

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.

Vrijdag 22/06/18

Week 4

Maandag 25/06/18

Webpagina

Er is gewerkt aan de webpagina.

Dinsdag 26/06/18

Woensdag 27/06/18

Donderdag 28/06/18

Vrijdag 29/06/18