

UNIVERZITA KARLOVA

Přírodovědecká fakulta

Studijní obor: Sociální geografie a geoinformatika



Ester Kaliská

Úvod do programování

Skúškové zadanie 2

Praha 2026

1. Zadanie

Cieľom úlohy je overiť normalitu dátového súboru obsahujúceho 100 hodnôt, ktoré sú načítané z externého textového súboru. Na testovanie normality som použila Kolmogorov-Smirnov test, keďže je vhodný pre menšie a stredne veľké výbery, je bežne používaný v praxi a je prehľadnejší na ručné naprogramovanie v jazyku python.

2. Popis a rozbor problému

Overenie normality rozdelenia je dôležitým krokom v štatistickej analýze, pretože veľa parametrických metód predpokladá, že dáta majú normálne rozdelenie. Ak tento predpoklad nie je splnený, výsledky ďalších testov môžu byť skreslené alebo nesprávne.

V tejto úlohe som sa zamerala na Kolmogorov-Smirnov test, ktorý slúži na porovnanie empirickej distribučnej funkcie načítaného výberu s teoretickou distribučnou funkciou normálneho rozdelenia. Podstatou testu je zoradenie hodnôt a následné hľadanie maximálnej absolútnej odchýlky (označovanej ako štatistika D) medzi reálnymi dátami a ich ideálnym matematickým modelom.

Rozhodovanie o výsledku prebieha na základe porovnania vypočítanej štatistiky D s kritickou hodnotou zistenou pre daný rozsah súboru (v tomto prípade $n = 100$) pri zvolenej hladine významnosti 0,05. Ak je vypočítaná hodnota D menšia ako kritická hodnota 0,134, nulovú hypotézu o normálnom rozdelení dát nezamietame. V opačnom prípade, ak odchýlka presiahne túto hranicu, predpokladáme, že dáta normálne rozdelenie nemajú.

$$KS = \max_x |F_1(x) - F_2(x)|$$

V praxi je potrebné brať do úvahy, že samotný test normality nemusí vždy poskytnúť jednoznačný záver. Pri väčších výberoch môže byť Kolmogorov-Smirnov test citlivý aj na malé odchýlky od normálneho rozdelenia, ktoré nemusia mať významný vplyv na ďalšiu analýzu. Z tohto dôvodu sa testovanie normality často kombinuje s grafickými metódami, ako sú histogramy alebo Q-Q grafy, ktoré umožňujú vizuálne posúdenie tvaru rozdelenia dát.

3. Použitý algoritmus

Program je rozdelený do niekoľkých logických krokov. Najskôr sa rieši správne načítanie vstupného súboru, následne samotná štatistická analýza a nakoniec vyhodnotenie výsledkov.

Postup algoritmu je nasledovný:

- program zistí absolútну cestu k priečinku, v ktorom sa nachádza skript, aby sa predišlo problémom s relatívnymi cestami
- používateľ zadá názov vstupného súboru

- súbor sa otvorí a číta po jednotlivých riadkoch
- každý riadok sa očistí od bielych znakov a prevedie na typ float
- v prípade chyby (neexistujúci súbor alebo nečíselná hodnota) sa vypíše chybové hlásenie
- hodnoty v dátach sa zoradia (podmienka testu)
- po načítaní dát sa skontroluje, či výber obsahuje hodnoty
- program z dát vypočíta aritmetický priemer a smerodajnú odchýlku, ktoré sú potrebné na štandardizáciu hodnôt do podoby Z-skóre
- pre každý bod sa vypočíta teoretická kumulatívna pravdepodobnosť pomocou chybovej funkcie math.erf a empirická kumulatívna pravdepodobnosť na hornej aj dolnej hranici aktuálneho poradia
- porovnávajú sa rozdiely medzi týmito distribúciami a ukladá sa najväčšia nájdená absolútна odchýlka pod názvom max_difference
- vypočítaná štatistika sa porovná s kritickou hodnotou 0,134 a vypíše sa slovný záver

Zvolený postup algoritmu je navrhnutý tak, aby bolo možné jednotlivé kroky jednoducho kontrolovať a prípadne upravovať. Oddelenie načítania dát, ich kontroly a samotnej štatistickej analýzy zvyšuje prehľadnosť kódu a uľahčuje jeho údržbu. Tento prístup je vhodný najmä pri práci s externými dátami, kde môže dôjsť k rôznym neočakávaným chybám vo vstupoch.

4. Ošetrenie problematických situácií

Pri tvorbe programu som sa snažila myslieť aj na situácie, ktoré by mohli viesť k chybe počas behu programu.

Konkrétnie ide o:

- nesprávny pracovný adresár, ktorý je riešený pomocou os.path.abspath(__file__)
- neexistujúci vstupný súbor, zachytený výnimkou FileNotFoundError
- nečíselné alebo prázdne riadky v súbore, ktoré sú ošetrené pomocou ValueError a podmienky na preskočenie prázdnych riadkov
- nedostatočný počet dát, keďže Kolmogorov-Smirnov test vyžaduje minimálne jednu hodnotu

Vďaka týmto kontrolám program nespadne, ale používateľ dostane zrozumiteľnú informáciu o probléme.

5. Vstupné dáta

Vstupom je externý textový súbor s príponou .txt, ktorý obsahuje jednu číselnú hodnotu na každom riadku. Desatinné čísla musia byť zapísané pomocou bodky. Odporúčaný počet hodnôt je približne 100, čo je postačujúca veľkosť výberu na testovanie normality.

Pre potreby zadania bol vygenerovaný pomocou umelej inteligencie súbor so 100 náhodnými hodnotami, nad ktorým bol program testovaný. Súbor data.txt je uložený v priečinku spolu s kódom.

6. Výstupné dátá

Výstupom programu je:

1. vypočítaná D štatistika zaokrúhlená na štyri desatinné miesta
2. slovná interpretácia výsledku, teda či dátá majú alebo nemajú normálne rozdelenie

Okrem samotnej D štatistiky je dôležité správne interpretovať výsledok v kontexte riešeného problému. Ak program vyhodnotí, že dátá majú normálne rozdelenie, je možné v ďalších krokoch použiť parametrické štatistické metódy, ako napríklad t-test alebo analýzu rozptylu. V prípade, že normalita potvrdená nie je, je vhodné zvážiť použitie neparametrických metód alebo transformáciu dát. Výstup programu tak slúži ako podklad pre ďalšie štatistické spracovanie dát a nie ako konečný cieľ analýzy.

7. Programová dokumentácia

Program je implementovaný ako trieda "DataAnalysis", ktorá zabezpečuje prehľadnosť a zapuzdrenie dát. Na výpočty je použitá knižnica "math" a odtiaľ funkcia sqrt. Metóda "normal_distribution" je vypočítaná na základe vzorca pre Kolmogorov-Smirnov test.

Knižnica "os" slúži na prácu so súborovým systémom. Pri testovaní bol problém s nájdením súboru a cesty k súboru data.txt a využitie knižnice "os" tento problém efektívne rieši za pomoci funkcie "os.path.dirname" a "os.path.join". Spracovanie výnimiek pomocou blokov "try-except" zvyšuje stabilitu programu a zabráňuje jeho pádu pri chybnom vstupe. Ošetruje chybu hodnoty (ValueError), nenájdený súbor (FileNotFoundException) a všeobecnú chybu (Exception). V bloku with open je taktiež ošetrená chyba, ak nie je v riadku žiadna hodnota. V tomto prípade sa riadok preskočí a program nespadne.

8. Možné rozšírenia riešenia

Program splňa zadanie, avšak do budúcnosti bolo možné ho ďalej rozšíriť. Medzi možné vylepšenia patrí:

- doplnenie vizuálnej analýzy (histogram, Q-Q graf),
- zavedenie logovania chýb do externého súboru,
- podpora vstupných súborov vo formáte CSV,
- možnosť volby medzi viacerými testami normality.

Takto navrhnutý program je možné využiť aj v iných úlohách, kde je potrebné rýchlo overiť základné štatistické vlastnosti dát. Môže slúžiť napríklad ako podporný nástroj pri spracovaní výsledkov meraní, dotazníkových prieskumov alebo experimentálnych dát. Vďaka

modulárnej štruktúre je možné program jednoducho prispôsobiť aj iným typom vstupných dát.

Zoznam literatúry:

Social science statistic: <https://www.socscistatistics.com/tests/kolmogorov/> (10.2.2026).

Economipedia: Prueba de Kolmogorov – Smirnoff (K-S),
<https://economipedia.com/definiciones/prueba-de-kolmogorov-smirnoff-k-s.html> (10.2.2026)

Python Software Foundation: Mathematical functions: `math.erf()`,
<https://docs.python.org/3/library/math.html> (10.2.2026)

Vlastné poznámky a prezentácie Úvod do programovaní