

Úvod

- Moje semestrální práce se zaměřuje na to, co dělá knihu populární.
- Analyzovala jsem data z Goodreads a hledala souvislosti mezi vlastnostmi knihy a počtem pozitivních recenzí.
- Cílem bylo vytvořit predikční model popularity a identifikovat klíčové faktory.

Dataset

- Metadata knih (800 000+) a uživatelská hodnocení (350 000+).
- Obsahují autora, rok vydání, jazyk, počet stran a textová hodnocení.
- Tato kombinace umožňuje kvantifikovat popularitu knih z více pohledů.

Čištění a zpracování

- Odstranění duplicit a prázdných hodnot.
- Normalizace jazyků a nakladatelství.
- Kategorizace délky knihy pomáhá odhalit vztahy bez vlivu extrémů.

2. Data enrichment and filtering:

- Missing values in 'Publisher' and 'Language' columns are filled with default labels.
- 'PublishMonthName' is derived from 'PublishMonth' using the calendar module.
- 'LengthCategory' is added based on the number of pages (short/medium/long).
- Rare languages (less than 500 books) are filtered out to focus on major languages.

First 5 rows of the dataframe:

	pagesNumber	LengthCategory	Language	PublishMonthName
0	652.0	long	eng	Unknown
1	870.0	long	eng	January
2	309.0	long	eng	January
3	352.0	long	eng	January
4	435.0	long	eng	January

Most common languages in the dataset:

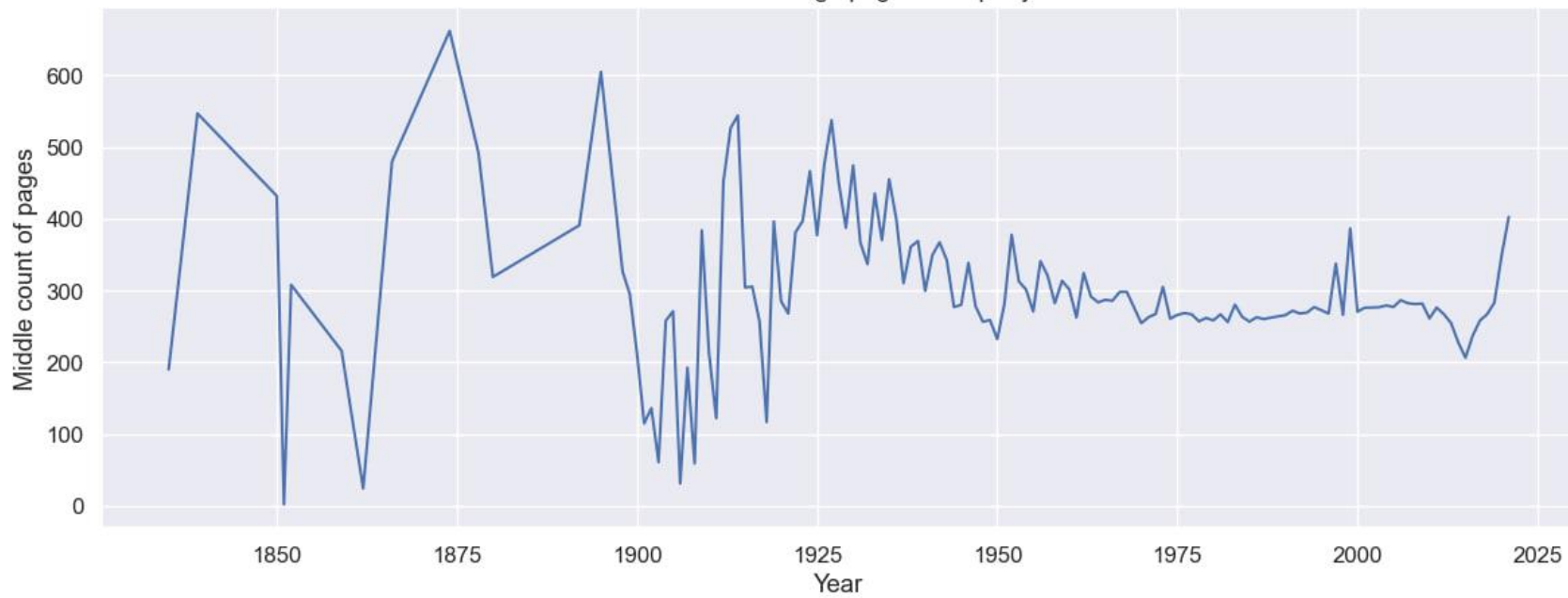
Language	
Unknown Language	669928
eng	119371
en-US	16399
fre	6519
en-GB	5035

Name: count, dtype: int64

Vývoj počtu stran v čase

- Délka knih se v moderní době ustálila okolo 300 stran.
- Stabilita naznačuje standardizaci i preference čtenářů.

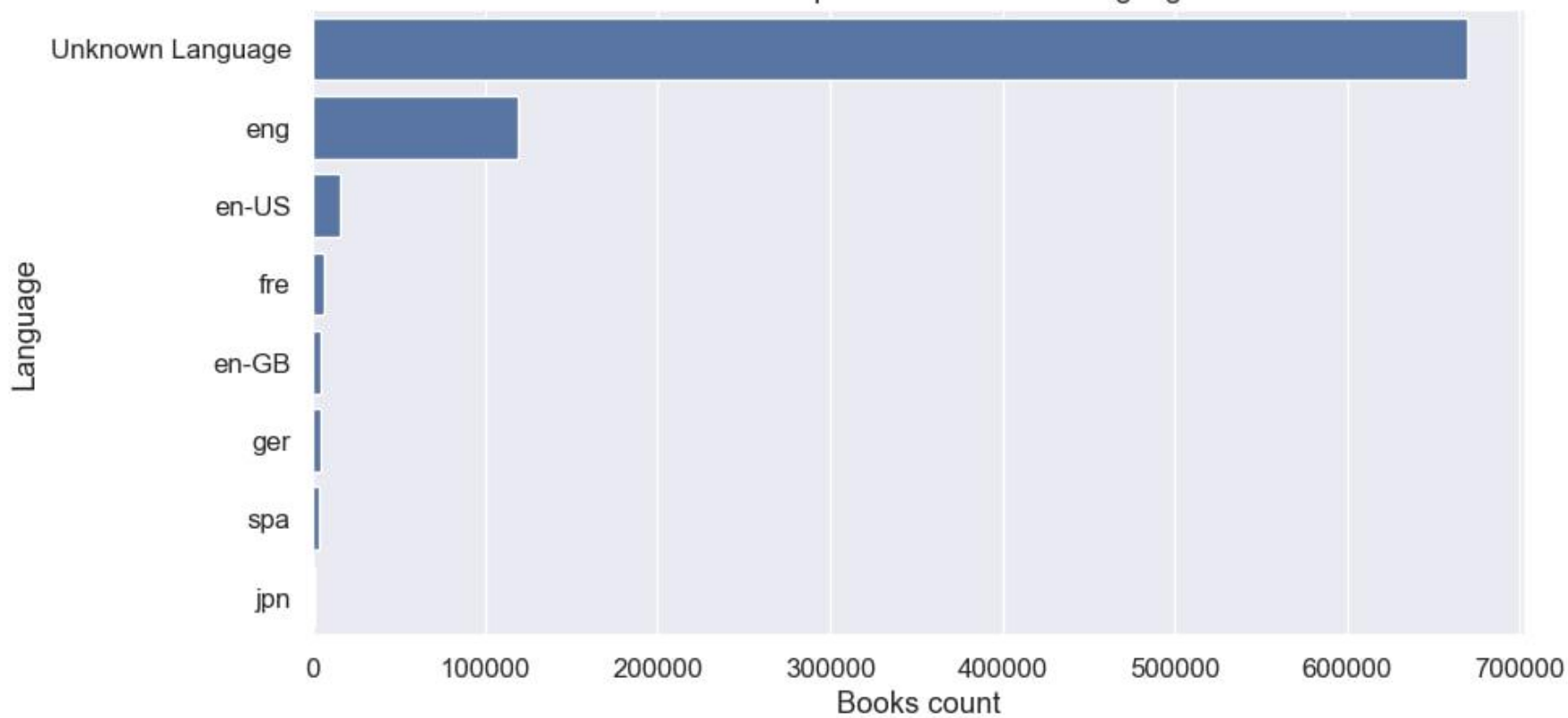
Line chart of average page count per year



Jazyky knih

- Angličtina dominuje. Ostatní jazyky sloučeny do skupiny 'Unknown'.
- Zjednodušuje model a odstraňuje šum z málo častých jazyků.

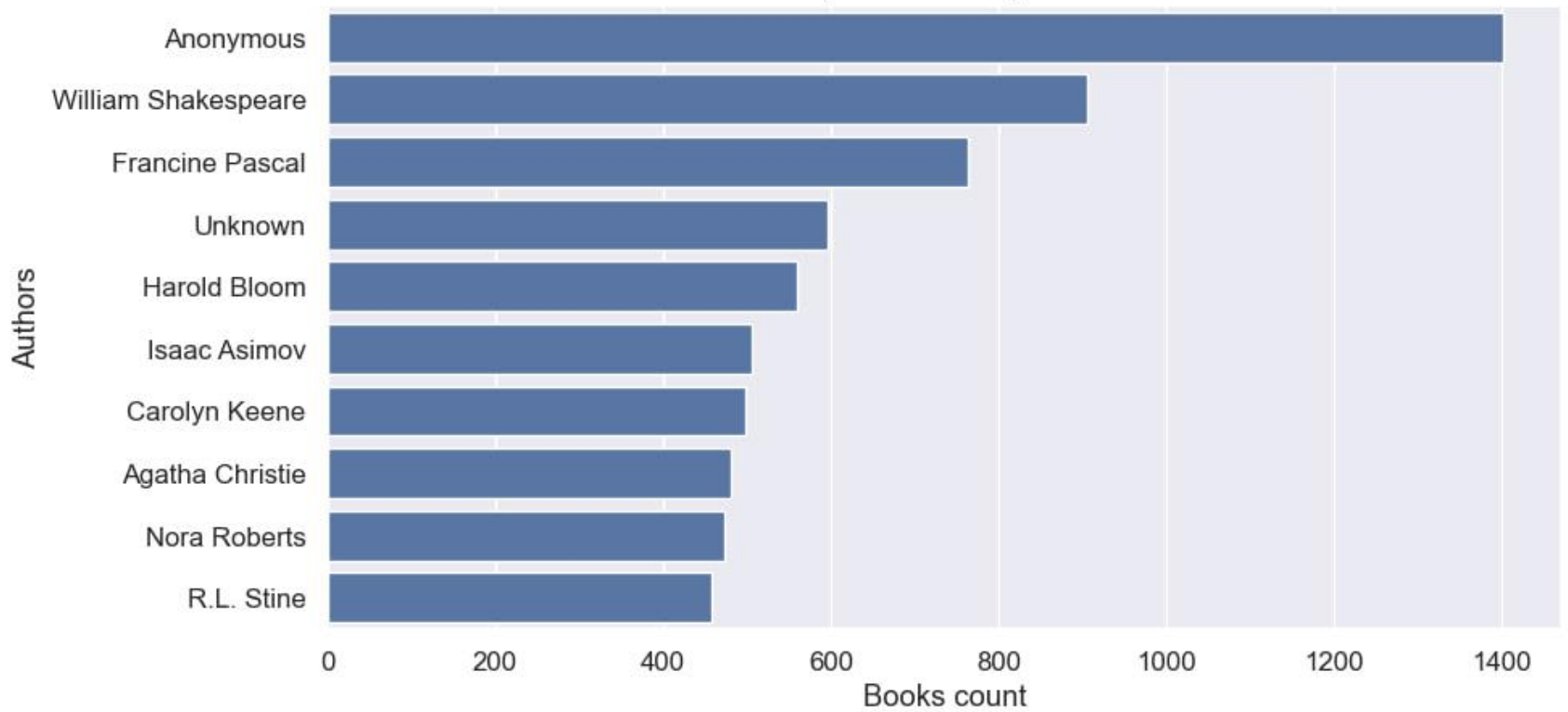
Bar chart of top 10 most common languages



Top autoři

- J.K. Rowling, George Orwell a další dominují hodnocením.
- Odráží kombinaci komerčního úspěchu a literární hodnoty.

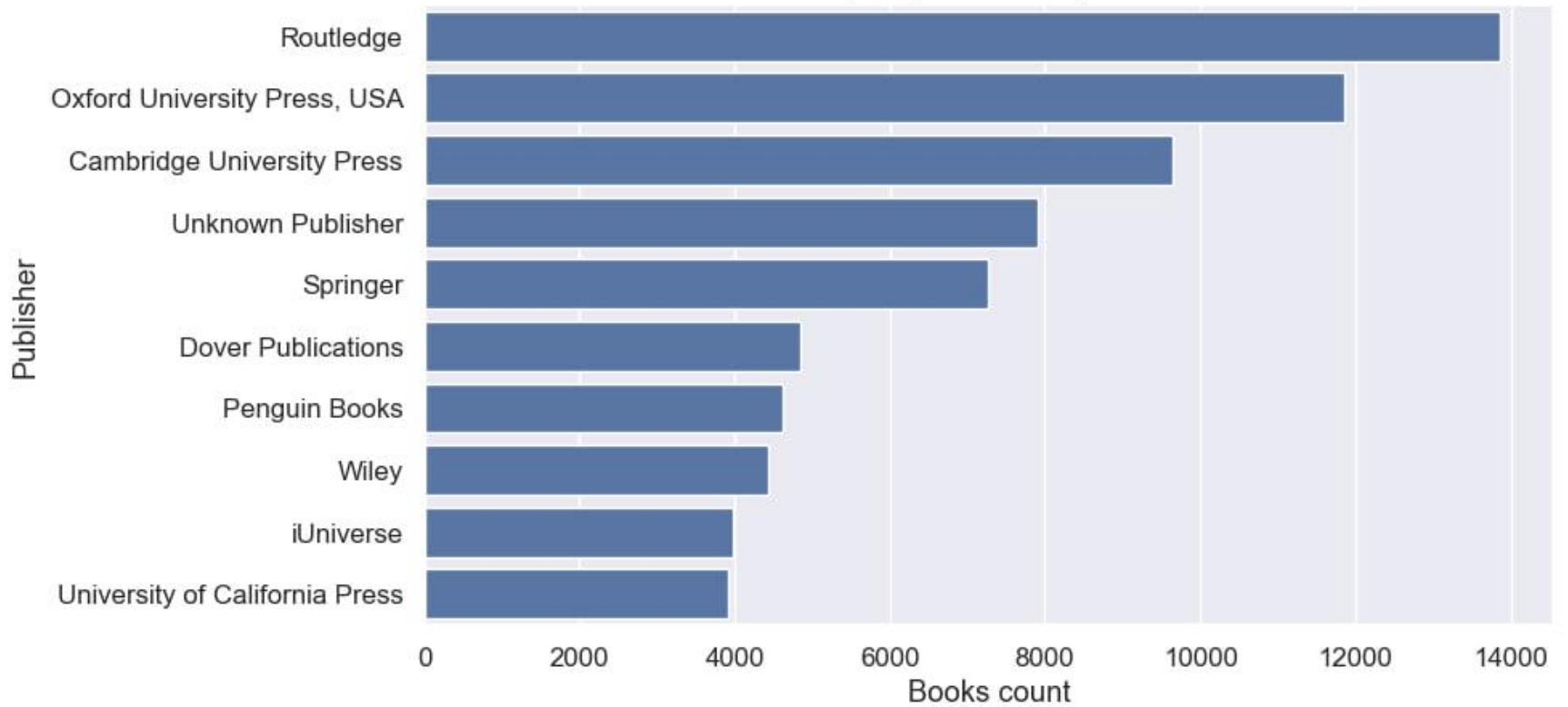
Bar chart of top 10 authors by number of books



Nejčastější vydavatelé

- Routledge a Oxford UP vydávají nejvíce knih – převážně odborné.
- Ovlivňuje charakter celého datasetu.

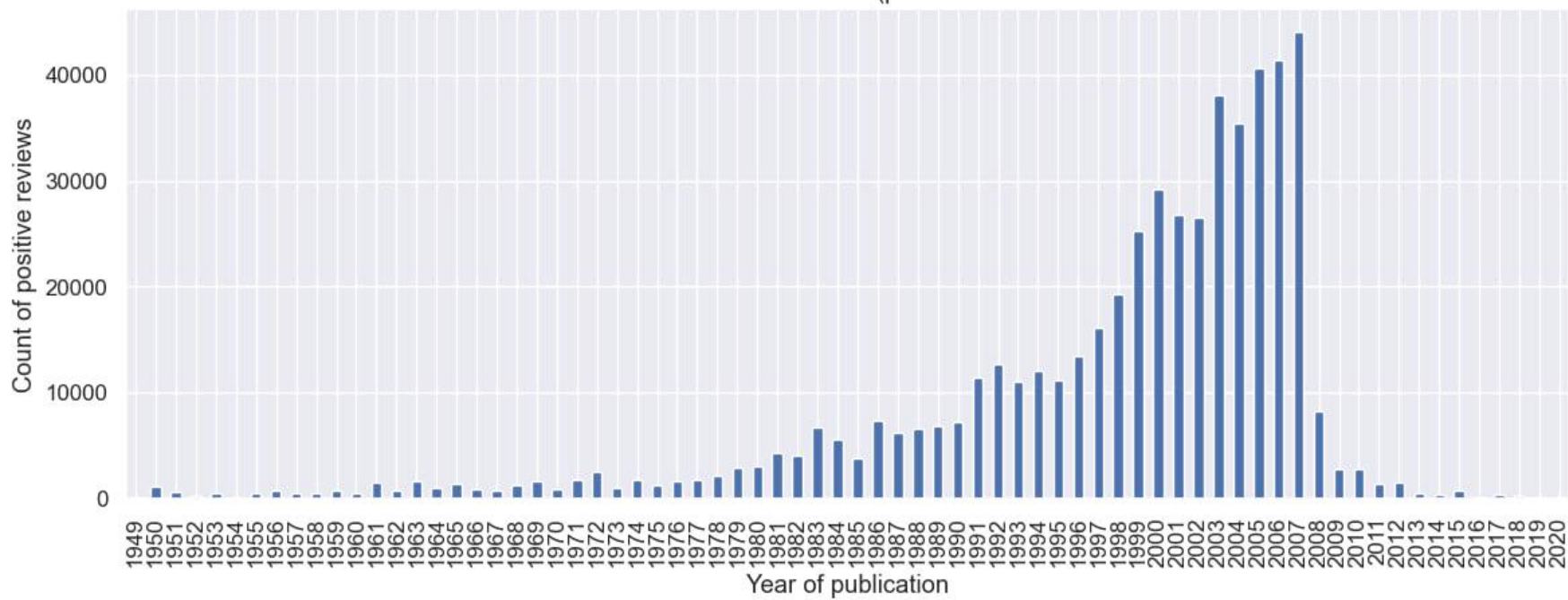
Bar chart of top 10 publishers by number of books



Recenze podle roku

- Od 90. let roste počet pozitivních recenzí.
- Důvod: nástup internetu a populárních sérií jako Harry Potter.

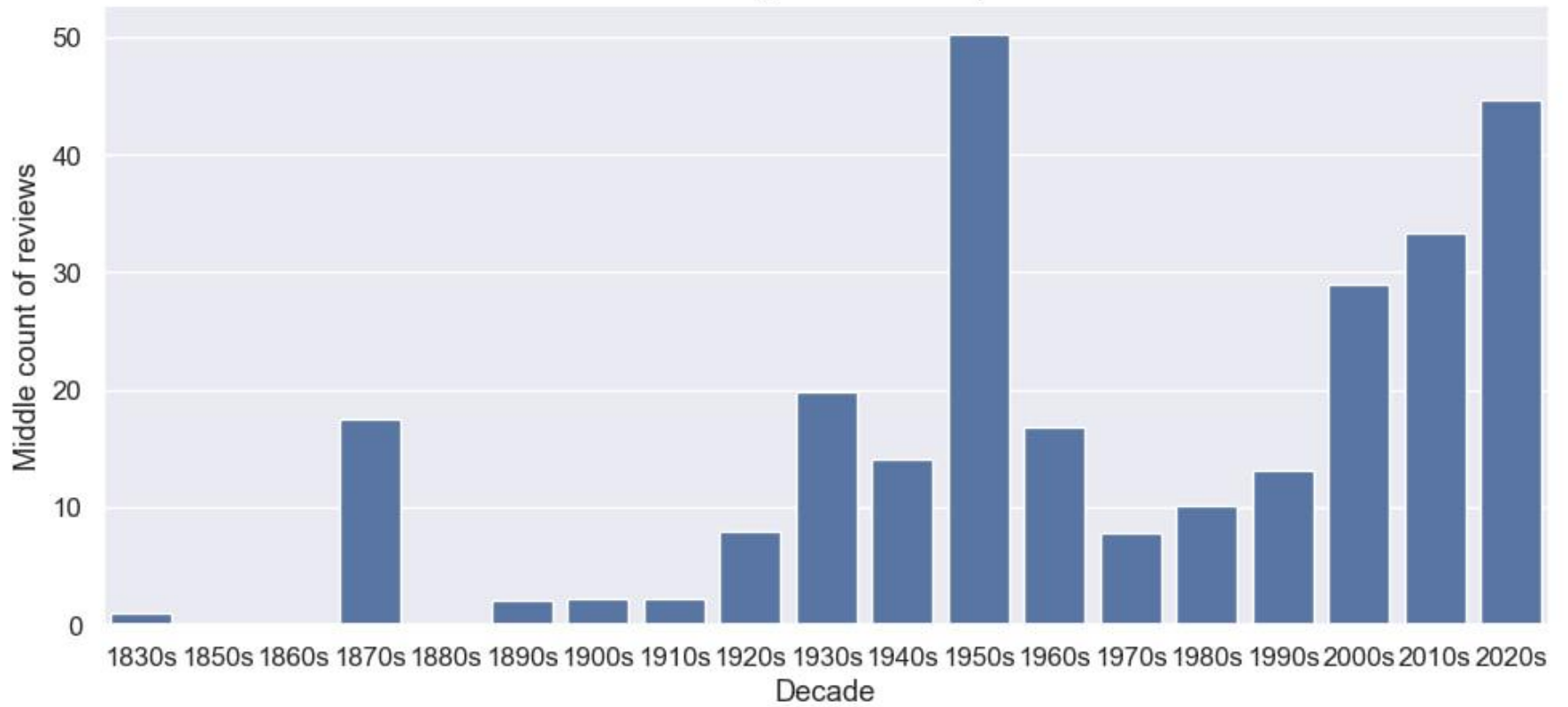
Positive reviews on books (published from 1949)



Recenze podle dekády

- Nejvíce recenzí mají knihy z let 2000–2009.
- Rozmach digitální dostupnosti a online katalogů.

Bar chart of average review count per decade



Korelační matice

- Ukazuje vztahy mezi proměnnými: počet stran, rok vydání apod.
- Pomáhá vybrat vhodné prediktory pro modelování.

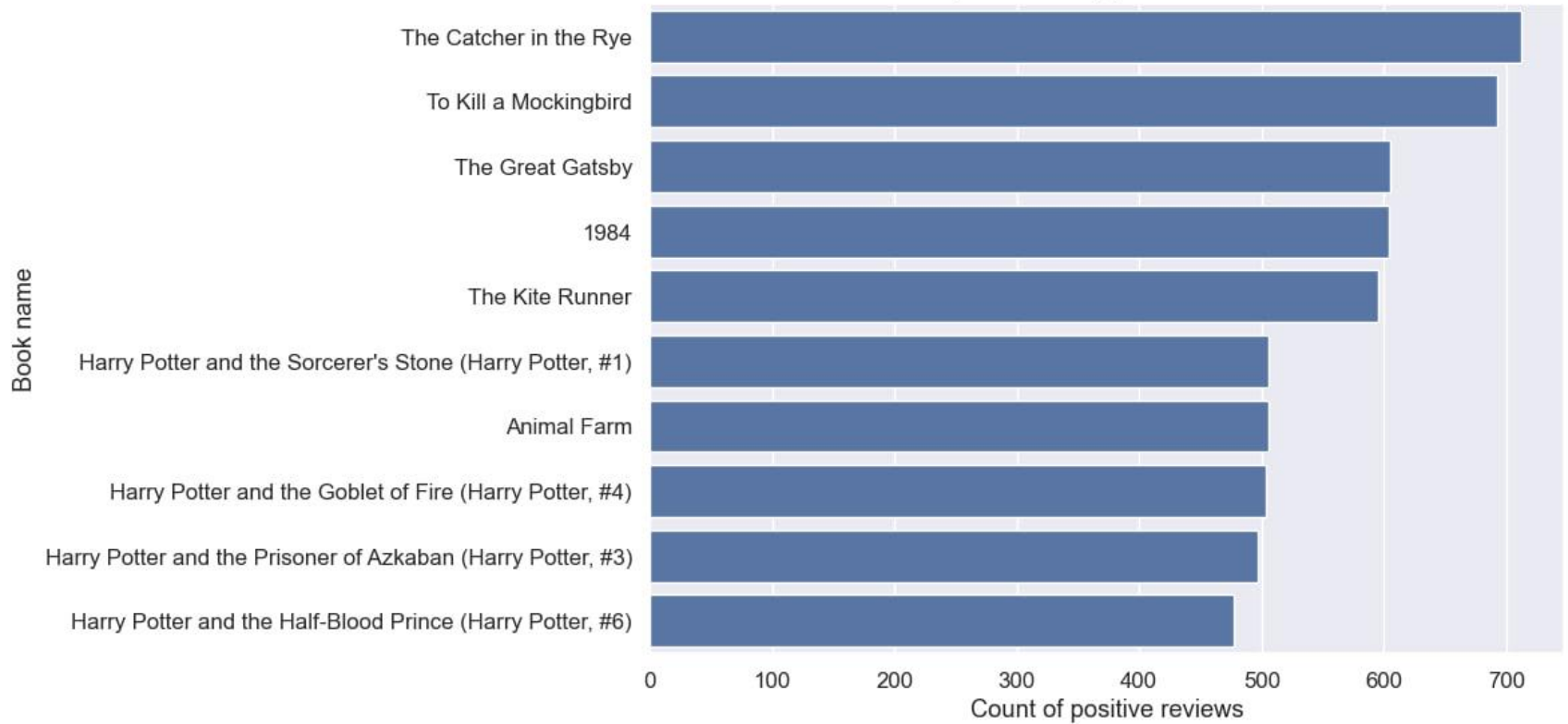
Jazyk a recenze

- Boxplot: anglické knihy mají více recenzí.
- Jazyk je významný prediktor popularity.

Autoři a hodnocení

- Známi autoři mají vyšší průměrné hodnocení.
- Čtenáři reagují na reputaci a mediální vliv.

Top-10 books by positive reviews



Výstup modelu

- Koeficienty modelu ukazují sílu vlivu každé proměnné.
- Model má stabilní chování bez výrazných chyb.

6. Analyzing positive reviews by publication year (from 1949 onward):

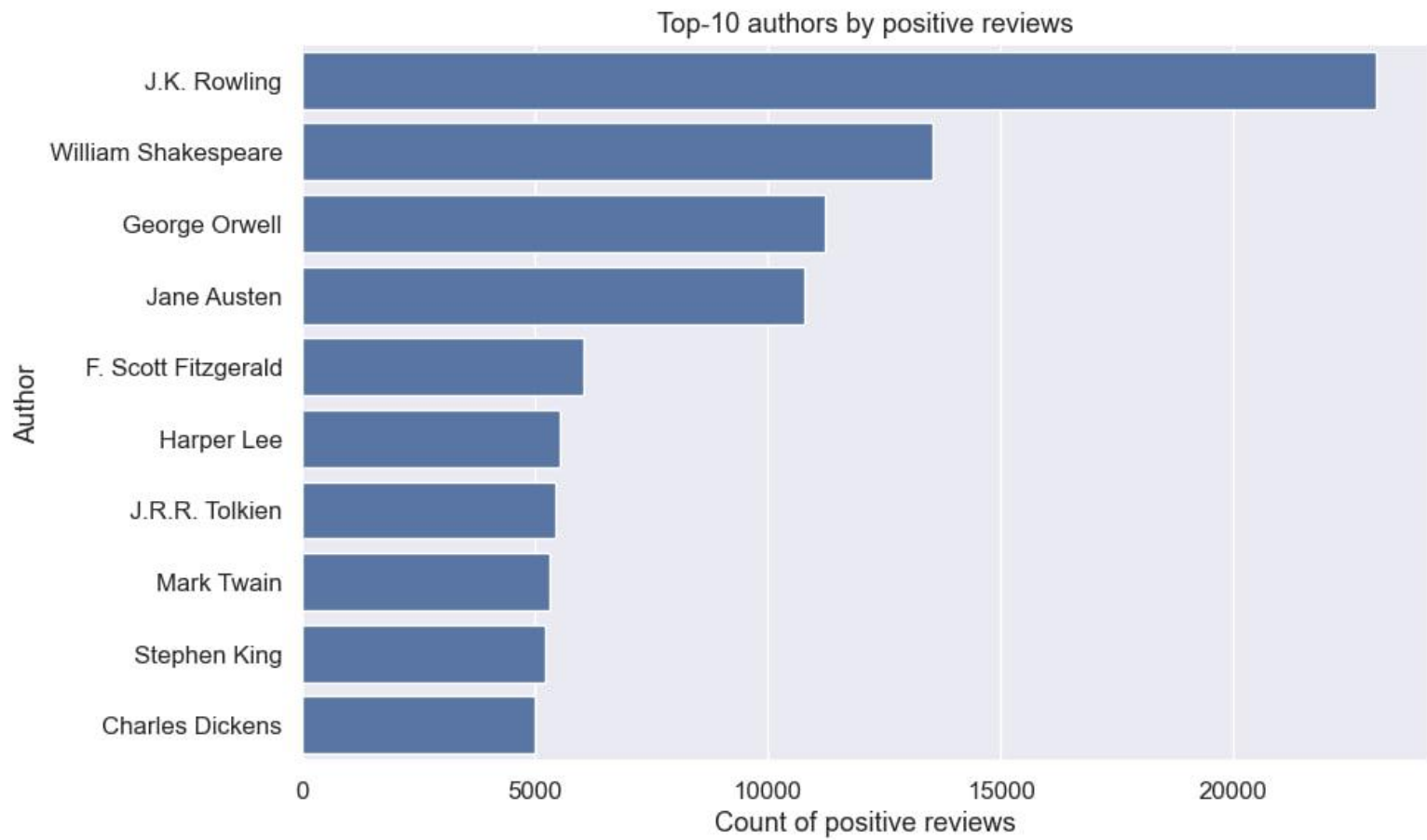
- Merges the ratings dataset with the main books dataset on book name
- Filters to include only positive reviews ('it was amazing', 'really liked it')
- Considers only books published in or after 1949
- Plots a bar chart showing the number of positive reviews by publication year

	Name	Rating \
0	Agile Web Development with Rails: A Pragmatic ...	it was amazing
1	The Restaurant at the End of the Universe (Hit...	it was amazing
2	The Restaurant at the End of the Universe (Hit...	it was amazing
3	Siddhartha	it was amazing
4	Siddhartha	it was amazing

	PublishYear	Authors
0	2005	Dave Thomas
1	2005	Douglas Adams
2	1981	Douglas Adams
3	2004	Hermann Hesse
4	2000	Hermann Hesse

Validace modelu

- Predikce a rezidua potvrzují přiměřenou přesnost.
- Model dobře vystihuje trendy popularity.



Shrnutí metody

- Lineární regrese s log-transformací počtu recenzí.
- Log pomáhá stabilizovat rozdělení a omezit vliv extrémů.
- $R^2 = 0.22$ znamená, že model vysvětluje 22 % variability.

Závěry

- Angličtina, renomovaný vydavatel a moderní knihy vedou k vyšší oblíbenosti.
- Krátké nebo velmi dlouhé knihy a neznámé jazyky mají méně recenzí.
- Model lze rozšířit o NLP, žánrové prvky nebo demografii.

Děkuji

- Děkuji za pozornost.
- Ráda zodpovím dotazy. Analýza mi ukázala sílu dat i důležitost jejich interpretace.