

# Posudek oponenta závěrečné práce

České vysoké učení technické v Praze

Fakulta informačních technologií

**Student:** Bc. Matuš Tóth  
**Oponent práce:** Ing. Martin Kopp  
**Název práce:** Analýza leteckých dat a hledání anomálních pasažérů  
**Obor:** Znalostní inženýrství

**Datum vytvoření:** 6. 6. 2017

<b>Hodnotící kritérium:</b>	<b>Způsob hodnocení - následující škálou 1 až 5:</b>
<b>1. Náročnost a další komentář k zadání</b>	1=mimořádně náročné zadání, 2=náročnější zadání, <b>3=průměrně náročné zadání,</b> 4=lehčí, ale ještě dostatečně náročné zadání, 5=nedostatečně náročné zadání
<b>Popis kritéria:</b> Podrobněji charakterizujte diplomovou (bakalářskou) práci a její případné návaznosti na předchozí nebo běžící projekty. Dále posuďte, čím je zadání této ZP náročné. (U obtížnější ZP lze dále tolerovat některé nedostatky, které by u ZP standardní obtížnosti tolerovány nebyly; a naopak u jednoduché ZP mohou být zjištěné nedostatky hodnoceny přísněji.)	
<b>Komentář:</b> Zadání je velice zajímavé a vzhledem k tématu i velmi důležité. Rozsahem odpovídá průměrně diplomové práci.	
<b>Hodnotící kritérium:</b>	<b>Způsob hodnocení - následující škálou 1 až 4:</b>
<b>2. Splnění zadání</b>	1=zadání splněno, 2=zadání splněno s menšími výhradami, <b>3=zadání splněno s většími výhradami,</b> 4=zadání nesplněno
<b>Popis kritéria:</b> Posuďte, zda předložená ZP splňuje zadání. V komentáři uveďte body zadání, které nebyly zcela splněny, případně rozšíření ZP oproti původnímu zadání. Nebylo-li zadání zcela splněno, pokuste se posoudit závažnost, dopady a případně i příčiny jednotlivých nedostatků.	
<b>Komentář:</b> Zadání bylo splněno jen částečně. První bod zadání byl nastudovat metody detekce anomálií z grafových dat. Jedina zmínka o grafech je v sekci 1.16.2, kde student píše, že tato oblast je neprozkoumána. Toto tvrzení rozporuje například: "Akoglu (2015). Graph based anomaly detection and description: a survey", což je aktuální přehled metod detekce anomálií z grafových dat, který má 68 stran a obsahuje citace na desítky relevantních článků.  Hlavním výstupem práce měla být odpověď na analytické otázky zformulované společně s policií CR. V závěru místo jasných odpovědí student popisuje pouze kroky potřebné k jejich zodpovězení. Některé otázky jsou sice zodpovězeny v experimentální části, ale bohužel ne všechny a analýza nutná pro jejich zodpovězení je nejlepšími nedostatky a neprůkazná, zejména vezmeme-li v potaz závažnost tématu.	
<b>Hodnotící kritérium:</b>	<b>Způsob hodnocení - následující škálou 1 až 4:</b>
<b>3. Rozsah písemné zprávy</b>	1=splňuje požadavky, 2=splňuje požadavky s menšími výhradami, <b>3=splňuje požadavky s většími výhradami,</b> 4=nesplňuje požadavky
<b>Popis kritéria:</b> Zhodnoťte přiměřenost rozsahu předložené ZP vzhledem k obsahu, tj. zda všechny části ZP jsou informačně bohaté a ZP neobsahuje zbytečné části.	
<b>Komentář:</b> Práce má sice 75 číslovaných stran, bohužel relevantního obsahu je zde minimum. Kapitoly 2,3,4 a polovina kapitoly 5 popisuje stále dokola totéž a i další části se opakují. Hlavní tézou práce, kterým by měla být vlastní analýza dat, je naopak nedostatečná.	
<b>Hodnotící kritérium:</b>	<b>Způsob hodnocení - bodové hodnocení 0 až 100 bodů (známka A až F):</b>
<b>4. Věcná a logická úroveň práce</b>	20 (F)
<b>Popis kritéria:</b> Posuďte, zda předložená ZP je po věcné stránce v pořádku, případně vyskytují-li se v práci věcné chyby nebo nepřesnosti. Zhodnoťte dále logickou strukturu ZP, návaznosti jednotlivých kapitol a pochopitelnost textu pro čtenáře.	

#### Komentář:

Prace je naprosto nelogicky strukturovana, viz cela prvni kapitola. Student zde naprosto jasne ukazuje, nejen ze nerozumi dane problematice, ale ani vztahum mezi jednotlivymi obory do kterych zasahuje.

Prace se velmi casto opakuje cimz komplikuje citelnost.

Dale prace obsahuje velke mnozstvi nepresnosti viz. napriklad kapitola o neparametrickech metodach, u ktere nasleduje popis parametru ktere je potreba optimalizovat. Prace casto rozporuje sama sebe a to i v ramci jedne kapitoly, viz 1.8.3 vs 1.8.3.2 kde nejprve student tvrdi, ze pro detekci anomalii je tezi sehnat labelovana normalni data, nacez tvrdi ze normalni data lze sehnat snadneji nez data anomalni.

Dalsi nepresnosti napr: str.21 student tvrdi, ze regresni model je statisticke rozdeleni; jako nevyhody neuronovych siti a nahodnych lesu uvadi, ze vraci primo tridy (anomalni/normalni) namisto miry anomalie, nicmene oba modely se daji pouzit i tak aby vraceli miru anomalie; u metody zalozene na hledani asociacnich pravidel navrhuje pouzit min support threshold, coz je standardni technika pro klasifikaci nebo doporučovaci systemy ktera odstrani malo caste vzory z dat, coz je naopak naprosto nezadouci pro detekci anomalii.

Hodnotící kritérium:

Způsob hodnocení - bodové hodnocení 0 až 100 bodů (známka A až F):

### 5. Formální úroveň práce

60 (D)

Popis kritéria:

Posuďte správnost používání formálních zápisů obsažených v práci. Posuďte typografickou a jazykovou stránku ZP, viz Směrnice děkana č. 14/2015, článek 3.

#### Komentář:

I po formální stránce má práce několik zásadních nedostatků jako je nekonzistentní značení vzorcu, které student jednoduše odněkud opsal a aniž by se nad nimi přemyslel. Kapitola 1.12.1.1 je jasným příkladem. Student zde zavedl dva znaky pro průměrnou hodnotu, či směrodatnou odchylku, pro hodnoty používá malé  $x$  i velké  $X$  a podobně. Naopak některé znaky jako velké  $S$  nezavedl vůbec. Notace je nekonzistentní v rámci celé práce, velká písmena někde značí množiny, jinde hodnoty atp.

Další nekozistence, tektokrat typografická je napr. v sekci 1.15.3. Obvykle student popisuje výhody a nevýhody jako číslovane seznamy, zde jsou v textu odstavce.

Jazykovou úroveň bohužel nemohu přilic posoudit, jelikož je práce psána slovensky. Nicméně mi přijde, že je práce psána poněkud hovorově. Dale práce obsahuje překlepy a opakující se slovesa.

Hodnotící kritérium:

Způsob hodnocení - bodové hodnocení 0 až 100 bodů (známka A až F):

### 6. Práce se zdroji

40 (F)

Popis kritéria:

Vyjádřete se k aktivitě studenta při získávání a využívání studijních materiálů k řešení ZP. Charakterizujte výběr studijních pramenů. Posuďte, zda student využil všechny relevantní zdroje nebo zda se pokoušel řešit již vyřešené problémy. Ověřte, zda jsou všechny převzaté prvky řádně odlišeny od vlastních výsledků a úvah, zda nedošlo k porušení citační etiky a zda jsou bibliografické citace úplné a v souladu s citačními zvyklostmi a normami.

#### Komentář:

Cela sekce 1.7.2 je opsána z článku: "Chandola, Banerjee and Kumar 2009 Anomaly detection: A survey" bez jediného odkazu na původní práci. Článek je sice v seznamu literatury, nicméně v této kapitole citace uvedena není.

Obrazky a tabulky nejsou obvykle v textu referovány vůbec.

Citovaná literatura je k obsahu práce relevantní, bohužel faktické chyby v textu naznačují, že jí student ve skutečnosti necetl.

Dle zadání se měl student věnovat detekci anomalii na grafových datech, kde tvrdí, že oblast není dostatečně prozkoumána, což není pravda viz bod 2.

Hodnotící kritérium:

Způsob hodnocení - bodové hodnocení 0 až 100 bodů (známka A až F):

### 7. Hodnocení výsledků, publikační výstupy a ocenění

40 (F)

Popis kritéria:

Vyjádřete se k úrovni dosažených hlavních výsledků ZP, např. k úrovni teoretických výsledků, nebo k úrovni a funkčnosti technického nebo programového vytvořeného řešení, apod. Případně také zhodnoťte, zda software nebo zdrojové texty, které vytvořil sám student, byly v ZP použity v souladu s licenčními podmínkami a autorským právem. Popište případnou publikační činnost a získaná ocenění související s řešením této ZP.

#### Komentář:

Výstupy práce jsou nedostatečné a nepřesné. V kapitole 6.2 student obhájí svůj binární klasifikátor s přesností 46%, což je horší než náhodný klasifikátor. Nasleduje porovnání s rozhodovacím stromem, který na téměř všechny vstupy vrací třídu 1 a s neuronovou sítí, která se rozhoduje na základě jednoho vstupního atributu obr. 6.2 str. 58.

U modelu nejsou typicky uvedeny žádné parametry ani nastavení. Například pro SVM je volba jádrové funkce naprosto zásadní.

Následná analýza nebezpečných letišť a letů, více méně zachraňuje alespoň částečné splnění zadání.

Hodnotící kritérium:

Způsob hodnocení - nehodnotí se

## 8. Komentář o využitelnosti výsledků

Popis kritéria:

Uvedte, zda hlavní výsledky ZP rozšiřují již publikované známé výsledky a/nebo přinášející zcela nové poznatky. Uvedte možnosti využití výsledků ZP v praxi.

*Komentář:*

Výsledky ohledně nebezpečných letišť a nebezpečných letů jsou vcelku průkazně založeny na dostupných datech. Jejich analýza spočívala v prostém vypočtení frekvence nebezpečných pasáží.

Výsledky, které se týkají označení anomálních pasáží jsou naopak naprosto nedostatečné a jejich využití v praxi je nereálné, viz předchozí bod hodnocení.

Hodnotící kritérium:

Způsob hodnocení - nehodnotí se

## 9. Otázky k obhajobě

Popis kritéria:

Uvedte případné dotazy, které by měl student zodpovědět při obhajobě ZP před komisí (body oddělte odřádkami).

*Otázky:*

V kapitole 7.2 uveďte, že by bylo možné data dále analyzovat algoritmy detekce anomálií. Měl jste na mysli nějaké konkrétní algoritmy? Jaké a proč?

Proč jste z důvodu anonymizace vynechal i parametr typu dokladu?

Hodnotící kritérium:

Způsob hodnocení - bodové hodnocení 0 až 100 bodů (známka A až F):

## 10. Celkové hodnocení

40 (F)

Popis kritéria:

Shrňte stránky ZP studenta, které nejvíce ovlivnily Vaše celkové hodnocení. Celkové hodnocení **nemusí** být aritmetickým průměrem či jinou hodnotou vypočtenou z hodnocení v předchozích jednotlivých kritériích 1 až 9.

*Text hodnocení:*

Práce obsahuje velké množství logických, formálních i typografických chyb.

Analýza je nedostatečná a prezentované modely pro detekci anomálií (ve skutečnosti binární klasifikatory) nejsou použitelné ani teoreticky natož pak v reálné situaci.

Navzdory názvu i zadání, práce ve skutečnosti nepoužívá metody detekce anomálií, kterým se věnuje v teoretické části, ale spolehně se na binární klasifikaci.

Student neuvádí zdroj pro sekci 1.7.2, která je opsána včetně formátování.

Vzhledem k tomu, že zde uvedené připomínky jsou jen špička ledovce a vzhledem k závaznosti tématu hodnotím práci známkou F a navrhuji aby student práci přepracoval.

Podpis oponenta práce: