

Katedra informatiky
Přírodovědecká fakulta
Univerzita Palackého v Olomouci

BAKALÁŘSKÁ PRÁCE

Obrázkový dataset s atributy



2023

Vedoucí práce:
Mgr. Tomáš Mikula

Lukáš Jiříček

Studijní program: Aplikovaná informatika,
prezenční forma

Bibliografické údaje

Autor: Lukáš Jiříček
Název práce: Obrázkový dataset s atributy
Typ práce: bakalářská práce
Pracoviště: Katedra informatiky, Přírodovědecká fakulta, Univerzita Palackého v Olomouci
Rok obhajoby: 2023
Studijní program: Aplikovaná informatika, prezenční forma
Vedoucí práce: Mgr. Tomáš Mikula
Počet stran: 42
Přílohy: 1 DVD
Jazyk práce: český

Bibliographic info

Author: Lukáš Jiříček
Title: Image dataset with attributes
Thesis type: bachelor thesis
Department: Department of Computer Science, Faculty of Science, Palacký University Olomouc
Year of defense: 2023
Study program: Applied Computer Science, full-time form
Supervisor: Mgr. Tomáš Mikula
Page count: 42
Supplements: 1 DVD
Thesis language: Czech

Anotace

Práce popisuje postup vytvoření obrázkového datasetu s atributy, návrh a princip princip fungování webové aplikace, použité ke sběru dat pro tento dataset. Je kladen důraz na jednotlivé kroky nutné k přípravě tohoto datasetu. Na začátku práce je popsán postup výběru domény, volbu velikosti, obecnosti a formátů výsledného datasetu. Následně je rozepsán postup sestavení univerza atributů, způsobu sběru a prezentace dat. V poslední části práce popisuje fungování a implementaci webové aplikace, použité k hlavnímu sběru dat.

Synopsis

The thesis describes the process of creating an image dataset with attributes, the design and principle of operation of the web application used to collect data for this dataset. Emphasis is placed on the individual steps necessary to create this dataset. At the beginning of the work, the domain selection procedure, the choice of size, generality and formats of the resulting dataset are described. Subsequently, the procedure for creating a universe of attributes and the method of data collection and presentation is detailed. The last part of the thesis describes the functioning and implementation of the web application used for the main data collection.

Klíčová slova: dataset; webová aplikace; sběr dat

Keywords: dataset; web application; data collection

Děkuji vedoucímu své bakalářské práce panu Mgr. Tomáši Mikulovi, za jeho vedení a odborné rady. Dále chci poděkovat rodičům a všem, kteří mě podporovali po celou dobu studia.

Místopřísežně prohlašuji, že jsem celou práci včetně příloh vypracoval/a samostatně a za použití pouze zdrojů citovaných v textu práce a uvedených v seznamu literatury.

datum odevzdání práce

podpis autora

Obsah

1 Úvod	7
1.1 Motivace a základní koncept	8
1.2 Struktura datasetu	9
2 Vytvoření datasetu	11
2.1 Výběr domény datasetu	11
2.1.1 Proces výběru domény	12
2.1.2 Zvolená doména	14
2.2 Výběr vhodných objektů	15
2.2.1 Licence k použití obrázků	15
2.2.2 Způsob hledání obrázků	16
2.2.3 Vybrané objekty	16
2.3 Sestavení univerza atributů	17
2.3.1 Sběr atributů	18
2.4 Stanovení počtu objektů a atributů	21
2.4.1 Počet atributů	21
2.4.2 Počet objektů	21
2.5 Hlavní sběr dat	22
2.5.1 Návrh dotazníku a formuláře	23
2.5.2 Průběh sběru	25
2.6 Výsledky a prezentace datasetu	26
2.6.1 JSON	27
2.6.2 CSV	28
2.6.3 Kompletní dataset a obrázky objektů	28
3 Webová aplikace	29
3.1 Programovací část	29
3.1.1 Struktura aplikace	30
3.1.2 Databáze	31
3.1.3 Publikace	33
3.2 Uživatelská část	33
3.2.1 Vzhled aplikace	33
3.2.2 Stránky aplikace	34
Závěr	38
Conclusions	39
A Obsah přiloženého DVD	40
Literatura	41

Seznam obrázků

1	Název, popis a první otázka ve formuláři dotazníku o doménách	12
2	Počet unikátních objektů	13
3	Příklad objektu <i>stůl</i> [7]	17
4	Příklad objektu <i>židle</i> [8]	17
5	Příklad objektu <i>postel</i> [9]	17
6	Příklad objektu <i>skříň</i> [10]	17
7	Formulář dotazníku ke sběru atributů	19
8	Popis hlavního sběru dat	23
9	Část formuláře s obrázkem v hlavním sběru dat	24
10	Objekt a atributy ve formuláři hlavního sběru dat	25
11	Stránka webové aplikace s odkazy ke stažení datasetu	26
12	Schéma databáze	31
13	Barevná paleta webové aplikace a HEX kódy použitých barev	33
14	Úvod webové aplikace	34
15	Stránka s popisem dotazníku hlavního sběru dat	35
16	Dotazník webové aplikace	36
17	Stažení dat na webové aplikaci	37

Seznam tabulek

1	Příklad jednoduchého datasetu	7
2	Příklad záznamu v ukázkovém datasetu	9
3	Příklad záznamu v datasetu, rozšířeném o hodnoty	9
4	Výsledky dotazníku o doménách	13
5	Univerzum atributů	20
6	Příklad výřezu části datasetu ve formátu CSV	28

1 Úvod

Cílem této bakalářské práce bylo vytvoření obrázkového datasetu s atributy. Práce je rozdělena do dvou hlavních částí. První část, popisující kroky návrhu datasetu, jako jsou výběr domény¹ obrázků, sestavení univerza atributů², způsobu sběru dat a další. Druhá část práce, popisuje funkčnost a principy webové aplikace, vytvořené za účelem přípravy univerza atributů a hlavního sběru dat.

Za dataset můžeme považovat jakoukoli sadu dat, například databázi zdravotních záznamů, seznam potravin a jejich ingrediencí nebo tabulky dat o bankovních účtech a klientech banky.

Dataset je kolekce dat a často představuje právě tabulku, kde každý řádek obsahuje sadu dat odpovídající jedné položce datasetu a sloupce představují jednotlivé hodnoty atributů k dané položce.

<i>Jméno</i>	<i>Pohlaví</i>	<i>Věk</i>	<i>Výška</i>	<i>Váha</i>
Pavel	muž	27	182	76
David	muž	34	178	80
Ivana	žena	24	174	64
Petr	muž	48	169	92
Zuzana	žena	31	165	54
John	muž	24	191	86

Tabulka 1: Příklad jednoduchého datasetu

Za *obrázkový dataset* si můžeme představit kolekci dat (tabulku), kde každá sada dat (řádek tabulky) obsahuje obrázek a další informace k tomuto obrázku.

Časté novodobé využití datasetů se nachází ve strojovém učení a vývoji umělé inteligence. Některé metody strojového učení využívají trénovací data (dataset) pomocí kterých by měl algoritmus předpovídat výstupní hodnoty. Jako jednoduchý příklad si můžeme představit algoritmus, který má za vstupní hodnotu: výšku a váhu, výstupní hodnotu: zvíře, nebo člověk. Takový algoritmus, s dostatečně velkým množstvím trénovacích dat, by měl být schopný s určitou přesností rozlišovat na základě vstupních hodnot, zda se jedná o zvíře, nebo osobu.

¹Doména - téma objektů v datasetu, množina do které patří všechny objekty v datasetu

²Univerzum atributů - množina všech relevantních prvků, v tomto případě se jedná o množinu všech atributů v datasetu

1.1 Motivace a základní koncept

Na internetu je veřejně dostupné velké množství různých obrázkových datasetů, použitelných pro spoustu různých účelů. Motivace k vytvoření obrázkového datasetu v této práci souvisí s několika základními rozdíly, které ho odlišují od většiny ostatních:

- **Použitý jazyk** - velká většina veřejných, volně použitelných datasetů jsou dostupné v anglickém jazyce. Překlad některého takového datasetu do českého jazyka je možný, ale vždy se určitá část informací ztratí v překladu. Dataset v této práci byl navržen a vytvořen v českém jazyce, za účelem možnosti využití širší škály atributů a kvůli sběru dat od česky mluvících respondentů.
- **Způsob hlavního sběru dat** - dataset může být vytvořen různými způsoby. Mezi nejčastější volbu sběru dat, pro obrázkové datasety podobného typu patří manuální popsání obrázků malou vybranou skupinou lidí, například samotných tvůrců datasetu. Způsob zvolený u datasetu v této práci je založený na sběru dat od širší náhodné skupiny lidí, za pomoci webové aplikace.
- **Proces návrhu datasetu** - samotný návrh, vybrání domény datasetu a sestavení univerza atributů bylo uskutečněno pomocí sběru dat z dotazníků a průběžně upravován na základě sběru relevantních dat.

Dataset by měl splňovat určité požadavky. Různé datasety si mohou stanovit rozdílné požadavky, záležící na použití, typu, doméně, velikosti a dalších vlastnostech specifického datasetu. V následující části je rozepsáno jak by obrázkový dataset s atributy měl vypadat, a požadavky, které by měl splňovat, z čehož na několik jsme už narazili v předešlé části.

Každý z předešlých bodů již značí určitý požadavek na dataset a proces jeho vytvoření. Dataset musí být vytvořen v českém jazyce. Hlavní sběr musí být proveden za pomoci webové aplikace, vytvořené k tomuto účelu. A k návrhu datasetu budou použita data sesbírána od respondentů. Součástí těchto požadavků je i struktura a specifická podoba uložených dat.

Typ obrázkového datasetu s atributy, který je předmětem této práce, by měl mít základní datovou strukturu obsahující seznam jednotlivých záznamů, kde každý záznam obsahuje:

- **Obrázek s určitým objektem**, typ objektu bude záležet na zvolené doméně.
- **Atributy**, neboli vlastnosti, **popisující objekt vyobrazený na obrázku**.
- **Dodatečné informace**, jako zdroj obrázku, typ objektu a další relevantní informace.

1.2 Struktura datasetu

Následující tabulka zobrazuje příklad jednoho záznamu z ukázkového datasetu. Takový dataset zobrazuje pouze informace o atributech, které objekt na obrázku popisují.

Obrázek	Atributy
	dřevěný čtyřnohý kulatý matný nalakovaný obyčejný ornamentální hnědý

Tabulka 2: Příklad záznamu v ukázkovém datasetu

Výše uvedený příklad znázorňuje nejčastější strukturu u obrázkového datasetu. Tento přístup je jednoduchý a efektivní. Hlavní výhodou je jednoduchost sběru dat. Pokud je cílem datasetu pouze seznam atributů přiřazených k objektu na obrázku, není třeba velké skupiny respondentů, nebo je možné takový dataset vytvořit manuálně bez respondentů. Další výhodou může být větší univerzálnost použití. Na druhou stranu, očividný nedostatek této struktury pro tuto práci je malá informační hodnota a nedostatečné zameření na psychologii respondentů. Data v datasetu z této práce by měli ukázat nejen, jestli atribut popisuje, či nepopisuje objekt, ale s jakou mírou by průměrný respondent popsal daný objekt tímto atributem.

Z těchto důvodů dataset v této práci používá záznam vypadající následujícím způsobem:

Obrázek	Atributy	Hodnoty
	dřevěný	1.0
	čtyřnohý	1.0
	kulatý	0.85
	matný	0.75
	nalakovaný	0.8
	obyčejný	0.6
	škaredý	0.2
	růžový	0.0

Tabulka 3: Příklad záznamu v datasetu, rozšířeném o hodnoty

Jak je znázorněné na příkladě, dataset v této práci přidává ke každému atributu *hodnotu*. Hodnota u každého atributu značí nakolik daný atribut popisuje objekt. Taková hodnota je získána díky sběru dat od více tazatelů a znázorňuje podíl kladným odpovědí a všech odpovědí. Například u atributu *dřevěný* je hodnota 1.0, to značí, že každý tazatel, který měl rozhodnout, zda tento atribut popisuje objekt na obrázku, odpověděl kladně. U atributu *škaredý* je hodnota 0.2, značící, že 20% tazatelů odpovědělo kladně. A u atributu *růžový* kladně neodpověděl žádný tazatel. Specifický způsob sběru dat je rozepsán v sekci 2.5.

Hlavní výhoda datasetu tohoto typu je větší informační hodnota. Díky této hodnotě můžeme v datasetu použít i subjektivní atributy, neboli vlastnosti objektu, u kterých záleží na dojmu každého respondenta. Jestli je objekt na obrázku určité barvy, je možné objektivně určit, ale zda je objekt na obrázku například *škaredý*, *obyčejný*, nebo *ornamentální* záleží na subjektivních zkušenostech a představě každého respondenta. Nevýhoda tohoto přístupu je potřeba většího množství respondentů na každý atribut a více atributů na každý objekt.

Stanovení konceptu a požadavků na strukturu datasetu byl důležitý krok. Bez jasně stanovené výsledné podoby datasetu by nebylo možné pokračovat se samotným návrhem a vytvořením datasetu. Tento přístup předchází případným problémům a komplikacím v dalších částí práce.

2 Vytvoření datasetu

Po stanovení požadavků a struktury, jak by měl dataset v této práci vypadat, je možné začít s procesem vytváření tohoto datasetu. Tento proces je rozdělen do několika hlavních částí: *výběr domény datasetu*, *výběr vhodných objektů*, *sestavení univerza atributů*, *stanovení počtu objektů a atributů*, *hlavní sběr dat*, *prezentace dat*.

2.1 Výběr domény datasetu

Jelikož je tento dataset vytvořen za pomoci sběru dat od respondentů, výběr vhodné domény je jedna z nejdůležitějších částí vytvoření tohoto datasetu. Pro vybrání vhodné domény je stanoveno několik požadavků:

1. **Objekty z domény musí být jednoduše popisovatelné.** Objekty by měli být jednoznačně rozpoznatelné a respondent by měl být schopen jednoduše poznat z obrázku, jaké vlastnosti objekt má.
2. **Objekty z domény musí sdílet určité charakteristiky.** Různé objekty z domény musí sdílet dostatečné množství stejných kategorií atributů³. Ideálně by měla určitá skupina kategorií atributů popisovat všechny vybrané objekty z domény. Například domény jako *zvíře* nebo *rostlina* obsahují příliš rozdílných typů objektů, které nesdílejí dostatek charakteristik. Tento požadavek souvisí s obecností domény, ta by neměla být příliš velká, ale stále dostatečně obecná, aby bylo možné vybrat zajímavé objekty.
3. **Průměrný respondent musí mít dostatek zkušeností s vybranou doménou.** Pokud chceme aby respondent mohl popisovat objektivní i subjektivní vlastnosti objektů, je třeba aby měl s těmito objekty reálnou zkušenost. Například objekty z domén jako *letadla* a *počítačové komponenty* by průměrný respondent nemohl subjektivně popisovat, pokud s nimi nemá dostatek zkušeností.
4. **Doména by měla mít veřejně nasdílené a použitelné obrázky různých objektů.** Aby doména byla vhodný kandidát, musí být možné najít a vybrat dostatek vhodných objektů. O způsobu výběru objektů je psáno v sekci 2.2.

³Kategorie atributů - skupina atributů, popisující stejný aspekt objektu, například *barva*, *materiál*, *velikost*

2.1.1 Proces výběru domény

Po stanovení požadavků, bylo vybráno několik domén, jako vhodných kandidátů: *Budovy, Nábytek, Dopravní prostředky, Domácí potřeby, Kuchyňské potřeby a nádobí*.

Pro vybrání správné domény a získání více informací byl vytvořen dotazník pomocí platformy Google Forms.

Google Forms

Google Forms[1] je software pro správu průzkumů. Tento software umožňuje uživatelům vytvářet a editovat online dotazníky. Je možné jednoduše a rychle vytvořit formuláře a sesbírané data je možné automaticky uložit v tabulkovém procesoru. Bohužel tento software je vhodný pouze pro základní sběr dat a proto byl použit pouze u dotazníku pro pomoc s výběrem domény.

Návrh dotazníku

Pro lepší porozumění jakou má průměrný respondent představu o vybraných doménách byl vytvořen jednoduchý dotazník. Cílem tohoto dotazníku bylo zjistit kolik unikátních objektů dokáže respondent vyjmenovat, které z těchto objektů se budou opakovat a zda ty nejčastější objekty budou sdílet charakteristiky, související s požadavky na začátku této sekce. Tyto data pomůžou při výběru vhodné domény, která splňuje zadané požadavky.

Kategorie a jejich objekty

Tento krátký průzkum se věnuje přiřazování věcí do jejich kategorií. U každé kategorie vyjmenuj co nejvíce objektů, které by měla obsahovat.

[Přihlaste se do Googlu](#), abyste mohli uložit dosavadní postup. [Další informace](#)

***Povinné pole**

Vyjmenujte co nejvíce objektů, které patří do kategorie „Domácí potřeby“.
(oddělujte čárkou) *

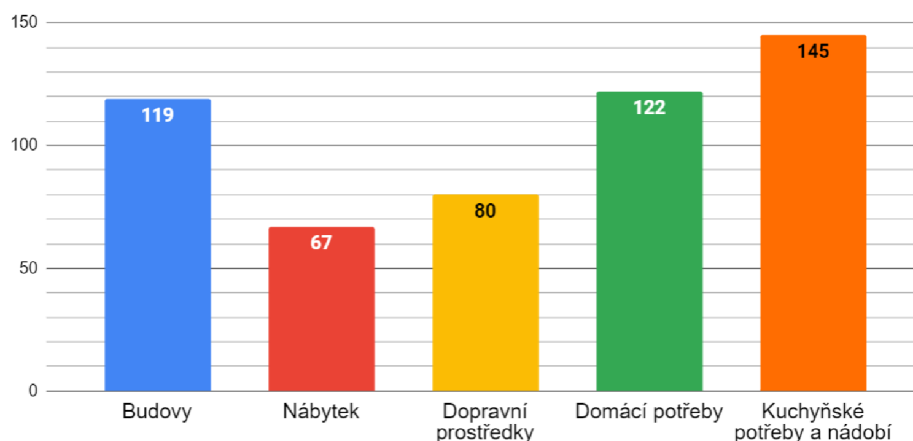
Vaše odpověď

Obrázek 1: Název, popis a první otázka ve formuláři dotazníku o doménách

V tomto dotazníku měl respondent za úkol vyjmenovat co nejvíce objektů, které patří do každé kategorie (domény). Respondent neměl žádný časový limit, nebo jiné podmínky k odeslání odpovědí.

Výsledek dotazníku

Konečný počet odpovědí byl 35. V této sekci je popis výsledků tohoto datasetu a několik grafů a tabulek se zpracovanými výsledky. Následující graf ukazuje počet unikátních objektů z těchto odpovědí u každé kategorie a tabulka ukazuje nejčastěji vyjmenované objekty s počtem výskytu v odpovědích.



Obrázek 2: Počet unikátních objektů

Budovy		Nábytek		Dopravní prostředky	
Panelový dům	28	Židle	35	Vlak	34
Dům	20	Skříň	35	Auto	34
Škola	19	Stůl	30	Autobus	33
Mrakodrap	19	Postel	30	Letadlo	32
Nemocnice	17	Křeslo	28	Kolo	31
Rodinný dům	12	Police	22	Tramvaj	27
Bytový dům	12	Komoda	17	Loď	25
Zámek	10	Gauč	17	Koloběžka	22

Domácí potřeby		Kuchyňské potřeby a nádobí	
Smeták	14	Nůž	35
Hadr	13	Talíř	33
Mop	12	Hrnec	31
Koš	12	Sklenice	28
Vysavač	11	Naběračka	27
Utěrka	9	Hrnek	27
Vědro	7	Pánev	24
Košťe	7	Vařečka	23

Tabulka 4: Výsledky dotazníku o doménách

Z těchto výsledků a následné analýze nejčastěji vyjmenovaných objektů a jak vhodné jsou tyto domény pro dataset v této práci, se došlo k následujícím poznatkům:

- **Budovy** - U domény *Budovy* respondenti vyjmenovali větší počet unikátních objektů. Respondenti se neshodli na nejčastějších objektech a už druhý nejčastěji zmíněný objekt byl vyjmenován pouze 20×. Objekty, které respondenti vyjmenali by nespĺňovali požadavky stanovené na doménu a jejich popis by byl obtížný.
- **Nábytek** - U doména *Nábytek* respondenti vyjmenovali nejmenší počet unikátních objektů. Respondenti se shodli na nejčastějších objektech a první čtyři objekty byly zmíněné velkou většinou respondentů. Objekty zmíněné, jsou lehce popisovatelné, sdílejí charakteristiky a splňují ostatní požadavky na doménu.
- **Dopravní prostředky** - Stejně jako u domény *Nábytek*, u doména *Dopravní prostředky* respondenti vyjmenovali menší počet rozdílných objektů a shodli se na nejčastějších objektech. I když by některé vyjmenované objekty splňovali stanovené požadavky, objekty z této domény jsou komplikované a jejich popis by pro případného respondenta mohl být zmatečný.
- **Domácí potřeby** - Z výsledků dotazníku lze vyčíst, že doména *Domácí potřeby* není vhodná. Respondenti se často neshodli na objektech, kvůli velmi obecnému názvu této domény. Z vyjmenovaných objektů, by bylo nemožné získat objekty rozdílné, ale s podobnou charakteristikou, které by byly vhodné pro tuto práci.
- **Kuchyňské potřeby a nádobí** - U poslední domény *Kuchyňské potřeby a nádobí* bylo vyjmenováno největší množství unikátních objektů a vysoká také shoda u nejčastěji zmíněných objektů. Tyto objekty mají společné charakteristiky a jejich popis by byl pro respondenty jednoduchý.

2.1.2 Zvolená doména

Z výsledků předešlého dotazníku a stanovených požadavků na doménu se došlo k závěru, že doména **Nábytek** bude nejvhodnější pro dataset v této práci. Hlavní důvody tohoto výběru jsou:

- Respondenti se shodli na objektech, které patří do této kategorie.
- Objekty z této kategorie sdílejí charakteristiky, ale i tak existuje velké množství rozdílných objektů.
- Objekty z této domény je lehké popsat. Respondent má s těmito objekty zkušenosti a dokáže rozpoznat jejich vlastnosti.
- Objekty z této domény mají často veřejně nasdílené obrázky, použitelné pro tento dataset.

2.2 Výběr vhodných objektů

Výběr vhodných objektů znamená výběr vhodných obrázků objektů, použitelných pro dataset v této práci. Aby byl obrázek objektu použitelný a vhodných pro tuto práci. Musí splňovat několik požadavků.

- Obrázek musí legálně použitelný, například sdílený pod některou s licencí *Creative Commons* nebo licencí *Public Domain*.
- Na obrázku musí být vyobrazený objekt z vybraných kategorií z domény datasetu. Více o těchto kategoriích je rozepsáno v této sekci.
- Objekt na obrázku musí být jasně viditelný, nesplývat s pozadím.
- Objekty na obrázku by měli mít odlišné vlastnosti a vzhled, ale sdílet charakteristiky.

Vybrané kategorie objektů

Aby objekty měly podobné charakteristiky, nebudou se vybírat všechny objekty spadající do domény datasetu, ale pouze objekty z několika vybraných kategorií. Díky výsledkům dotazníku v sekci 2.1 máme informace o nejčastějších typech objektu, které respondenti přiřčleňují k doméně. Tyto kategorie objektů jsou:

- židle,
- skříň,
- stůl,
- postel.

Tyto kategorie objektů budou zvolené jako reprezentativní kategorie objektů domény, pro dataset této práce. A všechny vybrané objekty budou patřit do těchto kategorií.

2.2.1 Licence k použití obrázků

Aby obrázky byli použitelné pro tuto práci, musí být sdílené pod licenci[2], která legálně umožňuje jejich použití. Licence je právní termín, který znamená povolení nebo oprávnění k určité činnosti, zvláštní smlouva opravňující k nakládání s autorským dílem. Kromě vlastnoručního nafocení objektů je možné použít veřejně sdílené obrázky, ale pouze, pokud jsou sdílené pod některou s veřejnou licenci. Veřejná licence umožňuje užití autorskoprávně chráněných děl. Takové užití, může být podmíněno, nebo částečně omezeno, záležící na typu licence.

V této práci jsou obrázky vybrané s licencí *Creative Commons*, nebo jako *Public Domain*. A jejich použití v této práci je chráněné a povoleno v rámci těchto licencí. Každý použitý materiál uvádí originální název a autora, i pokud to daná licence nevyžadovala.

Creative Commons

Creative Commons[3] je americká nezisková organizace nabízející škálu různých licencí. Všechny Creative Commons licence umožňují vystavovat, sdělovat, rozmnožovat a rozšiřovat díla a z něj odvozená díla při uvedení autora. Rozdílné Creative Commons licence mají navíc určité omezení vůči komerčnímu použití a dalším.

U veškerých cizích, převzatých obrázků v této práci, šířených pod jednou s licencí Creative Commons, je jejich použití v této práci chráněné pod touto licencí.

Public Domain

Public Domain[4], neboli *Volné dílo* je autorské dílo, jehož autorská práca nejsou chráněna. V kontextu této práce, se jedná o obrázky, kde se autor rozhodl, že dovolí svoje obrázky volně užívat, bez žádných podmínek nebo nároky na ochranu.

2.2.2 Způsob hledání obrázků

Pro samotné vyhledávání vhodných obrázků bylo použito několik způsobů.

Flickr

Flickr[5] je komunitní web pro sdílení fotografií. Tento web obsahuje obrovské množství fotografií všeho druhu a umožňuje vyhledávání podle názvu i licence pod kterou je fotografie sdílena.

Wikimedia Commons

Wikimedia Commons[6] je úložiště mediálních souborů, kde jsou veškeré sdílené materiály volně licencované. Autoři mohou volně sdílet jejich mediální soubory pod licencí jako *Creative Commons*, nebo úplně bez podmínek jako volné dílo. Tento web také umožňuje vyhledávání podle názvu a licence pod kterou je fotografie sdílena.

Ostatní

K vyhledávání vhodných obrázků je možné použít jakýkoli vyhledávač jako *Google Images*, který umožňuje filtrovat podle licence. Bohužel u vyhledávačů podobného typu, není toto filtrování bez chyb a je zapotřebí ověřit, zda je licence uvedena správně.

2.2.3 Vybrané objekty

I když finální počet objektů v datasetu ještě nebyl stanoven. Bylo vybráno dostatečné množství okolo 60 objektů, které budou použité i k sestavení univerza

atributů v sekci 2.3. Toto množství může být dále rozšířeno, nebo omezeno. Finální seznam obrázků, list zdrojů, názvů a autorů každého obrázku je obsažen v příloženém datovém médiu.

Na následujících obrázcích je znázorněn reprezentující objekt z každé kategorie objektů. Tyto objekty splňují všechny podmínky, které jsme si stanovili a jsou součástí datasetu.



Obrázek 3: Příklad objektu *stůl*[7]



Obrázek 4: Příklad objektu *židle*[8]



Obrázek 5: Příklad objektu *postel*[9]



Obrázek 6: Příklad objektu *skříň*[10]

2.3 Sestavení univerza atributů

Sestavení univerza atributů je velmi důležitou částí vypracování datasetu. S pojmem *univerzum atributů* jsme se už v této práci setkali a byl vysvětlen jako *množina všech relevantních prvků*, v kontextu této práce se jedná o *množinu*

všech atributů v datasetu. Proč ale takovou množinu potřebujeme, proč nemůžeme práci zjednodušit a v hlavním sběru dat nechat respondenty ručně psát vlastnosti objektů?

Proč vytvářet univerzum atributů

Vytvoření množiny atributů před hlavním sběrem dat není podmínkou pro každý obrázkový dataset. V případě datasetu v této práci je několik hlavních důvodů za tímto rozhodnutím.

Jeden z hlavních cílů hlavního sběru dat je zautomatizování celého sběru bez nutnosti zásahu v průběhu nebo při prezentaci dat. Toho lze docílit pouze, pokud atributy, na které respondenti odpovídají, jsou předem ověřené a relevantní. Pokud by respondenti sami vypisovali atributy objektů, tato automatizace by nebyla možná.

Předem ověřený seznam atributů zaručí relevantnost, ideální popis vlastností objektů a pokrytí charakteristik rozdílných typů objektů v doméně.

Požadavky na univerzum atributů

Před vytvořením této množiny je zapotřebí si stanovit cíle a požadavky:

- Sestavení univerza atributů musí být uskutečněno za pomoci sběru dat od respondentů.
- Atributy by měli popisovat větší množství rozdílných charakteristik.
- Cílem by měl být seznam atributů, rozdělených do kategorií podle charakteristik. Z většiny vytvořen sběrem dat od respondentů a ručně doplněn.

2.3.1 Sběr atributů

Pro samotný sběr atributů bude použitý výběr objektů, vytvořený v sekci 2.2 a jako platforma bude použita webová aplikace, jejíž návrh a funkce je popsána v sekci 3.

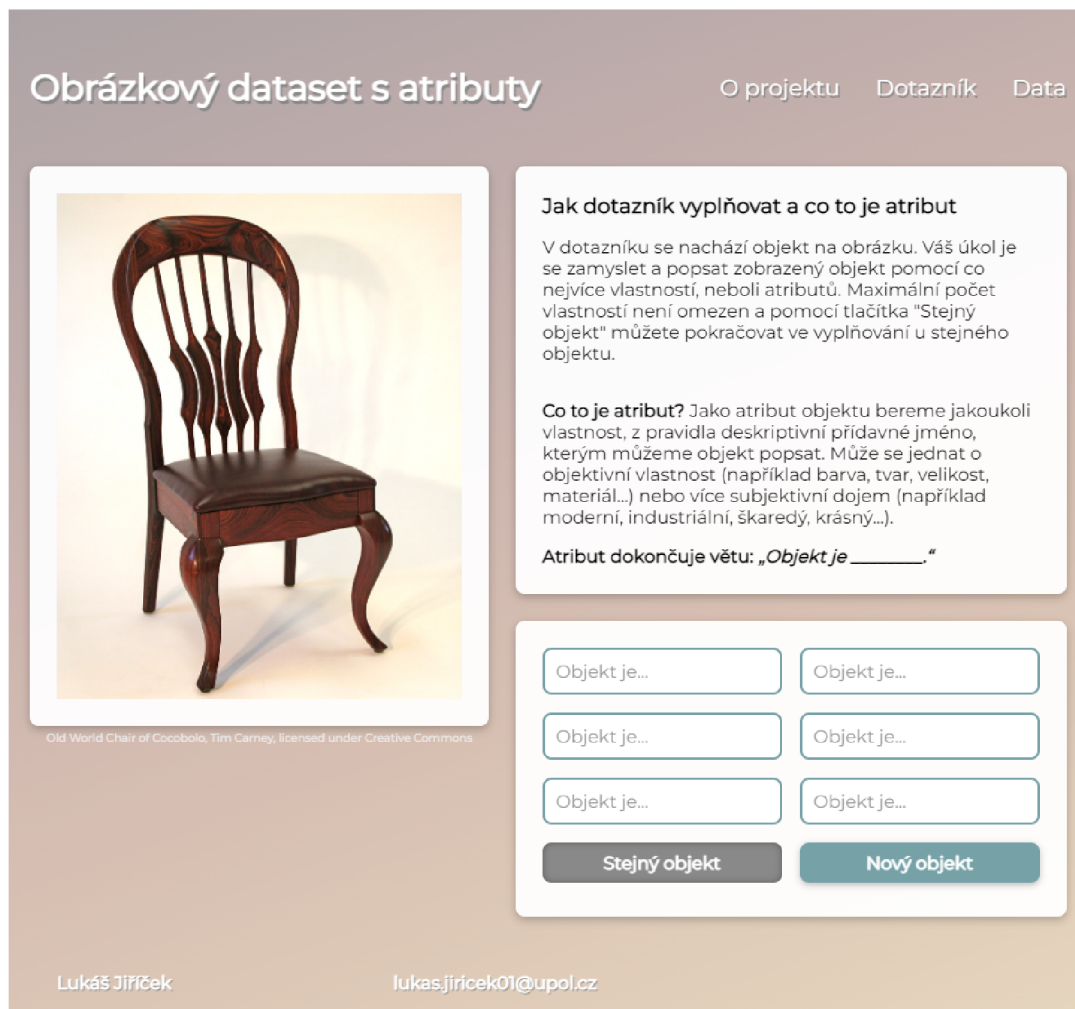
Návrh dotazníku

Ke sběru dat je vytvořen formulář. V tomto formuláři je zobrazen náhodný objekt u z datasetu s nejmenším počtem odpovědí, popis jak formulář vyplnit a několik prázdných textových boxů k vepsání atributů. U každého obrázku je uveden originální název a autor.

Atribut je ve formuláři vysvětlen jako, z pravidla deskriptivní přídavné jméno, kterým se může objekt popsat. Může se jednat o objektivní i subjektivní vlastnost.

Respondentu je vysvětleno, že jeho úkolem je popsat zobrazený objekt pomocí co nejvíce vlastností, neboli atributů. Po popsání objektu až 6 atributy

formulář odešle a má na výběr, zda popisovat nový objekt, nebo dále popisovat stejný objekt. Respondent nemá žádný časový limit nebo limit kolik vyplněných formulářů může odeslat.



Obrázkový dataset s atributy O projektu Dotazník Data

Jak dotazník vyplňovat a co to je atribut

V dotazníku se nachází objekt na obrázku. Váš úkol je se zamyslet a popsat zobrazený objekt pomocí co nejvíce vlastností, neboli atributů. Maximální počet vlastností není omezen a pomocí tlačítka "Stejný objekt" můžete pokračovat ve vyplňování u stejného objektu.

Co to je atribut? Jako atribut objektu bereme jakoukoli vlastnost, z pravidla deskriptivní přídavné jméno, kterým můžeme objekt popsat. Může se jednat o objektivní vlastnost (například barva, tvar, velikost, materiál...) nebo více subjektivní dojem (například moderní, industriální, škaredý, krásný...).

Atribut dokončuje větu: „Objekt je _____.“

Old World Chair of Cocobolo, Tim Carney, licensed under Creative Commons

Lukáš Jiríček lukas.jiricek01@upol.cz

Obrázek 7: Formulář dotazníku ke sběru atributů

Výsledky dotazníku

Průběh dotazníku dva týdny a celkem bylo sesbíráno 90 odpovědí s celkovým počtem 426 atributů. Tyto výsledky bylo nutné manuálně zpracovat následujícími základními kroky:

1. Odstranění nerelevantních odpovědí.
2. Odstranění duplicitních hodnot.
3. Oprava diakritiky a nespisovných tvarů.

Výsledný seznam bylo nutné dále zpracovat více komplikovanými kroky. První krok byl pročištění od atributů, které jsou příliš specifické k jednomu objektu,

ale nepoužitelné k žádnému jinému. Seznam atributů by měl být specifický pro vybrané kategorie objektů, ale dostatečně obecný pro všechny teoretické objekty z těchto kategorií, nejen pár vybraných objektů v tomto datasetu.

Tento výsledek ale netvoří celkové univerzum atributů. Jako finální krok je rozdělení atributů do kategorií a doplnění těchto chybějících atributů v každé kategorii. Toto doplnění je nutné, aby bylo případně možné v budoucnu popisovat i objekty, které nebyli součástí tohoto sběru, ale spadají do vybraných kategorií objektů v tomto datasetu.

Následující tabulka představuje celé univerzum, celkem 90 atributů. Zvýrazněné atributy jsou převzaty přímo z výsledků dotazníku, nezvýrazněné byly manuálně doplněné.

<i>Barva</i>	<i>Styl</i>	<i>Dojem</i>	<i>Provedení</i>	<i>Počet nohou</i>
Běžový	Elegantní	Detailní	Drsný	Jednonohý
Bílý	Exteriérový	Drahý	Lesklý	Dvounohý
Černý	Interiérový	Krásný	Matný	Třínohý
Červený	Luxusní	Levný	Nalakovaný	Čtyřnohý
Fialový	Minimalistický	Nepohodlný	Natřený	Pětinohý
Hnědý	Moderní	Nevkusný	Průhledný	
Modrý	Ornamentální	Obyčejný	Pruhovaný	
Oranžový	Přírodní	Originální	Vyřezávaný	
Růžový	Rustikální	Ošklivý	Vykládaný	
Šedý	Staromódní	Pohodlný		
Zelený	Umělecký	Zvláštní		
Zlatý	Vyzdobený			
Žlutý				

<i>Materiál</i>	<i>Umístění</i>	<i>Stav</i>	<i>Velikost</i>	<i>Tvar</i>	<i>Funkčnost</i>
Dřevěný	Dětský	Čistý	Dlouhý	Hranatý	Skládací
Kamenný	Kancelářský	Nový	Krátký	Kulatý	Nastavitelný
Kovový	Koupelnový	Odřený	Malý	Oblý	
Mramorový	Kuchyňský	Požkozený	Nízký	Ostrý	
Plastový	Nemocniční	Rezavý	Široký		
Plechový	Restaurační	Roztrhaný	Úzký		
Polstrovaný	Vojenský	Starý	Velký		
Skleněný	Zahradní	Špinavý	Vysoký		
Vyplétaný	Zámecký				

Tabulka 5: Univerzum atributů

2.4 Stanovení počtu objektů a atributů

Po sestavení univerza atributů je nutné rozmyslet a stanovit počet objektů a atributů ve finálním sběru. Tyto počty budou záviset na několika aspektech:

- Jelikož hodnoty atributů v datasetu nejsou pouze binární, je potřeba několik odpovědí na každý atribut u každého objektu.
- Jedna z podmínek této práce je zaměření na kognitivní psychologii, z toho důvodu je celý dataset zaměřen na sběr subjektivních i objektivních dojmů z objektu na obrázku a také vypracování v českém jazyce. Český jazyk má různorodé přídavné jména vyjadřující podobné dojmy. Právě tyto poznatky jsou důvodem, proč by dataset v této práci měl obsahovat spíše větší množství atributů, na úkor objektů.
- Návrh datasetu je postaven na vytvoření uceleného univerza atributů před hlavním sběrem dat, ale umožňující další přidávání objektů do datasetu i v průběhu sběru dat.
- Je potřeba určité minimum odpovědí na každý atribut, aby výsledné hodnoty byly vypovídající.

Hlavní kritérium tohoto výběru závisí na *celkovém počtu potřebných odpovědí*. Aby hodnoty atributů byly vypovídající, minimum odpovědí pro každý atribut, u každého objektu, je stanoven na 4. Způsob jak se tyto odpovědi sbírají bude popsán v sekci 2.5, ale v principu by se měl každý atribut zobrazit u každého objektu minimálně 4×. Každou takovou datovou hodnotu, budeme označovat jako *datový bod*⁴.

2.4.1 Počet atributů

Sestavení univerza atributů záleželo na doplnění dat sesbíraných od respondentů, aby vybrané univerzum dostatečně popisovalo objekty. Současný počet 90 atributů se může zdát příliš velký, ale toto bylo odůvodněno na začátku sekce 2.4. Z těchto důvodů a zvážení konceptu celého datasetu, bude počet atributů stanoven na 90.

2.4.2 Počet objektů

S počtem atributů pevně stanoven, počet objektů bude variabilní a bude záviset na způsobu hlavního sběru dat. Objektů, by mělo být co nejvíce, ale je potřeba, aby bylo možné sesbírat dostatečný počet odpovědí.

Způsob hlavního sběru dat je rozepsán v sekci 2.5, ale tento způsob je už relevantní při stanovení počtu objektů. Za jeden vyplněný formulář, neboli *odpověď*, v tomto sběru dat považujeme vyobrazený objekt a seznam atributů (v základním

⁴Datový bod - jedna odpověď, zda daný atribut popisuje objekt, u jednoho objektu

provedení se jedná o 10 atributů), kde respondent vybere zda každý z atributů buď popisuje objekt, nebo nepopisuje.

Pokud za jeden formulář o 10 attributech získáme 10 datových bodů, to je 1 datový bod pro celé univerzum atributů u daného objektu, po 9 odeslaných formulářích. Pokud minimální množství datových bodů u atributu, aby data byla odpovídající, jsou 4, je zapotřebí minimálně 36 vyplněných formulářů na jeden objekt.

Z výše uvedených důvodů je stanoven počet objektů, aby minimální počet vyplněných formulářů byl okolo 2000, neboli okolo 55 objektů, finální počet objektů byl stanoven na 57. Toto číslo bere v potaz, že respondenti budou vyplňovat několik formulářů a hlavní sběr je navržen způsobem, aby bylo možné počet objektů navýšit.

2.5 Hlavní sběr dat

Hlavní sběr dat je nejspíš ten nejdůležitější krok při tvorbě datasetu v této práci. Tento krok staví na všech přechozích přípravách a podkladech. Následující shrnutí ve zkratce představuje všechny důležité provedené kroky, které umožňují začít proces hlavního sběru dat.

Jako první bylo potřeba navrhnout strukturu a koncept datasetu. V závislosti na konceptu, následovala analýza a výběr vhodné domény. S tímto připraveným podkladem, bylo možné vybrat vhodné obrázky objektů a sestavit množinu všech atributů, pomocí kterých budou objekty popsány. A jako finální krok, před zahájením sběru, bylo potřeba stanovit počty v těchto množinách objektů a atributů, v závislosti na způsobu finálního sběru dat.

Požadavky a koncept hlavního sběru dat

Před zahájením finálního sběru dat je posledním krokem stanovení požadavků na tento sběr a navržení formuláře a procesu, jak tento sběr bude probíhat. Pro hlavní sběr dat byli stanovené tyto požadavky:

- Hlavní sběr dat musí zajistit sběr dat ke každému atributu, u každého objektu v datasetu.
- Formulář v tomto sběru by měl zobrazit jeden z objektů a list vybraných atributů.
- Respondent nebude nijak časově omezen, ani omezen v počtu odeslání formulářů.
- Výsledná data budou automaticky zpracována bez nutnosti manuálního zásahu.

Některé z těchto požadavků souvisí s vnitřním procesem webové aplikace, na které je tento sběr prováděn. Jak tyto procesy fungují je popsáno v sekci 3. Návrh formuláře a další popis samotného sběru je v této sekci.

Dataset byl navržen takovým způsobem, že čím větší množství odpovědí k hodnotám atributů, tím lepší přesnost těchto hodnot na škále ukazující jak moc by průměrný respondent daným atributem popsal objekt. Toto je důvod, proč trvání hlavního sběru dat není nijak stanoveno.

V závislosti na tomto konceptu, respondent není nijak omezen v počtu odeslaných formulářů. Kombinací objektů a atributů je velké množství a jelikož, webová aplikace vybírá zobrazené objekty a atributy podle počtu uložených odpovědí, situace při které by se respondent setkal se stejným objektem a stejnými atributy je minimální.

2.5.1 Návrh dotazníku a formuláře

Formulář hlavního sběru je sestaven ze dvou webových stránek. Úvod s popisem dotazníku a samotného formuláře, skládajícího se ze dvou částí.

Popis dotazníku

První částí, která je celým obsahem první stránky, je popis dotazníku, kde je respondentu popsáno jak formulář vypadá a způsob jak dotazník vyplňovat. Dále je zdůrazněno, že počet odpovědí ve formuláři není omezen a je možné ho vyplňovat opakovaně. Dotazník je anonymní, žádné osobní údaje nejsou ukládány a pokračováním k dotazníkům dává respondent souhlas se zapojením do sběru dat v rámci této práce.



Obrázek 8: Popis hlavního sběru dat

Obrázek objektu

Na stránce s dotazníkem se nachází samotný formulář, skládající se ze dvou částí, první částí formuláře je obrázek objektu. V každém novém vygenerovaném formuláři, je vybrán objekt s nejmenším počtem odpovědí. To zaručuje, že se data v datasetu budou plnit rovnoměrně. Pod každým obrázkem je uveden jeho originální název a autor v souladu s podmínkou Creative Commons licence.



Obrázek 9: Část formuláře s obrázkem v hlavním sběru dat

Seznam atributů

Poslední částí formuláře je vybraný seznam atributů. Aby formulář nebyl příliš dlouhý k vyplnění, ale zároveň poskytoval dostatečné množství dat, počet atributů v každém formuláři je stanoven na 12. Atributy v tomto seznamu jsou vybrané podle počtu uložených odpovědí u právě zobrazeného objektu. Jsou vždy

vybrány atributu, které se u daného objektu mají nejméně odpovědí. Tento způsob zaručí, že se data u atributů sbírají průběžně a nenastávají velké rozdíly v počtu odpovědí u jednoho atributu, na úkor jiného.

U každého atributu je možnost vybrat značku \times a \checkmark k označení, zda atribut objekt popisuje, nebo ne. Na konci seznamu je tlačítko k odeslání formuláře, které se zpřístupní po zakliknutí odpovědi u všech atributů.

„Je objekt na obrázku _____?“	
Plechový	\times \checkmark
Nízký	\times \checkmark
Široký	\times \checkmark
Rustikální	\times \checkmark
Skleněný	\times \checkmark
Ostrý	\times \checkmark
Minimalistický	\times \checkmark
Kuchyňský	\times \checkmark
Levný	\times \checkmark
Pruhovaný	\times \checkmark
Červený	\times \checkmark
Modrý	\times \checkmark

Odeslat

Obrázek 10: Objekt a atributy ve formuláři hlavního sběru dat

2.5.2 Průběh sběru

Po dokončení formuláře a funkcionalit webové aplikace s tím spojených. Byl hlavní sběr dat zahájen na webové adrese:

<http://obrazkovy-dataset.herokuapp.com/form2>

Odkaz na tuto webovou aplikaci byl sdílený na různých místech a webových stránkách. Jak bylo zmíněno v požadavcích pro tento sběr, žádné ukončení sběru není plánováno. S celkovým počtem 57 objektů, 90 atributů a 12 attributech v každém formuláři, je zapotřebí aspoň 1710 odeslaných formulářů, pro minimum 4 datových bodů u každého atributu.

Kromě uložení samotných dat o attributech a času odeslání formuláře, webová aplikace neukládala žádné osobní, nebo identifikační údaje respondenta.

Samotný sběr probíhal bez problémů a žádný změn v průběhu. Informace o výsledcích a prezentaci datasetu je uvedeno v sekci 2.6.

2.6 Výsledky a prezentace datasetu

V době psaní této práce je počet odpovědí přes 3000. Každá odpověď obsahuje několik informací: *Identifikátor objektu, Seznam atributů a odpovědí, Čas uložení formuláře*

Z těchto odpovědí je možné identifikovat a přiřadit data ke správnému objektu. Výsledky hlavního sběru dat jsou hned v průběhu automaticky ukládány a zpracovány. Prezentován je až samotný dataset v několik rozdílných formátech. Tyto soubory jsou dynamicky a pravidelně generovány ze všech aktuálně uložených odpovědí.

Všechny tyto formáty, složka s obrázky objektů a kompletní ZIP obrázků a všech formátů, jsou součástí přiloženého datového média a dostupné na adrese webová aplikace:

<http://obrazkovy-dataset.herokuapp.com/data>

Součástí každého formátu je i soubor `readme.txt`, který popisuje formátu a co každá položka znamená.

Obrázkový dataset s atributy

O projektu Dotazník Data

O datasetu

Dataset je zaměřený na okruh objektů z kategorie "Nábytek". Obsahuje 57 objektů z 4 kategorií (Postel, Židle, Stůl, Skříň) a 90 atributů z 11 kategorií, popisující objektivní a subjektivní vlastnosti. Všechny obrázky objektů byly publikovány pod licencí Creative Commons, zdroje jsou uvedeny při stažení kompletního datasetu.

Tazateli je ukázán objekt a 12 atributů. Jeho úkolem je vybrat, zda atribut popisuje objekt. Tazatel není nijak omezen v počtu vyplnění. Výsledný dataset obsahuje objekty a seznam atributů s daty z odpovědí, které atributy objekt popisují a které ne. Data v tomto datasetu představují odpovědi od různých tazatelů, kteří můžou odpovídat rozdílně. S větším množstvím odpovědí je cílem těchto dat ukázat jak by průměrný člověk objekty popsal.

Obrázky objektů Kompletní dataset

JSON

Dataset ve formátu JSON je dostupný ve dvou provedení. Kompletní data obsahují všechny informace o objektech, zdroje a odpovědi z datasetu. Zjednodušená data obsahují objekt a finální data z odpovědí. Detailní popisy obou provedení jsou součástí staženého souboru.

Stáhnout JSON

CSV

Formát CSV je dostupný ve zjednodušené formě, vhodné pro tabulkové zobrazení, třídění a řazení. Detailní popis dat součástí staženého souboru.

Stáhnout CSV

Lukáš Jiříček lukas.jiricek01@upol.cz

Obrázek 11: Stránka webové aplikace s odkazy ke stažení datasetu

2.6.1 JSON

Datový formát JSON[11] je způsob zápisu dat nezávislý na počítačové platformě. Jedná se o jeden z nejpoužívanějších formátů pro datasety. JSON umí pojmout pole hodnot, objekty a hodnoty, kterými mohou být čísla, řetězce, nebo speciální hodnoty.

Následující příklad ukazuje příklad jednoho objektu z pole všech objektů, s jeho uloženými daty o atributech ve formátu JSON:

```
1 {
2   "link": "obrazkovy-dataset.herokuapp.com/ ... .jpg",
3   "id": 1,
4   "source": "https://www.flickr.com/photos/ ... ",
5   "author": " MyBriefCase3000",
6   "type": "chair",
7   "title": "Simple Chairs. $5",
8   "filename": "Chair 1.jpg",
9   "attributes": {
10    "Bily": {
11      "Value": 0.0,
12      "true": 0,
13      "false": 4
14    },
15    "Sedy": {
16      "Value": 0.0,
17      "true": 0,
18      "false": 4
19    },
20    "Cerny": {
21      "Value": 1.0,
22      "true": 4,
23      "false": 0
24    },
25    .
26    .
27    .
28 }
```

Toto je příklad jednoho uloženého objektu s pár prvními atributy, tento list atributů obsahuje všech 90 atributů z univerza. Tento příklad je ve formátu JSON, který je jedním z dostupných formátů, ve kterém je možné dataset prohlížet.

Kromě informací o obrázku objektu a listu atributů, jsou u každého atributu uvedeny 3 hodnoty.

- **Value** - hodnota Value označuje hodnotu atributu, jak byla vysvětlená v sekci 1.2.

- **true** - označuje kolikrát respondenti označili, že tento atribut popisuje objekt.
- **false** - označuje kolikrát respondenti označili, že tento atribut nepopisuje objekt.

Dataset v tomto formátu obsahuje všechny data z datasetu, ale počet řádek celého datasetu v této podobě je přes 26000. Z toho důvodu je k dispozici i kompaktní soubor ve formátu JSON, kde nejsou uvedené všechny dodatečné informace o obrázku, objektu a hodnoty *true* a *false*. Místo nich je uvedena pouze číselná hodnota každého atributu. Tento soubor má okolo 5400 řádků.

2.6.2 CSV

Další z formátů, ve kterém je dataset dostupný je CSV[12]. Tento formát je souborový formát určený pro výměnu tabulkových dat. Takový soubor obsahuje řádky, ve kterých jsou jednotlivé položky odděleny čárkou. Díky jednoduchosti a čitelnosti je tento formát vhodný pro kompaktní zobrazení hodnot atributů u všech objektů a je ideálním formátem pro zobrazení datasetu.

Následující tabulka představuje pohled na 4 objekty a atributy z datasetu ve formátu CSV, zobrazeném jako tabulka. Zde číselné hodnoty u atributů značí jejich *hodnotu*.

<i>id</i>	<i>filename</i>	<i>link</i>	<i>Matný</i>	<i>Natřený</i>	<i>Nalakovaný</i>	<i>Nízký</i>
5	Chair 5.jpg	https://jpg	0.5	1.0	1.0	0.0
21	Table 2.jpg	https://jpg	1.0	0.0	1.0	1.0
37	Bed 4.jpg	https://jpg	1.0	1.0	1.0	0.5
52	Cabinet 8.jpg	https://jpg	1.0	1.0	1.0	1.0

Tabulka 6: Příklad výřezu části datasetu ve formátu CSV

2.6.3 Kompletní dataset a obrázky objektů

Dataset je možné taky stáhnout kompletně ve všech třech formátech najednou, kompletní a kompaktní JSON soubor a CSV soubor. K datasetu patří také složka se všemi obrázky objektů. Tato složka je ke stažení na stejné webové stránce. Kromě toho jsou tyto obrázky taky publikované na stejné webové aplikaci a jednotlivé odkazy jsou součástí JSON a CSV souborů.

3 Webová aplikace

V této sekci je rozepsáno o návrhu a implementaci webové aplikace použité ke dvěma sběrům dat a prezentaci datasetu v této práci. První sběr byl k vytvoření univerza atributů, popsaném v sekci 2.3. Druhý sběr, byl hlavní sběr dat k hodnotám atributů u každého objektu, popsaném v sekci 2.5.

Požadavky na jednotlivé části aplikace byly rozepsány ve výše uvedených sekcích. Programovací část této sekce je zaměřená na popis technologií a funkcionalit webové aplikace. Uživatelská část popisuje vzhled a funkčnost ze strany uživatele.

3.1 Programovací část

Tato sekce se bude zaměřovat na použité technologie, popis databáze a způsob publikování aplikace. Následující výčet technologií popisuje, které byly vybrány při programování této aplikace.

Flask

K vytvoření webové aplikace v této práci je použit webový framework Flask[13], napsaný v jazyce Python. Je klasifikován jako mikro webový framework, protože nevyžaduje konkrétní nástroje, neobsahuje vlastní ORM, nebo další vnitřní knihovny. Mezi jeho přednosti patří rychlý vývoj, usnadnění vývoje a volnost ve vybrání typu databáze.

Tento framework byl nainstalován pomocí návodu na oficiálních stránkách s dokumentací. Byla použita verze 2.0.1: <https://flask.palletsprojects.com/en/2.0.x/>

Jinja2

Jinja2[14] je šablonovací systém pro zjednodušení vytváření HTML souborů a nahrazování textových řetězců, při generování HTML souboru. Tento systém byl použit k nasazení dat z frameworku do formulářů.

PostgreSQL

Webový framework Flask umožňuje použít jakoukoli databázi. PostgreSQL je open-source, volně použitelný objektově-relační databázový systém, který je jeden z nejčastěji používaných ve Flask frameworku. Z těchto důvodů byl zvolen i pro tuto aplikaci.

SQLAlchemy

SQLAlchemy[15] je open-source SQL objektově-relační mapovač, neboli ORM. Jedná se o programovou vrstvu mezi relační databází a objektovými typy, která umožní pracovat s databází v aplikaci.

GitHub

GitHub[17] je webová platforma a pro hostování kódu a správy verzí softwarových projektů. Tato služba je bezplatná a používá verzovací nástroj *Git*. V tomto projektu je tato služba využita k hostování kódu, aktualizaci nových verzí aplikace v průběhu vývoje a následnou publikaci na platformu Heroku.

Heroku

Heroku[16] je cloudová platforma a webová aplikace, umožňující publikace a správu webových aplikací. Tato platforma byla zvolena ke zveřejnění a publikace webová aplikace. Tato platforma také umožňuje hostování databáze, kterou aplikace používá. Od druhé poloviny roku 2022 je použití Heroku služby zpoplatněné i pro menší projekty. Pro studenty nadále existuje možnost využití služby zdarma (více informací na <https://www.heroku.com/github-students>).

3.1.1 Struktura aplikace

Následující část popisuje jednotlivé složky a soubory webové aplikace. Webová aplikace běží ve virtuální prostředí, jak je doporučeno v instalačních instrukcích Flask frameworku: <https://flask.palletsprojects.com/en/2.1.x/installation/>

Kromě níže uvedených složek a souborů, aplikace obsahuje několik souborů `__pycache__`, `.git`, `Procfile` a `requirements.txt`. Tyto soubory jsou vygenerovány automaticky Python systémem, napojením na GitHub a publikací na Heroku službu.

download

Složka `download` obsahuje všechny stáhnutelné soubory v aplikaci. Dataset v několika formátech, ZIP soubory, sloužku s obrázky, instrukce k datům.

static

Složka `static` obsahuje všechny statické soubory v aplikaci, kromě HTML šablon. Obsahuje složku `css` s CSS soubory, `img` s obrázky všech objektů, `website_images` s ikonou aplikace.

templates

Složka `templates` obsahuje jednotlivé HTML soubory, šablony.

- `base.html` - základní HTML šablona se základní strukturou, hlavičkou, navigací, footerem a místem ke vkládání vnitřní sekce, z ostatních šablon.
- `data.html` - HTML šablona stránky `/data`, která obsahuje popis datasetu a odkazy ke stažení různých formátů datasetu.

- `form.html` - HTML šablona stránky `/form`, která obsahuje formulář použitý ke sběru atributů.
- `form2.html` - HTML šablona stránky `/form2`, která obsahuje formulář k hlavnímu sběru dat.
- `index.html` - HTML šablona uvodní stránky, která obsahuje popis projektu, popis dotazníku a kde je možné dataset stáhnout.

venv

Složka `venv` je automaticky vygenerovaná složka se soubory k Python virtuálnímu prostředí. Ke spuštění tohoto prostředí se používá skript `venv/scripts/activate`

app.py

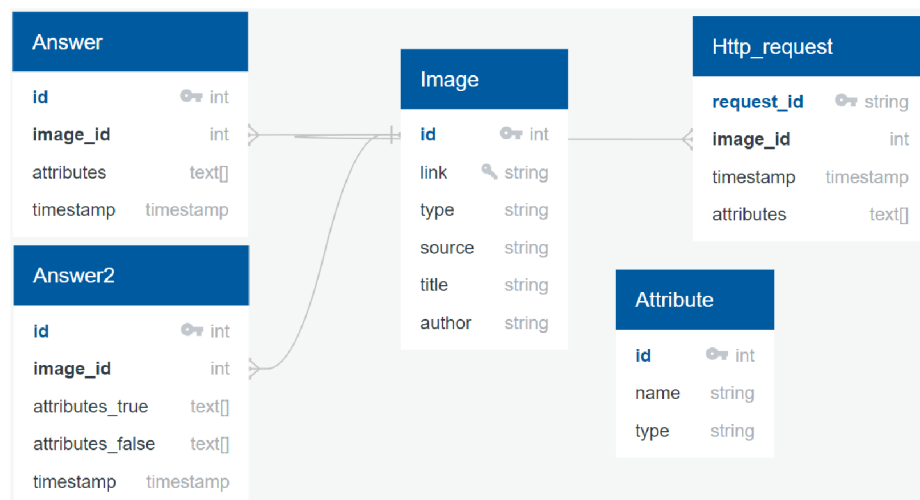
Soubor `app.py` je hlavním souborem aplikace, obsahující zdrojový kód celé aplikace. Tento soubor je rozdělen do několika částí pomocí komentářů v tomto souboru. Obsahuje všechny funkce ke každému formuláři, funkce ke vygenerování různých formátů datasetu. Routy ke zpracování a vrácení správné URL se správnými daty ke každé GET a POST metodě.

models.py

Soubor `models.py` obsahuje Python třídy, použité ke komunikaci s PostgreSQL databází.

3.1.2 Databáze

V této sekci je rozepsáno o jednotlivých tabulkách databáze této aplikace, atributech a použití.



Obrázek 12: Schéma databáze

- **image** - obrázky objektů
 - id - identifikátor objektu
 - link - odkaz na obrázek ve webové aplikaci
 - type - typ objektu
 - source - zdroj obrázku
 - title - originální název obrázku
 - author - autor obrázku
- **answer** - odpovědi ze sběru k vytvoření univerza atributů
 - id - identifikátor odpovědi
 - image_id - indentifikátor obrázku ve formuláři
 - attributes - seznam atributů
 - timestamp - čas uložení odpovědi
- **answer2** - odpovědi z hlavního sběru dat
 - id - identifikátor odpovědi
 - image_id - indentifikátor obrázku ve formuláři
 - attributes_true - seznam atributů, označených kladně
 - attributes_false - seznam atributů, označených negativně
 - timestamp - čas uložení odpovědi
- **http_request** - seznam aktivních, vygenerovaných formulářů, čekajících na vyplnění a uložení. Sloužící k ověření zda přichozí data odpovídají odeslanému formuláři.
 - request_id - identifikátor formuláře
 - image_id - indentifikátor obrázku ve formuláři
 - timestamp - čas uložení odpovědi
 - attributes - seznam atributů ve formuláři
- **attribute** - atributy v datasetu
 - id - identifikátor atributu
 - name - atribut
 - type - typ atributu

3.1.3 Publikace

Po naprogramování webové aplikace, bylo nutné tuto aplikaci veřejně publikovat. K tomuto byla by použita služba Heroku. Na této službě je i spuštěna PostgreSQL databáze.

Kód webové aplikace je hostován na webové platformě GitHub, odkaz na repositáři v této službě je https://github.com/LukeJiricek/bakalarska_prace. Tento kód je potom nahrán na cloud server služby Heroku, kde je tato aplikace publikována. Instrukce, specifické k Python projektu, k použití této služby se nachází na této adrese: <https://devcenter.heroku.com/articles/getting-started-with-python>

Po dokončení potřebných kroků ke spuštění, je webová aplikace publikována na následující URL adrese:

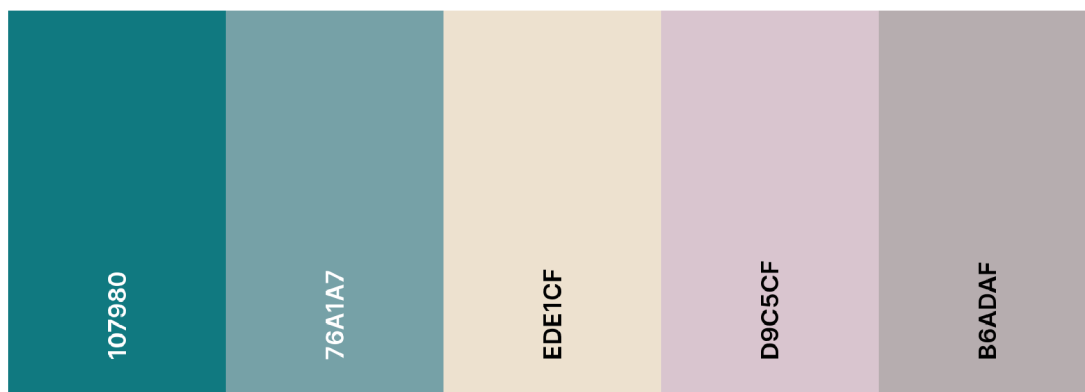
`https://obrazkovy-dataset.herokuapp.com/`

3.2 Uživatelská část

V této části je rozepsáno jak fungují jednotlivé části aplikace a jak s aplikací uživatel pracuje.

3.2.1 Vzhled aplikace

Vzhled aplikace byl navržen aby byla uživatelsky přívětivá. Aplikace je responzivní a je plně funkční na mobilních zařízeních. Při návrhu designu bylo rozhodnuto, že je potřeba použít přívětivou barevnou paletu, aby byla stránka respondentu příjemná na pohled. Bylo zvoleno následující barevné schéma.



Obrázek 13: Barevná paleta webové aplikace a HEX kódy použitých barev

Z této barevné palety byl vytvořen barevný gradient primárních barev béžové a šedé, který byl použit pro pozadí aplikace. Tato kombinace barev by neměla odvádět pozornost od objektů, ale zároveň navodit přívětivý dojem. Sekundární barvy byly použity k dodatečným prvkům, jako jsou tlačítka, podtržení a zvýraznění některých prvků.

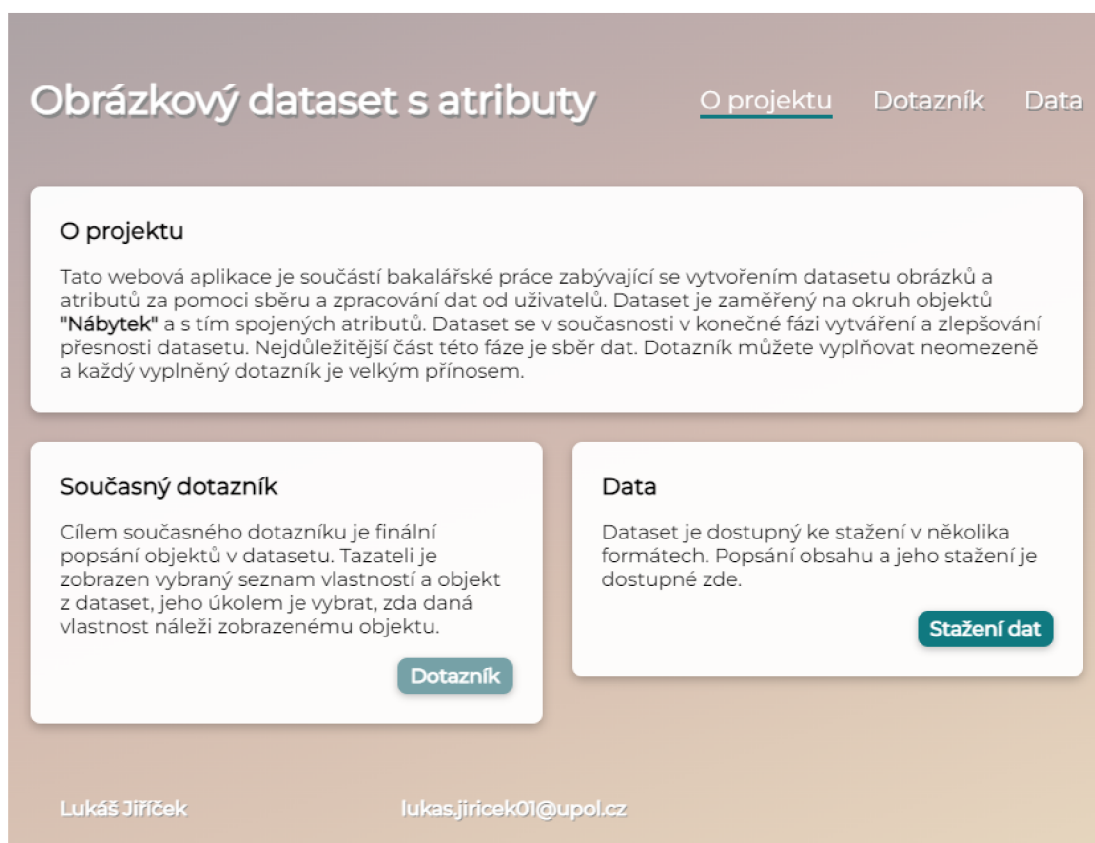
V samotných formulářích je součástí interaktivní zaklikávání odpovědí u každého atributu. Tlačítka a odesílání formuláře je navrženo responsivně.

Hlavní koncept celého vzhledu je založen na minimalismu, který je příjemný na pohled. To celé za účelem zaměření pozornosti respondenta na objektu a udržení uživatele na stránce co nejdéle, pro vyplnění co nejvyššího počtu formulářů.

3.2.2 Stránky aplikace

Úvod

Na úvodní stránce aplikace je uživatel srozuměn s projektem a popisem aktuálního dotazníku. Z této stránky má možnost se dostat k dotazníku, nebo ke stažení datasetu.



Obrázek 14: Úvod webové aplikace

Popis dotazníku

Po kliknutí na tlačítko k přesměrování na dotazník, je respondent zobrazena stránka s popisem dotazníku. Na této stránce je respondent seznámem s dotazníkem, jak bylo posáno v sekci 2.5.1. Tlačítkem pod tímto popisem se může uživatel dostat k formuláři hlavního sběru.



Obrázek 15: Stránka s popisem dotazníku hlavního sběru dat

Dotazník

Při přesměrování (server zaznamená GET metodu) na stránku s dotazníkem, je uživateli vygenerován nový formulář. Data, jako identifikační číslo, objekt a seznam atributů, která byla vygenerována se uloží v databázi, v tabulce `http_request`.

Uživatel může kliknutím označit svoji odpověď u každého atributu. Až po označení odpovědi u každého atributu, je tlačítko *Odeslat* zpřístupněno.

Po kliknutí tlačítka *Odeslat* (server obdrží POST metodu), je formulář a jeho data odeslána. Server ověří, zda tento formulář je stále aktivní, neboli, zda nebyl skrytý časový limit překročen (databáze maže vygenerované formuláře po 30 minutách). Dále také ověří, zda identifikační číslo, objekt a atributy odpovídají uložené kopii v databázi.

Tato kontrola zaručí, že nelze formulář odeslat duplicitně, manuálně zaměnit objekt nebo atributy. Pokud některá z těchto kontrol selže, formulář nebude uložen a tato chyba bude uživateli sdělena. Pokud kontrola proběhne v pořádku, stránka se aktualizuje a uživateli je vygenerován další formulář.

Obrázkový dataset s atributy

O projektu [Dotazník](#) [Data](#)

„Je objekt na obrázku _____?“

<input checked="" type="checkbox"/> Odřený	<input checked="" type="checkbox"/> Vykládaný
<input checked="" type="checkbox"/> Skleněný	<input checked="" type="checkbox"/> Pruhovaný
<input checked="" type="checkbox"/> Zelený	<input checked="" type="checkbox"/> Nízký
<input checked="" type="checkbox"/> Restaurační	<input checked="" type="checkbox"/> Nový
<input checked="" type="checkbox"/> Kamenný	<input checked="" type="checkbox"/> Nepohodlný
<input checked="" type="checkbox"/> Zámecký	<input checked="" type="checkbox"/> Zahradní

Odeslat

SOLD: Really old wooden chair, TheLivingRoominKenmore, licensed under Creative Commons

Lukáš Jiríček lukas.jiricek01@upol.cz

Obrázek 16: Dotazník webové aplikace

Data

Finální stránka webové aplikace nabízí popis datasetu, stručný popis způsobu hlavního sběru dat a především tlačítka umožňující stažení složky s obrázky objektů a datasetu v rozdílných formátech.

Všechny formáty datasetu jsou periodicky generovány, každých 15 minut, z právě aktuálních dat. Tento proces probíhá souběžně s hlavním sběrem dat.

Obrázkový dataset s atributy O projektu Dotazník Data

O datasetu

Dataset je zaměřený na okruh objektů z kategorie "**Nábytek**". Obsahuje 57 objektů z 4 kategorií (Postel, Židle, Stůl, Skříň) a 90 atributů z 11 kategorií, popisující objektivní a subjektivní vlastnosti. Všechny obrázky objektů byly publikovány pod licencí Creative Commons, zdroje jsou uvedeny při stažení kompletního datasetu.

Tazateli je ukázán objekt a 12 atributů. Jeho úkolem je vybrat, zda atribut popisuje objekt. Tazatel není nijak omezen v počtu vyplnění. Výsledný dataset obsahuje objekty a seznam atributů s daty z odpovědí, které atributy objekt popisují a které ne. Data v tomto datasetu představují odpovědi od různých tazatelů, kteří můžou odpovídat rozdílně. S větším množstvím odpovědí je cílem těchto dat ukázat jak by průměrný člověk objekty popsal.

Obrázky objektů Kompletní dataset

JSON

Dataset ve formátu JSON je dostupný ve dvou provedení. Kompletní data obsahují všechny informace o objektech, zdroje a odpovědi z datasetu. Zjednodušená data obsahují objekt a finální data z odpovědí. Detailní popisy obou provedení jsou součástí staženého souboru.

Stáhnout JSON

CSV

Formát CSV je dostupný ve zjednodušené formě, vhodné pro tabulkové zobrazení, třídění a řazení. Detailní popis dat součástí staženého souboru.

Stáhnout CSV

Lukáš Jiríček lukas.jiricek01@upol.cz

Obrázek 17: Stažení dat na webové aplikaci

Závěr

V této práci byl popsán proces vytvoření obrázkového datasetu s atributy. Na začátku práce byl čtenář seznámen s motivací, strukturou a požadavky na tento dataset. Následně byl proces vytvoření rozdělen do 6 kroků které byly podrobně popsány. V jednotlivých krocích byl kladen důraz na principy kognitivní psychologie a celý návrh byl ovlivněn těmito principy a souběžným sběrem dat od respondentů.

Pro vytvoření množiny atributů a hlavního sběru dat, byla vytvořena a popsána webová aplikace s důrazem na responzivní design a jednoduchost použití. V práci bylo rozepsáno o technologiích, vzhledu aplikace a postupu jak ji uživatelé mohou použít.

Conclusions

In this thesis, the process of creating an image dataset with attributes was described. At the beginning of the work, the reader was introduced to the motivation, structure and requirements for this dataset. Subsequently, the creation process was divided into 6 steps that were described in detail. In individual steps, emphasis was placed on the principles of cognitive psychology, and the entire design was influenced by these principles and the simultaneous collection of data from respondents.

To create a set of attributes and the main data collection, a web application was created and described with an emphasis on responsive design and ease of use. The work described the technologies, the appearance of the application and how users can use it.

A Obsah příloženého DVD

flask/

Kompletní zdrojový kód webové aplikace. Struktura této aplikace je popsána v sekci [3.1.1](#).

dataset.zip

Kompletní dataset, ve všech dostupných formátech. Obsahující složku všech obrázků objektů v datasetu `img.zip`. A `readme.txt` popisující jednotlivé formáty datasetu.

doc/

Text práce ve formátu PDF, vytvořený s použitím závazného stylu KI PřF UP v Olomouci pro závěrečné práce, včetně všech příloh, a všechny soubory potřebné pro bezproblémové vygenerování PDF dokumentu textu (v ZIP archivu), tj. zdrojový text textu, vložené obrázky, apod.

U veškerých cizích převzatých materiálů obsažených na médiu jejich zahrnutí dovoluují podmínky pro jejich šíření nebo přiložený souhlas držitele copyrightu. Pro všechny použité (a citované) materiály, u kterých toto není splněno a nejsou tak obsaženy na médiu, je uveden jejich zdroj (např. webová adresa) v bibliografii nebo textu práce nebo v souboru `readme.txt`.

Literatura

- [1] GOOGLE FORMS,About [online] [cit. 20.06.2022] Dostupné z: <https://www.google.com/forms/about/>
- [2] WIKIPEDIA, The Free Encyclopedia: Licence [online] [cit. 27.06.2022] Dostupné z: <https://cs.wikipedia.org/wiki/Licence>
- [3] CREATIVE COMMONS, O licencích [online] [cit. 27.06.2022] Dostupné z: <https://creativecommons.org/licenses/?lang=cs>
- [4] WIKIPEDIA, The Free Encyclopedia: Public domain [online] [cit. 27.06.2022] Dostupné z: https://en.wikipedia.org/wiki/Public_domain
- [5] FLICKR, About Flickr [online] [cit. 27.06.2022] Dostupné z: <https://www.flickr.com/about>
- [6] WIKIMEDIA COMMONS, Welcome [online] [cit. 27.06.2022] Dostupné z: <https://commons.wikimedia.org/wiki/Commons:Welcome>
- [7] FLICKR, Cherry Gate Leg Table, Tim Carney [online] [cit. 27.06.2022] Dostupné z: <https://www.flickr.com/photos/23622770@N00/4024999320>
- [8] WIKIMEDIA COMMONS, Armchair (fauteuil), Metropolitan Museum of Art [online] [cit. 27.06.2022] Dostupné z: [https://commons.wikimedia.org/wiki/File:Armchair_\(fauteuil\)_MET_DP276258.jpg](https://commons.wikimedia.org/wiki/File:Armchair_(fauteuil)_MET_DP276258.jpg)
- [9] PIXABAY, Bunkbed Metal Mattress, Jazella [online] [cit. 27.06.2022] Dostupné z: <https://pixabay.com/images/id-5004237/>
- [10] FLICKR, kitchen ikea cabinet, Rrin Williamson [online] [cit. 27.06.2022] Dostupné z: <https://www.flickr.com/photos/erinwilliamson/3254054404/>
- [11] JSON, Introducing JSON [online] [cit. 29.06.2022] Dostupné z: <https://www.json.org/json-en.html>
- [12] WIKIPEDIA, The Free Encyclopedia: CSV [online] [cit. 29.06.2022] Dostupné z: <https://cs.wikipedia.org/wiki/CSV>
- [13] PYTHON TUTORIAL, What is Flask Python: Flask [online] [cit. 29.06.2022] Dostupné z: <https://cs.wikipedia.org/wiki/Flask>
- [14] JINJA [online] [cit. 30.06.2022] Dostupné z: <https://jinja.palletsprojects.com/en/3.1.x/>
- [15] SQLALCHEMY, Home [online] [cit. 30.06.2022] Dostupné z: <https://www.sqlalchemy.org/>

- [16] HEROKU, About [online] [cit. 30.06.2022] Dostupné z: <https://www.heroku.com/about>
- [17] GITHUB, Hello World, About [online] [cit. 19.07.2022] Dostupné z: <https://docs.github.com/en/get-started/quickstart/hello-world>