

Skóre rizikovosti a sčítání lidu

Michal Půlpán

15. 11. 2022

1 Úvod

V poslední několika letech jsme v České republice mohli zaznamenat několik zajímavých mikroekonomických událostí. Výrazně narostla kupní síla, obyvatelstvo více utrácelo, ale i se znatelně zadlužilo. To se mimo jiné projevilo i na růstu e-commerce a maloobchodu. Další nevídanou situací byl i znatelný růst cen nemovitostí napříč českými městy i suburbii způsobený mimojiné nezastavitelnou poptávkou po novém bydlení, investicím a nebo jen touze získat více místa či zahradu.

To vše, vzhledem k negativnímu vývoji ekonomické situace, má vliv na rekordní zadlužení obyvatelstva. Každá instituce poskytující úvěry fyzickým osobám se tak musí chovat opatrněji a věnovat více pozornosti, komu peníze poskytuje. V tomto projektu se podíváme na to, jak lze využít moderních technologií, veřejně dostupných dat a data science ke svému prospěchu.

Pro finanční instituce poskytující úvěry, či obchodníky nabízející své zboží nebo služby s odloženou splatností, by tak mohlo být zajímavé, jaké faktory ovlivňují schopnost splácení klientů. To by samo o sobě mohlo napomoci k větší obezřetnosti prodejce, přesnějšímu marketingovému cílení a nebo naopak k uzpůsobení služeb jednotlivým skupinám obyvatelstva.

Podíváme se tak na vztah mezi agregovanými daty ze sčítání lidu a podílu exekucí a bankrotů v jednotlivých obcích a pokusíme se z nich odvodit, jaké sociodemografických údaje ovlivňují, zda jsou obyvatelé schopni splácet své dluhy či nikoli.

2 Data

Využijeme vybrané výsledky ze sčítání lidu z roku 2011 a 2021 a data o exekucích a bankrotech za jednotlivé obce České republiky.

2.1 Původ dat

Celá datová sada za jednotlivé obce České republiky obsahuje:

- vybrané výsledky sčítání lidu 2011 a 2021
- počet obyvatel celkem v jednotlivých sociodemografických skupinách
- výskyt vybraných trestných činů
- podíl osob v bankrotu či exekuci

Byla stažena ze stránek nextcloud.profinet.eu, pravděpodobně je ale původ neagregovaných datasetů od Českého statistického úřadu.

2.2 Stručný popis

Pojďme si představit používané datasety a jejich strukturu. Oba datasety jsou díky dělení na obce velmi granuloované a tím i relativně dlouhé.

2.2.1 Sčítání lidu (Census)

Data ze sčítání lidu z let 2011 a 2021 (`geodata/census11_21/data_obce_vyhl_nevyhl.csv`) jsou obsáhlá v jednom souboru a obsahují informace o jednotlivých obcích. Každá obec má svůj řádek a sloupce obsahují informace o jednotlivých sociodemografických skupinách. Výsledky jsou však agregované, takže neobsahují informace o jednotlivých obyvatelích, ale o celkovém počtu obyvatel v daných skupinách. V datech tak máme například pro každou obec v České republice informaci o počtu dětí, vzdělání, věku, rodinného stavu, občanství, ale i o vlastníkovi domů (družstvo, fyz. osoba, obec/stát, podílové) či stavebního materiálu domu, ve kterém obyvatelé žijí a mnoho více.

Dataset je složený z 6246 pozorování o 160 proměnných. Díky tomu je dataset velmi široký a obsahuje velkou škálu sociodemografických údajů. Každé pozorování je vytvořeno pro obec, kterých je celkem 6246.

2.2.2 Exekuce a bankroty

Data o exekucích a bankrotech (`geodata/score/obce_skore_rizikovosti.csv`) jsou dělená dle obcí České republiky a obsahují informace o podílu exekucí (klouzavý průměr za roky 2018-2021) a bankrotů (klouzavý průměr za roky 2018-2021) v jednotlivých obcích. Výsledky jsou agregované, takže neobsahují informace o jednotlivých obyvatelích, ale o celkovém počtu exekucí a bankrotů v dané obci. Mimo jiné obsahují i tzv. skóre obce, které je výsledkem normovaného průměru počtu exekucí a bankrotů v obci. Toto skóre je v rozmezí 0 až 1, kde 0 znamená, že v obci nebyly žádné exekuce ani bankroty a 1 znamená, že v obci byly exekuce a bankroty v průměru v každém roce.

Dataset je složený z 6254 pozorování o 10 proměnných. Každé pozorování je vytvořeno pro obec, kterých je celkem 6254.

2.3 Kvalita dat

Data ze sčítání lidu z nějakého důvodu neobsahují data ze 14 obcí. Dle [struktury území ČR mezi roky 2013 a 2022](#) od ČSÚ je v České republice (od 1.1.2016) celkem 6260 obcí. V našich datech je však “pouze” 6246 resp. 6254 obcí pro dataset exekucí. Těžko říct, zda je to chyba v našich datech, nebo v původních datech od ČSÚ, ale i mezi datasety je rozdíl v 8 obcích, které ve výsledcích ze sčítání lidu chybějí. Konkrétně se jedná o obce:

- Krhová v okrese Vsetín (500062)
 - 2024 obyvatel
- Poličná v okrese Vsetín (500071)
 - 1745 obyvatel
- Bražec v okrese Karlovy Vary (500101)
 - 221 obyvatel
- Doupovské Hradiště v okrese Karlovy Vary (500127)
 - 160 obyvatel
- Kozlov v okrese Olomouc (500135)
 - 270 obyvatel
- Luboměř pod Strážnou v okrese Přerov (500151)

- 122 obyvatel
- Město Libavá v okrese Olomouc (500160)
 - 0 obyvatel
- Polná na Šumavě v okrese Český Krumlov (500194)
 - 202 obyvatel

Ač se jedná o obce (pro nás) v relativně zajímavém místě (s předpokládanou horší ekonomickou situací), tak se jedná o relativně malé obce, takže bychom je nemuseli brát v úvahu.

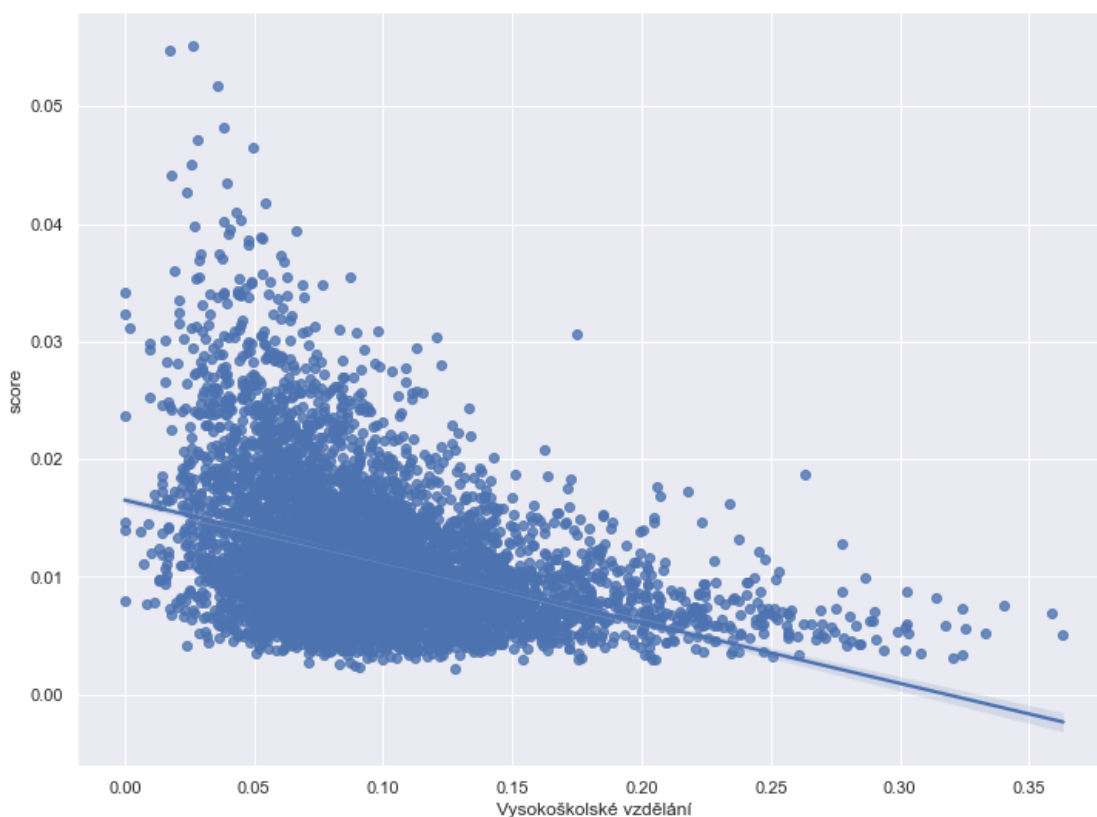
Dle výstupů z nástroje **pandas-profiling** žádný ze zmíněných datasetů nemá problém s chybějícími daty ve sloupcích. Jediné, co je vhodné upravit jsou názvy sloupců, ale ty jsou vždy vysvětleny v příložených textových souborech.

2.4 Příběhy

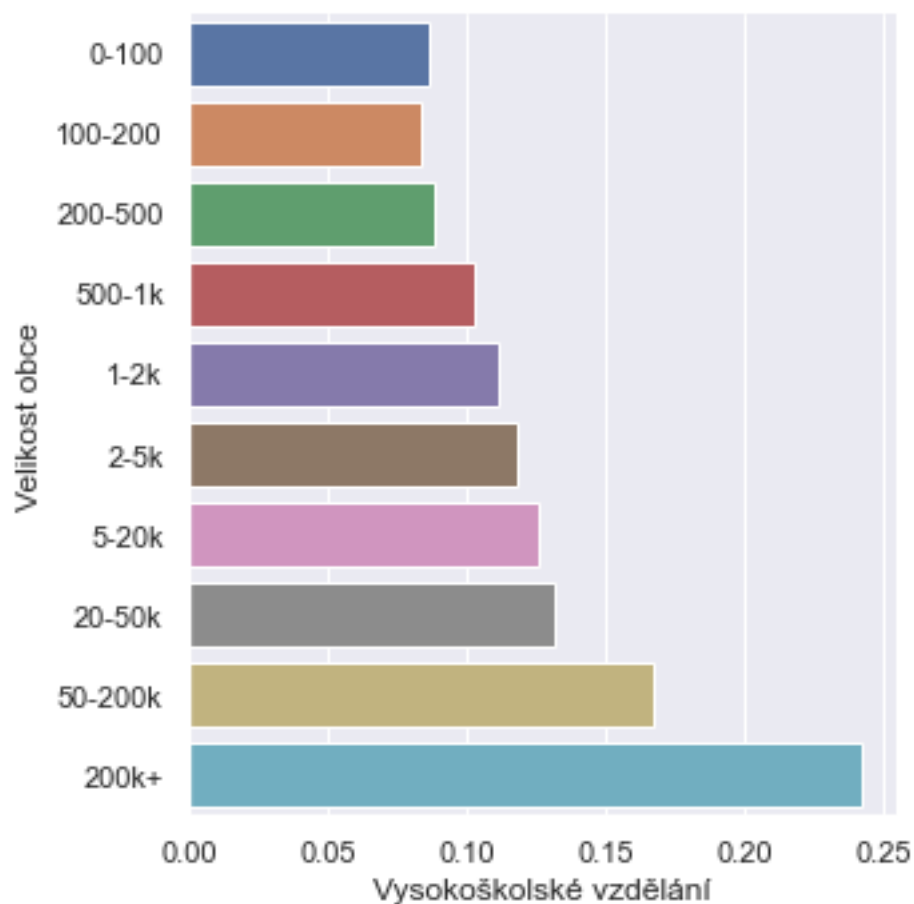
Poskytnuté datasety jsou velmi obsáhlé a obsahují veliké množství informací. Na jejich základě by mělo být možné vytvořit (a podložit) veliké množství příběhů. Ukážeme si jeden z nich.

2.4.1 Vzdělání a dluhy

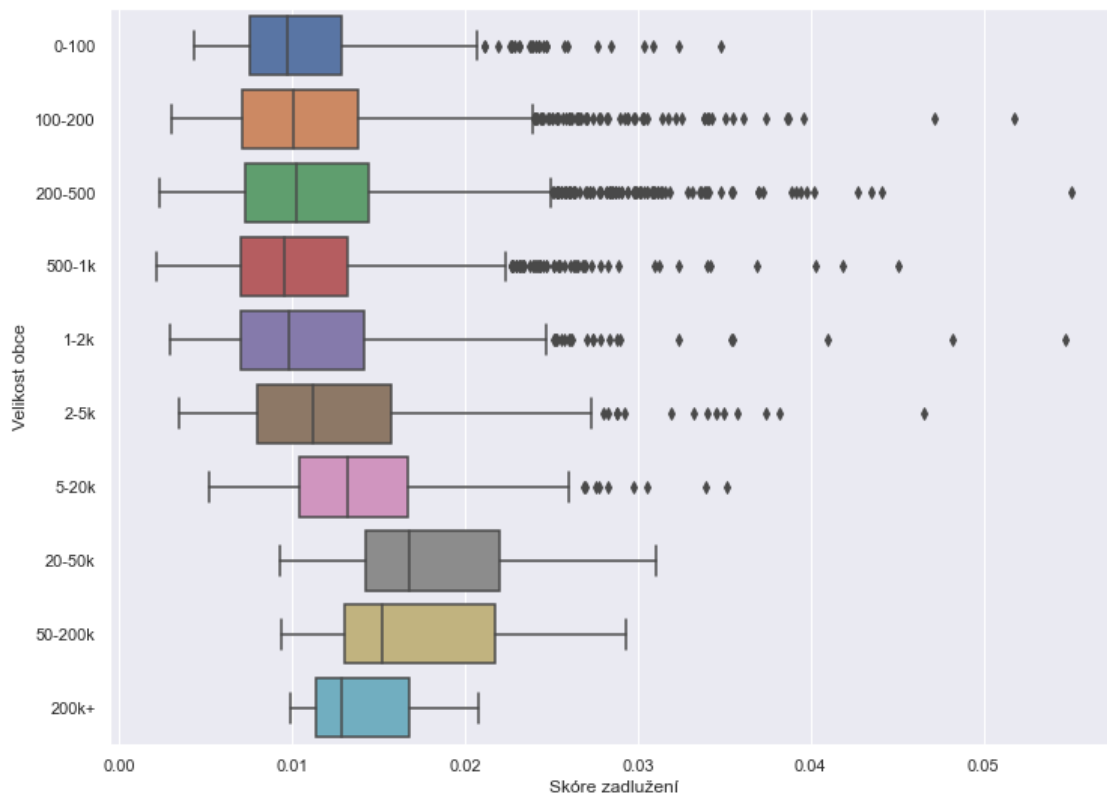
V tomto příběhu se podíváme na vztah vzdělání a zadlužení obyvatel českých obcí. Můžeme nějak využít informace o počtu vysokoškoláků zastoupených v obci?



Dle grafu výše se zdá, že s přibývajícím počtem vysokoškoláků v obci klesá i skóre zadlužení obce. To znamená, že obce s vysokoškoláky mají méně dluhů. To je ale jen první pohled na data. Co když se podíváme na poměr mezi počtem vysokoškoláků a počtem obyvatel v obci? Výsledek je následující:



Jak můžeme vidět, s velikostí obce nám roste i počet vysokoškoláků. Jaký vztah tedy má velikost a zadlužení obce?



Zde můžeme vidět, že velikost obce až takový vliv na její skóre zadlužení nemá (až na několik výjimek). Minimálně tedy na základě grafického zobrazení dat se může zdát, že vysokoškoláci nejsou zodpovědní za zadlužení v obcích.