

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

Fakulta elektrotechniky  
a komunikačních technologií

DIPLOMOVÁ PRÁCE

Brno, 2019

Bc. Kateřina Gregorová



# VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

## FAKULTA ELEKTROTECHNIKY A KOMUNIKAČNÍCH TECHNOLOGIÍ

FACULTY OF ELECTRICAL ENGINEERING AND COMMUNICATION

## ÚSTAV BIOMEDICÍNSKÉHO INŽENÝRSTVÍ

DEPARTMENT OF BIOMEDICAL ENGINEERING

## STATISTICKÁ ANALÝZA ANOMÁLIÍ V SENZOROVÝCH DATECH

STATISTICAL ANALYSIS OF ANOMALIES IN SENSOR DATA

### DIPLOMOVÁ PRÁCE

MASTER'S THESIS

### AUTOR PRÁCE

AUTHOR

**Bc. Kateřina Gregorová**

### VEDOUCÍ PRÁCE

SUPERVISOR

**Ing. Jiří Sekora**

**BRNO 2019**



# Diplomová práce

magisterský navazující studijní obor **Biomedicínské a ekologické inženýrství**  
Ústav biomedicínského inženýrství

**Studentka:** Bc. Kateřina Gregorová

**ID:** 155573

**Ročník:** 2

**Akademický rok:** 2018/19

## NÁZEV TÉMATU:

### Statistická analýza anomálií v senzorových datech

#### POKYNY PRO VYPRACOVÁNÍ:

1) Prostudujte problematiku konstrukčního řešení senzorů pro měření tlaku a teploty kapalin používaných u leteckých motorů. 2) Proveďte rešerši diagnostických přístupů pro detekci poruch senzorů. 3) Na základě dostupných informací navrhnete algoritmus pro rozlišení chybných měření a anomálií senzorových dat. 4) Proveďte statistické vyhodnocení výsledků. 5) V rámci diskuze validujte navržený algoritmus na dostupných datech. Zadáni je řešeno ve spolupráci s firmou Honeywell.

#### DOPORUČENÁ LITERATURA:

[1] BRANDT, Siegmund. Data analysis: statistical and computational methods for scientists and engineers. 4th ed. Přeložil Glen COWAN. Cham: Springer, c2014. ISBN 978-3-319-03761-5.

[2] FROST, Roger. Datalogging and Control: The IT in Science book of. Cambridge: IT in Science, 2002. ISBN 0-9520257-1-X.

**Termín zadání:** 4.2.2019

**Termín odevzdání:** 17.5.2019

**Vedoucí práce:** Ing. Jiří Sekora

**Konzultant:**

**prof. Ing. Ivo Provazník, Ph.D.**  
*předseda oborové rady*

#### UPOZORNĚNÍ:

Autor diplomové práce nesmí při vytváření diplomové práce porušit autorská práva třetích osob, zejména nesmí zasahovat nedovoleným způsobem do cizích autorských práv osobnostních a musí si být plně vědom následků porušení ustanovení § 11 a následujících autorského zákona č. 121/2000 Sb., včetně možných trestněprávních důsledků vyplývajících z ustanovení části druhé, hlavy VI. díl 4 Trestního zákoníku č.40/2009 Sb.

## ABSTRAKT

Tato diplomová práce se zabývá problematikou detekce poruchových stavů u leteckých motorů. Hlavním přístupem detekce je hledání anomálií v datech snímaných pomocí senzorů. Pro získání komplexní představy o systému a jednotlivých senzorech, je v úvodu této práce uveden popis celého systému a to konkrétně motoru typu HTF7000 a také popis senzorů. Pro samotnou detekci anomálií je zde uveden návrh algoritmu na základě tří různých detekčních metod, které jsou popsány ve druhé kapitole. Jedná se o metody SVM (Support Vector Machine), K-means a ARIMA (Autoregressive Integrated Moving Average). Implementace algoritmu je popsána v další části práce, včetně návrhu grafického uživatelského rozhraní. V závěru práce je pak statistická analýza získaných výsledků, srovnání účinnosti jednotlivých modelů a diskuze výstupů z navrženého algoritmu.

## KLÍČOVÁ SLOVA

SVM, K-means, ARIMA, letecký motor, senzor, HTF7000, detekce anomálií

## ABSTRACT

This thesis deals with the failure mode detection of aircraft engines. The main approach to the detection is searching for anomalies in the sensor data. In order to get a comprehensive idea of the system and the particular sensors, the description of the whole system, namely the aircraft engine HTF7000 as well as the description of the sensors, are dealt with at the beginning of the thesis. A proposal of the anomaly detection algorithm based on three different detection methods is discussed in the second chapter. The above-mentioned methods are SVM (Support Vector Machine), K-means a ARIMA (Autoregressive Integrated Moving Average). The implementation of the algorithm including graphical user interface proposal are elaborated on in the next part of the thesis. Finally, statistical analysis of the results, the comparison of efficiency particular models and the discussion of outputs of the proposed algorithm can be found at the end of the thesis.

## KEYWORDS

SVM, K-means, ARIMA, aircraft engines, sensor, HTF7000, anomaly detection

GREGOROVÁ, Kateřina. *Statistická analýza anomálií v senzorových datech*. Brno, Rok, 65 s. Diplomová práce. Vysoké učení technické v Brně, Fakulta elektrotechniky a komunikačních technologií, Ústav biomedicínského inženýrství. Vedoucí práce: Ing. Jiří Sekora

## PROHLÁŠENÍ

Prohlašuji, že svou diplomovou práci na téma „Statistická analýza anomálií v senzorových datech“ jsem vypracovala samostatně pod vedením vedoucího diplomové práce a s použitím odborné literatury a dalších informačních zdrojů, které jsou všechny citovány v práci a uvedeny v seznamu literatury na konci práce.

Jako autorka uvedené diplomové práce dále prohlašuji, že v souvislosti s vytvořením této diplomové práce jsem neporušila autorská práva třetích osob, zejména jsem nezasáhla nedovoleným způsobem do cizích autorských práv osobnostních a/nebo majetkových a jsem si plně vědoma následků porušení ustanovení § 11 a následujících autorského zákona č. 121/2000 Sb., o právu autorském, o právech souvisejících s právem autorským a o změně některých zákonů (autorský zákon), ve znění pozdějších předpisů, včetně možných trestněprávních důsledků vyplývajících z ustanovení části druhé, hlavy VI. díl 4 Trestního zákoníku č. 40/2009 Sb.

Brno .....

.....

podpis autorky

## PODĚKOVÁNÍ

Ráda bych poděkovala vedoucímu diplomové práce panu Ing. Jiřímu Sekorovi, za odborné vedení, konzultace, trpělivost, shovívavost a podnětné návrhy k práci. Také celému pracovnímu kolektivu za jejich podporu a pochopení a stejně tak mé rodině a blízkým.

Brno .....

.....

podpis autorky

# Obsah

Úvod	10
<b>1 Popis systému</b>	<b>11</b>
1.1 HTF7000 konstrukce motoru . . . . .	11
1.1.1 Základní koncept řízení motoru . . . . .	12
1.1.2 Monitorování stavu motoru . . . . .	12
<b>2 Detekce anomálií</b>	<b>18</b>
2.1 Co jsou anomálie? . . . . .	19
2.2 Typy anomálií . . . . .	19
2.2.1 Bodové . . . . .	19
2.2.2 Kontextové . . . . .	19
2.2.3 Kolektivní . . . . .	20
2.3 Metody detekce anomálií . . . . .	21
2.3.1 Klasifikační metody . . . . .	21
2.3.2 Metody založené na hledání nejbližšího souseda . . . . .	22
2.3.3 Statistické metody . . . . .	22
<b>3 Detekce chyb senzorů</b>	<b>24</b>
3.1 Support Vector Machines . . . . .	24
3.1.1 Rozdělení SVM . . . . .	25
3.1.2 One-Class SVM . . . . .	32
3.1.3 Analýza hlavních komponent . . . . .	33
3.2 K-means . . . . .	36
3.2.1 Detekce anomálií pomocí K-Means . . . . .	36
3.3 ARIMA . . . . .	39
3.3.1 Zlepšení modelu ARIMA . . . . .	41
<b>4 Návrh algoritmu</b>	<b>44</b>
4.1 Popis dat . . . . .	44
4.2 Obecný přístup k návrhu algoritmu pro monitorování stavu systému . . . . .	47
4.2.1 Pracovní postup pro vývoj algoritmů . . . . .	48
4.3 Vlastní návrh algoritmu . . . . .	50
4.3.1 Předzpracování vstupních dat . . . . .	51
4.3.2 Hlavní funkce algoritmu . . . . .	51
4.3.3 SVM . . . . .	52
4.3.4 K-means . . . . .	52
4.3.5 ARIMA . . . . .	54

<b>5</b>	<b>Vyhodnocení detekce</b>	<b>56</b>
5.1	Vyhodnocení SVM . . . . .	56
5.2	Vyhodnocení K-MEANS . . . . .	56
5.3	Vyhodnocení ARIMA . . . . .	57
5.4	Shrnutí . . . . .	58
<b>6</b>	<b>Závěr</b>	<b>59</b>
	<b>Literatura</b>	<b>60</b>
	<b>Seznam symbolů, veličin a zkratk</b>	<b>65</b>



# Seznam obrázků

1.1	Schematický diagram proudového motoru s vysokým poměrem ob- toku vzduchu. . . . .	11
1.2	Středně velká obchodní letadla . . . . .	12
1.3	Typické senzory používané pro řízení motoru. . . . .	13
1.4	Schéma termočlánku typu K (chromel–alumel) ve standardním ter- močláňkovém měřícím zapojení. . . . .	15
1.5	Schéma řízení paliva pomocí hydromechanické/elektronické jednotky. Převzato z [10]. . . . .	15
1.6	Schéma senzoru LVDT . . . . .	16
2.1	Jednoduchý příklad anomálií ve dvourozměrném souboru dat. . . . .	18
2.2	Příklad kontextové anomálie. . . . .	20
2.3	Příklad kolektivní anomálie. . . . .	21
3.1	Princip klasifikace metodou SVM. . . . .	25
3.2	Lineární dělicí nadroviny pro separabilní případ. . . . .	27
3.3	Lineární dělicí nadroviny pro neseperabilní případ. . . . .	30
3.4	Znázornění nelineární SVM. . . . .	32
3.5	Příklad procentuálního zastoupení informace v každé z hlavních kom- ponent. . . . .	35
3.6	Klasifikace pro $K = 2$ . . . . .	38
3.7	Detekce odlehlých bodů . . . . .	38
3.8	Kombinace klasifikace a detekce odlehlých bodů . . . . .	39
3.9	ARIMA model . . . . .	40
3.10	Metoda exponenciálně váženého klouzavého průměru . . . . .	41
3.11	Vylepšený model detekce anomálií, založený na modelu ARIMA . . . . .	43
4.1	Příklad vzorů chování pro vybrané senzory. . . . .	46
4.2	Detekované poruchy dle vybraných vzorů chování. . . . .	46
4.3	Fotografie reálného motoru s dírou v plášti kompresoru s vykresleným záznamem dat. . . . .	46
4.4	Závislosti mezi senzory u datového souboru s anomáliemi. . . . .	47
4.5	Blokové schéma vývoje detekčního algoritmu. . . . .	48
4.6	Data pro vybrané senzory v závislosti na čase . . . . .	51
4.7	Metoda SVM - příklad modelu pro trénovací data . . . . .	53
4.8	Příklad detekce anomálií pomocí metody K-means pro kombinaci N2 - ITT . . . . .	54
4.9	Příklad parciální autokorelační funkce - metoda ARIMA . . . . .	55
5.1	Příklad detekce anomálií pro jeden motor metodou SVM . . . . .	57

# Seznam tabulek

1.1	Přehled senzorů . . . . .	14
4.1	Popis testovaných senzorů. . . . .	44
4.2	Statistické hodnoty použitých senzorů pro TKO data. . . . .	45
4.3	Statistické hodnoty použitých senzorů pro CRU data. . . . .	45
5.1	Počet motorů, u kterých byla detekována anomálie metodou K-means pro daný senzor. . . . .	57

# Úvod

Senzory v komerčních leteckých motorech pracují v těžkých podmínkách, a proto jsou náchylné k poruchám. Jakékoliv nedetekované poruchy snímačů mohou mít za následek katastrofální následky. Proto jsou detekční a prediktivní algoritmy nedílnou součástí monitorování systému a mají rozhodující vliv na zvýšení spolehlivosti, účinnosti a bezpečnosti motoru. Pro návrh těchto algoritmů je možné využít celou řadu metod, avšak nejvíce používaným způsobem [29], [30] je detekce anomálií. Tato metoda je důležitým nástrojem pro monitorování, diagnostiku a určení prognózy daného systému. Pokud dojde k chybě systému, můžeme předpokládat, že dojde ke změnám měřených hodnot. Včasně odhalení anomálie může pomoci k detekci chybného chodu systému. To umožní provedení údržby dříve, než porucha způsobí sekundární poškození, které vede k vyřazení celého systému z provozu.

Tato diplomová práce je vypracována ve spolupráci s firmou Honeywell a jejím cílem je navrhnout spolehlivé metody pro identifikaci normálního pracovního režimu a poruchového režimu ve velkém množství dat snímaných pomocí senzorů.

Úvodní část práce se bude věnovat obecnému popisu systému, kterým je v tomto případě letecký motor typu HTF7000. Dále pak detailnějšímu popisu senzorů a jejich technických parametrů. Ve druhé kapitole jsou popsány typy anomálií a obecné metody detekce. Následně pak budou diskutovány vybrané metody, uveden jejich popis a důvody volby těchto metod. Zmíněná obecná část bude následována detailním popisem navrženého algoritmu a vyhodnocením detekce.



který obsahuje pouze dvě sběrné jímky oleje. Obě tyto jímky jsou uloženy v chladné části motoru, aby nedošlo k jeho vznícení. Spalovací komora je prstencová s přímým vstřikováním. Pro snížení hluku a zvýšení účinnosti se na výstupu z motoru míchá vzduch obtékající jádro s plyny z jádra motoru. K míchání dochází před vstupem do výfukové trysky [1].



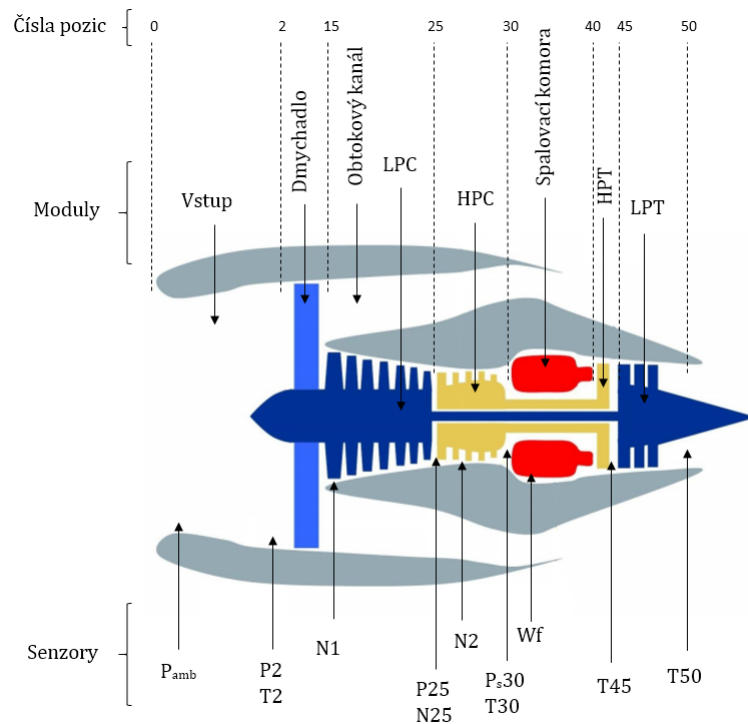
Obr. 1.2: Letadla, která jsou osazena motory typu HTF7000 od firmy Honeywell (Challenger 300, Gulfstream G280, Embraer Legacy, Cessna Citation Longitude).

### 1.1.1 Základní koncept řízení motoru

Motor je řízen dvoukanálovým systémem digitálního elektronického řízení s plnou autoritou (FADEC). Tento systém má k dispozici dvě samostatné kontrolní jednotky (ECU). Tento systém přijímá několik vstupních proměnných aktuálního letového stavu, včetně hustoty vzduchu, polohy páky škrticí klapky, teploty motoru, tlaku motoru a mnoha dalších parametrů. Vstupy jsou přijímány a analyzovány až sedmdesátkrát za sekundu. Provozní parametry motoru, jako je průtok paliva, poloha satorové lopatky, poloha odvzdušňovacího ventilu a další, jsou vypočteny z dostupných údajů. FADEC také řídí startování a restartování motoru. Základním účelem FADEC je poskytnout optimální účinnost motoru pro daný letový stav [6].

### 1.1.2 Monitorování stavu motoru

Monitorování stavu motoru a diagnostika poruch se vyvíjí stejně jako probíhá vývoj samotného motoru. Nejdůležitějšími parametry pro diagnostiku motoru je monitorování teploty výfukových plynů (EGT), měření rychlosti otáčení nízkotlakého kompresoru (N1) a vysokotlakého kompresoru (N2), měření průtoku paliva ( $W_f$ ), měření celkové teploty vzduchu (TAT), vstupní teploty motoru (T2), okolního tlaku



Obr. 1.3: Typické senzory používané pro řízení motoru.

( $P_{amb}$ ) a celkového vstupního tlaku ( $P_2$ ), vstupní teploty a tlaku vysokotlakého kompresoru (HPC) ( $T_{25}$ ,  $P_{25}$ ), výstupní teploty a tlaku HPC ( $T_{30}$ ,  $P_{s30}$ ) a další viz obrázek 1.3 [8].  $N_1$  představuje rychlost otáčení nízkotlakého kompresoru a  $N_2$  je rychlost otáčení vysokotlakého kompresoru. Hodnoty jsou měřeny v otáčkách za minutu (RPM). Po startu je rychlost nízkotlakého kompresoru řízena turbínovým kolem  $N_1$ , které je připojeno k nízkotlakému kompresoru přes soustředný hřídel a obdobně je vysokotlaký kompresor řízen turbínovým kolem  $N_2$ , které je připojeno k vysokotlakému kompresoru.

### Monitorování otáček vysokotlakého a nízkotlakého kompresoru

Principiálně se jedná o induktivní snímače. Současné induktivní snímače se skládají většinou z tyčového magnetu s magneticky měkkým pólovým nástavcem, na kterém je umístěna indukční cívka se dvěma vývody. Otáčí-li se před tímto snímačem feromagnetické ozubené kolo nebo jiný obdobně konstruovaný rotor, indukuje se v cívkce přibližně sinusové napětí. Při vysokých rychlostech se elektrický puls generovaný z každého cílového zubu kombinuje a objevuje se jako nepřetržitá sinusová vlna. Na rotoru může být umístěna jedna nebo více vztažných značek [7].

Tab. 1.1: Přehled senzorů

<b>Senzor</b>	<b>Popis senzoru a umístění</b>
N1	rychlost otáčení nízkotlaké turbíny
N2	rychlost otáčení vysokotlaké turbíny
P0	okolní tlak
P15	statický tlak na obtokovém vedení
P2	vstupní tlak dmyhadla
P25	vstupní tlak vysokotlakého kompresoru
Ps30	výstupní statický tlak vysokotlakého kompresoru
P40	vstupní tlak vysokotlaké turbíny
P48	vstupní tlak nízkotlaké turbíny
P50	výstupní tlak nízkotlaké turbíny
T0	okolní teplota
T2	vstupní teplota dmyhadla
T21	výstupní teplota dmyhadla
T24	výstupní teplota nízkotlakého kompresoru
T30	výstupní teplota vysokotlakého kompresoru
T40	vstupní teplota vysokotlaké turbíny
T48	vstupní teplota nízkotlaké turbíny
T50	výstupní teplota nízkotlaké turbíny
Wf	průtok paliva
W21	průtok vzduchu na vstupu dmyhadla
W22	průtok vzduchu na vstupu vysokotlakého kompresoru
W36	průtok vzduchu na výstupu vysokotlakého kompresoru
W40	průtok vzduchu na vstupu vysokotlaké turbíny
W48	průtok vzduchu na vstupu nízkotlaké turbíny
W50	průtok vzduchu na výstupu nízkotlaké turbíny

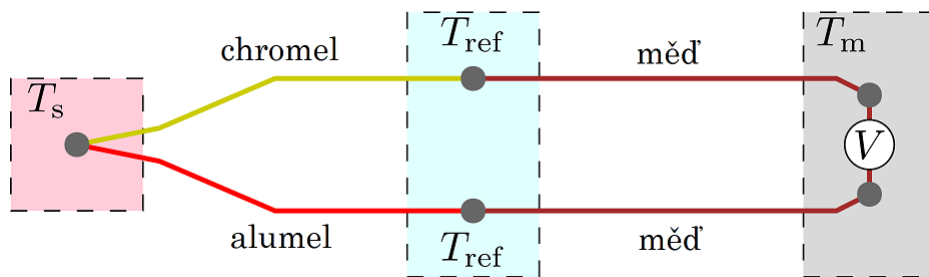
### **Měření teploty**

Teplota je v různých částech motoru měřena obvykle pomocí termočlánků viz obrázek 1.4. Termočlánky v motorech HTF7000 jsou typu K tedy chromel-alumelové jejichž teplotní rozmezí je -200 až 1250 °C s citlivostí 40,8  $\mu$ V [9].

### **Palivový systém**

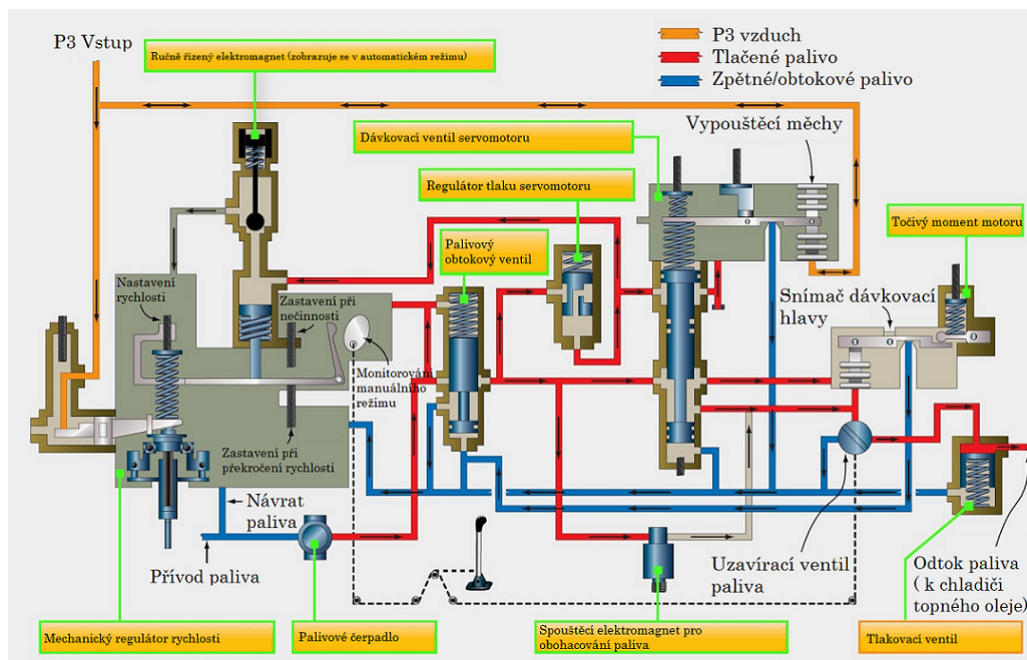
Řídící prvky palivového systému by se daly rozdělit do tří podjednotek:

1. hydromechanická jednotka,
2. hydromechanická/elektronická jednotka,



Obr. 1.4: Schéma termočlánku typu K (chromel–alumel) ve standardním termočlávkovém měřícím zapojení.

### 3. FADEC.



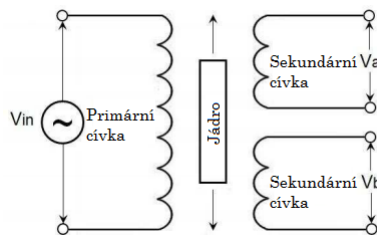
Obr. 1.5: Schéma řízení paliva pomocí hydromechanické/elektronické jednotky. Převzato z [10].

Hydromechanická/elektronická jednotka viz obrázek 1.5 je hybridní systém, ale může fungovat pouze na hydromechanickém principu. V duálním režimu jsou vstupy a výstupy elektronické a tok paliva je řízen servomotory. Třetí typ, FADEC, využívá pro své vstupy elektronická čidla a řídí tok elektronickými výstupy. Ovládání typu FADEC dává elektronickému ovladači úplnou kontrolu. Bez ohledu na typ, mají všechny ovládací prvky pro řízení paliva stejnou funkci a to řídit tok paliva tak, aby odpovídal výkonu požadovanému pilotem. Do řídicí jednotky vstupuje mnoho



parametrů, na jejichž základě je řízen tah motoru pro daný průtok paliva. Například otáčky vysokotlaké i nízkotlaké turbíny, vstupní tlak a teplota kompresoru, tlak ve spalovací komoře a mnoho dalších parametrů [10].

Pro měření průtoku paliva do spalovací komory se používají senzory, které snímají a odesílají data do ECU. Nemají řídicí funkci. Tento senzor je z fyzikálního hlediska hmotnostní průtokoměr. Hmotnostní průtok  $Q_m$  udává hmotnost tekutiny, která projde průtočným průřezem v určitém systému za jednotku času. Pro přímé měření existují dvě základní metody - průtokoměry založené na Coriolisově principu a tepelné hmotnostní průtokoměry. Coriolisovy průtokoměry využívají tzv. Coriolisovy síly, která vzniká ve vibrujících měřicích trubicích při průtoku hmotného média. Fázový posuv v rezonančním kmitání trubice, který vzniká jako důsledek působení Coriolisovy síly, je úměrný hmotnostnímu průtoku tekutiny trubicí a frekvence vlastních kmitů odpovídá hustotě tekutiny. Tepelný hmotnostní průtokoměr vychází ze závislosti výměny tepla mezi zdrojem a okolím, které tvoří proudící tekutina, na hmotnostním průtoku [11].



Obr. 1.6: Schéma senzoru LVDT. Proud je veden do primární cívky  $V_{in}$  a vzájemnou indukčností vzniká napětí na sekundárních cívkách  $V_a$  a  $V_b$ . Převzato z [12]

## Tlakové senzory

Pro měření tlaku jsou u leteckých motorů používané senzory typu LVDT (Linear Variable Differential Transformer) viz obrázek 1.6. Tento senzor je tvořen třemi cívkami, které jsou navinuté vedle sebe po celé délce nevodivé trubice, kterou prochází feromagnetické jádro, jež je při měření posouváno. Prostřední cívka je primární, zbylé dvě jsou sekundární. Primární cívka je napájena střídavým proudem. Sekundární cívky jsou propojeny tak, že jakýkoli proud indukovaný primární cívkou bude mezi dvěma sekundárními cívkami fázově posunutý o 180 stupňů a výsledný výstup bude nulový, pokud nedojde k posunutí. Pokud se jádro posune doleva nebo doprava, vzájemná indukčnost mezi primární a sekundárními cívkami se zvýší v jedné a sníží se ve druhé cívce, čímž se vytvoří diferenční napětí mezi oběma sekundárními

cívkami. Produkovaný výstup diferenčního napětí je přímo úměrný posunutí jádra [12].

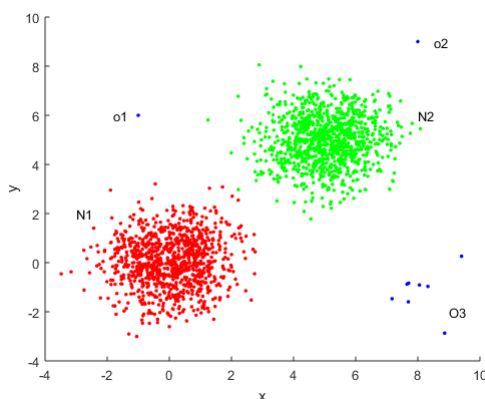
## 2 Detekce anomálií

Detekce anomálií je technika používaná pro identifikaci vzorku dat, který neodpovídá předpokládanému chování. Tyto nekonformní vzorce jsou často označovány jako anomálie, odlehlé hodnoty, výjimky, aberace nebo nesouhlasná pozorování v různých aplikačních oblastech. První dva uvedené výrazy jsou nejčastěji používané termíny v souvislosti s detekcí anomálií. Tato metoda se používá v široké škále aplikací, jako je detekce podvodů kreditních karet, pojistných podvodů, detekce narušení kybernetické bezpečnosti, poruch v zabezpečovacích systémech a další [5].

Význam detekce anomálií je důsledkem skutečnosti, že anomálie v datech jsou nositeli významných a často kritických informací v širokém spektru aplikačních domén. Například anomální provozní vzorec v počítačové síti může znamenat, že napadený počítač odesílá citlivá data do neoprávněného cíle [2]. Anomální obraz MRI může indikovat přítomnost maligních nádorů [3] a například anomálie v datech transakcí mohou znamenat krádež kreditní karty nebo krádež identity [4].

Detekce odchylek nebo anomálií v datech byla studována ve statistice již v 19. století. Postupně se v několika výzkumných komunitách vyvíjela celá řada detekčních technik. Některé techniky byly vyvinuté pro určité aplikační domény, jiné jsou obecnější [5].

V této kapitole jsou uvedeny různé typy anomálií dále jsou představeny skupiny detekčních metod a příklady konkrétních metod, které budou následovně diskutovány podrobněji s přihlédnutím na jejich praktickou realizaci.



Obr. 2.1: Jednoduchý příklad anomálií ve dvourozměrném souboru dat.

## 2.1 Co jsou anomálie?

Anomálie jsou určité vzory v datech, které neodpovídají dobře definovanému pojetí normálního chování. Na obrázku 2.1 je uveden příklad anomálií v dvourozměrném souboru dat. Data mají dvě normální oblasti N1 a N2. Jsou určeny tak, že většina dat leží v těchto oblastech. Body označené o1, o2 a množina bodů O3 vyskytující se v dostatečné vzdálenosti od množin N1 a N2, jsou anomálie. Anomálie se v datech vyskytují z různých důvodů, ale mají jednu společnou vlastnost. Jedná se o analyticky významné informace.

Detekce anomálií se stejně jako filtrace signálů zabývá odstraňováním nechtěného šumu v datech. Šum může být definován například jako jev, který nemá pro analýzu dat žádný význam, je nežádoucí. Je třeba jej před jakýmkoliv zpracováním dat odstranit. Dalším tématem vztahujícím se k detekci anomálií je detekce novotvarů. Cíl této metody spočívá v odhalení doposud nepozorovaných vzorců dat. Rozdíl mezi novotvary a anomáliemi spočívá v tom, že nové vzory jsou typicky začleněny do normálního modelu poté, co byly detekovány. Je nutno podotknout, že řešení těchto souvisejících problémů se často používají k detekci anomálií a naopak [5].

## 2.2 Typy anomálií

Důležitým aspektem použité techniky detekce je povaha dané anomálie. Můžeme je rozdělit do tří kategorií:

### 2.2.1 Bodové

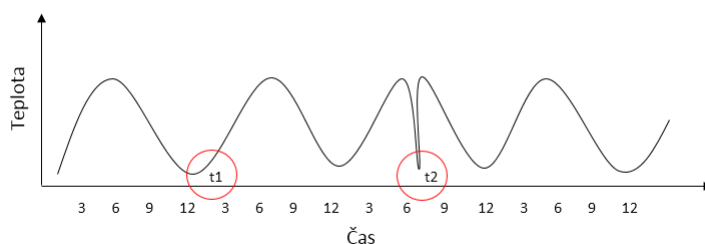
Odlišuje-li se jeden vzorek dat od ostatních, potom ho označíme jako bodovou anomálii. Jde o nejjednodušší typ anomálie, na který je zaměřena většina výzkumů. Vezmeme-li v úvahu obrázek 2.1, body o1 a o2 jsou bodové anomálie. Jako příklad z reálného života můžeme za bodovou anomálii považovat podvodnou bankovní transakci. Kvůli jednoduchosti předpokládejme, že data jsou definována pouze jedním znakem: hodnotou transakce. Transakce, jejíž hodnota je velmi vysoká ve srovnání s běžnými převody, bude bodovou anomálií [5].

### 2.2.2 Kontextové

Kontextová anomálie, jinak také nazývaná podmíněná anomálie, se vyskytuje, je-li jeden datový bod nebo skupina bodů anomální vzhledem ke kontextu. Pojem kontextu je dán strukturou datového souboru a musí být specifikován jako součást formulace problému. Každý datový soubor je definován pomocí následujících dvou sad atributů:

1. *Kontextové atributy.* Kontextové atributy se používají k určení kontextu pro daný soubor dat. Například v systému souřadnic je zeměpisná šířka a délka dané lokace kontextovým atributem. V časových řadách je čas kontextuální atribut, který určuje polohu datové instance v celé sekvenci.
2. *Atributy chování.* Atributy chování naopak definují ty vlastnosti, které nejsou definovány pomocí kontextových atribut. Například v již zmiňovaném systému souřadnic může být také popsán průměrné množství srážek na planetě. Množství srážek na libovolném místě je pak atribut chování.

Neobvyklé chování v datech se určuje pomocí atributů chování v konkrétním kontextu. Jedna datová instance může být kontextovou anomálií za určitých podmínek, pokud ale bude stejná datová instance podmíněna jiným způsobem, může být považována za normální. Kontextové anomálie byly nejčastěji zkoumány v časových řadách [13], [14] a v prostorových datech [15]. Příklad kontextové anomálie je uveden na obrázku 2.2, který zobrazuje vývoj teploty v posledních letech. Teplota v čase  $t_1$  tedy v zimě je normální, avšak stejná teplota naměřená v létě (v čase  $t_2$ ) je považována za anomálii [5].

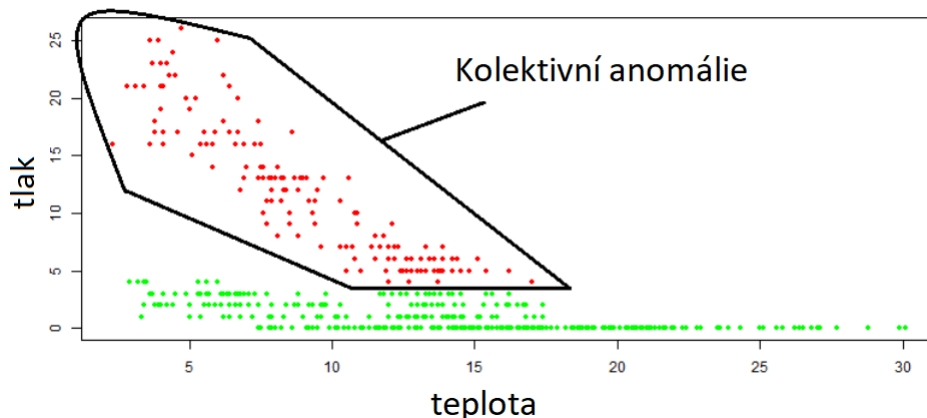


Obr. 2.2: Příklad kontextové anomálie  $t_2$  v časové řadě. Teplota v čase  $t_1$  je stejná jako teplota v čase  $t_2$ , ale vyskytuje se v jiném kontextu, tudíž není považována za anomálii.

### 2.2.3 Kolektivní

Vyskytuje-li se v určité časové řadě úsek, který je vzhledem k celé sadě dat neobvyklý, označuje se jako kolektivní anomálie. Jednotlivé datové body vyskytující se v úseku kolektivní anomálie nemusejí být samy o sobě anomáliemi, ale vzhledem ke společnému výskytu jsou takto označovány. Na obrázku 2.3 je záznam EKG, kde zvýrazněná část záznamu je kolektivní anomálie, z medicínského pohledu se jedná o předčasnou kontrakci síní. Kolektivní anomálie se mohou vyskytovat pouze v takových datových souborech, kde jednotlivé instance vykazují příbuznost. Na rozdíl od bodových anomálií, které mohou nastat v jakémkoliv datovém souboru. Naproti

tomu výskyt kontextových anomálií závisí na dostupnosti kontextových atributů v datech. Bodová anomálie nebo kolektivní anomálie mohou být kontextuální anomálií, pokud jsou analyzovány podmíněně s ohledem na daný kontext [5].



Obr. 2.3: Kolektivní anomálie uvedená na příkladu závislosti tlaku na teplotě.

## 2.3 Metody detekce anomálií

K detekci anomálií lze použít mnoho metod. V této kapitole budou představeny různé detekční algoritmy a tento přehled bude základem ke zvolení správného detekčního algoritmu pro praktickou část této diplomové práce. Důkladnější rozbor algoritmů je uveden například v [5].

### 2.3.1 Klasifikační metody

Klasifikační metody se snaží klasifikovat záznam jako normální nebo anomální. Příkladem klasifikátorů jsou umělé neuronové sítě, nejčastěji dopředné preceptronové sítě, Bayesovské sítě, Support Vector Machines <sup>1</sup> a klasifikátory založené na pravidlech, kdy pravidla mohou být nastavena pomocí strojového učení, nebo lidské odbornosti. Tyto klasifikátory bývají trénovány na množině dat, která obsahuje všechna možná normální chování, aniž by obsahovala anomálie. Jestliže nemohou být nové záznamy zařazeny do žádné z klasifikovaných tříd, jsou označeny jako anomálie. Klasifikátory lze taktéž natrénovat pomocí souboru dat sestávajícího z normálních a anomálních záznamů, které jsou označeny odpovídajícím způsobem za účelem vytvoření klasifikátoru schopného přesně klasifikovat budoucí záznamy buď

<sup>1</sup>Support vector machines (SVM) je metoda strojového učení s učitelem, sloužící zejména pro klasifikaci a také pro regresní analýzu.

jako normální, nebo anomální. Hlavním nedostatkem těchto metod je, že všechny s výjimkou metod založených na pravidlech určených lidskou odborností a metody jednotřídních SVM potřebují velké soubory značených dat. Vzhledem k tomu, že tyto množiny je třeba vytvořit manuálně, často neexistují [17].

### 2.3.2 Metody založené na hledání nejbližšího souseda

Metody detekce anomálií založené na analýze nejbližších sousedů se opírají o následující předpoklad: *data vykazující normální chování se vyskytují v hustých shlucích, zatímco anomálie jsou od těchto shluků značně vzdáleny*. Hlavním problémem této metody je určení vhodné míry vzdálenosti (podobnosti). Nejčastěji se používá jednoduchá Euklidovská vzdálenost. Tato metoda se skládá buď z analyzování vzdálenosti nejbližších sousedů nebo z analýzy počtu záznamů v určité vzdálenosti od testovaného vzorku, což je obvykle výpočetně velmi náročný úkol a závisí na měřené vzdálenosti [17].

### Shluková analýza

Metody detekce anomálií založené na shlukové analýze definují shluky pro běžné datové instance a s těmito shluky porovnávají nové záznamy. Existují dva hlavní přístupy pro detekci anomálií pomocí této metody. První přístup spočívá ve shlukování celé sady dat a analyzuje hustotu každého shluku. Předpokládá se, že data vykazující normální chování patří do velkých a hustých shluků na rozdíl od anomálií, které se obvykle vyskytují v menších a řidších skupinách. Druhý přístup je založen na předpokladu, že anomálie se vyskytují ve větší vzdálenosti od středu seskupení než normální data. Techniky shlukování často vyžadují stejné měření vzdálenosti jako analýza založená na hledání nejbližšího souseda. Hlavní rozdíl mezi těmito dvěma metodami je že při shlukové analýze jsou data porovnávána s celými shluky, ne s jednotlivými záznamy. Tím pádem je tato metoda ve fázi testování rychlejší [17].

### 2.3.3 Statistické metody

Klíčovým předpokladem pro jakoukoliv statistickou metodu detekce anomálií je: *Normální záznamy se vyskytují v oblasti stochastického modelu s vysokou pravděpodobností, zatímco anomálie se vyskytují v oblasti s nízkou pravděpodobností* [5]. Statistické metody detekce anomálií se skládají ze dvou kroků. V prvním kroku je vybrán statistický model. Ve druhém kroku jsou nové datové instance testovány oproti tomuto modelu, aby se zjistila pravděpodobnost výskytu daného záznamu. V časových řadách jsou použité statistické modely často založeny na regresi a pro detekci anomálií se využívá množství regresních modelů. Například Autoregresní

modely (AR) [19] [20], Autoregresní model klouzavého průměru (ARMA) a Autoregresní integrované modely klouzavých průměrů (ARIMA) [17]. Tyto metody detekce anomálií jsou nejjednodušším přístupem. Lze říci, že principem těchto metod je označení těch datových bodů, které se odchyľují od běžných statistických vlastností distribuce, zahrnující průměr, medián, modus a kvantily. Řekněme, že anomálie je takový bod, který je vzdálený od průměru o více než je standardní odchylka datového souboru, případně její násobek [18].



## 3 Detekce chyb senzorů

Detekce abnormálního chování či poruch v leteckých motorech je pro prognostické aplikace velmi náročná, jelikož data bývají obvykle částečně znehodnocena a to zejména šumem vznikajícím na snímacích senzorech. Mírné odchylky ve výrobním procesu mají za následek, že snímaná data se mohou motor od motoru lišit. Jednotlivé komponenty mají různou životnost a jejich údržba neprobíhá ve stejných intervalech, tudíž vznikají další odchylky. Proto data reálného leteckého motoru typicky obsahují odlehlé hodnoty, nejsou normálně distribuovány a obvykle obsahují šum. Pro zavedení detekčního algoritmu je třeba, aby splňoval jisté předpoklady. Zejména je vyžadováno malé množství falešně pozitivních výsledků, aby nedocházelo ke zbytečnému vysazení motoru z provozu, což by znamenalo nemalé finanční ztráty. Dále je důležité, aby i přes uvedené problémy se vstupními daty byl algoritmus natolik přesný, aby dokázal včas lokalizovat začátek problému a odhadnout zbývající životnost motoru.

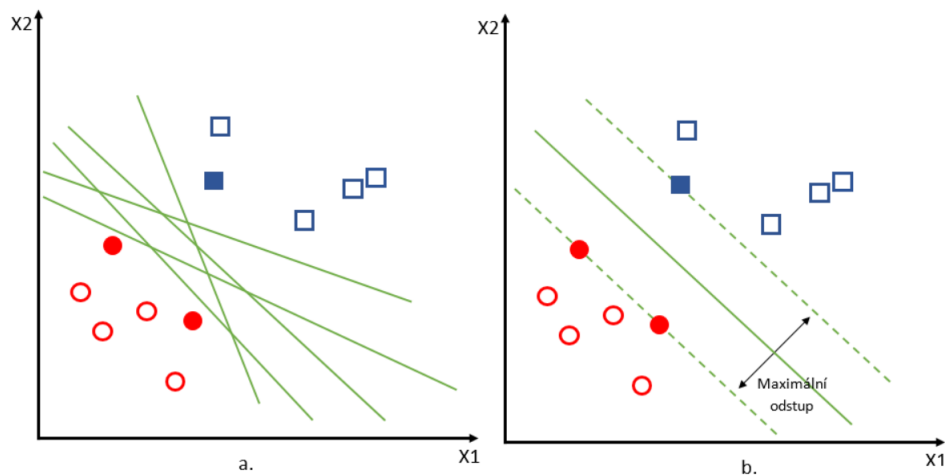
V této diplomové práci budou pro detekci anomálií v sensorových datech využity tři již zmíněné detekční metody. První z nich, metoda SVM, je jedinnou klasifikační metodou, kterou lze jako *1-Class SVM* využít pro detekci odlehlých hodnot. Druhá vybraná metoda, K-means, je z kategorie metod shlukové analýzy. Tato metoda je popsána v [40] a patří k relativně jednoduchým, avšak efektivním metodám, které byly již dříve využívány pro detekci anomálií. Poslední uvedenou metodou je metoda ARIMA. Opět je to hojně diskutovaná metoda, která je založena na odlišném principu než dvě předchozí metody. Inspiraci pro výběr této metody jsem našla v [42]. Model ARIMA se zaměřuje spíše na predikci dat, ale za využití dat historických je možné predikovanou hodnotu odečítat od reálně naměřené hodnoty a na základě nastavených hranic je pak možné definovat, zda se jedná o anomálii, či nikoliv.

### 3.1 Support Vector Machines

Support Vector Machines (SVM) neboli metoda podpůrných vektorů je relativně nová třída algoritmů strojového učení, které byly původně vyvinuty pro třídění dat do dvou různých tříd. Základy metody SVM položil v roce 1982 zejména Vladimir Vapnik [31]. Tato metoda se stala proslulou díky úspěchu v klasifikaci ručně psaných poštovních směrovacích čísel v Americe [32]. Jedná se o typický algoritmus strojového učení, hledající nadrovinu, která v prostoru příznaků optimálně rozděljuje vstupní data do dvou tříd. Algoritmus slouží k lineární separaci i takových dat, kde není separace možná. Požadavkem při hledání nadroviny je maximální vzdálenost mezi nadrovinou a nejbližším prvkem jednotlivých tříd 3.1. Cílem SVM je nalezení

jediného optimálního lineárního klasifikátoru. Takto nalezená optimální nadrovina je uprostřed co nejširšího pásma, které odděluje jednu třídu dat od třídy druhé. V tomto pruhu se nenachází žádný bod. Česky je tento pruh nazýván často pásmo necitlivosti či hraniční pásmo. K popisu nadroviny slouží nejbližší body, kterých je obvykle málo, nazývají se *podpůrné vektory*. Své jméno získala metoda právě podle těchto vektorů. Původní metoda SVM je binární, což znamená, že rozděluje data do dvou tříd [33].

Důležitou součástí techniky SVM je jádrová transformace (angl. kernel transformation) prostoru příznaků dat do prostoru transformovaných příznaků typicky vyšší dimenze. Tato jádrová transformace umožňuje převést původně lineárně neseparovatelnou úlohu na úlohu lineárně separovatelnou, na kterou lze dále aplikovat optimalizační algoritmus pro nalezení rozdělující nadroviny [33].



Obr. 3.1: Princip klasifikace metodou SVM. Na obrázku a. jsou všechny možné nadroviny a na obrázku b. je již nalezena optimální nadrovina.

### 3.1.1 Rozdělení SVM

SVM zahrnuje velké množství algoritmů, které můžeme rozdělit do skupin podle linearit a separace dat:

- lineárně separabilní případ,
- lineárně neseparabilní případ,
- nelineární SVM.

Lineární SVM je jednodušší varianta metody SVM, kdy není potřeba převod z původního prostoru příznaků a nedochází k žádné jádrové transformaci. Výsledkem je pak čistě lineární klasifikátor. Pokud jsou klasifikační třídy lineárně separabilní,

existuje nekonečně mnoho hranic, které dokáží rozdělit prostor tak, aby na jedné straně hranice byly pouze objekty z jedné třídy a na druhé straně hranice pouze objekty z druhé třídy viz obrázek 3.1a. Složitější variantou lineárních SVM je případ, kdy se snažíme lineárně oddělit data, která nejsou plně lineárně separovatelná. Například se může jednat o zašuměná data, kde se jednotlivé třídy částečně překrývají a není proto možné najít jednoznačnou hranici. V takovém případě chceme najít takovou rozdělovací nadrovinu, aby k chybné klasifikaci docházelo co nejméně [27].

### Lineárně separabilní případ

Obecně nejjednodušším případem použití SVM je lineární klasifikátor pro lineárně separabilní data, tedy taková, ve kterých lze obě třídy od sebe oddělit nadrovinou. Mějme tedy množinu trénovacích dat  $\{\mathbf{x}_i, y_i\}, i = 1, \dots, l, y_i \in \{-1, 1\}, \mathbf{x}_i \in \mathbf{R}^d$ . Předpokládejme, že máme nadrovinu, která odděluje pozitivní od negativních vzorků dat. Body  $\mathbf{x}$ , které vyhovují podmínce  $\mathbf{w} \cdot \mathbf{x} + b = 0$ , kde  $\mathbf{w}$  je normálový vektor k nadrovině,  $|b| / \|\mathbf{w}\|$  je vzdálenost od nadroviny k počátku a  $\|\mathbf{w}\|$  je Euklidovská norma  $\mathbf{w}$ . Necht  $d_+(d_-)$  je nejkratší vzdálenost od dělicí nadroviny k nejbližšímu pozitivnímu (negativnímu) bodu. Pásmo necitlivosti oddělovací nadroviny jsou pak definováno jako  $d_+ + d_-$ . Pro lineárně separovatelná data algoritmus jednoduše hledá oddělovací nadrovinu s nejširším pásmem necitlivosti. To lze formulovat následovně: předpokládejme, že všechna trénovací data splňují následující podmínku [34]:

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 \quad \text{pro } y_i = +1, \quad (3.1)$$

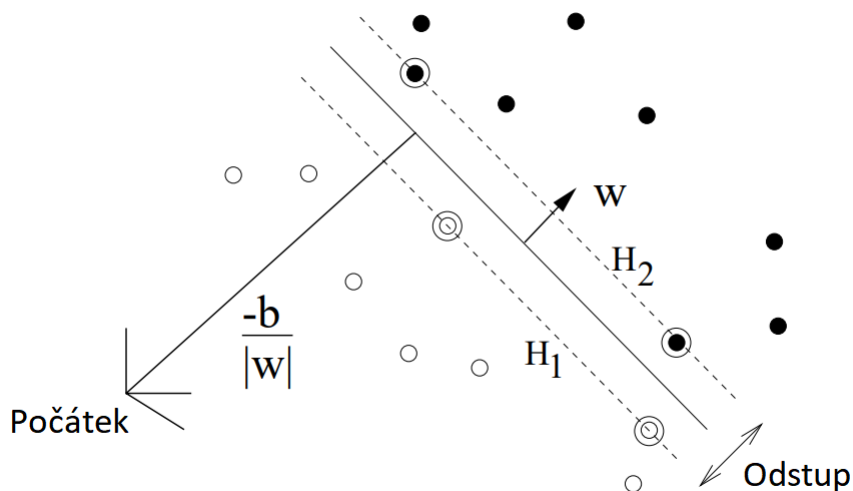
$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \quad \text{pro } y_i = -1, \quad (3.2)$$

což lze sloučit do jedné nerovnice

$$y_i (\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad \forall i. \quad (3.3)$$

Nyní uvažujme body, pro které platí rovnost v 3.1 (požadavek, že takový bod existuje, je ekvivalentní volbě rozsahu pro  $\mathbf{w}$  a  $b$ ). Tyto body leží na nadrovině  $H_1 : \mathbf{x}_i \cdot \mathbf{w} + b = 1$ . Podobně body, pro které platí rovnost v 3.2, leží v nadrovině  $H_2 : \mathbf{x}_i \cdot \mathbf{w} + b = -1$ . Nadroviny  $\mathbf{H}_1$  a  $\mathbf{H}_2$  jsou kolmé a mají stejný normálový vektor  $\mathbf{w}$  a jejich vzdálenost od rozdělovací nadroviny je  $1 / \|\mathbf{w}\|$ . Šířka hraničního pásma je tedy  $2 / \|\mathbf{w}\|$  a nadroviny s nejširším pásmem můžeme určit minimalizací  $\|\mathbf{w}\|^2$  vzhledem k podmínkám 3.3 [34].

Očekáváme tedy, že řešení pro typický dvourozměrný případ bude mít podobu znázorněnou na obrázku 3.2. Tyto trénovací body, pro které platí rovnost 3.3 a tedy ty, jejichž odstraněním by se změnilo nalezené řešení, se nazývají podpůrné vektory [34].



Obr. 3.2: Lineární dělicí nadroviny pro separabilní případ. Podpůrné vektory jsou zakroužkované. Převzato z [34].

Nyní přejdeme k Lagrangovské formulaci problému. Pro to existují dva důvody. První je, že omezení 3.3 budou nahrazena omezeními samotných Lagrangeových multiplikátorů, s nimiž bude mnohem snazší manipulovat. Druhým je, že v tomto přeformulování problému se trénovací data vyskytují ve všech výrazech jen ve skalárních součinech. Toto je klíčová vlastnost, která nám umožní zobecnit postup pro nelineární případ.

Zavedeme tedy nezáporné Lagrangeovy multiplikátory  $\alpha_i, i = 1, \dots, l$ , jeden pro každou nerovnost v 3.3. Lagrangeova funkce:

$$L_P = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l \alpha_i y_i (x_i \cdot \mathbf{w} + b) + \sum_{i=1}^l \alpha_i \quad (3.4)$$

Nyní musíme minimalizovat  $L_P$  vzhledem k  $\mathbf{w}$  a derivace  $L_P$  vzhledem ke všem  $\alpha_i$  musejí být nulové. Vše za omezení  $\alpha_i \geq 0$ . Jedná se o konvexní problém kvadratického programování. To znamená, že můžeme místo toho řešit duální problém - hledání maxima  $L_P$  za následujících podmínek: že gradient  $L_P$  je vzhledem k  $\mathbf{w}$  a  $b$  nulový, a že  $\alpha_i \geq 0$ . Tato konkrétní dvojité formulace problému se nazývá Wolfova dualita [35].

Aby byl gradient  $L_P$  podle  $\mathbf{w}$  a  $b$  nulový, musí platit:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \quad (3.5)$$

$$\sum_i \alpha_i y_i = 0. \quad (3.6)$$

Dosazením do 3.4 se získá duální Lagrangeova funkce

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \quad (3.7)$$

Nyní máme dvě různá označení Lagrangeových funkcí.  $L_P$  pro primární a  $L_D$  pro duální funkci. Je nutné zdůraznit, že tyto dvě formulace jsou odlišné. Vyplývají ze stejné objektivní funkce, ale mají různá omezení. Řešení je pak nalezeno minimalizací  $L_P$  nebo maximalizací  $L_D$ .

Trénování SVM spočívá v maximalizaci  $L_D$  vzhledem k  $\alpha_i$ , při omezení 3.7 a  $\alpha_i \geq 0$ , s řešením daným v 3.5. Pro každý trénovací bod máme jeden Lagrangeův multiplikátor  $\alpha_i$ , přičemž body, pro které platí  $\alpha_i > 0$ , jsou podpůrné vektory a leží v nadrovinách  $\mathbf{H}_1$  nebo  $\mathbf{H}_2$ , zatímco všechny ostatní trénovací body mají  $\alpha_i = 0$  a leží buď v jedné z nadrovin, nebo na té straně od nich, pro kterou platí ostrá nerovnost v 3.3. Podpůrné vektory jsou nejdůležitějšími body trénovací množiny - pokud by se všechny ostatní body odstranily, novým natrénováním by se získala tatáž dělicí nadrovina.

### Karushovy-Kuhnovy-Tuckerovy podmínky

Karushovy-Kuhnovy-Tuckerovy (KKT) podmínky jsou nutné podmínky pro optimální řešení v úloze nelineárního programování za předpokladu splnění určitých podmínek regularity. Pro konvexní problémy jsou zároveň i postačujícími podmínkami. SVM představují konvexní problém a podmínky regularity jsou pro ně splněné vždy. KKT podmínky jsou tím pádem nutné a postačující, aby  $\mathbf{w}$ ,  $\mathbf{b}$ ,  $\alpha$  byly řešením [34].

Pro primární problém KKT podmínky vypadají následovně:

$$\frac{\partial}{\partial w_v} L_P = w_v - \sum_i \alpha_i y_i x_{iv} = 0 \quad v = 1, \dots, d \quad (3.8)$$

$$\frac{\partial}{\partial w_v} L_P = w_v - \sum_i \alpha_i y_i = 0 \quad (3.9)$$

$$y_i (\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad i = 1, \dots, l \quad (3.10)$$

$$\alpha_i \geq 0 \quad \forall i \quad (3.11)$$

$$\alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0 \quad \forall i \quad (3.12)$$

KKT podmínky hrají důležitou roli při numerickém řešení SVM úlohy.

### Lineárně neseparabilní případ

Složitější variantou lineárních SVM je případ, kdy se snažíme lineárně oddělit data, která nejsou plně lineárně separovatelná. Například se může jednat o zašuměná

data, kde se jednotlivé třídy částečně překrývají a není proto možné najít jednoznačnou hranici. V takovémto případě chceme najít takovou rozdělující nadrovinu, aby k chybné klasifikaci docházelo co nejméně.

V takové situaci se provádí zmírnění podmínek 3.1 a 3.2 s přidanou cenou za jejich porušení, a to zavedením přídatné proměnné  $\xi, i = 1, \dots, l$ :

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 - \xi_i \quad \text{pro } y_i = +1, \quad (3.13)$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 + \xi_i \quad \text{pro } y_i = -1, \quad (3.14)$$

$$\xi_i \geq 0 \quad \forall i. \quad (3.15)$$

Aby nastala chyba, musí odpovídající  $\xi_i$  překročit jednotu, takže  $\sum_i \xi_i$  je horní hranice počtu chyb v trénovacích datech. Přirozeným způsobem jak přiřadit dostatečnou cenu za chyby je stanovení minimalizované funkce místo  $\|\mathbf{w}\|^2 / 2$  jako  $\|\mathbf{w}\|^2 / 2 + C(\sum_i \xi_i)^k$ , kde  $C$  je volitelný parametr, přičemž větší  $C$  znamená větší penalizaci chyb.

Duální problém pak je:

Maximalizovat

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \quad (3.16)$$

za podmínek

$$0 \leq \alpha_i \leq C, \quad (3.17)$$

$$\sum_i \alpha_i y_i = 0. \quad (3.18)$$

Řešení je opět

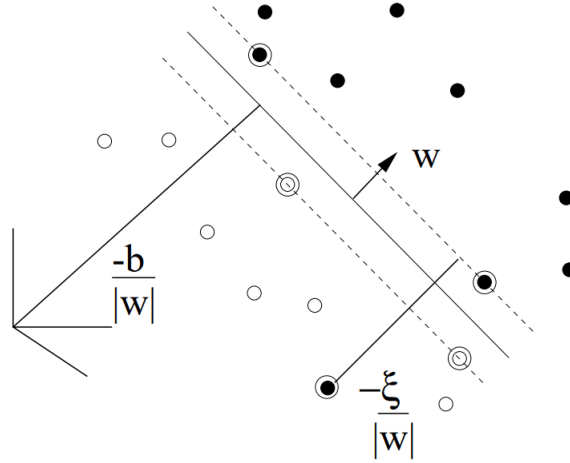
$$\mathbf{w} = \sum_{i=1}^{N_S} \alpha_i y_i \mathbf{x}_i, \quad (3.19)$$

kde  $N_S$  je počet podpůrných vektorů. Jediný rozdíl oproti předchozímu případu je tedy ten, že  $\alpha_i$  jsou nyní shora omezené  $C$ . Tato situace je schématicky znázorněna na obrázku 3.3.

Budeme potřebovat KKT podmínky pro primární problém. Primární Lagrangeova funkce je:

$$L_P = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i - \sum_i \alpha_i \{y_i (\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i\} - \sum_i \mu_i \xi_i, \quad (3.20)$$

kde  $\mu_i$  jsou Lagrangeovy multiplikátory zavedené pro zajištění  $\xi_i \geq 0$ . KKT podmínky pro primární problém jsou pak



Obr. 3.3: Lineární dělicí nadroviny pro neseparabilní případ. Převzato z [34].

$$\frac{\partial L_P}{\partial w_\nu} = w_\nu - \sum_i \alpha_i y_i x_{i\nu} = 0 \quad (3.21)$$

$$\frac{\partial L_P}{\partial b} = - \sum_i \alpha_i y_i = 0 \quad (3.22)$$

$$\frac{\partial L_P}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad (3.23)$$

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i \geq 0 \quad (3.24)$$

$$\xi_i \geq 0 \quad (3.25)$$

$$\alpha_i \geq 0 \quad (3.26)$$

$$\mu_i \geq 0 \quad (3.27)$$

$$\alpha_i \{y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i\} = 0 \quad (3.28)$$

$$\mu_i \xi_i = 0, \quad (3.29)$$

kde  $i$  jde od 1 do počtu trénovacích bodů,  $\nu$  jde od 1 do dimenze dat.

Stejně jako v předchozím případě můžeme použít KKT podmínky komplementarity, rov. 3.28 a rov. 3.29 pro stanovení prahu  $b$ . Všimněme si, že rov. 3.28 v kombinaci s rov. 3.29 ukazuje, že  $\xi_i = 0$  jestliže  $\alpha_i < C$ . Pak lze pro určení  $b$  jednoduše dosadit do rov. 3.28 kterýkoliv trénovací bod, pro který platí  $0 < \alpha_i < C$  (s  $\xi_i = 0$ ) [34].

## Nelineární SVM

Jak lze výše uvedené metody zobecnit na případ, kdy rozhodovací funkce <sup>1</sup> není lineární? Základní myšlenkou je použití tzv. jádrového triku (anglicky kernel trick), s jehož pomocí se provádí transformace dat z původního prostoru příznaků do prostoru vyšší dimenze, ve kterém jsou data lineárně separabilní. Jinými slovy, provádíme zobrazení trénovacích dat z původního prostoru do jiného vícerozměrného, eukleidovského prostoru  $\mathcal{H}$ , pomocí funkce:

$$\Phi : \mathbf{R}^d \mapsto \mathcal{H} \quad (3.30)$$

V tomto vícerozměrném prostoru lze též definovat novou optimální separační nadrovinu. Kdyby existovala jaderná funkce  $K$ , kde  $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$ , nemusíme vůbec znát zobrazení  $\Phi$  a k natrénování algoritmu by stačilo pouze  $K$ , pro které v tomto případě platí:

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2}. \quad (3.31)$$

V tomto konkrétním případě má  $\mathcal{H}$  nekonečný rozměr, takže by nebylo snadné explicitně pracovat s  $\Phi$ . Nicméně pokud nahradíme  $\mathbf{x}_i \cdot \mathbf{x}_j$  výrazem  $K(\mathbf{x}_i, \mathbf{x}_j)$  kdekoliv v trénovacím algoritmu, budou všechny předchozí úvahy platit, jelikož se bude stále jednat o lineární separaci, pouze v jiném prostoru. V testovací fázi lze použití  $\mathbf{w}$  také obejít - pro určení do které třídy bod spadá, počítáme funkci sign z výrazu

$$f(x) = \sum_{i=1}^{N_S} \alpha_i y_i \Phi(\mathbf{s}_i) \cdot \Phi(\mathbf{x}) + b = \sum_{i=1}^{N_S} \alpha_i y_i K(\mathbf{s}_i, \mathbf{x}) + b, \quad (3.32)$$

kde  $\mathbf{s}_i$  jsou podpůrné vektory. Takže se opět můžeme vyhnout explicitnímu výpočtu  $\Phi(\mathbf{x})$  a namísto toho použít  $K(\mathbf{s}_i, \mathbf{x}) = \Phi(\mathbf{s}_i) \cdot \Phi(\mathbf{x})$ . Znaménko  $f(x)$  opět reprezentuje příslušnost vektorů k jedné ze dvou tříd [34].

### Příklady nelineárních SVM

Problematika jádrových funkcí používaných v SVM je velmi rozsáhlá a stále se vyvíjí. Proto v následujícím odstavci budou uvedeny pouze tři nejčastěji používané funkce, ačkoli jako jádrové funkce mohou být použity všechny ty, které splňují tzv. Mercerovu podmínku [34].

Mezi často používané jádrové funkce patří například:

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^p \quad (3.33)$$

$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2} \quad (3.34)$$

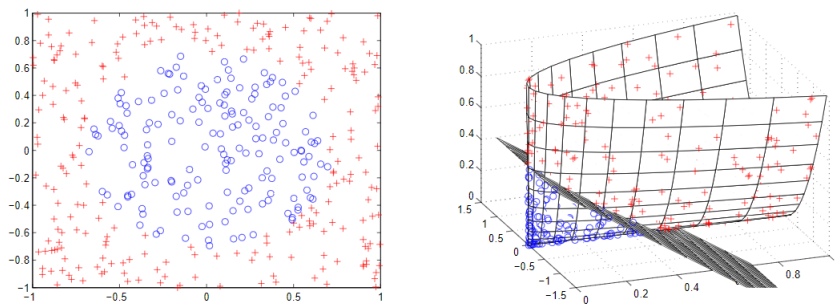
$$K(\mathbf{x}, \mathbf{y}) = \tanh(\kappa \mathbf{x} \cdot \mathbf{y} - \delta) \quad (3.35)$$

Rovnice 3.33 je polynom stupně  $p$ , rovnice 3.34 je Gausovská radiální bázová funkce a rovnice 3.35 je konkrétní typ dvouvrstvé sigmoidální neuronové sítě [34].

---

<sup>1</sup>„Rozhodovací funkcí“ se rozumí funkce  $f(x)$ , jejíž znaménko představuje třídu přiřazenou datovému bodu  $x$ .





Obr. 3.4: Znázornění nelineární SVM.

### 3.1.2 One-Class SVM

Metoda One-Class SVM byla navržena Schölkopfem a kol. [36] jako varianta původního SVM algoritmu. Cílem této techniky je právě detekce odlehlých bodů, která je taktéž primárním cílem této diplomové práce. Hlavní rozdíl mezi One-Class SVM a klasickou SVM metodou je ten, že první uvedená metoda vyžaduje pro trénování modelu pouze datové body normálního provozního režimu, zatímco metoda uvedená jako druhá vyžaduje pro natrénování modelu vzorky jak normálního provozního režimu, tak vzorky chybových dat, na jejichž základě je definován klasifikátor. V podstatě jde o to že, metoda One-Class SVM definuje hranici, která vyhovuje většině trénovacích bodů. Tato část algoritmu představuje právě trénovací fázi. Jestliže daný datový bod spadá do této hranice, je zařazen do kategorie normálních operačních dat, jinak je klasifikován jako odlehlý bod, tedy anomálie. Pro nelineární hranice je k mapování dat do prostoru s vyššími dimenzemi použita opět jádrová funkce.

Tato metoda je založena na dříve uvedených úvahách. Uvažujeme-li tedy trénovací množinu dat  $x_i \in \mathcal{R}^n, i = 1, \dots, m$  a každý bod  $x_i$  může být přiřazen k označení třídy  $y_i \in 1$ . Vzorové vektory jsou mapovány do prostoru s vyššími dimenzemi pomocí mapovací funkce  $\Phi$ , ve kterém je lineární nadrovina konstruována tak, aby její vzdálenost od počátku byla maximální. To může být definováno následovně: Je třeba minimalizovat  $\frac{1}{2} \| \mathbf{w} \|^2$  vzhledem k podmínce  $\mathbf{w} \cdot \Phi(x_i) + b \geq 0 \forall i$ . V prostoru příznaků však není vždy možné dokonale oddělit vzorové vektory od počátku. V takovém případě nebude optimalizační problém uvedený výše konvergovat s konečným řešením. Za tímto účelem definujeme  $\nu \in (0, 1)$  a proměnnou  $\xi_i$  k uvolnění optimalizačních podmínek. Proměnná  $\nu$  odpovídá maximální hodnotě zlomku chyb tréninkového souboru a  $b$  je vzdálenost od počátku v prostoru příznaků. Upravený primární problém optimalizace může být zapsán jako [37]:

$$\min_{w, \xi_i, b} \frac{1}{2} \| w \|^2 + \frac{1}{\nu m} \sum_{i=1}^m \xi_i - b \quad (3.36)$$

podléhající podmínce

$$w \cdot \Phi(x_i) \geq b - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, m. \quad (3.37)$$

Stejně jak bylo uvedeno u binárních metod SVM jsou zde zavedeny Lagrangeovy multiplikátory  $\alpha_i \geq 0$  a  $\mu_i \geq 0$ . Parciální derivace Lagrangeovy funkce jsou nastaveny na nulu následovně [37]:

$$L(w, \xi, b) = \frac{1}{2} \|w\|^2 + \frac{1}{\nu m} \sum_{i=1}^m \xi_i + b - \sum_{i=1}^m \alpha (w \cdot \Phi(x_i) + b + \xi_i) - \sum_{i=1}^m \nu_i \xi_i \quad (3.38)$$

$$\frac{\partial L}{\partial w} = 0 \rightarrow w = \sum_{i=1}^m \alpha_i \Phi(x_i) \quad (3.39)$$

$$\frac{\partial L}{\partial \xi_i} = 0 \rightarrow \alpha_i = \frac{1}{\nu m} - \mu_i, \quad \nu_i \in (0, 1) \quad (3.40)$$

$$\frac{\partial L}{\partial b} = 0 \rightarrow \sum_{i=1}^m \alpha_i = 1. \quad (3.41)$$

Je zvolena jádrová funkce  $K(x, y) = \Phi(x_i) \cdot \Phi(x_j)$  a duální forma Lagrangeovy funkce je pak zapsána jako:

$$\min_{\alpha_s} \alpha_i \alpha_j K(x_i, x_j), \quad (3.42)$$

vzhledem k podmínkám

$$0 \leq \alpha_i \leq 1, \quad \sum_{i=1}^m \alpha_i = \nu m \quad (3.43)$$

Řešením dvojího problému pro  $\alpha_i$  a rekonstrukcí  $b$  z primárního problému se získá rozhodovací funkce

$$f(x) = \sum_{i=1}^m \alpha_i k(x_i, x) - b. \quad (3.44)$$

Záporná hodnota  $f(x)$  značí, že bod  $x$  je odlehlá hodnota. Datové body, pro které platí  $0 < \alpha_i < \frac{1}{m\nu}$ , jsou podpůrné vektory.

### 3.1.3 Analýza hlavních komponent

Na základě dostupných informací je pro předzpracování dat využita metoda analýzy hlavních komponent (PCA) [37].

Analýza hlavních komponent (Principal Component Analysis, PCA) je metoda sloužící k dekorelaci dat. Často se používá ke snížení dimenze dat s co nejmenší ztrátou informace. Snížení počtu proměnných daného datového souboru je obvykle

na úkor přesnosti, avšak takové řešení je kompromisem přinášející zjednodušení pro další zpracování. Data lze pak mnohem snadněji prozkoumat, vizualizovat a analyzovat a algoritmy strojového učení konzumující takto předzpracovaná data jsou pak mnohem rychlejší. Vzhledem k tomu, že menší datové sady lze snadněji prozkoumat a vizualizovat je analýza dat pomocí algoritmů strojového učení mnohem jednodušší a rychlejší. Myšlenka PCA je tedy jednoduchá - snížení počtu proměnných dané sady dat při zachování co největšího počtu informací.

## Normalizace dat

Cílem tohoto kroku je normalizovat rozsah vstupních proměnných. Důvodem proč je nezbytné provést normalizaci dat před samotným užitím PCA je takový, že je tato metoda velmi citlivá na odchylky vstupních proměnných. To znamená, že pokud jsou mezi rozsahy vstupních proměnných velké rozdíly, budou proměnné s většími rozsahy dominovat nad proměnnými s malými rozsahy (například proměnná, která se pohybuje v rozsahu mezi 0 a 100, bude dominovat nad proměnnou, která se pohybuje mezi 0 a 1), což povede ke zkresleným výsledkům.

## Výpočet kovarianční matice

Cílem tohoto kroku je pochopit, jak se proměnné v sadě vstupních dat liší od sebe navzájem, nebo jinými slovy, zjistit, zda mezi nimi existuje nějaký vztah. Někdy jsou proměnné korelovány takovým způsobem, že obsahují nadbytečné informace. Abychom tedy identifikovali tyto korelace, spočítáme kovarianční matici.

Kovarianční matice je symetrická matice  $p \times p$  (kde  $p$  je počet dimenzí), která má na vstupu kovariance všech možných párů všech vstupních proměnných. Například pro třídimenziální datovou množinu s třemi proměnnými  $x$ ,  $y$  a  $z$  je kovarianční matice  $3 \times 3$  matice:

$$\begin{bmatrix} Cov(x, x) & Cov(x, y) & Cov(x, z) \\ Cov(y, x) & Cov(y, y) & Cov(y, z) \\ Cov(z, x) & Cov(z, y) & Cov(z, z) \end{bmatrix} \quad (3.45)$$

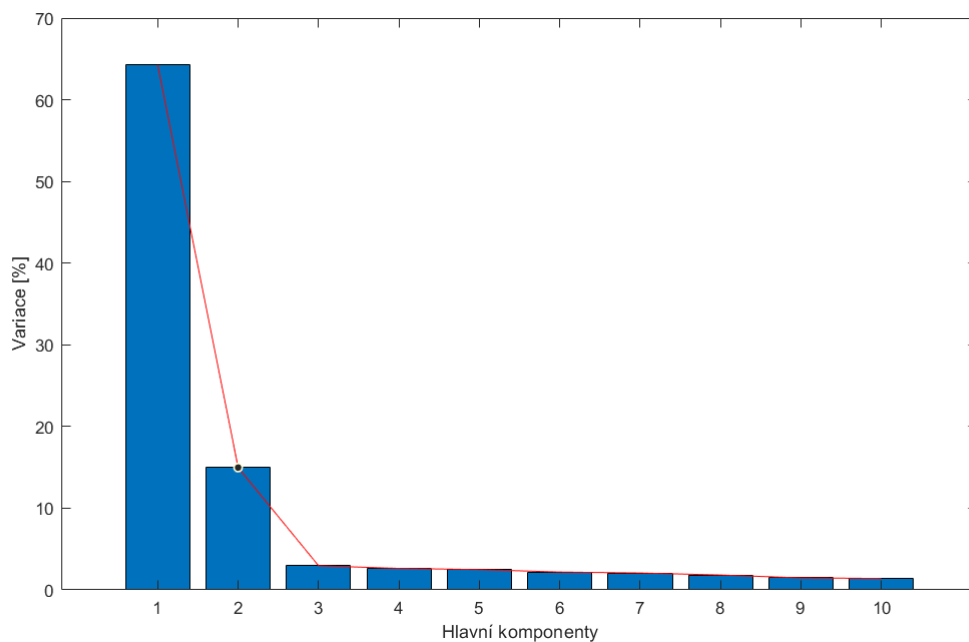
Jelikož kovariance proměnné samé se sebou je její variance ( $Cov(a, a) = Var(a)$ ), je pak na hlavní diagonále variance každé vstupní proměnné. A jelikož má kovariance komutativní vlastnost ( $Cov(a, b) = Cov(b, a)$ ) jedná se o symetrickou matici. Jednotlivé hodnoty v kovarianční matici nám říkají zda a jak moc jsou jednotlivé vstupní proměnné korelovány.

- Jestliže je hodnota pozitivní, potom jednotlivé proměnné rostou či klesají dohromady, tedy jsou korelovány.

- Pokud je hodnota negativní pak jedna proměnná roste, zatímco druhá klesá, tedy nejsou korelovány.

### Výpočet vlastních vektorů a vlastních hodnot pro identifikaci hlavních komponent

Vlastní vektory a vlastní hodnoty jsou pojmy lineární algebry, které je nutné vypočítat z kovarianční matice, aby bylo možné určit hlavní komponenty. Definice hlavních komponent je následující: Hlavní komponenty jsou nové proměnné jež jsou konstruovány jako lineární kombinace počátečních proměnných. Tyto kombinace se provádějí takovým způsobem, že nové proměnné (tj. hlavní komponenty) nejsou korelovány a většina informací uvnitř počátečních proměnných je komprimována do první složky. Myšlenka je taková, že desetirozměrná data dají na výstupu deset hlavních komponent, ale PCA se snaží dát maximální možné informace do první složky, zbývající informace do druhé složky a tak dále, jak je naznačeno na obrázku 3.5.



Obr. 3.5: Příklad procentuálního zastoupení informace v každé z hlavních komponent.

## 3.2 K-means

Pojem *K-means* poprvé použil James MacQueen v roce 1967 [39]. K-means je metoda nehierarchické shlukové analýzy, která seskupuje objekty na základě jejich hodnot do  $K$  nesouvisejících shluků, jejichž počet je předem dán. Objekty které jsou klasifikovány do stejného shluku mají podobné funkční hodnoty. Shluky jsou definovány svými centroidy, což jsou body ve stejném prostoru jako shlukované objekty. Objekty se zařazují do toho shluku, jehož centroidu jsou nejbližší. Algoritmus postupuje v následujících krocích [40]:

1. Definice počtu shluků.
2. Inicializace  $K$  centroidů a to tak, že všechny objekty budou libovolně rozděleny do  $K$  shluků, dále se definují jejich centroidy a dále je pak nutné ověřit, zda se jednotlivé centroidy liší. Alternativně mohou být centroidy inicializovány na  $K$  náhodně vybraných různých objektech.
3. Dále je pak provedena iterace přes všechny objekty a vypočítána vzdálenost každého objektu ke každému z definovaných centroidů. Objekty jsou pak přiřazeny do shluku s nejbližším centroidem. Přepočítání centroidů nově modifikovaných shluků.
4. Krok 3 je opakován, dokud dochází ke změnám centroidů po každé iteraci.

Pro výpočet vzdálenosti mezi dvěma objekty je definována funkce. Nejčastěji používaná distanční funkce je euklidovská, která je definována jako:

$$d(y, x) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}, \quad (3.46)$$

kde  $x = (x_1, \dots, x_m)$  a  $y = (y_1, \dots, y_m)$  jsou dva vstupní vektory s  $m$  kvantitativními znaky. V euklidovské vzdálenostní funkci, všechny příznaky přispívají stejně k výsledné hodnotě funkce. Avšak ve většině případů mají různá měření různé rozsahy, proto je nutné před použitím euklidovské vzdálenostní funkce tyto veličiny normalizovat.

Alternativou k Euklidovské vzdálenosti je Mahalanobisova vzdálenost, která používá inverzní kovarianční matici  $S^{-1}$  k určení korelace mezi jednotlivými veličinami.

$$d(y, x) = \sqrt{(x - y)^T S^{-1} (x - y)} \quad (3.47)$$

Výpočet a inverze kovarianční matice je však výpočetně náročná pro velký počet vícerozměrných veličin [40].

### 3.2.1 Detekce anomálií pomocí K-Means

Algoritmus K-means je aplikován na sadu dat, která může obsahovat jak normální, tak anomální data, aniž by byly předem označeny. Existuje předpoklad, že anomální

a normální data budou v prostoru příznaků tvořit různé shluky. Data však mohou také obsahovat odlehlé hodnoty, které nepatří do většího shluku. Pokud je počet takovýchto hodnot malý, proces shlukování není narušen.

K-means algoritmus rozděluje trénovací data do  $K$  shluků, ale neurčuje, zda daný shluk odráží normální či nenormální data. To musí být definováno manuálně nebo heuristicky. Může se stát, že clustery budou velmi blízko u sebe, což může mít několik důvodů. Jedním z nich je například, že trénovací data jsou příliš homogenní, to znamená že trénovací data neobsahují žádné anomální hodnoty, nebo anomální chování vypadá velmi podobně jako normální chování.

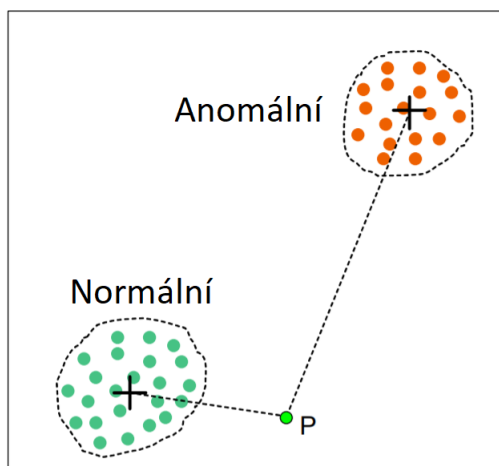
Základním problémem metody K-means je definování vhodného počtu shluků  $K$ . Jako počáteční hodnota byla zvolena  $K = 2$ , jelikož předpokládáme, že normální a anomální trénovací data tvoří dva různé shluky. Proces shlukování vede k vytvoření centroidů pro normální a anomální data, které mohou být použity pro detekci anomálií v nově přichozích datech. Nové záznamy musí být předem zpracovány a transformovány, aby bylo docíleno stejných podmínek jako u trénovacích dat. Pro účely detekce anomálií byly použity dvě metody založené na měření vzdálenosti, které je možné použít samostatně nebo kombinovaně [40].

## Klasifikace

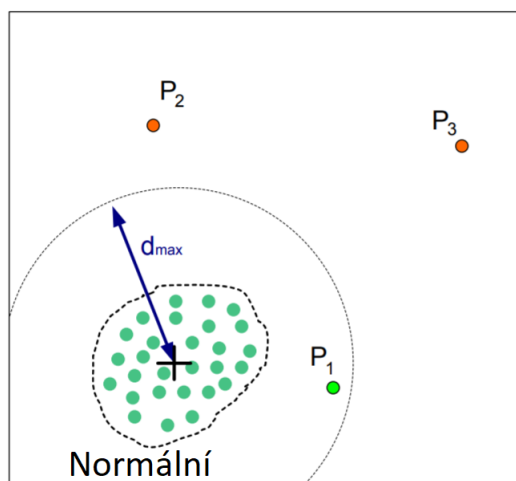
Vzdálenost k centroidu odpovídajícímu shluku je vypočítána pomocí vážené Euklidovské funkce. Daný datový bod je klasifikován jako normální nachází-li se blíže k centroidu jež byl definován jako normální než k anomálnímu a naopak. To je znázorněno na obrázku 3.6 s dvojrozměrným příznakovým prostorem. Bod  $P$  je blíže k centroidu, který byl určen jako normální, proto klasifikujeme  $P$  jako normální. Tato klasifikace na základě vzdálenosti umožňuje detekovat známé typy anomálií, tj. anomální data s podobnými charakteristikami jako v souboru trénovacích dat [40].

## Detekce odlehlých bodů

Odlehlý bod je takový datový bod, který se výrazně liší od ostatních. Proto jej lze považovat za anomálii. Pro detekci odlehlých bodů se vypočítá pouze vzdálenost k příslušnému centroidu normálního shluku. Pokud je vzdálenost mezi objektem a centroidem větší než předdefinovaná prahová hodnota  $d_{max}$  je bod považován za anomálii. To je znázorněno na obrázku 3.7, kde bod  $P_2$  a  $P_3$  leží mimo okruh  $d_{max}$ . Na rozdíl od klasifikační metody, detekce odlehlých bodů nevyužívá centroidy shluků, jež byly označeny jako anomální. To znamená, že může být méně přesná pro detekci známých druhů anomálií. Na druhou stranu umožňuje detekci nových anomálií, které se nevyskytovaly v souboru trénovacích dat.



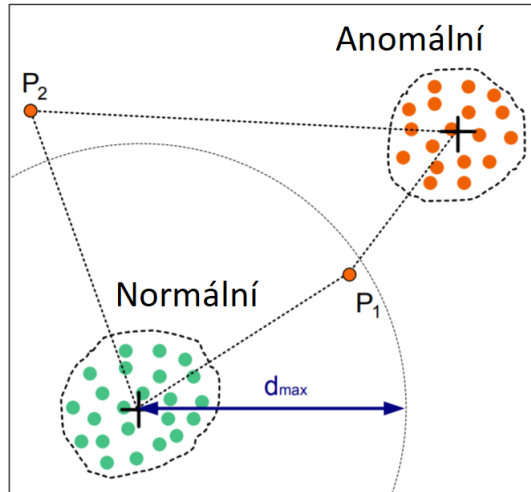
Obr. 3.6: Klasifikace pro  $K = 2$ . Převzato z [40].



Obr. 3.7: Detekce odlehlých bodů. Převzato z [40].

### Kombinovaná metoda klasifikace a detekce odlehlých bodů

Za účelem překonání omezení každé jednotlivé metody lze použít klasifikaci a detekci odlehlých hodnot kombinovaným způsobem. Pokud jsou tyto dvě metody aplikovány současně, je s objektem zacházeno jako s anomálií, pokud je blíže k centroidu anomálního shluku než k normálnímu, nebo pokud je jeho vzdálenost k centroidu normálního shluku větší než předdefinovaná prahová hodnota. Na obrázku 3.8 jsou oba body  $P_1$  a  $P_2$  definovány jako anomálie. Bod  $P_1$  je blíže anomálnímu shluku a vzdálenost bodu  $P_2$  od centroidu normálního shluku je větší než  $d_{max}$ .



Obr. 3.8: Kombinace klasifikace a detekce odlehlých bodů. Převzato z [40].

### 3.3 ARIMA

Metoda ARIMA (zkratka anglického AutoRegressive Integrated Moving Average, „autoregresní integrovaný klouzavý průměr“) je metoda sloužící k pochopení vlastností časových řad a k predikci jejich chování. Metodu navrhli v sedmdesátých letech George Box a Gwilym Jenkins [41]. Model ARIMA se skládá ze tří komponent, kdy každá z nich pomáhá modelovat určitý typ vzoru [42]:

- **Autoregresní (AR)**, která vyjadřuje, že současná hodnota časové řady se dá vysvětlit jako lineární kombinace předchozích hodnot. Řád této komponenty se označuje  $p$  a označuje počet předchozích hodnot, ze kterých je tvořena současná hodnota časové řady. Například  $p = 2$  znamená, že současná hodnota časové řady je tvořena maximálně dvěma předcházejícími hodnotami.
- **Integrační (I)** komponenta, znamená diferenci časové řady. Řád integrační složky se značí  $d$  a udává počet po sobě se opakujících diferencí.
- Komponenta **klouzavého průměru (MA)**, která je lépe chápána jako termín označující chybovou zpětnou vazbu, měří přizpůsobení nových prognóz chybám předchozích prognóz. Řád MA složky se označuje  $q$  a (podobně jako u AR parametru  $p$ ) vyjadřuje z kolika časových intervalů v minulosti se chyby v modelu uplatní [42].

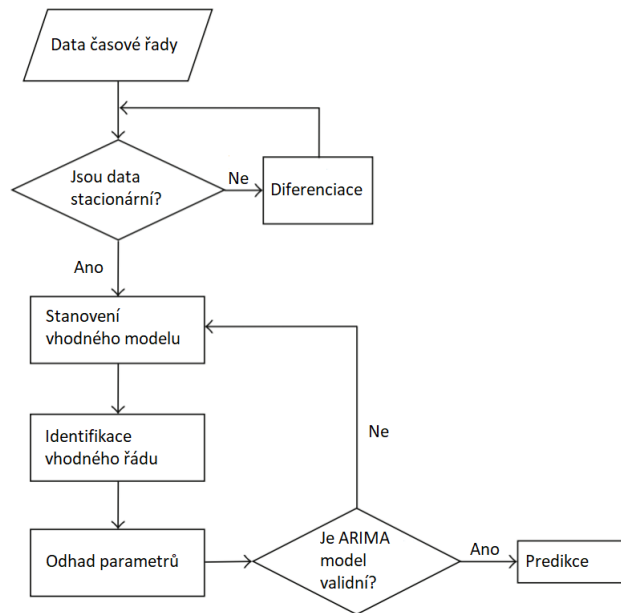
Obecně pro časové řady  $\{X_t, t = 1, 2, 3, \dots, n\}$  je střední hodnota vyjádřena jako  $E(X_t) = \mu$  a model ARIMA( $p, q, d$ ) lze vyjádřit jako:

$$\phi(B) (\nabla^d X_t - \mu) = \Theta(B) \epsilon_t, \quad t > d, \quad (3.48)$$

kde  $B, BX_t = X_{t-1}$  je operátor zpětného posunutí a  $\nabla^d \nabla^d X_t = (1 - B)^d X_t$  je  $d$ -tá



diference modelovaného procesu.



Obr. 3.9: ARIMA model. Převzato z [42].

Polynomy operátoru zpětného posunu jsou pak definovány jako:

$$\phi(B) = -\phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p \quad (3.49)$$

$$\theta(B) = -\theta_1 B - \theta_2 B^2 - \dots - \theta_p B^p \quad (3.50)$$

Vstupní data pro modely ARIMA jsou obecně nestacionární. Nejdříve je nutné použít vhodný proces pro úpravu dat, tak aby byla daná časová řada stacionární. Poté je nutné stanovit vhodný model a provést další kroky, jako je například identifikace řádu modelu, odhad parametrů a testování hypotéz. Celý postup návrhu ARIMA modelu je uveden na obrázku 3.9 [42].

1. V prvním kroku je nutné určit, zda je sekvence stacionární pomocí vhodného testu stacionarity. Pokud tento test selže je nutné provádět diferenciaci, dokud není sekvence stacionární.
2. Dále je vytvořen vhodný model na základě identifikačních pravidel modelu ARIMA.
3. Identifikace optimálního řádu pro stanovení modelu pomocí AIC (Akaike Information Criterion) nebo BIC (Bayesian Information Criterion) <sup>2</sup>.

<sup>2</sup>Ve statistice, Bayesovské informační kritérium (BIC) nebo Schwarzovo kritérium (také SBC, SBIC) je kritérium pro výběr modelu mezi konečnou sadu modelů. Je částečně založeno na funkci pravděpodobnosti a úzce souvisí s informačním kritériem Akaike (AIC) [43].

4. Odhad neznámých parametrů v modelu pomocí odhadu korelačního momentu, odhadu nejmenších čtverců nebo metody odhadu maximální pravděpodobnosti.
5. Identifikace reziduí. Zda je reziduum pouze sekvence bílého šumu, jestliže není je nutné vrátit se ke kroku 2.
6. Provedení prediktivní analýzy na základě zavedeného ARIMA modelu [42].



Obr. 3.10: Metoda exponenciálně váženého klouzavého průměru. Převzato z [42].

### Výhody modelu ARIMA

1. Model ARIMA je založen na Markovově náhodném procesu, který by mohl odrážet dynamické charakteristiky modelu ARIMA. Model ARIMA zahrnuje model AR, modely MA a model ARMA, který plně absorbuje výhody regresní analýzy a posiluje dobré vlastnosti klouzavého průměru.
2. ARIMA model může být aplikován na nestacionární časové řady, což je opravdu vhodné pro nestacionární charakteristiky většiny sensorových dat.
3. Model ARIMA může být použitý pro predikci provozu a detekci anomálií z důvodu dostupné výpočetní složitosti [42].

### 3.3.1 Zlepšení modelu ARIMA

Pro zlepšení výstupu za použití modelu ARIMA je vhodné například použití klouzavého okna pro určení historických dat. Model ARIMA používá k modelování a předvídání budoucích dat vzorkovaná historická data. Aby bylo zajištěno rychlé a přesné modelování, je navržena pevná velikost posuvného okna. Velikost okna by měla být co nejmenší v rámci předpokladu relativně vysoké přesnosti modelování. To nejenže

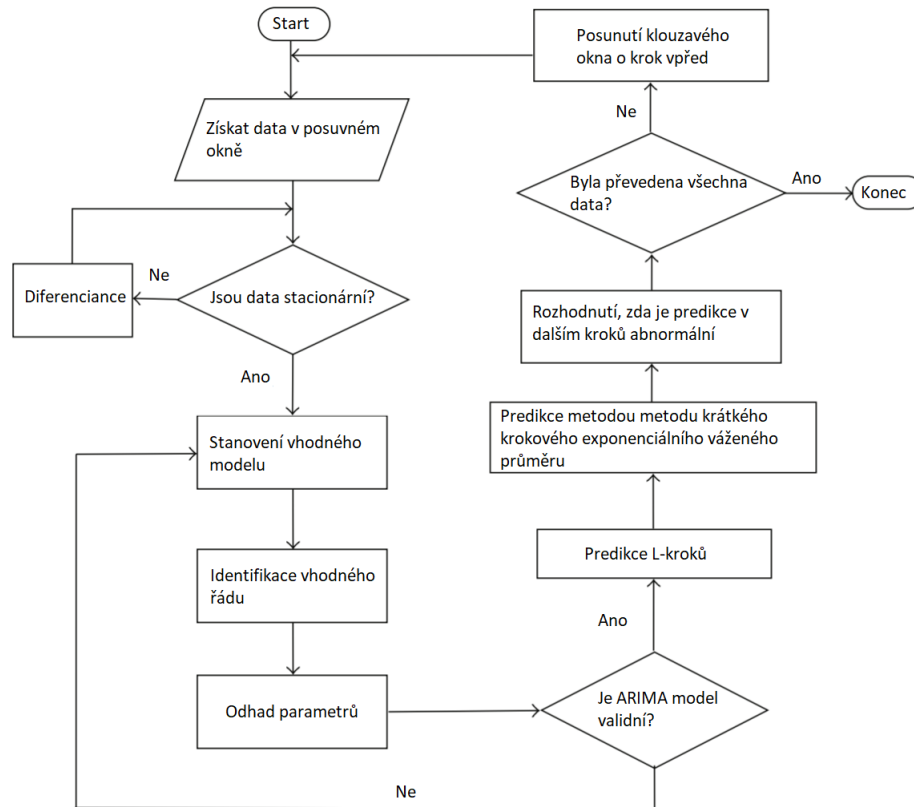
zajistí rychlé modelování, ale také zaručí použití nejnovějších historických dat, takže můžeme získat přesnější predikovanou hodnotu. Dále je nutné po každém posunutí okna aktualizovat model. Jelikož jsou sensorová data obecně nestacionární, používá se v každém okně zavedený model ARIMA. I přes to je nutné znovu definovat řád a odhady parametrů v závislosti na datech v novém okně a poté vypočítat predikovanou hodnotu. Včasně aktualizovaný nový model tak zajišťuje přesnost predikované hodnoty [42].

Jako další a klíčový krok je predikce sensorových dat metodou "Short Step Exponential Weighted Average", tedy metodu krátkého krokového exponenciálního váženého průměru. Tato metoda do jisté míry snižuje přesnost predikce, ale zvyšuje schopnost hodnocení abnormálního chování dat. Postup je obvykle rozdělen do dvou kroků: Predikce krátkého kroku a predikce hodnoty dat pomocí váženého průměru. Data mohou být korelována v různých časech. Jestliže je interval mezi korelacemi kratší je závislost těchto dat větší a naopak, pokud je interval mezi nimi delší, jejich korelace je menší. Jestliže pro model ARIMA platí  $L < L_0$  je jeho predikovaná hodnota platná a čím kratší je krok predikce, tím přesnější je predikovaná hodnota.

Na základě výše uvedené analýzy tedy platí, je-li  $L = 1$  je predikovaná hodnota nejpřesnější a nejúčinnější, ale není vhodná k detekci abnormálního vývoje dat. Vyskytne-li se v novém časovém okně anomálie, model se rychle přizpůsobí a tato predikovaná hodnota bude označena jako nová *přesná* abnormální hodnota. Dojde-li skutečně k abnormálnímu vývoji dat, může být tato hodnota označena jako normální. Pro návrh algoritmu detekce anomálií je proto někdy nutné dosáhnout *ne-přesné* predikční hodnoty, aby bylo možné provést přesnou detekci. Z tohoto důvodu je zavedena určitá *setrvačnost* k predikované hodnotě. Jestliže nastane anomální vývoj dat, parametry modelu nemohou být tak snadno přizpůsobeny a je tak snadnější získat predikovanou hodnotu, kterou lze snadno použít k detekci anomálií. Teoretickým základem k tomuto způsobu je to, že vývoj snímaných dat je v určitém ustáleném stavu a má jistou provozní *setrvačnost*, takže jakákoli náhlá změna je způsobena chybnou funkcí snímacích sensorů, nebo poruchou systému. Pro snadnější detekci anomálií je tedy do normálního provozu zavedena již dříve zmíněná setrvačnost a to metodou krátkého krokového exponenciálního váženého průměru, která je znázorněna na obrázku 3.11 a popsána v následujícím textu.

1. Nejprve jsou k modelování použita data v klouzavém okně, je provedeno  $L$ -kroků predikce a odpovídající výstupní hodnoty jsou uloženy na odpovídající pozici v daném časovém okně.
2. Dále je pak ve stejném časovém okně na  $L$ -hodnotách vytvořen exponenciálně vážený průměr. Výstupní hodnota je reprezentována jako

$$D = \sum_{i=1}^L \lambda_i D_i, \quad (3.51)$$



Obr. 3.11: Vylepšený model detekce anomálií, založený na modelu ARIMA. Převzato z [42].

kde  $\lambda_i$  je vážená hodnota  $D_i$ .

3. Posledním krokem je posouzení zda je vývoj v následujícím momentu abnormální či nikoliv. K tomu je vypočítán rozdíl mezi predikovanou hodnotou a reálnou hodnotou dat. Pro detekci anomálií jsou nastaveny různé prahy v závislosti na testovaném senzoru. Jestliže relativní chyba překročí prahovou hodnotu, předpokládáme, že se v daném čase vyskytuje anomálie [42].

## 4 Návrh algoritmu

Algoritmus byl navržen v prostředí Matlab na základě vybraných metod, tedy metody SVM, K-means a ARIMA. Algoritmus je koncipován tak, aby byl snáze aplikovatelný i pro budoucí analýzu dat ve firmě Honeywell.

### 4.1 Popis dat

Vstupní data pochází z databáze firmy Honeywell a jedná se o záznamy dvanácti různých senzorů (viz tabulka 4.1) pro 1641 motorů (pro TKO fázi a 1752 pro CRU fázi) typu HTF7000, které jsou v rámci analytického týmu ve firmě Honeywell zpracovávány a mají pro vybrané období dostatek dostupných záznamů. V případě uložené datové struktury, kde jsou data dle legislativních podmínek firmy částečně anonymizována, se jedná o data z období od 1. 1. 2016 do 30. 4. 2019. Pro případné další využití algoritmu, je možné nastavit časové období v podstatě libovolné, avšak nejstarší data uložená v databázi pochází z roku 2000.

Tab. 4.1: Popis testovaných senzorů.

Senzor	Popis	Jednotka
poil	tlak oleje	psia
oilt	teplota oleje	°C
N1	otáčky nízkotlaké turbíny	rpm
N2	otáčky vysokotlaké turbíny	rpm
ITT	teplota vzduchu vstupujícího do turbíny	°C
TT2	teplota vzduchu na vstupu do motoru	°C
PS	statický vstupní tlak měřený	psia
P3	tlak vzduchu vstupující do vysokotlakého kompresoru	psia
FFM	průtok paliva	PPH
EGT	teplota výfukových plynů	°C
FUEL_USAGE	spotřeba paliva	PPH
WFT	teplota paliva	°C

Sběr dat probíhá během celého letu, ale významné záznamy jsou především z TKO fáze a CRU fáze, proto jsou v praktické části mé diplomové práce zkoumána data právě z těchto uvedených letových fází. Z pohledu statistických funkcí jsou data popsána v tabulce 4.2 pro TKO a v tabulce 4.3 pro CRU.

Z předchozích analýz zkoumaných dat a závislostí, jsou známá určitá fakta. Například pro některá sensorová data jsou definovány vzory, říkající uživateli, která porucha je predikována a je tak možné včas zabránit rozsáhlejší poruše či kolapsu celého systému. Příklady zmíněných vzorů chování jsou uvedeny na obrázku 4.1. Na základě těchto vzorů jsou stanoveny možné poruchy systému jako je uvedeno na obrázku 4.2. Kde je na prvním řádku možné vidět vzor pro poruchu vysokotlaké turbíny

Tab. 4.2: Statistické hodnoty použitých senzorů pro TKO data.

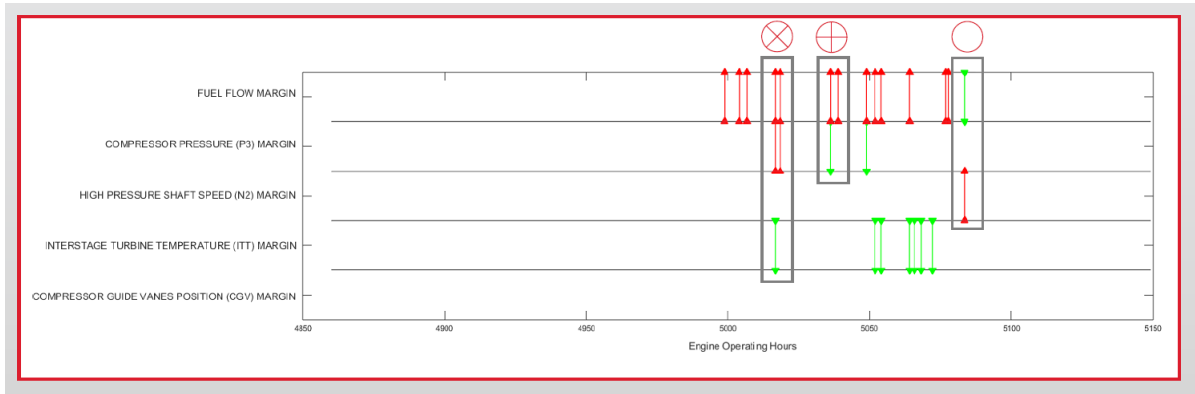
Statistická funkce	Aritmetický průměr	Medián	Modus	Směrodatná odchylka	Min	Max
poil	120.51	120.55	137.19	6.14	-50.00	137.25
oilt	50.86	51.35	55.18	7.84	13.98	84.11
N1	8910.29	8916.00	8927.00	214.16	742.00	9633.00
N2	26192.95	26265.00	26509.00	494.85	4934.00	27286.00
ITT	836.73	843.37	858.96	34.65	416.46	932.09
TT2	19.44	20.93	28.44	10.31	-33.51	52.93
PS	14.29	14.53	14.72	0.75	1.68	15.15
P3	323.89	324.03	322.07	12.99	15.59	363.72
FFM	2881.74	2868.91	2838.94	155.97	-1000.00	3489.02
EGT	596.46	601.47	613.67	28.67	408.80	813.12
FUEL_USAGE	84.87	70.11	59.68	105.21	0.00	5811.95
WFT	44.45	44.96	46.86	8.16	-150.00	78.89

Tab. 4.3: Statistické hodnoty použitých senzorů pro CRU data.

Statistická funkce	Aritmetický průměr	Medián	Modus	Směrodatná odchylka	Min	Max
poil	99.61	100.23	101.00	6.97	-50.00	129.89
oilt	71.55	71.61	68.81	7.33	-150.00	101.81
N1	8238.	8381.00	8640.00	578.58	4144.00	9662.00
N2	24768	24861.00	24961.00	679.75	19879.00	27108.00
ITT	728.3	735.28	727.50	57.75	395.84	910.07
TT2	-17.6	-20.11	-26.00	13.74	-54.34	29.90
PS	3.85	3.47	2.14	1.42	1.98	6.75
P3	116.9	114.90	103.81	28.37	32.23	346.46
FFM	971.7	954.25	878.00	231.52	-1000.00	2164.16
EGT	511.5	516.22	0.00	47.00	0.00	642.44
FUEL_USAGE	1594.	1293.18	0.00	1033.88	0.00	7298.23
WFT	65.13	65.38	62.38	7.96	-150.00	96.00

(HPT Issue), dále pak vzor pro poruchu kompresoru (Compressor Issue), problémy s lopatkami na přívodu vzduchu (Inlet Guide Vane Issue) a poruchu toku vzduchu (Bleed Issue). Reálný problém je vyobrazen na obrázku 4.3, kdy došlo k protržení pláště kompresoru a v datech je tak možné vidět snížení otáček vysokotlaké turbíny, která není doprovázena změnou teploty. Tento vývoj značí právě poruchu na kompresoru, která se v tomto případě potvrdila. Uvedený obrázek pochází z předchozích analýz, nicméně může posloužit jako pravdivě pozitivní příklad výskytu anomálie v datech.

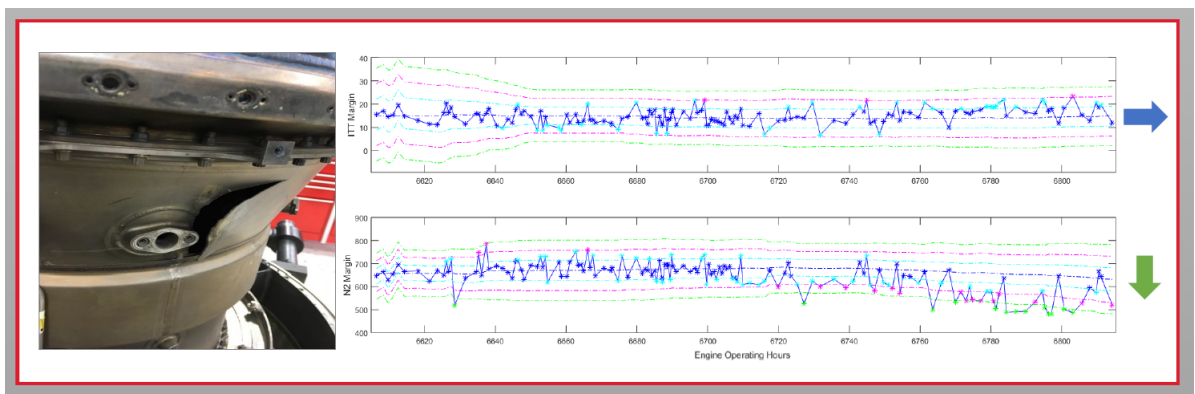
Jako příklady pravdivě pozitivních výskytu anomálií byl použitý soubor pěti motorů o celkovém počtu 8362 záznamů. Anomálie jsou v tomto datovém souboru detekovány u senzorů ITT, N2, FFM a P3. Závislosti mezi jednotlivými senzory jsou vykresleny na obrázku 4.4.



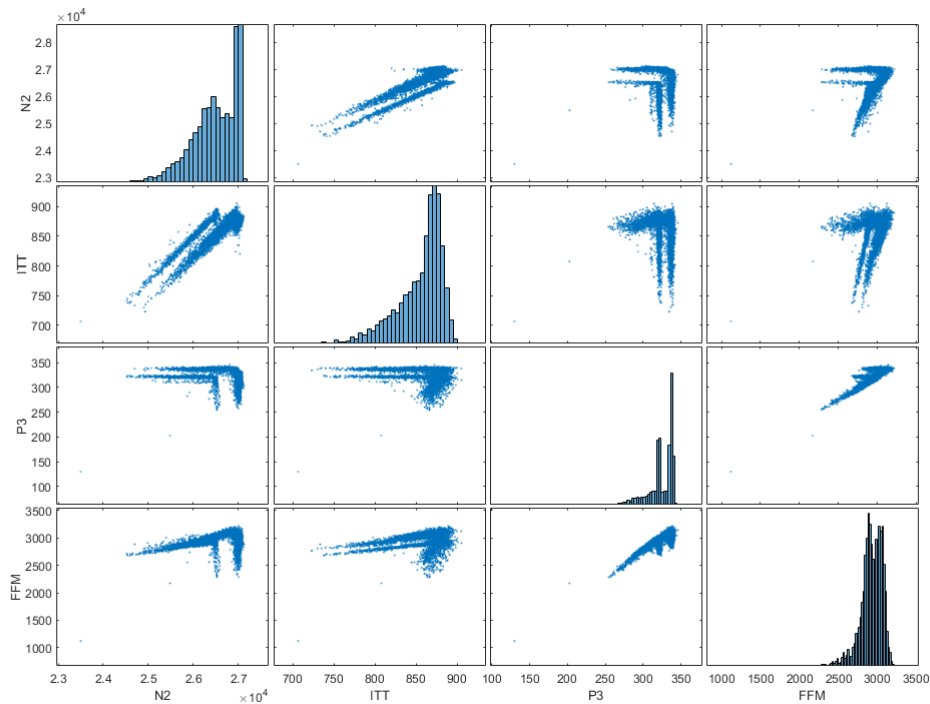
Obr. 4.1: Příklad vzorů chování pro vybrané senzory.

		FUEL FLOW MARGIN	N2 MARGIN	P3 MARGIN	ITT MARGIN
HPT Issue	⊗	UP ↑	UP ↑	DOWN ↓	DOWN ↓
Compressor Issue		UP ↑	FLAT →	FLAT →	DOWN ↓
Inlet Guide Vane Issue 1		UP ↑	DOWN ↓	FLAT →	FLAT →
Inlet Guide Vane Issue 2	○	DOWN ↓	UP ↑	FLAT →	FLAT →
Bleed Issue	⊕	UP ↑	FLAT →	DOWN ↓	DOWN ↓

Obr. 4.2: Detekované poruchy dle vybraných vzorů chování.



Obr. 4.3: Fotografie reálného motoru s dírou v plášti kompresoru s vykresleným záznamem dat.



Obr. 4.4: Závislosti mezi senzory u datového souboru s anomáliemi.

## 4.2 Obecný přístup k návrhu algoritmu pro monitorování stavu systému

Jak již bylo dříve uvedeno, prediktivní algoritmy a monitorovací algoritmy umožňují uživatelům a výrobcům zařízení posoudit provozní stav strojů, diagnostikovat poruchy nebo odhadnout kdy pravděpodobně dojde k selhání zařízení. Díky včasné diagnostice nebo predikci selhání, je možné předem naplánovat údržbu, lépe spravovat zásoby, zkrátit prostoje a tak zvýšit provozní efektivitu [28].

Vývoj těchto algoritmů vyžaduje dobře navrženou strategii pro včasné zhodnocení provozního stavu zařízení a včasné odhalení vznikajících chyb. To vyžaduje využití jak sensorových měření, tak znalost systému. Je nutné vzít v úvahu mnoho faktorů, včetně:

- Pozorované zdroje poruch a jejich relativní četnost.
- Dostupnost měření procesu pomocí senzorů. Počet, typ a umístění senzorů a jejich spolehlivost a redundanci ovlivňují jak vývoj algoritmů, tak náklady.
- Jak se různé zdroje chyb projevují na pozorovaných symptomech. Taková analýza příčin a následků může vyžadovat rozsáhlé zpracování dat z dostupných



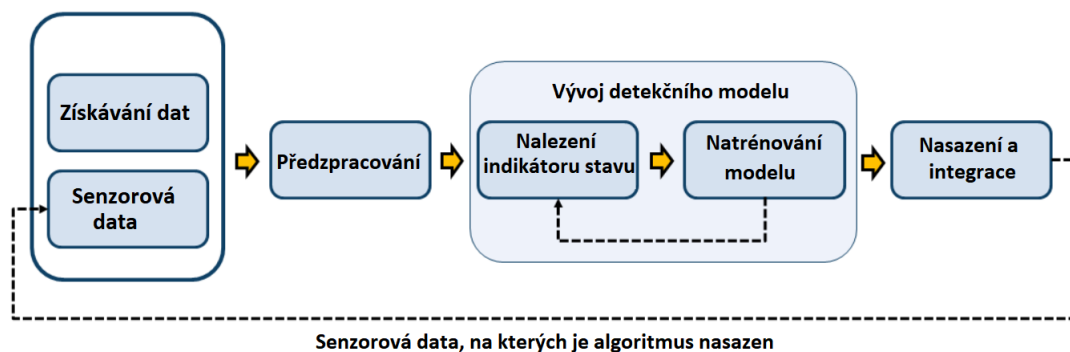
senzorů.

- Fyzikální znalosti o dynamice systému. Tyto znalosti mohou pocházet z matematického modelování systému a jeho chyb a z poznatků odborníků v oboru.

*Monitorování stavu* využívá sensorová data k vyhodnocení aktuálního stavu konkrétního zařízení a k detekci nebo diagnostice poruch. Vstupními daty mohou být například záznamy z měření teploty, tlaku, napětí, hluku nebo vibrací, shromážděná pomocí konkrétních senzorů. Algoritmus monitorování stavu odvozuje metriky z dat označovaných jako indikátory stavu. Indikátor stavu je jakákoli funkce systémových dat, jejichž chování se mění v předvídatelném stavu, v souvislosti s degradací systému. Algoritmus monitorování stavu může tedy provádět detekci chyb nebo diagnózu systému porovnáním nových dat s již existujícím označeným anomálním vývojem dat [28].

### 4.2.1 Pracovní postup pro vývoj algoritmů

Následující obrázek 4.5 ukazuje postup pro vývoj detekčního či predikčního algoritmu.



Obr. 4.5: Blokové schéma vývoje detekčního algoritmu. Převzato z [28].

#### Sběr dat

Vývoj samotného algoritmu začíná sběrem dat. Často se jedná o správu a zpracování velkého objemu dat. Obvykle má vývojář či uživatel přístup k jednomu nebo více z následujících typů dat.

- Reálná data z normálně pracujícího systému.
- Reálná data ze systému pracujícího v chybném stavu.
- Reálná data ze selhání systému.

Avšak reálná data ze selhání systému nejsou ve většině případů k dispozici a to z důvodu pravidelné údržby a relativní vzácnosti takových incidentů. V tomto případě mohou být data generována pomocí modelů, které simulují provoz systému v různých poruchových stavech [28].

### **Předzpracování dat**

Předběžné zpracování je často nezbytné pro převod dat do podoby, ze které lze snadno odvodit stavové indikátory. Předzpracování dat zahrnuje jednoduché techniky, například odstraňování odlehlých a chybějících hodnot, a pokročilé techniky zpracování signálů, jako je krátkodobá Fourierova transformace. K určení vhodné techniky pro předzpracování dat je nutné mít dostatečnou znalost systému a charakteristiky dat [28].

### **Identifikace indikátorů stavu**

Klíčovým krokem ve vývoji detekčních algoritmů je identifikace indikátorů stavu, tzn. funkcí systému, jejichž chování je možné předvídat v závislosti na degradaci systému. Indikátor stavu může být jakákoliv funkce, která je schopna rozlišit normální stav od poruchového, nebo je užitečná pro predikci zbývající životnosti. Mezi techniky, které se běžně používají pro extrahování stavových indikátorů, patří:

- řádová analýza,
- modální analýza,
- spektrální analýza,
- spektrální analýza obálky,
- analýza únavy materiálu<sup>1</sup>,
- nelineární analýza časových řad,
- analýza založená na modelu.

### **Trénování modelu**

Hlavní komponentou algoritmu je model detekce nebo predikce. Tento model analyzuje extrahované stavové indikátory k určení aktuálního stavu systému (detekce a diagnostika chyb) nebo předpovídá jeho budoucí stav (predikce životnosti systému). Detekce a diagnostika poruch závisí na použití jedné nebo více hodnot indikátoru stavu pro rozlišení mezi zdravým a chybovým provozem a mezi různými typy poruch. Jednoduchý model detekce poruch je prahová hodnota indikátoru stavu, která indikuje stav poruchy při jejím překročení. Jiný model může porovnávat aktuální hodnoty s jejich statistickým rozdělením. Je tak možné určit pravděpodobnost

---

<sup>1</sup>Fatigue analysis, neboli analýza únavy je prováděna buď deterministickým přístupem pouhým použitím Minerova pravidla nebo spektrálního přístupu (ve frekvenční oblasti).

konkrétního poruchového stavu. Komplexnějším přístupem k diagnostice je použití klasifikátoru, který porovnává aktuální hodnotu jednoho nebo více indikátorů s hodnotami spojenými s chybovými stavy a je tak možné určit pravděpodobnost výskytu konkrétního chybového stavu [28].

### Nasazení a integrace

Nasazení nebo integrace detekčního algoritmu je obvykle poslední fází pracovního postupu vývoje algoritmu. Již v dřívějších fázích je nutné brát v úvahu, jakým způsobem bude algoritmus nakonec nasazen. Vliv na požadavky a další aspekty návrhu systému má například to, zda algoritmus běží na vestavěném hardwaru, jako samostatný spustitelný soubor nebo jako webová aplikace [28].

## 4.3 Vlastní návrh algoritmu

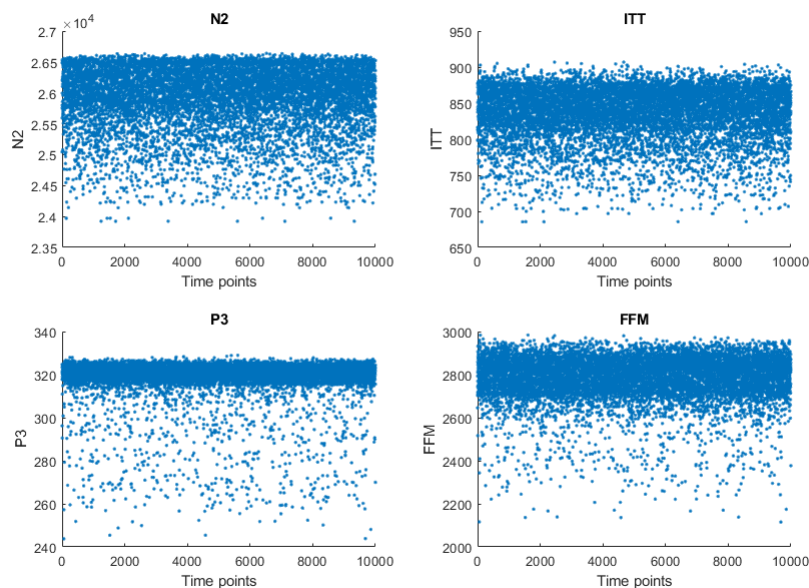
Algoritmus je vytvořen v programu MATLAB 2015a. Tato verze je zvolena záměrně proto, aby bylo případně možné algoritmus využívat i nadále ve firmě Honeywell. Ze stejného důvodu je koncipován způsobem uvedeným v následujícím textu.

Hlavním spouštěcím skriptem je *DP\_Gregorova\_Runner.m*, ve kterém je zapotřebí definovat všechny vstupní parametry a to včetně zkoumané fáze letu, sledovaného časového období, použité metody, sledovaných senzorů a dalších. V tomto skriptu je rovněž implementována funkce *getSQLdata.m* která stahuje data z databáze firmy Honeywell na základě definovaných vstupních parametrů jako je sériové číslo motoru, požadované časové období a požadované výstupní proměnné. Na výstupu této funkce dochází k anonymizaci dat, na základě legislativních podmínek firmy a dále jsou sériová čísla jednotlivých motorů označena jako ESN\_1 až ESN\_1641 pro vzletovou fázi a ESN\_1 až ESN\_1752 pro CRU fázi. Pro testování a spouštění skriptu bez přístupu do databáze je nutné nastavit parametr *connection\_type = extract*. V tomto případě budou data načítána z uložených struktur pro TKO a CRU.

Součástí spouštěcího skriptu je rovněž funkce *statDescription.m*, která poskytuje statistický popis všech dostupných dat a v rámci této funkce jsou dostupné statistické výstupy ukládány do tabulek ve formátu excel. Spouštění této funkce je možné nastavením parametru *statistics = 0* vypnout. Další funkce, která je součástí tohoto skriptu je funkce *dataPreProcessing*. zde dochází k předzpracování dat pro další analýzy. Detailnější popis předzpracování dat je uveden v následující kapitole. Ze spouštěcího skriptu vstupuje algoritmus do funkce *DP\_Gregorova\_Deployment.m*. Jak již bylo zmíněno, koncepce celého algoritmu je vytvářena na základě zvyklostí u ostatních analytických algoritmů používaných k analýze dat ve firmě Honeywell.

### 4.3.1 Předzpracování vstupních dat

Jak již bylo zmíněno, předzpracování dat je důležitou součástí celého algoritmu. Data jsou zde rozdělena na trénovací a testovací množinu a výstupem jsou dvě tabulky, kde každý sloupec jak trénovacích tak testovacích dat představuje záznamy pro určitý senzor. Vstupním parametrem je kromě samotné datové struktury taktéž datové pole s definovanými testovanými sensory a velikost testovací a trénovací matice. Součástí *dataPreProcessing* funkce je grafický výstup sensorových dat v čase, viz obrázek 4.6.



Obr. 4.6: Data pro vybrané senzory v závislosti na čase

### 4.3.2 Hlavní funkce algoritmu

Koncept hlavní funkce *DP\_Gregorova\_Deployment.m* je přizpůsoben tak aby bylo možné použít algoritmus třetí stranou. Jelikož všechny vyvíjené algoritmy jsou poskytovány zákazníkům, jako jsou společnosti zajišťující údržbu letadel koncových zákazníků.

Součástí této funkce je například také funkce *getConfig.m*, která načítá data z příloženého csv souboru. Parametry které jsou načítány z externího souboru jsou konfigurovatelné a tak snáze přizpůsobitelné aktuálním požadavkům zákazníka a to bez zásahu do samotného algoritmu. Struktura s těmito parametry je načítána téměř do všech následujících funkcí, které s těmito parametry pracují. Vstupy do všech aplikovaných funkcí, jsou nastavovány právě v *DP\_Gregorova\_Deployment.m* a tato funkce tak trojí v podstatě kostru celého algoritmu.

### 4.3.3 SVM

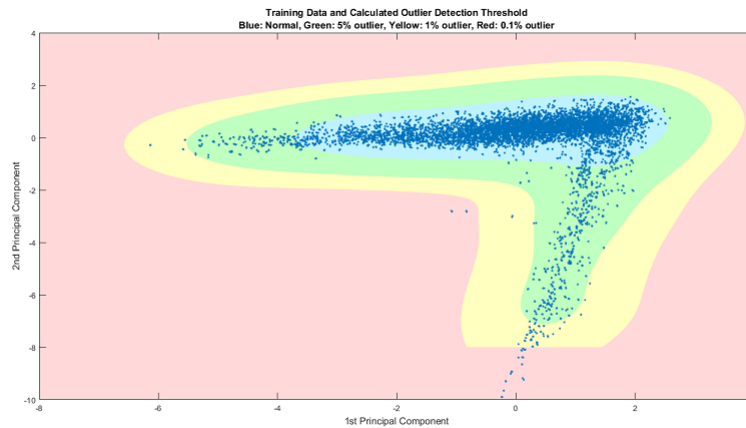
Do funkce *SVM.m* se algoritmus dostává jak již bylo uvedeno přes *DP\_Gregorova\_Deployment.m*. Stejně tak do funkce *PCA.m*, která bezprostředně předchází samotnému SVM algoritmu, jelikož je to pro něj nezbytný proces. V PCA funkci dojde k normalizaci dat pomocí metody *z-score* a následně je na takto upravená data aplikována původní funkce *pca*, která je součástí Matlabu. SVM pak pracuje z výstupy z funkce *pca* jako jsou například koeficienty hlavních komponent, které byly vráceny jako  $p \times p$  matice, kde každý sloupec obsahuje koeficienty pro jednu hlavní komponentu, nebo matice skóre.

Následně dojde k natrénování SVM modelu, a jeho aplikaci na testovací data. Pro trénování je třeba nastavit typ jádrové funkce. Pro OneClass SVM je vhodné použít funkci RBF, kterou jsem zvolila. Dále je nutné definovat vektor, přiřazující všem vstupním datům hodnotu 1, to znamená, že jsou všechna vstupní trénovací data označena jako normální a také jsem definovala hranice oddělující 5%, 1% a 0.1% anomálií. Výstup pro trénovací sadu dat je znázorněn na obrázku 4.7. Pro vykreslování barevně oddělených hranic je zde použita funkce *plotContourmatrix.m*. Po natrénování modelu je možné jej aplikovat na testovací množinu dat a následně na celý soubor motorů, kdy se za pomoci této metody snažím detekovat anomálie pro různé kombinace senzorů u jednotlivých motorů. K tomuto účelu slouží u metody SVM funkce *fleetRunTestingSVM.m*. Testování probíhá tak, že jsou data pro jednotlivé motory rozdělena do dvou skupin a to na data aktuální a historická. Toto dělení je typické pro používané algoritmy ve firmě Honeywell, proto jsem se ho takto snažila ponechat. Dále jsou aktuální data za pomoci výstupu z funkce *PCA.m* normalizována a je na ně aplikován definovaný SVM model. Tímto způsobem je možné najít anomálie jak vizuálně tak pak pomocí metriky, která určuje, které body ze vstupní množiny dat, spadají za hranici normálního chování.

Pro jednodušší posouzení a definování v jakém časovém okně se vyskytla anomálie jsou testovaná data vykreslena v závislosti na čase včetně hodnoty klouzavého průměru a hranic, které jsou definovány jako součet průměru a směrodatné odchylky, nebo jejího násobku. Tyto proměnné jsou výstupem z funkce *calcAverageSigma.m*. Anomálie jsou zvýrazněny viz obrázek.

### 4.3.4 K-means

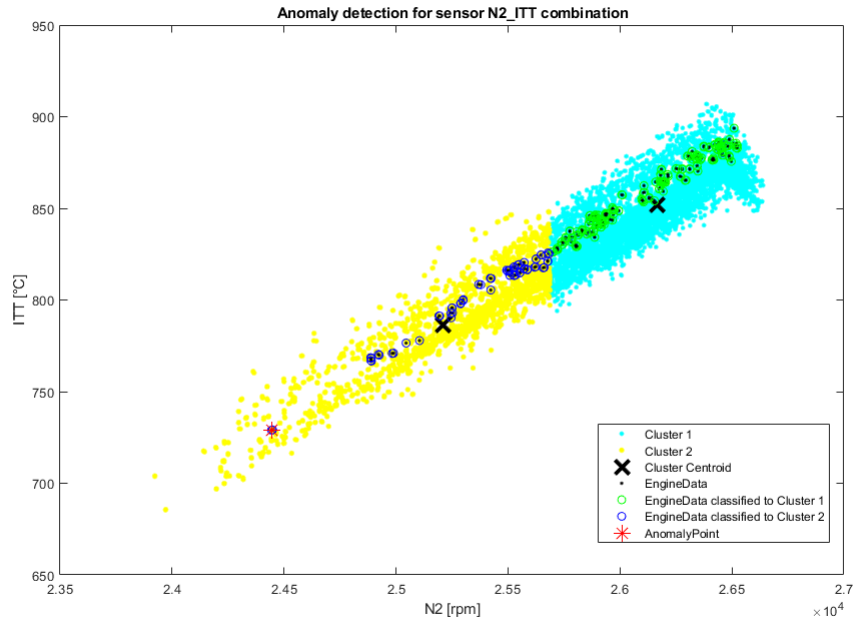
Algoritmus pro detekování anomálií metodou Kmeans je naprogramován ve stejnojmenné funkci *k\_mean.m*. Na vstupu této funkce je datová struktura, trénovací data, pole s testovanými senzory a konfigurační struktura. Metodou Kmeans jsou procházeny všechny možné kombinace testovaných senzorů a jsou tak hledány anomální hodnoty pro jednotlivé senzory či skupinu senzorů.



Obr. 4.7: Metoda SVM - příklad modelu pro trénovací data

V prvním kroku je vytvořena matice aktuálně testované kombinace senzorů  $all-SensorData(x,y)$ , kde hodnota  $x$  je délka sensorových dat pro daný motor a  $y = 2$ , jelikož je vždy testována pouze kombinace dvou senzorů a to z důvodu zobrazení a následné vizuální kontroly. Následně jsou pak vypočteny pozice centroidů definující dva shluky. Jelikož je nutné předem určit počet shluků, je tato hodnota definována opět v konfiguračním souboru. V tomto případě byly zvoleny shluky dva, jelikož se jedná o závislost mezi dvěma různými senzory, tudíž se dá předpokládat, že data nebudou tvořit pouze jeden kompaktní shluk. Konfigurovatelná je také metoda pro měření vzdálenosti. K tomuto účelu jsem zvolila Euklidovskou vzdálenost. Dále jsou pro stejnou kombinaci senzorů, avšak pro data pocházející pouze z jednoho motoru vypočítány vzdálenosti od již definovaných centroidů a určena příslušnost k jednomu či druhému shluku. Anomálie jsou pak takové body, jejichž vzdálenost od jednoho či druhého shluku je větší než součet průměrné vzdálenosti a trojnásobku směrodatné odchylky viz obrázek 4.8. Takto označené anomální body jsou platné pouze pro aktuální kombinaci testovaných senzorů a jsou tak mezikrokem k detekci takových odlehlých hodnot pro které platí, že se vyskytují na stejné pozici ve všech testovaných kombinacích. Výstupem K-mean algoritmu je struktura s hodnotami pro každý testovaný motor a každý testovaný senzor.

Pro testování všech motorů je použita funkce *fleetRunTestingKmean.m*, které je obdobná jako dříve uvedená funkce pro testování metody SVM. Výstupem je opět struktura, ve které je uvedeno zda daný motor a daný testovaný senzor má nějaká anomální data či nikoliv a také je možné tento výstup graficky znázornit.



Obr. 4.8: Příklad detekce anomálií pomocí metody K-means pro kombinaci N2 - ITT

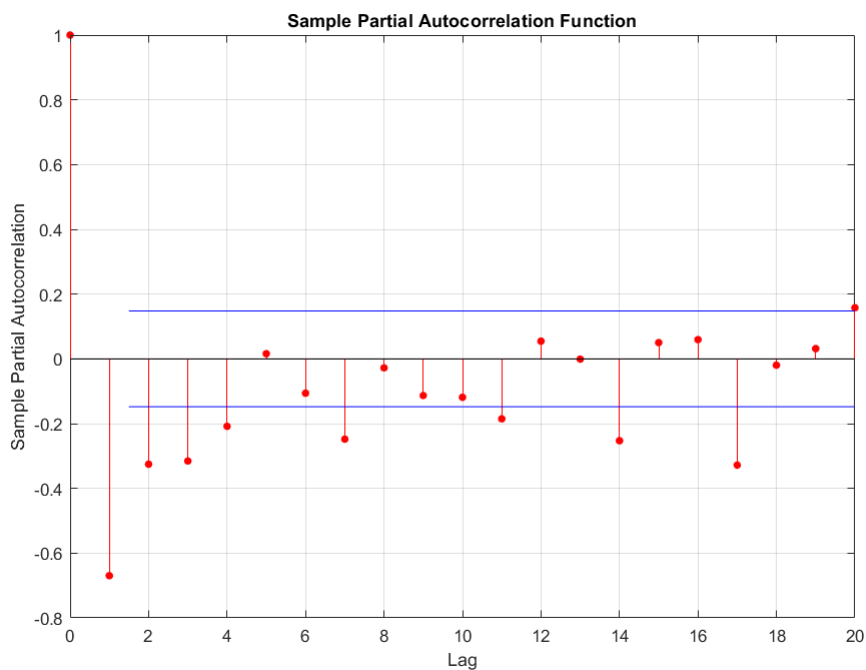
### 4.3.5 ARIMA

Metoda ARIMA je od začátku postavena na jiném konceptu, než předchozí dvě uvedené metody. Na data je v tomto případě nahlíženo jako na data časových řad, proto je celý algoritmus tomuto přizpůsoben.

V první řadě jsou vybrána data pro testovaný motor a senzor a jsou převedena do formátu *timetable* kdy první sloupec udává časovou posloupnost a ve druhém sloupci jsou příslušná data pro testovaný senzor. Jelikož chceme testovat pouze data označena jako aktuální, jsou proto tato data vybrána v daném časovém intervalu a s těmito daty se následně pracuje jako s novým datovým polem. V dalším kroku je nutné identifikovat, zda jsou data stacionární či nikoliv a to pomocí Dickey-Fullerova testu, který je v matlabu implementován jako funkce *adftest*. Jestliže je výstupní hodnota jedna, test je validní a data jsou stacionární. Následně se pomocí autokorelace a parciální korelace (obrázek 4.9) určí zda je model vhodný a definuje se významné zpoždění MA. Na základě výše uvedené analýzy postačuje pro transformaci modelu na stacionární diferenciaci 1.řádu.

V dalším kroku je pomocí funkce *arima* definován model a určeny predikované hodnoty. Následně je ověřena vhodnost modelu, případně jsou definovány nové parametry a znovu vypočtena predikovaná hodnota. Výstupem této detekce je hodnota rezidua. Jestliže je rozdíl mezi predikovanou hodnotou a hodnotou rezidua v abso-

lutní hodnotě větší než trojnásobek směrodatné odchylky vypočtené z originálních dat, je tato hodnota označena jako anomálie.



Obr. 4.9: Příklad parciální autokorelační funkce - metoda ARIMA



## 5 Vyhodnocení detekce

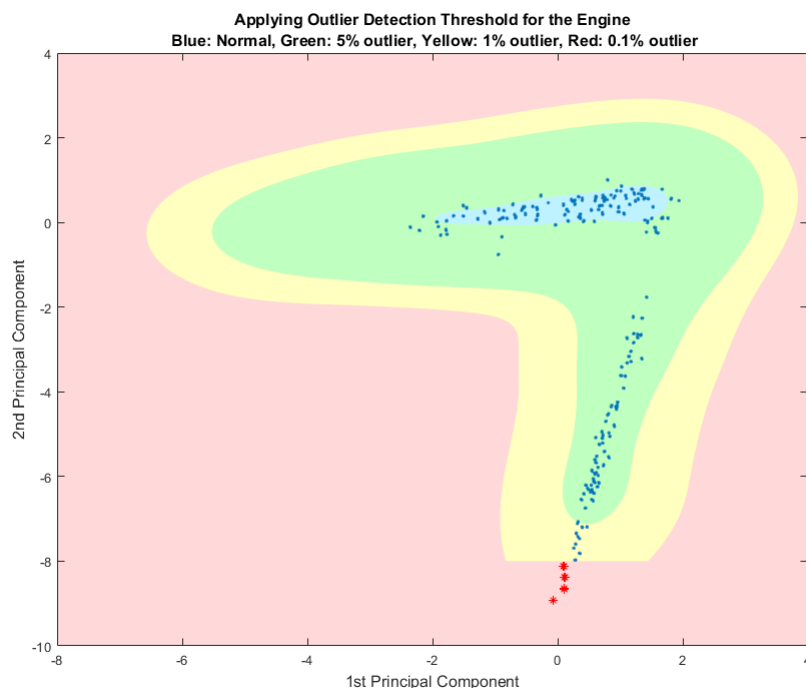
Algoritmy navržené v této diplomové práci se zprvu jevily jako velmi vhodné pro detekci anomálií v sensorových datech, nicméně ani jedna metoda nespĺňuje přísná kritéria pro detekci anomálií pro letecký průmysl. Zřejmě by bylo vhodné tyto metody více propracovat a lépe definovat vstupní parametry, což může být v budoucnu předmětem práce analytického týmu firmy Honeywell.

### 5.1 Vyhodnocení SVM

V testovací fázi jsem zvolila čtveřici senzorů, pro které známe závislosti, jak bylo uvedeno dříve. Jelikož tato metoda pracuje právě s kombinacemi senzorů, je možné podle znalostí systému hledat anomálie v datech na základě jejich kombinací. Pro zvolené senzory (N2, ITT, P3 a FFM) z celkového počtu 1641 testovaných motorů v TKO fázi jsem metodou SVM detekovala anomálie u z 974 nich, což je 59,35 % motorů. Tato hodnota je významná pro detekci bodových anomálií a podle vizuální kontroly viz obrázek 5.1 je do jisté míry vyhovující. Avšak pro detekci kontextuálních anomálií, které jsou pro sensorová data signifikantní je tato metoda ve fázi vývoje v jaké se nachází nevyhovující. Bylo by však vhodné pokračovat ve vývoji i nadále, protože jak se zdá, má tato metoda potenciál. Po vhodném zvolení kombinace senzorů je možné odhalit poruchy, které závisí na vícero vstupních parametrech, ne pouze na datech z jednoho senzoru.

### 5.2 Vyhodnocení K-MEANS

Metoda K-means má podobný výstup jako předchozí metoda SVM. Avšak tímto způsobem jsem testovala všechny dostupné senzory a to tak že výstupem je anomálie, vyskytující se v kombinaci všech senzorů. Stejně tak jako u SVM se vždy jedná pouze o bodové anomálie a pro lepší funkčnost této metody je také potřeba dále pracovat na vývoji. Metodou K-means jsem detekovala u každého motoru minimálně jeden anomální bod alespoň pro jeden senzor. Výstup je uložený ve výstupním souboru *K-MEANS results.csv*. Celkem bylo detekováno pro každý senzor relativně velké množství anomálií. V tabulce 5.1 jsou shrnuté počty výskytu anomálií pro jednotlivé senzory. Průměrně tak byla detekována anomálie u 61 % motorů z celkového počtu 1641 pro každý senzor.



Obr. 5.1: Příklad detekce anomálií pro jeden motor metodou SVM.

Tab. 5.1: Počet motorů, u kterých byla detekována anomálie metodou K-means pro daný senzor.

Statistická funkce	Počet motorů, pro které byly detekovány anomálie	Procentuální výskyt anomálie [%]
poil	1102	67,15
oilt	1428	87,02
N1	1228	74,83
N2	498	30,35
ITT	1158	70,57
TT2	1294	78,85
PS	982	59,84
P3	954	58,14
FFM	826	50,34
EGT	1117	68,07
FUEL_USAGE	0	0
WFT	1456	88,73

### 5.3 Vyhodnocení ARIMA

Metoda ARIMA má s počtem 856 detekovaných motorů z celkového počtu 1641 motorů nejlepší výsledky, ze všech testovaných metod, ale pro běžnou praxi je to stále nedostačující. Stele je míra detekovaných motorů více než 50%, konkrétně 52,16 %. A

to ze souboru, kde je pravděpodobnost výskytu anomálie za dané období cca. 10 %.

Tato metoda je nejvíce vyhovující metodou z pohledu kontextuálních anomálií, které se vyskytují v sensorových datech nejčastěji, jelikož pracuje přímo s daty časových řad. Nicméně je nutné na vývoji této detekční metody pro použití v praxi dále pracovat. Existuje nesčetně možností jak metodu přizpůsobit datům tak, aby vyhovovala potřebám zákazníků. Je zřejmě ze všech uvedených metod nejvíce variabilní, proto si trůfám říci, že je nejvhodnější metodou pro další vývoj.

## 5.4 Shrnutí

Obecně vzato všechny metody mají více než 50 % pozitivních detekcí. Je to zřejmě proto, že definovaná hranice pro indikaci anomálie je nastavena jednoduše jako součet průměrné hodnoty vstupních dat a trojnásobku směrodatné odchylky, případně u metody ARIMA musí být reziduum větší než je trojnásobek směrodatné odchylky. Nicméně je tato diplomová práce dobrým základem k dalším analýzám. V běžné praxi analytika trvá vývoj jediného algoritmu někdy i více než rok. V souvislosti s tímto faktem lze očekávat, že uvedené algoritmy je potřeba více promyslet a propracovat a snad někdy v budoucnu budou využitelné k pečlivým analýzám, díky kterým bude možné předcházet katastrofám, které nejsou v leteckém průmyslu výjimkou.

## 6 Závěr

V této práci jsou shrnuty obecné poznatky z oblasti detekce anomálií. Konkrétně je zaměřena na detekci anomálií v datech, které jsou snímány různými typy senzorů z leteckých motorů typu HTF7000. Všechny tyto senzory jsou v úvodu práce popsány, včetně jejich fyzikálních principů. Stejně tak je popsán motor, který v tomto případě složí jako systém poskytující data. Tento obecný popis technických parametrů je následován popisem typů anomálií a detailním popisem metod vybraných pro praktickou část.

Cílem zmíněné praktické části, bylo navržení algoritmu pro detekci anomálií. Byly zvoleny tři detekční metody a to metoda SVM, K-means a ARIMA. Všechny tyto metody bohužel nedosahují tak kvalitních výsledků, aby je bylo možné okamžitě používat v praxi. U všech metod byly detekovány anomálie u více jak 50% motorů, což lze vyhodnotit jako nevyhovující množství, jelikož z praxe víme, že obvyklá a přípustná míra je cca. 10% motorů na období, které jsem zvolila. To bylo pro data označená jako aktuální od začátku roku 2017 do konce roku 2018.

Je však nutné zamyslet se důkladněji nad reálnými anomáliemi vyskytujícími se v datech tohoto typu, tedy datech pocházejících z leteckého motoru. Z pohledu použitých metod nemusí být výstup natolik špatný, hledáme-li opravdu pouze anomální bod. Letecký motor je však natolik komplexní systém, že je zřejmě nutné nahlížet na něj více komplexněji. Je tedy možné že tyto metody se v budoucnu stanou základním kamenem pro návrh komplexního algoritmu detekujícího chyby tak dobře, že bude možné zabránit katastrofám, které jak již bylo zmíněno, nejsou v leteckém průmyslu výjimkou. K tomuto cíli však vede ještě dlouhá cesta.

## Literatura

- [1] WINSTANLEY, D., *HTF7000 Engine Design, Development and Uses* [online]. SAE Int. J. Aerosp. 6(2):545-554, 2013,. [cit. 28. 11. 2018]. Dostupné z URL: <[https://www.davidwinstanley.com/uploads/7/2/0/0/72002645/development\\_of\\_the\\_htf7000\\_gas\\_turbine\\_engine\\_for\\_business\\_aviation\\_aircraft.pdf](https://www.davidwinstanley.com/uploads/7/2/0/0/72002645/development_of_the_htf7000_gas_turbine_engine_for_business_aviation_aircraft.pdf)>.
- [2] KUMAR, V.: *Parallel and distributed computing for cybersecurity* [online]. IEEE Distributed Systems Online, 6(10). [cit. 28. 11. 2018]. Dostupné z URL: <<https://doi.org/10.1109/MDSO.2005.53>>
- [3] SPENCE, C., PARRA, L., SAJDA, P.: *Detection, Synthesis and Compression in Mammographic Image Analysis with a Hierarchical Image Probability Model. Proceedings of the Workshop on Mathematical Methods in Biomedical Image Analysis(2001)*, 10.1109/MMBIA.2001.991693. [cit. 28. 11. 2018].
- [4] ALESKEROV, E., FREISLEBEN, B., AND RAO, B.: *A neural network based database mining system for credit card fraud detection*. In Proceedings of the IEEE Conference on Computational Intelligence for Financial Engineering. 1997. Cardwatch. 220–226. [cit. 5. 12. 2018].
- [5] CHANDOLA, V., BANERJEE, A., KUMAR, *Anomaly detection: A survey*. V. 2009. ACM Comput. Surv. 41, 3, Article 15 (July 2009), 58 pages. [cit. 6. 12. 2018]. Dostupné z URL: <<http://doi.acm.org/10.1145/1541880.1541882>>
- [6] FEDERAL AVIATION ADMINISTRATION (FAA) *FAA-H-8083-3B Airplane Flying Handbook Handbook (PDF)*. Federal Aviation Administration. [cit. 12.12.2018]. Dostupné z URL: <[http://www.faa.gov/regulations\\_policies/handbooks\\_manuals/aviation/pilot\\_handbook/media/PHAK-Chapter06.pdf](http://www.faa.gov/regulations_policies/handbooks_manuals/aviation/pilot_handbook/media/PHAK-Chapter06.pdf)>
- [7] AMETEK AEROSPACE [online katalogový list]. *Speed Sensors* [cit. 12.12.2018]. Dostupné z URL: <<https://www.ameteksfms.com/-/media/ameteksensors/documents/sensor%20data%20sheets/speed-sensors.pdf?la=en>>
- [8] VOLPONI A.J. *Gas Turbine Engine Health Management: Past, Present, and Future Trends*. ASME. J. Eng. Gas Turbines Power. 2014;136(5):051201-051201-20. doi:10.1115/1.4026126. [cit. 17.12.2018]. Dostupné z URL: <<http://gasturbinespower.asmedigitalcollection.asme.org/article.aspx?articleid=1787593>>

- [9] Thermocouples: Types, What It Is & How It Works | Omega. OMEGA Engineering | Thermocouples, Pressure Transducers, Flow Meters, PID Controllers [online]. Copyright © [cit. 14.12.2018]. Dostupné z URL: <<https://www.omega.com/prodinfo/thermocouples.html>>
- [10] Aircraft Systems: Aircraft Turbine Engine Fuel System Requirements. Aircraft Systems [online]. [cit. 17.12.2018]. Dostupné z URL: <<https://okigihan.blogspot.com/p/turbine-engine-fuel-systemgeneral.html>>
- [11] ORLÍKOVÁ, Soňa. *Měření průtoku tekutin - principy průtokoměrů* [online]. [cit. 17.12.2018]. Dostupné z URL: <[http://www.elektrorevue.cz/clanky/01049/index.html#\\_Tepel%C3%BD\\_hmotnostn%C3%AD\\_pr%C5%AFtokom%C4%9Br](http://www.elektrorevue.cz/clanky/01049/index.html#_Tepel%C3%BD_hmotnostn%C3%AD_pr%C5%AFtokom%C4%9Br)>
- [12] LVDT PRINCIPLES OF OPERATION. Technical Paper. Metrolog Measurement Control Systems [online]. Copyright © [cit. 18.12.2018]. Dostupné z URL: <[https://www.metrolog.net/files/lvdt\\_principles\\_en\\_metrolog.pdf](https://www.metrolog.net/files/lvdt_principles_en_metrolog.pdf)>
- [13] WEIGEND, A. S., MANGEAS, M., AND SRIVASTAVA, A. N. 1995. *Nonlinear gated experts for time-series: Discovering regimes and avoiding overfitting*. Int. J. Neural Syst. 6, 4, 373–399.[cit. 19.12.2018]. Dostupné z URL: <[https://scholar.colorado.edu/cgi/viewcontent.cgi?referer=https://search.yahoo.com/&httpsredir=1&article=1749&context=csci\\_techreports](https://scholar.colorado.edu/cgi/viewcontent.cgi?referer=https://search.yahoo.com/&httpsredir=1&article=1749&context=csci_techreports)>
- [14] SALVADOR S., CHAN P., BRODIE J. (2004). *Learning States and Rules for Time Series Anomaly Detection*. Proceedings of the Seventeenth International Florida Artificial Intelligence Research Society Conference, FLAIRS 2004. 1. [cit. 19.12.2018]. Dostupné z URL: <[https://www.researchgate.net/publication/221438301\\_Learning\\_States\\_and\\_Rules\\_for\\_Time\\_Series\\_Anomaly\\_Detection](https://www.researchgate.net/publication/221438301_Learning_States_and_Rules_for_Time_Series_Anomaly_Detection)>
- [15] BOBBIA M., MISITI M., MISITI Y., POGGI M., PORTIER B., *Spatial outlier detection in the PM10 monitoring network of Normandy (France)*, Atmospheric Pollution Research, Volume 6, Issue 3, 2015, Pages 476-483, ISSN 1309-1042. [cit. 19.12.2018]. Dostupné z URL: <<http://www.sciencedirect.com/science/article/pii/S1309104215302178>>
- [16] Wikipedie: Otevřená encyklopedie: Dvouproudový motor [online]. c2018 [cit. 19. 12. 2018]. Dostupné z URL: <[https://cs.wikipedia.org/w/index.php?title=Dvouproudov%C3%BD\\_motor&oldid=16772402](https://cs.wikipedia.org/w/index.php?title=Dvouproudov%C3%BD_motor&oldid=16772402)>

- [17] JOHAN FLORBÄCK. Anomaly Detection in Logged Sensor Data. Goteborg, 2015. Master's thesis. ISSN 1652-8557. Department of Applied Mechanics Division of Vehicle Engineering and Autonomous Systems Chalmers University of Technology. Dostupné z URL: <<http://publications.lib.chalmers.se/records/fulltext/223871/223871.pdf>>
- [18] Introduction to Anomaly Detection. DataScience.com | Enterprise Data Science Platform Provider [online]. Copyright © 2018, Oracle and [cit. 28.12.2018]. Dostupné z URL: <<https://www.datascience.com/blog/python-anomaly-detection>>
- [19] F. J. ANSCOMBE, *Rejection of outliers*, Technometrics, vol. 2, no. 2, pp. 123–146, 1960. [cit. 28.12.2018].
- [20] A. J. FOX, *Outliers in time series*, Journal of the Royal Statistical Society. Series B (Methodological), pp. 350–363, 1972. [cit. 28.12.2018] Dostupné z URL: <[https://www.jstor.org/stable/2985071?origin=JSTOR-pdf&seq=1#page\\_scan\\_tab\\_contents](https://www.jstor.org/stable/2985071?origin=JSTOR-pdf&seq=1#page_scan_tab_contents)>
- [21] HU, Xiao, EKLUND, Neil, GOEBEL, Kai, CHEETHAM, William. (2007). *Hybrid Change Detection for Aircraft Engine Fault Diagnostics*. Proc. IEEE Aerosp. Conf. 1 - 10. 10.1109/AERO.2007.352848. [cit. 29.12.2018]
- [22] R. FISHER, *The logic of inductive inference*, Journal of the Royal Statistical Society, 98, pp.39-54, 1935. [cit. 29.12.2018]
- [23] EKLUND, Neil, GOEBEL, Kai,. *Using Neural Networks and the Rank Permutation Transformation to Detect Abnormal Conditions in Aircraft Engines*. Computing and Decision Sciences. General Electric Global Research. One Research Circle, Niskayuna, NY 12309 USA. [cit. 30.12.2018]
- [24] David V. Hinkley, *Inference about the change-point in a sequence of random variables* Biometrika, Volume 57, Issue 1, 1 April 1970, Pages 1–17, <https://doi.org/10.1093/biomet/57.1.1>. [cit. 30.12.2018]
- [25] SIEGMUND, D., *Boundary crossing probabilities and statistical applications*. Ann. Statist. 14, 361-404, 1986. [cit. 30.12.2018]
- [26] FRANCOURT, C.L., PRINCIPE, J.C. , *On the Use of Neural Networks in the generalizaed likelihood ratio test for detecting abrupt changes in signals*, Proceedings of the IEEE-NNNS-ENNS International Joint Conference on Neural Networks, Vol. 2, pp. 243-248, 2000. [cit. 30.12.2018]

- [27] Matematická biologie učebnice: Úvod do klasifikace pomocí hranic [online]. [cit. 15.04.2019] Dostupné z URL: <<http://portal.matematickabiologie.cz/>>
- [28] *Designing Algorithms for Condition Monitoring and Predictive Maintenance - MathWorks - MATLAB & Simulink* [online]. Copyright © 1994 [cit. 18.04.2019]. Dostupné z URL: <<https://nl.mathworks.com>>
- [29] YU, Q., JIBIN, L., JIANG, L. (2016). An Improved ARIMA-Based Traffic Anomaly Detection Algorithm for Wireless Sensor Networks. *International Journal of Distributed Sensor Networks*. [cit. 19.04.2019] Dostupné z URL: <<https://doi.org/10.1155/2016/9653230>>
- [30] LI, L., GARIEL, M., HANSMAN, R., PALACIOS, R. (2011). *Anomaly detection in onboard-recorded flight data using cluster analysis*. [cit. 19.04.2019] Dostupné z URL: <[https://www.researchgate.net/publication/254010406\\_Anomaly\\_detection\\_in\\_onboard-recorded\\_flight\\_data\\_using\\_cluster\\_analysis](https://www.researchgate.net/publication/254010406_Anomaly_detection_in_onboard-recorded_flight_data_using_cluster_analysis)>
- [31] VAPNIK, Vladimir Naumovich. Estimation of dependences based on empirical data. New York: Springer-Verlag, c1982. ISBN 0387907335. [cit. 19.04.2019]
- [32] BOTTOU, L, CORTES, C, DENKER, JS, DRUCKER, H, GUYIN, I, JACKEL, LD, LECUN, Y, MULLER, UA, SACKINGER, E, SIMARD, P, VAPNIK, V 1994, *Comparison of classifier methods: A case study in handwritten digit recognition*. in IAPR (ed.), Proceedings of the International Conference on Pattern Recognition, Jerusalem, October 1994. vol. II, IEEE, pp. 77-82. [cit. 20.04.2019]
- [33] STEINWART, I.; CHRISTMANN, A.: *Support vector machines*. New York : Springer, 2008. 601 s. ISBN 9780387772424. [cit. 21.04.2019]
- [34] BURGESS, C. *Tutorial on Support Vector Machines for Pattern Recognition*. In *Data Mining and Knowledge Discovery*, 1998. s. 121 – 167. [Online]. [cit. 22.04.2019]. Dostupné z URL: <<http://research.microsoft.com/pubs/67119/svmtutorial.pdf>>
- [35] R. FLETCHER, *Practical Methods of Optimization*, John Wiley and Sons, Inc., 2nd edition, 1987. [cit. 24.4.2019]
- [36] B. SCHÖLKOPF, J. PLATT, J. SHAWE-TAYLOR, A. SMOLA, R. WILLIAMSON, *Estimating the support of a high-dimensional distribution*, *Neural Computation* 13 (2003) 1443–1471. [cit. 25.4.2019] Dostupné z URL: <<http://users.cecs.anu.edu.au/~williams/papers/P132.pdf>>



- [37] MAHADEVAN, S., SHAH, S. (2009). *Fault Detection and Diagnosis in Process Data Using One-Class Support Vector Machines*. Journal of Process Control. 19. 1627-1639. 10.1016/j.jprocont.2009.07.011. [cit. 25.4.2019] Dostupné z URL: <file:///C:/Users/H226838/Downloads/j.jprocont.2009.07.011.pdf>
- [38] Outlier Detection with One-Class SVMs: An Application to Melanoma Prognosis. National [cit. 25.4.2019] Dostupné z URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3041295/>
- [39] MacQUEEN, J. *Some methods for classification and analysis of multivariate observations*. Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics, 281–297, University of California Press, Berkeley, Calif., 1967. [cit. 26.4.2019] Dostupné z URL: <https://projecteuclid.org/euclid.bsmsp/1200512992>
- [40] MÜNZ, G., LI, S., CARLE, G. (2007). *Traffic Anomaly Detection Using K-Means Clustering*. [cit. 1.5.2019] Dostupné z URL: <https://pdfs.semanticscholar.org/634e/2f1a20755e7ab18e8e8094f48e140a32dacad.pdf>
- [41] BOX, G.E.P., JENKINS, G.M. (1976). *Time Series Analysis: Forecasting and Control*, Revised Edition, San Francisco: Holden Day. [cit. 1.5.2019]
- [42] YU, Q., JIBIN, L., JIANG, L. (2016). An Improved ARIMA-Based Traffic Anomaly Detection Algorithm for Wireless Sensor Networks. International Journal of Distributed Sensor Networks. [cit. 4.5.2019] Dostupné z URL: <https://doi.org/10.1155/2016/9653230>
- [43] Bayesian information criterion. [cit. 6.5.2019] Dostupné z URL: <https://www.immagic.com/eLibrary/ARCHIVES/GENERAL/WIKIPEDI/W120607B.pdf>

## Seznam symbolů, veličin a zkratek

<b>MRI</b>	Magnetická rezonance ( z anglického „Magnetic Resonance Imaging“)
<b>BPR</b>	Bypass ratio
<b>FF</b>	Fuel flow
<b>FADEC</b>	Full Authority Digital Engine Control
<b>ECU</b>	Electronic Control Unit
<b>RPM</b>	Revolutions Per Minute
<b>LVDT</b>	Linear Variable Differential Transformer
<b>SVM</b>	Support Vector Machine (v českém překladu "Metoda podpůrných vektorů")
<b>KKT</b>	Karushovy-Kuhnovy-Tuckerovy podmínky
<b>ARIMA</b>	Zkratka anglického AutoRegressive Integrated Moving Average, „autoregresní integrovaný klouzavý průměr)
<b>TKO</b>	Takeoff - vzletová fáze letu
<b>CRU</b>	Cruise - fáze letu, kdy letadlo letí delší dobu v určité výšce, to znamená, že ani neklesá, ani nestoupá.
<b>RBF</b>	Radial Base Function