

MUNI
FI

Etické nakládání s daty

Mgr. Tomáš Foltýnek, Ph.D.

foltynek@fi.muni.cz

Osnova dnešní přednášky

- Opakování: Základní pojmy; Morální dilema
- Dnešní téma: Etické otázky sběru, uchovávání a využívání dat
- Etika internetového výzkumu
- GDPR
- Základní etické normy
- Výzkum nad daty vs. výzkum na lidech

Opakování: Základní pojmy

- Etika
 - Filosofická disciplína zkoumající morálku a morální hodnoty
 - Rozprava o tom, co je **správné**
 - Zkoumá principy rozhodování v situacích ovlivnitelných pomocí svobodné vůle
- Morálka
 - Individuální morálka = souhrn individuálních přesvědčení a postojů **o tom, co je správné** (v souladu se svědomím)
 - Skupinová morálka = společenský konsenzus **o tom, co je správné**
- Právo = Řád, soubor pravidel a nařízení
 - Nezbytné dobro k zachování funkce formální organizace

Opakování: Jak poznat, co je morální?

- Deontologie
 - Stanovení nepřekročitelných pravidel a povinností
 - Teorie božích přikázání, Kantův kategorický imperativ
- Konsekvencialismus
 - Při posuzování situace jsou relevantní pouze důsledky činů
 - Čisté blaho a jeho distribuce ve společnosti
- Teorie společenské smlouvy
 - Pravidla fungování v rámci společnosti, práva a povinnosti
 - Absolutní a omezená práva; Pozitivní a negativní práva
- Etika ctností
 - Morální charakter jednotlivce

Opakování: Morální dilema

- Je situace vyžadující rozhodnutí mezi možnostmi, z nichž každá znamená jednání v rozporu s morálními hodnotami
- Typicky jsou možnosti v rozporu s různými hodnotami
 - Řešení dilematu vyžaduje zvážení relativní důležitosti těchto hodnot

Etika internetového výzkumu



Výzkum nad daty Dropboxu

- Článek “[Was it Ethical for Dropbox to Share Customer Data with Scientists?](#)”
 - Anonymizovaná data 16 000 výzkumníků z 1000 kateder/ústavů/oddělení
 - Informace o složkách (sdílení, využívání)
 - Informace o vědeckém výkonu (počet článků, počet citací)
- Cíl [studie](#): Zjistit, jak spolupracují nejlepší týmy
 - Pracují v malých týmech (průměr 2,3 členů, vs. 3,0 u nejhorších)
 - Pracují na projektech déle (průměr 172 dnů vs. 130 u nejhorších)
 - Udržují stabilní týmy (průměr 5 projektů vs. 3,5 u nejhorších)
 - Rovnoměrné rozdělení práce mezi členy týmu (nikdo se „neveze“)
 - Zkušenější pracují více (63 % vs. 58 % u nejhorších)
- Bylo etické data sdílet?

Tastes, Ties and Time (3T) Dataset

- 2008: [Výzkumníci z Harvardu](#) publikovali [dataset](#) 1700 Facebookových profilů
 - Profily patřily studentům z Harvardu; anonymizované, kódované
 - Ukazuje vývoj sociálních sítí studentů v průběhu studia
 - Veřejně dostupná jen kódová kniha; dataset na vyžádání
- Obrovský potenciál pro sociální vědce
- Problémy
 - Nikdo ze studentů nedal ke sběru dat souhlas
 - Demografická data, obor a další údaje umožňovaly deanonymizaci
 - A tedy zveřejnění informací určených jen pro přátele
 - Některé údaje byly přístupné jen uživatelům ze stejné školy
 - Jejich zveřejnění by tedy porušilo uživatelská nastavení
- Dataset nakonec nebyl zveřejněn
- Podrobný rozbor viz Zimmer, M. (2010). “But the data is already public”: on the ethics of research in Facebook. *Ethics and Information Technology*, 12(4), 313–325. <https://doi.org/10.1007/s10676-010-9227-5>

OkCupid Dataset

- 2016 dánští výzkumníci [Kirkegaard & Bjerrekaer publikovali dataset](#) uživatelů seznamovací stránky OkCupid
 - Data získali pomocí web scraperu
 - Uživatelské jméno, věk, pohlaví, sexuální orientace, město
 - Odpovědi na otázky o náboženství, politických názorech, drogy, nevěra, svazování během sexu...
- Všechny údaje byly dobrovolně vloženy uživateli s vědomím, že budou veřejně dostupné
 - Výzkumníci je poskytli v „užitečnější formě“
- Výzkumníci neměli ke sběru dat souhlas platformy, uživatelů, ani etické komise univerzity

Otázky do diskuse

- Bylo **dobré** zveřejnit data?
- Jaké hodnoty stojí proti sobě?

- Různá pojetí soukromí
 - Způsobí zveřejnění informací nějakou **újmu**?
 - Ohrozí zveřejnění informací **důstojnost**?

- Soukromí (a další hodnoty) **závisí na kontextu**

GDPR

- **G**eneral **D**ata **P**rotection **R**egulation
- Právní rámec EU pro ochranu osobních údajů
- Osobní údaj = každý údaj týkající se **identifikované nebo identifikovatelné** osoby
- Pseudonymizace (kódování) nemusí zamezit identifikovatelnosti!
- **Anonymní** informace pro statistické nebo výzkumné účely nepodléhají GDPR

Zásady zpracování osobních údajů

- Zákonnost – musí existovat důvod pro zpracovávání dat
- Transparentnost – subjekt musí vědět, jaká data se uchovávají
- Omezení účelu – data nesmí být využívána v rozporu s tímto účelem
- Omezení uložení – uložení ve formě umožňující identifikaci
- Minimalizace – uložení jen nezbytných dat
- Integrita, přesnost, důvěrnost

”Anonymizovaná” geolokace

- Studie (2020) ukázala, že z anonymizovaných dat o poloze je možné získat údaje o bydlišti a zaměstnání s 98% přesností
 - Rada: Šifrovat ID uživatelů pomocí soli, která se několikrát denně mění → Rozbití sekvencí polohových dat
- 2013: New York City Taxi & Limousine Commission
 - Zveřejnění datasetu o 173 milionech jednotlivých jízd
 - Čas a místo začátku a konce, jízdné, spropitné
 - ID taxikářů hashována
 - Anonymita datasetu byla jen zdánlivá (viz článek v Big Data & Society)
 - Deanonymizace taxikářů → Informace o jejich příjmech
 - Spojení s dalšími zdroji (blogy, články,...) → Sledování celebrit
 - Odvození adres lidí navštěvujících strip kluby

Etické otázky pro nakládání s daty

- 10 otázek podle Data Science Process Alliance
- Jaké zákony a předpisy se vztahují na daný projekt?
 - Jaké důsledky může mít jejich nedodržení?
- Kdo je zodpovědný za minimalizaci etických rizik?
 - Proaktivní vyhodnocování možných rizik
- Jak se využití dat dotkne práv jednotlivců?
 - Jaká práva se vztahují na data (jejich sběr, shromažďování a využití)
- Jak se propojování dat dotkne soukromí a anonymity jednotlivců?
 - Jak bude zajištěno soukromí při přenosu a uložení dat?
 - Jak zabránit deanonymizaci při propojení s jinými zdroji?

Etické otázky pro nakládání s daty

- Jak víme, že je etické data využít k danému účelu?
 - Sběr dat a přístup k nim ještě neznamena, že je etické je využít
 - Kdo vlastní data? Jaká jsou jejich práva a očekávání?
 - Jsou data využívána v souladu s očekáváním uživatelů?
- Jak víme, že data jsou k danému účelu vhodná?
 - Zejména pokud byla sesbírána za jiným účelem
 - Doplnování chybějících hodnot, čištění dat... → ovlivnění výsledků
- Jak rozpoznat a snížit zkreslení v datech?
 - Strojové učení ze zkreslených dat → Zkreslený model

Etické otázky pro nakládání s daty

- Jak zmírnit subjektivitu při návrhu systému?
 - Subjektivní rozhodnutí: Jakou použít metriku? Jaký algoritmus?
Z jakých zdrojů získat data?
- Jak velká transparentnost je potřeba a jak ji dosáhnout?
 - Hrozí znevýhodňování určité skupiny osob?
 - Bude potřeba vysvětlovat rozhodnutí systému (např. půjčky)
- Jak pravděpodobný je mylný výklad výsledků a jak mu zabránit?
 - Většina modelů je z podstaty pravděpodobnostní → Nedává záruku
 - Ukazují na vyšší pravděpodobnost určitého jevu
 - Korelace vs. kauzalita

Etika internetového výzkumu

- **Je výzkum nad daty sociálních sítí výzkumem na lidech?**
- Živý jedinec se stává subjektem výzkumu, pokud výzkumník
 - Shromažďuje informace nebo biologické vzorky pomocí intervence či interakce
 - Interakce zahrnuje komunikaci nebo jiný kontakt mezi výzkumníkem a subjektem výzkumu
 - Získává, využívá, studuje, analyzuje či generuje identifikovatelné soukromé informace či identifikovatelné biologické vzorky
- **Soukromá informace**
 - lze důvodně očekávat, že nebyla vůbec zaznamenána, nebo
 - byla poskytnuta pro stanovený účel a lze důvodně očekávat, že nebude zveřejněna
- **Různé výklady „důvodného očekávání“, „stanoveného účelu“ a „zveřejnění“**
 - Více viz např. <https://plato.stanford.edu/entries/ethics-internet-research/>
- **Etika internetového výzkumu je relativně mladá**
 - Nejsou konkrétní závazné předpisy; Aplikují se obecnější normy

Norimberský kodex (1947)

- První kodifikace **etických zásad výzkumu na lidech**
 - Reakce na výzkumy prováděné nacisty během 2. světové války
1. **Dobrovolný souhlas** lidské osoby je absolutně nevyhnutelný.
 2. Pokus by měl být takový, aby přinesl **plodné výsledky** pro dobro společnosti, nebyl zjistitelný jinými metodami nebo studijními prostředky a neměl by být ve své podstatě náhodný a zbytečný.
 3. Pokus by měl být navržen a založen na výsledcích experimentů na zvířatech a na znalosti přirozené historie onemocnění nebo jiného studovaného problému tak, aby **očekávané výsledky ospravedlňovaly provádění pokusu**.
 4. Pokus by měl být prováděn tak, aby **se zabránilo veškerému zbytečnému tělesnému a duševnímu utrpení a zranění**.
 5. Neměl by být prováděn **žádný pokus, pokud lze předpokládat, že dojde ke smrti nebo poranění**. Snad kromě pokusů, kde experimentující lékaři slouží jako pokusné objekty.
 6. **Stupeň rizika**, které je třeba podstoupit, by nikdy neměl překročit míru určenou **humanitárním významem problému**, který má být pokusem vyřešen.
 7. Měly by být provedeny vhodné přípravy a zajištěno náležité **vybavení na ochranu účastníků** pokusu proti sebemenší možnosti zranění, zdravotního postižení nebo smrti.
 8. Pokus by měly provádět pouze **vědecky kvalifikované osoby**. Měl by být vyžadován nejvyšší stupeň dovedností a péče při všech fázích pokusu u všech účastníků pokusu.
 9. Během pokusu by měl mít účastník **možnost pokus ukončit**, pokud se dostal do takového tělesného nebo duševního stavu, kdy mu pokračování v pokusu přijde nemožné.
 10. Odpovědný **vědec musí být připravený ukončit pokus** v kterékoli fázi, má-li, na základě svých nejlepších zkušeností, důvod věřit, že by pokračování v pokusu mohlo mít za následek zranění, postižení nebo smrt pokusného objektu.

Belmontská zpráva (1978)

- Reakce na „Tuskegee Syphylis Study v USA“ (1972)
 - Zkoumání průběhu syfilitidy na afroameričanech, kterým nebyl podán penicilin
- Základní principy
 - **Úcta k lidem**: Respekt lidské autonomie, důstojnost účastníků výzkumu
 - Anonymita, důvěrnost
 - **Užitečnost**: Maximalizovat přínosy výzkumu, minimalizovat rizika pro účastníky, požadavek neškodit
 - **Spravedlivost**: férovost při výběru subjektů, rovnost
- Vodítko pro posuzování výzkumu etickými komisemi
- Základní principy převzaté i do evropských a dalších předpisů
 - Ženevská deklarace (1947), Helsinská deklarace (1964),...

Dilemma Game

Kontrola dat

Pro svoji závěrečnou práci jste sesbírali data od respondentů, kterým jste slíbili anonymitu. Data máte uložená na fakultním serveru, osobní údaje respondentů máte ve svém počítači. Na jiné fakultě se objevil případ fabrikace dat v diplomové práci, což vyústilo v plošnou kontrolu dat ve studentských pracích na celé univerzitě. Byli jste požádáni o informace o vašich respondentech, aby se potvrdilo, že jste skutečně sesbírali data od reálných osob.

- A. Anonymita respondentů je klíčová. Žádné informace nikomu neposkytnu.
- B. Dodám informace o identitě respondentů způsobem, který neumožní spárování respondentů s daty o jejich osobě.
- C. Poskytnu plný přístup ke svým datům pod podmínkou, že kontrolující osoba podepíše závazek mlčenlivosti.
- D. Poskytnu plný přístup ke svým datům. Institucionální kontrola má vyšší prioritu než mnou daný slib anonymity.

Sdílení dat

Jsem začínající výzkumník na Fakultě informatiky, který si pečlivě sesbíral velké množství dat ke své práci. Můj první článek, kde využívám tato data, byl právě přijat k publikaci v odborném časopisu. Starší kolega z mé katedry mě požádal, abych mu poskytl získaná data. Daný kolega má velké slovo v mém dalším působení na katedře. Co mám dělat?

- A. Kolegovi pošlu zmiňovaná data.
- B. Kolegovi řeknu, že data dostane k dispozici ihned, jakmile odešlu poslední práci, která bude daná data využívat. To však klidně může trvat i dva roky.
- C. Kolegovi řeknu, že mu nechci poskytnout speciální zacházení.
- D. Kolegovi řeknu, že mu rád pošlu všechna data s podmínkou, že u všech publikací a článků, které budou využívat zmiňovaná data, budu uveden jako spoluautor.

Úkoly na příště

- Příští téma: Sociální sítě
- Přečíst si článek
 - [How Filter Bubbles Distort Reality: Everything You Need to Know](#)
- Seznámit se s pojmem Kapitalismus dohledu
 - Možné zdroje: [video](#) nebo [článek](#)