aan te pakken (bijv. bestendiging van stereotypen, bevestigingsvooroordelen, onevenwichtige klassen of weggelaten beïnvloedende variabelen)?</p> <p>Dit is het onderwerp van dit onderzoek: het vaststellen van bias en mogelijkheden om deze te reduceren ten gunste van studenten, het onderwijs en het onderwijsbeleid van De HHs.</p>
C.3 Eerlijke representatie	<p>Zijn onze visualisaties, samenvattende statistieken en rapporten ontworpen om de onderliggende gegevens eerlijk weer te geven?</p> <p>De visualisaties, samenvattende statistieken en rapporten zijn erop gericht om bias juist in kaart te brengen met inbegrip van de effecten van aanpassingen. Dit is een integraal onderdeel van de gebruikte software: de DALEX en <code>fairmodels</code> packages.</p>
C.4 Privacy in de analyse	<p>Hebben we ervoor gezorgd dat gegevens met PII niet worden gebruikt of weergegeven tenzij dit noodzakelijk is voor de analyse?</p> <p>Ja. Zie het antwoord op A.3.</p>
C.5 Controleerbaarheid	<p>Is het proces voor het genereren van de analyse goed gedocumenteerd en reproduceerbaar als we in de toekomst problemen ontdekken?</p> <p>Ja. Zie hiervoor het hoofdstuk over Reproduceerbaarheid</p>

8.4 D. Modelleren

Aandachtspunt	Overwegingen
D.1 Proxy discriminatie	<p>Hebben we ervoor gezorgd dat het model niet berust op variabelen of benaderingen voor variabelen die oneerlijk discriminerend zijn?</p> <p>Dit is het onderwerp van dit onderzoek.</p>

Aandachtspunt	Overwegingen
D.2 Eerlijkheid tussen groepen	<p>Hebben we de modelresultaten getest op eerlijkheid met betrekking tot verschillende getroffen groepen (bijv. getest op ongelijke foutpercentages)?</p> <p>Dit is het onderwerp van dit onderzoek.</p>
D.3 Variabele selectie	<p>Hebben we de effecten van het optimaliseren voor onze gedefinieerde variabelen overwogen en hebben we aanvullende variabelen overwogen?</p> <p>De optimalisatie van de variabelen en hun effect op eventuele voorspellingen is onderwerp van onderzoek (via ceteris paribus analyses en bias analyses).</p> <p>Eventueel aanvullende variabelen zijn deels overwogen en - vanwege het gebrek daaraan - afgevalen. Daarnaast verwachten we dat gaandeweg het onderzoek mogelijk aanvullende ideeën over variabelen zullen ontstaan, die gerapporteerd zullen worden.</p>
D.4 Uitlegbaarheid	<p>Kunnen we in begrijpelijke termen een beslissing uitleggen die het model heeft genomen in gevallen waarin een rechtvaardiging nodig is?</p> <p>Ja. De DALEX en fairmodels packages zijn erop gericht de interne werking van analyses en modellen exact uit te splitsen ten gunste van 'glass box' modellen.</p>
D.5 Communicatie over bias	<p>Hebben we de tekortkomingen, beperkingen en vooroordelen van het model gecommuniceerd aan relevante belanghebbenden op een manier die algemeen begrepen kan worden?</p> <p>De beperkingen van het onderzoek zullen worden meegenomen in de rapportages/artikelen en verschillen vormen waarin we de uitkomsten van het onderzoek zullen bekend maken. Zie Verwachte resultaten</p>

8.5 E. Inzet

Aandachtspunt	Overwegingen
E.1 Monitoring en evaluatie	<p>Hoe zijn we van plan het model en de effecten ervan te bewaken nadat het is ingevoerd (bijv. prestatiebewaking, regelmatige controle van steekproefvoorspellingen, menselijke beoordeling van beslissingen die veel op het spel staan, beoordeling van downstreameffecten van fouten of beslissingen met een lage betrouwbaarheid, testen op concept drift)?</p> <p>Het onderzoek levert naar verwachting inzicht in variabelen waarop bias bestaat in De HHs. Afhankelijk van de uitkomsten, kunnen deze worden overgenomen in het onderwijs en onderwijsbeleid. Dit kan onderdeel worden van het programma voor de implementatie van het Instellingsplan. Hierover zullen we in gesprek gaan met het CvB en de Directeur Strategie.</p>
E.2 Het repareren en voorkomen van schade	<p>Hebben we met onze organisatie een plan van aanpak besproken voor het geval gebruikers schade ondervinden van de resultaten (bijv. hoe evalueert het datascience team deze gevallen en hoe werkt het analyses en modellen bij om schade in de toekomst te voorkomen)?</p> <p>In de ervaring van de lector is goede communicatie over de uitkomsten van belang, met name met de pers. Bij de presentatie en bespreking van de resultaten zullen meerdere perspectieven op eerlijkheid behandeld worden om te voorkomen dat er een dominant wordt. Daarnaast de monitoring van het gebruik van deze inzichten. De invloed hierop is – zoals bij elk wetenschappelijk onderzoek – beperkt.</p>
E.3 Terugdraaien	<p>Is er een manier om het model in productie uit te schakelen of terug te draaien als dat nodig is?</p> <p>Het onderzoek ontwikkelt geen modellen die in productie gaan. Dit kan onderdeel zijn van een vervolgproject en zal dan geadresseerd worden.</p>
E.4 Onbedoeld gebruik	<p>Hebben we stappen ondernomen om onbedoeld gebruik en misbruik van het model te identificeren en te voorkomen en hebben we een plan om dit te controleren zodra het model is ingevoerd?</p> <p>Zie het antwoord bij vraag E.2</p>

Bron: Data Science Ethics Checklist, gemaakt met [deon](#) en vertaald door T. Bakker met behulp van DeepL; zie de site van deon voor de [originele Engelse versie](#).