

# Process Book

Jardenna, Janne, Jonne, Julius

June 2018

## Week 3

### Maandag 18/06/18

#### Clustering

Er is een eigen versie van het K-means algoritme geïmplementeerd om het clusteren van data mogelijk te maken. De gegeven versie van SK-Learn is overwogen. Echter aangezien dit algoritme minder goed omgaat met missing values is er voor gekozen om de eigen versie van K-means te gebruiken. Namelijk bij het clusteren werden soms producten vergeleken over verschillende tijdsspannen waardoor er dus als het ware missing values ontstonden. In onze versie van K-means worden missing values simpelweg genegeerd. Dit betekent dat deze dus ook niet worden meegeteld als datapunten bij het gemiddelde waardoor het gemiddelde een betere representatie geeft van de beschikbare datapunten.

Als Exploratory Data Analysis hebben het K-means algoritme toegepast op de producten van India en Oekraïne zodat clusters ontstonden. We hebben deze landen gekozen omdat deze een grote variatie aan producten hebben waarvan we de classificaties zelf konden beoordelen.

### Dinsdag 19/06/18

#### Bedenken van deelvragen

Naast de drie deelvragen gegeven op canvas hebben we nog een aantal andere vragen bedacht die ons interessant leken om te onderzoeken. Hiermee hebben we gelijk ook bepaald welke invalshoek we zullen hebben bij onder andere het geven van de presentatie. Twee van onze deelvragen gaan over de olieprijsen omdat wij verwachten dat er een correlatie bestaat tussen de voedselprijzen en olieprijsen. Namelijk voor het vervoeren van voedsel is brandstof nodig. Als de brandstof duurder wordt, wordt het ook duurder om een product te leveren. Vandaar dat wanneer de olieprijsen toenemen ook de voedselprijzen toenemen.

De andere deelvragen gaan over het verband tussen de voedselprijzen en het sterftcijfer in een land en de vluchtelingen-stromen. Dit leek ons interessant

omdat we verwachten dat als de voedselprijzen stijgen dat in arme landen dan ook het sterftecijfer toeneemt. Bovendien kan het stijgen van de voedselprijzen ook een oorzaak zijn voor het toenemen van de vluchtelingen-stromen.

### **Extra datasets**

Naast de extra datasets van het bruto nationaal product en valuta die in week 1 zijn uitgekozen is er naar aanleiding van de bedachte deelvragen gezocht naar een paar andere datasets. Namelijk die van vluchtelingenstromen en de sterftecijfers binnen een land.

## **Woensdag 20/06/18**

### **Correlatie**

Voor het beantwoorden van de deelvragen is het handig om de mate van correlatie tussen twee factoren te kunnen bepalen. Dit is gedaan met behulp van regressie. De methodes die zijn overwogen zijn Pearson en de Spearman's rank. De Spearman correlatie coëfficiënt is gebaseerd op de geordende waarden voor elke variabele in plaats van de ruwe data. Vandaar dat de Spearman's rank vooral wordt gebruikt wanneer er sprake is van ordinale variabelen. Een variabele bij onze dataset is tijd, wat een continue variabele. Ook zijn we vooral geïntereiseerd in de correlatie tussen de prijzen van twee producten. Vandaar dat Spearman voor ons niet de juiste keuze was. De Pearson correlatie evalueert de relatie tussen twee continue variabelen. Vandaar dat er is gekozen om Pearson te gebruiken.

## **Donderdag 21/06/18**

### **Visuele analyse**

Er is een begin gemaakt aan het beantwoorden van de deelvragen. We gebruiken hiervoor de bokeh plots, Pearson en k-means. Allereerst analyseren we de bokeh plots en kijken we waar veel data punten zijn zodat de resultaten representatief zijn. Daarna bereken we de correlatie coëfficiënt en kijken we met behulp van K-means welke clusters gevormd worden.

## **Vrijdag 22/06/18**

### **Technisch rapport**

Vandaag is vooral gewerkt aan een opzet van het technische rapport. Er wordt in LaTeX gewerkt omdat dit het verslag op een overzichtelijke universele manier weergeeft. Voor het technische rapport is er voor gekozen om voor alle deelvragen een aparte hypothese, opzet en verwachting te maken. Ook de methode resultaten en discussie zal geschreven worden aan de hand van de gescheiden deelvragen. Dit hebben we gedaan omdat er zo het best overzicht gehouden kan worden. Ook is zo duidelijk is welke stappen we hebben ondernomen om

bepaalde resultaten te verkrijgen. Bovendien is verder gegaan met het beantwoorden van deelvragen dit is wederom gedaan met behulp van K-means en de correlatie coëfficiënt en bokeh-plots.

## **Week 4**

### **Maandag 25/06/18**

#### **Webpagina**

Er is gewerkt aan de visualisatie en webpagina. De belangrijkste beslissingen die zijn gemaakt zijn hoe we ons dashboard eruit willen laten zien. Er is gekozen voor een interactief dashboard zodat bij de presentatie elke deelvraag ondersteund kan worden met behulp van de website. Hierbij kan de data van elk land in een bepaalde periode worden opgevraagd. Ook kan er worden geselecteerd op regio en op product en wordt dit vervolgens op een overzichtelijke manier gevisualiseerd.

### **Dinsdag 26/06/18**

Vandaag is voornamelijk gewerkt aan openstaande processen waaronder het beantwoorden van de deelvragen, het technisch rapport en de website.

### **Woensdag 27/06/18**

Vandaag zijn keuzes gemaakt voor het visualiseren van data. Er is overwogen om een tabel te gebruiken om clusters te visualiseren of een grafiek waarbij de lijnen die bij een cluster behoren dezelfde kleur hebben. Echter is elke datum bij ons een dimensie. T-SNE is een goede methode om *high dimensional data* te visualiseren. Vandaar dat we ervoor hebben gekozen om dit gebruiken.

### **Donderdag 28/06/18**

Vandaag is een groot deel van de presentatie voorbereid.

### **Vrijdag 29/06/18**

Vandaag