

Examen Data Scientist - Data analyseren in Python

Introductie

Voor dit examen is het belangrijk dat je volgende zaken zeker doet:

- Maak een Virtual Environment aan & exporteer een `.yaml` (of soortgelijke) file aan die je toevoegt aan je codebase.
- Maak een github pagina aan met een bijhorende `.gitignore` en `readme.md`. Voeg een link naar de github repo toe aan je readme file.
- Zorg dat je op het einde de tijd neemt om je code op te schonen zodat deze begrijpelijk is voor derden.
- Wanneer je klaar bent, maak een `.zip` bestand van je volledige codebase aan en laad deze op in de leeromgeving.

In deze opgave bekijken we een gesimuleerde dataset over de geboortes in België van 2019. De data is te vinden in `.csv` bestanden die gelabeld zijn met de echte geboortedag van de personen in het bestand. Laat ons bijvoorbeeld even kijken naar de eerste 5 lijnen van het bestand `2019-1-1.csv`:

gemeente	naam	geslacht	verwachte datum
Hooglede	Elias	Mannelijk	01/14/2019
Sint-Niklaas (Sint-Niklaas)	Pauline	Vrouwelijk	01/05/2019
Wijnegem	Anita	Vrouwelijk	01/05/2019
Grâce-Hollogne	Jean-Paul	Mannelijk	01/13/2019
Boussu	François	Mannelijk	01/26/2019

Hier hebben we de kolommen:

- **gemeente:** De gemeente waar de persoon geboren is.
- **naam:** De voornaam van de persoon.
- **geslacht:** Het geslacht van de persoon (Mannelijk of Vrouwelijk).
- **verwachte datum:** De verwachte geboortedatum van de persoon.

Dus hier hebben we dan bijvoorbeeld op de eerste lijn de mannelijke persoon Elias die geboren is in Hooglede waarvan men verwachtte dat hij zou geboren worden op `01/14/2019` (Amerikaans formaat) maar in feite geboren is op 1 januari 2019.

Doorheen deze opgave gaan we de data van dichterbij inspecteren! Je kan deze opgaves alternatief lezen in de bijgeleverde notebook `examen.ipynb` waar ook de (meeste) afbeeldingen & resultaten zichtbaar zijn.

Stap 1: data inlezen

Maak een lus over alle bestanden in de map `data/geboortes`, lees elk `.csv` bestand in en voeg deze samen tot 1 grote DataFrame. Merk op dat de datum verwerkt zit in de filename, voeg deze toe aan de kleine DataFrames voor je de dataframes samenvoegt met een `pd.concat(dfs)`.

Voeg een extra kolom toe met een "dag van het jaar" getal, dus 1 januari 2019 is 1, 2 januari 2019 is 2, enzovoort.

Opmerking: Er is ook data voor 29 februari 2019, terwijl dit geen schrikkeljaar is! Plaats deze geboortes in een afzonderlijke DataFrame `df_wrong`. Zet in `df_wrong` een extra kolom met een referentie naar waarom deze data foutief is.

Stap 2: Dagelijks aantal geboortes bekijken

Vraag 1: Maak een plot van het aantal geboortes per dag van het jaar.

Vraag 2: Outliers vinden & behandelen

- We klassificeren een datum als een outlier als het aantal geboortes meer dan 50% afwijkt van de gemiddelde waarde overheen de volledige dataset. Zoek alle outliers in de dataset.
- Ik zal je vertellen; de outliers op 1 januari en 1 juli komen doordat alle geboortes waarvoor de datum niet correct genoteerd is, de datum 1 januari (voor de eerste helft van het jaar) en op 1 juli (voor de tweede helft van het jaar) geplaatst worden. Om dit op te lossen, halen we in de grote geboortes DataFrame (met 1 lijn per geboorte) deze twee datums er volledig uit. Dus we gaan **alle** geboortes van 1 januari en 1 juli verwijderen uit deze dataframe en toevoegen aan de `df_wrong` dataframe met als reden `2019-01-01` en `2019-07-01`.
- Maak dezelfde plot nu opnieuw met het aantal geboortes per dag.
- We zien dat er nog steeds redelijk wat uitschieters zijn, kan je code schrijven om de 5 extreme uitschieters in de 2e helft van het jaar identificeren (in totaal 8 datums)?

Stap 3: Onderzoeksvragen

Onderzoek 1: Unisex namen

Een aantal statistieken

Sommige namen zijn unisex, dat wil zeggen dat zij hetzelfde zijn voor mannen en vrouwen. Beantwoord volgende 3 vragen over unisex namen in onze dataset:

- Hoeveel unisex namen zijn er?
- Wat is de meest voorkomende unisex naam; hiervoor zoek ik 3 namen (en getallen):
 - De meest populaire unisex naam bij mannen.
 - De meest populaire unisex naam bij vrouwen.
 - De meest populaire unisex naam.

Tip: Neem hiervoor terug de volledige dataset in acht (dus met de foutief geklasseerde namen).

We zien dat er een vrij groot aantal namen veel meer voorkomen bij mannen dan vrouwen (en omgekeerd). Dit vinden we eigenlijk geen "echte" unisex namen. We noemen een naam "echt unisex" als de naam niet meer dan 50% meer voorkomt bij 1 van de 2 geslachten. Dus als er `x` mannen en `y` vrouwen zijn genaamd Chris, dan noemen we Chris "echt unisex" als $x \leq 1.5 * y$ en $y \leq x * 1.5$. Filter je unisex dataframe tot een `df_real_unisex` waarbij je enkel de "echte" unisex namen weerhoudt.

- Beantwoord nu bovenstaande 3 vragen opnieuw voor deze dataset. Hierbij werk je best met een functie zodat je geen code moet herhalen!
- Zijn de echte unisex namen populairder bij mannen of vrouwen, of is de populariteit hetzelfde? Deze vraag kan je beantwoorden door te kijken naar het percentage mannen/vrouwen met een echte unisex naam.

Visualizatie

Maak een visualizatie die alle echte unisex namen toont en de relatieve voorkomens bij mannen en vrouwen.

Onderzoek 2: de accuraatheid van de geschatte bevallingsdatum

Evolutie vergelijken

Maak een afbeelding van het totaal aantal geboortes per dag en het totaal aantal verwachte geboortes.

Bonusvraag: Gewoon om over na te denken: waarom zien de uiteindes van je plot er wat gek uit? Schrijf het antwoord neer in een strategisch geplaatste markdown cel.

Verbanden bekijken

Maak een histogram van het aantal dagen dat babies te vroeg geboren werden in 2019 en een scatterplot die het verband toont tussen de verwachte en effectieve geboortedatum.

Onderzoek 3: Aantal namen versus aantal babies

Naarmate dat er meer babies geboren worden, gaan deze ook meer verschillende namen krijgen. Zolang er weinig babies geboren worden lijkt het logisch dat dit verband min of meer lineair is, maar naarmate dat er meer en meer babies geboren worden lijkt het niet meer logisch dat dit verband lineair blijft. We verwachten bijvoorbeeld niet dat het aantal unieke namen gegeven aan 2 miljoen babies het dubbel is van het aantal unieke namen gegeven aan 1 miljoen babies. Kan je op basis van deze dataset het verband tussen het aantal unieke namen in functie van het aantal babies eens onderzoeken? Hiervoor krijg je geen voorstel tot visualizatie maar moet je zelf mij proberen te overtuigen van je antwoord.