



UNIVERZITA KOMENSKÉHO, BRATISLAVA  
FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

# VIZUALIZÁCIA VERIFIKÁCIE PREDPOVEDNÝCH MODELOV POČASIA

Diplomová práca

Bratislava, 2015

Bc. Marek Kružliak



UNIVERZITA KOMENSKÉHO, BRATISLAVA  
FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

# VIZUALIZÁCIA VERIFIKÁCIE PREDPOVEDNÝCH MODELOV POČASIA

Diplomová práca

Študijný program:	Aplikovaná informatika
Študijný odbor:	2511 Aplikovaná informatika
Školiace pracovisko:	Katedra aplikovanej informatiky
Školiteľ:	RNDr. Andrej Lúčny, PhD.

**Bratislava, 2015**

**Bc. Marek Kružliak**

Tu bude zadanie

Čestne prehlasujem, že som túto diplomovú prácu vypracoval samostatne s použitím citovaných zdrojov.

.....

TODO - tu bude podakovanie

# Abstrakt

TODO

**Kľúčové slová:** vizualizácia informácií, verifikácia predpovedí počasia

# Abstract

TODO.

**Keywords:** information visualization, verification of weather forecasts

# Obsah

<b>1</b>	<b>Úvod</b>	<b>1</b>
<b>2</b>	<b>Verifikácia predpovedných modelov počasia</b>	<b>2</b>
2.1	Predpovedný model počasia . . . . .	2
2.1.1	WRF model . . . . .	4
2.2	Dáta . . . . .	4
2.2.1	Predpovedané dáta . . . . .	5
2.2.2	Pozorované dáta . . . . .	6
2.2.3	Párovanie dát . . . . .	6
2.3	Meranie chyby predpovede . . . . .	8
2.3.1	Stredná chyba predpovede . . . . .	8
2.3.2	Stredná absolútna chyba . . . . .	9
2.3.3	Stredná kvadratická chyba . . . . .	9
2.3.4	Všeobecná kumulovaná chyba . . . . .	9
2.3.5	Medián absolútnych chýb . . . . .	11
<b>3</b>	<b>Predchádzajúce riešenia</b>	<b>12</b>
3.1	Verifikačný softvér . . . . .	12
3.1.1	Štatistický softvér . . . . .	12
3.1.2	Špecializovaný softvér . . . . .	15
3.2	Techniky vizualizácie vo verifikácii . . . . .	15
3.2.1	Bodový graf . . . . .	15
3.2.2	Histogram . . . . .	16
3.2.3	Krabicový diagram . . . . .	16



3.2.4	Time series plot . . . . .	21
<b>4</b>	<b>Návrh riešenia</b>	<b>22</b>
4.1	Návrh systému . . . . .	22
4.2	Návrh vizualizácie . . . . .	22
4.2.1	Charakteristika dát . . . . .	22
4.2.2	Špecifikácia požiadaviek na vizualizáciu . . . . .	22
4.2.3	Návrh rozloženia prvkov vizualizácie . . . . .	22
4.2.4	Návrh vizualizácie štatistík verifikácie . . . . .	22
4.2.5	Návrh vizualizácie distribúcie chýb . . . . .	22
4.2.6	Návrh farebnej palety . . . . .	30
<b>5</b>	<b>Implementácia</b>	<b>31</b>
5.1	Použité technológie . . . . .	31
5.1.1	Java . . . . .	31
5.1.2	JavaScript . . . . .	31
5.1.3	d3.js . . . . .	31
<b>6</b>	<b>Výsledky</b>	<b>32</b>

# Zoznam obrázkov

2.1	Flowchart systému predpovedného model počasia od edukačného programu The COMET [LE11]. Na obrázku je zvýraznená časť, ktorej sa venujeme v tejto práci. . . . .	3
2.2	Vizuálne znázornenie dvoch bežne používaných metód na získavanie hodnôt z mriežky . . . . .	7
3.1	Pôvodný návrh krabicového diagramu, ako bol prezentovaný v práci <i>Exploratory Data Analysis</i> (1977) [Tuk77] . . . . .	18
3.2	a) Klasický krabicový diagram b-f) Vizuálne variácie krabicového diagramu b-c) 2 variácie pre kvartilový graf [Tuf83] c) Skrátený krabicový diagram [PKR07] e) Range-bar chart [Spe52] f) Farebná variácia [Car94] . . . . .	19
3.3	Krubicové diagramy s pridanou informáciou a) Krabicový diagram s variabilnou šírkou [MTL78] b) Vrúbkovaný krabicový diagram [MTL78] c) Krabicový diagram s informáciou o šikmosti dát [CM05] . . . . .	20
4.1	Obrázky z článku <i>Functional Boxplots</i> [SG11] a) Funkcie meraní teploty hladiny mora b) Funkčný krabicový diagram c) Rozšírený Funkčný krabicový diagram o centrálné regióny $C_{0.25}$ a $C_{0.75}$ . . . . .	29

# Zoznam tabuliek

# Zoznam skratiek

WRF	Weather Research and Forecasting
NCEP	National Centers for Environmental Prediction
NCAR	National Center for Atmospheric Research

# Kapitola 1

## Úvod

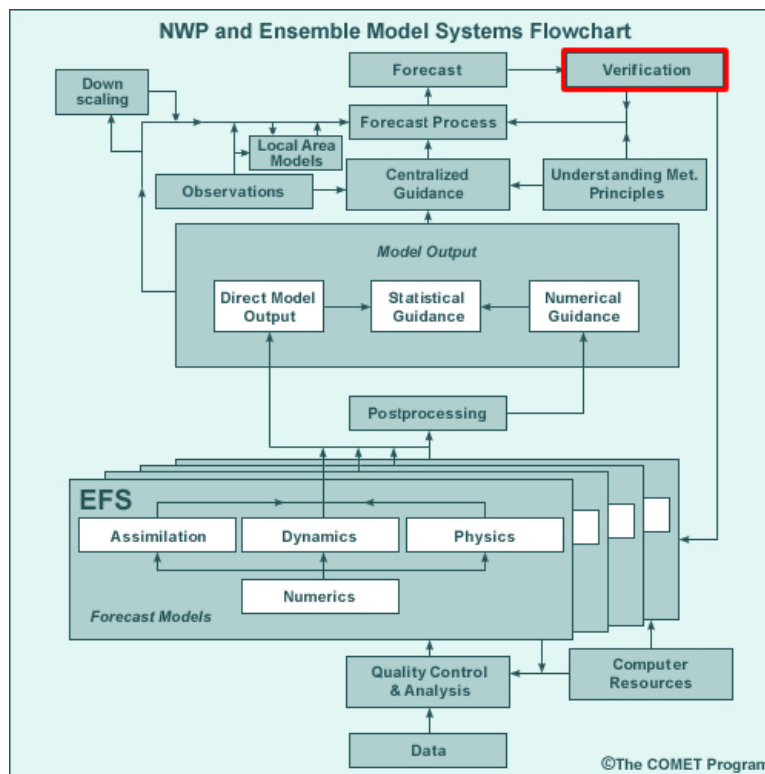
# Kapitola 2

## Verifikácia predpovedných modelov počasia

Verifikácia je proces, ktorý má overiť správnosť fungovania predpovedného modelu počasia. Z tohto dôvodu je nepostrádateľnou súčasťou meteorologického výskumu a taktiež celkového procesu predpovedania počasia. [CWS<sup>+</sup>08] Ciele verifikácie môžeme rozdeliť do troch skupín: *administratívne*, *vedecké* a *ekonomické*. Medzi *administratívne ciele* patrí monitorovanie úspešnosti predpovedania modelu a nasmerovanie užívateľov na jeho správnu konfiguráciu alebo voľbu iného modelu. *Vedekými cieľmi* sú identifikovanie a oprava slabín modelu a taktiež vylepšovanie predpovedí. *Ekonomickými cieľmi* sú rozhodovanie, kam majú smerovať investície do výskumu a iné závažné ekonomické rozhodnutia. [FJB12]

### 2.1 Predpovedný model počasia

Už v 19. storočí vývoj termodynamiky na základe Newtonovskej fyziky vyvrcholil v ucelení množiny fundamentálnych princípov, ktoré riadia prúdenie plynov v atmosfére. Začiatkom 20. storočia sa o matematický prístup k predpovedaniu počasia najviac zaslúžili osobnosti ako Vilhelm Bjerknes alebo Lewis F. Richardson. Avšak na ďalší úspech, v tejto oblasti, sa muselo čakať až na vynájdenie prvých počítačov počas 2. svetovej vojny ako bol IAS alebo ENIAC. [Lyn07] Prvá úspešná predpoveď bola vykonaná v 50. rokoch minulého storočia a to hlavne vďaka práci Jula Charneyho. Následný vývoj vo výpočtovej sile počítačov,



Obr. 2.1: Flowchart systému predpovedného model počasia od edukačného programu The COMET [LE11]. Na obrázku je zvýraznená časť, ktorej sa venujeme v tejto práci.

používanie satelitných pozorovaní a vývoj samotnej meteorológie ako vedy zapríčinil, že je numerická predpoveď počasia (NWP) dnes najúspešnejším prístupom ako predpovedať počasie. [Gol]

Odvtedy vzniklo veľké množstvo modelov, ako sú napríklad GFS, NAM, RUC, WRF, SREF, GEFS, ECMWF, ALADIN a mnoho ďalších. Naša práca sa zameriava konkrétne na verifikáciu modelu *WRF*. Taktiež pokračuje neustály vývoj aj vďaka novým modelovacím technikám, novým parametrizáciám, a zvyšovaniu výkonu výpočtových zdrojov.

Ako môžeme vidieť na obrázku 2.1 proces predpovedania počasia má okrem numerického modelu, ktorý je jej jadrom, aj iné časti. Ako príklad môžeme spomenúť získavanie vstupných dát, ich predspracovanie, postprocesing, spracovanie výstupu a následne poskladanie samotnej predpovede. Cieľom nášho záujmu, celého procesu predpovedania, je *verifikácia*. Ako môžeme vidieť z obrázka, verifikácia vplýva na vyladenie parametrov modelu, avšak tento proces sa nedeje automaticky, ale vyžaduje prácu meteorológov a ich

chápanie základných meteorologických princípov.

### 2.1.1 WRF model

Ako sme už spomenuli *The Weather Research and Forecasting* (WRF) model je *numerická predpoveď počasia* (NWP) a systém atmosferickej simulácie.

WRF je podporovaný, ako bežný nástroj pre univerzity, výskum a operačné komunity, pričom sa usiluje o splnenie požiadaviek ich všetkých súčasne. Vývoj WRF modelu bol snahou mnohých spoločností ako napríklad *The National Center for Atmospheric Research's* (NCAR), *Mesoscale and Microscale Meteorology* (MMM), *The National Oceanic and Atmospheric Administration's* (NOAA) *National Centers for Environmental Prediction* (NCEP) a *Earth System Research Laboratory* (ESRL), oddelenie ministerstva obrany *Air Force Weather Agency* (AFWA) a *Naval Research Laboratory* (NRL), *The Center for Analysis and Prediction of Storms* (CAPS) [WCSDG<sup>+</sup>08].

WRF model je vhodný pre širokú škálu aplikácií od *metódy vzdušných vírov* (Large Eddy Simulation - LES) až po globálne simulácie počasia. Takéto aplikácie vyžadujú numerické predpovede v reálnom čase, vývoj a štúdium asimilácie dát, výskum parametrizovanej fyziky, výskum parametrizovanej fyziky, modelovanie kvality ovzdušia, idealizované simulácie, čo všetko WRF model spĺňa.

V roku 2008 evidovala WRF viac ako 6000 užívateľov, no dnes (2014) eviduje viac ako 25000 užívateľov vo viac ako 130 krajinách sveta. Tieto fakty poukazujú na to, že WRF model má nie len veľkú základňu užívateľov, ale aj vývojárov a má v budúcnosti istotne svoje miesto a preto si myslíme, že sa oplatí investovať čas a úsilie do verifikácie tohto modelu.

## 2.2 Dáta

Na správne zhodnotenie úspešnosti modelu potrebujeme dva druhy dát. V prvom rade sa jedná o dáta, ktoré sú výstupom z daného predpovedného modelu počasia, teda **predpovedané dáta**. Tieto umelo získané dáta chceme konfrontovať s realitou, aby sme si mohli vytvoriť obraz o správnom fungovaní celého modelu. Realitu v našom prípade predsta-

vujú dáta namerané špecializovanými meteorologickými senzormi, ktoré označujeme ako **pozorované dáta** alebo skrátene *pozorovania*.

### 2.2.1 Predpovedané dáta

Predpovedané dáta z modelu WRF sa ukladajú vo formáte **GRIB**, čo je skratka pre *GRIdded Binary* [WMO94] alebo na iných miestach uvádzané ako *General Regularly-distributed Information in Binary form* [WMO03]. Tento formát je štandardom Svetovej meteorologickej organizácie teda *World Meteorological Organization* (WMO). Jedná sa o pomerne rozšírený formát, používaný pri veľkom množstve meteorologických aplikácií a je taktiež používaný ako výstupný formát pre iné predpovedné modely ako WRF, či už ECMWF, GFS, NAM, SREF alebo mnohé iné [NCE14].

Doteraz boli vyvinuté 3 verzie tohoto formátu od 0 po 2. Verzia 0 bola určená pre malé projekty typu TOGA a to iba s limitovaným použitím a dnes sa táto verzia už vôbec nepoužíva. Verzia grib 1 [WMO94], grib 2 [WMO03] sú dnes bežne používané väčšinou meteorologických centier.

Medzi verziami 1 a 2 nie sú žiadne rozdiely v obsahovej filozofii, preto popis obsahu gribovského formátu, ktorý tu uvádzame je spoločný pre obe tieto verzie. *Gribovský súbor* (ďalej iba *Grib*) pozostáva z viacerých *Gribovských záznamov*, pričom jeden záznam môže existovať ako samostatný Grib. Vďaka tomu je možné ľahko spájať Griby, a to tiež v ľubovoľnom poradí, bez toho, aby sme ich nejako poškodili. Samozrejme musí byť zachovaná homogenita, čo sa týka verzií Gribov, teda verziu 1 nemožno miešať s verziou 2 a naopak. Už samotný názov *Gridded Binary* nám napovedá, že dáta sú usporiadané v pravidelnej mriežke. Každý Gribovský záznam obsahuje dvojrozmernú mriežku (zemepisná šírka x zemepisná dĺžka) hodnôt v určitom čase a vertikálnej hladine. Taktiež v hlavičke záznamu sa nachádzajú metainformácie, ktoré nám hovoria o aké dáta ide, teda o akú premennú sa jedná, čas predpovede, výškovú hladinu a podobne. Grib je zvyčajne z tohto dôvodu 2 až 5 rozmerná dátová štruktúra s veľkým množstvom veličín ako je napríklad teplota, tlak, relatívna vlhkosť, rosný bod,  $u$  a  $v$  súradnice vetra a ďalšie, ktoré sú definované v rôznych hladinách. Taktiež je dôležité povedať, že Grib zriedkakedy zachytáva povrch celej planéty, ale iba vymedzenú skúmanú oblasť - *doménu*.



### 2.2.2 Pozorované dáta

Pozorovania sa získavajú meraním priamo v teréne pomocou špecializovaných meracích zariadení, ktoré sú súčasťou meteo staníc. Každá stanica môže obsahovať iné vybavenie, ku príkladu teplomer, zrážkomer, barometer, vetromer a im podobné [Vas98], ktorými môžeme zachytávať informácie o rôznych skúmaných veličinách.

Majoritná časť meraní sa deje pri povrchu zeme priamo na meteo staniciach a nazývajú sa *surface* merania. Tieto merania najlepšie popisujú dianie v oblasti najväčšieho záujmu (biosfére), avšak neobsahujú informáciu o dianí v iných výškových hladinách. Pozorovania týchto hladín sa dejú pomocou *radiosondy*, ktorá je pripojená k meteo balónu alebo vypustená z lietadla smerom k zemi. Takéto pozorovania sa nazývajú *upper air* merania.

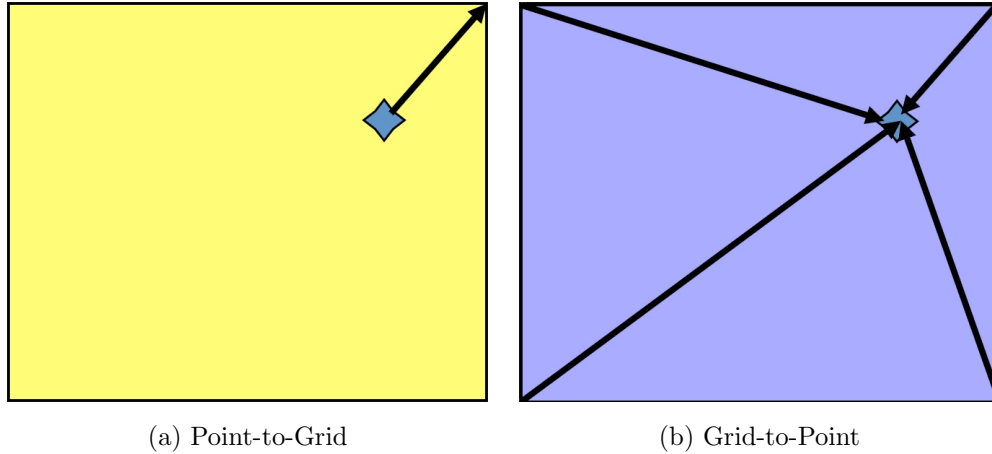
Narozdiel od predpovedaných dát, pozorované dáta nemajú štandardizovaný formát a zvyčajne sa ukladajú do databázy. Aby sme zhrnuli charakteristiku týchto dát, jedná sa o niekoľko meraných veličín, nameraných v konštantných časových krokoch - napríklad každú minútu alebo každú hodinu - v jednom konkrétnom geografickom bode a zvyčajne pri povrchu zeme, teda ak sa nejedná o *upper air* merania, ktoré sa uskutočňujú v štandardných výškových hladinách, ktoré sa merajú v hPa.

### 2.2.3 Párovanie dát

Z predpovedného modelu a rovnako aj z merania získame veľké množstvo hodnôt. Aby sme mohli korektne porovnať predpovede s pozorovaniami, je nevyhnutné nájsť správne párovanie týchto hodnôt, teda zistiť, ktorú hodnotu porovnať s ktorou, aby sme získali zmysluplný výsledok.

Vždy sa snažíme nájsť správnu predpoveď pre pozorovanie a nie naopak. Dôvodom je, že chceme skúmať vzťah predpovede s realitou a preto v párovaní musí byť zahrnutých čo najviac **meraných** hodnôt, ak nie všetky.

Každá pozorovaná hodnota, ktorú chceme spárovať má štyri kľúče podľa ktorých hľadáme pár: *meraná veličina* (napríklad teplota), *čas merania*, *výšková hladina* a *geografická poloha*. Nájsť všetky hodnoty podľa kľúča meranej veličiny v Grike je ľahké, keďže sa jedná o kategorickú premennú, teda môže nadobúdať iba určitý konečný počet hodnôt.



Obr. 2.2: Vizuálne znázornenie dvoch bežne používaných metód na získavanie hodnôt z mriežky

Toto sa však nedá povedať o čase, hladine a polohe, ktoré sú spojitými premennými.

Pre čas pozorovania, čas predpovede a výškovú hladinu existujú štandardy, ktoré určujú v akých časoch resp hladinách sa robia merania a predpovede, čo nám uľahčuje prácu. Ak sa napriek tomu čas alebo hladina v Grike nevyskytuje, tak pár vyhadzujeme z párovania.

V prípade polohy zo samozrejmych dôvodov neexistuje žiaden štandard a hustota mriežky v Grike nemôže byť nikdy tak veľká, aby poloha našej stanice vždy dopadla na presný bod mriežky. Z tohoto dôvodu získavame hodnoty z mriežky z okolitých bodov a to dvoma metódami *Point-to-Grid* a *Grid-to-Point* [FJB12], ktoré sú znázornené na obrázku 2.2. Jedná sa vlastne o dve interpolačné metódy. Point-to-Grid predstavuje metódu *najbližší sused* (*Nearest Neighbour*) a Grid-to-Point *bilineárnu interpolačnú metódu*.

Výber správnej metódy môže značne ovplyvniť výsledok. Dôvodom je, že môžu byť veľké rozdiely hodnôt v okolitých mrežových bodoch a tak, ak pomocou Point-To-Grid metódy získame nízku hodnotu, tak pomocou Grid-To-Point môžeme získať hodnotu omnoho väčšiu, vplyvom zvyšných troch bodov, ktoré vstúpili do interpolácie. Nemožno však jednoznačne povedať, ktorá z metód je lepšia, keďže obe môžu v istých prípadoch dávať lepšie výsledky.

## 2.3 Meranie chyby predpovede

Výsledkom procesu párovania je  $n$  párov (predpoveď, pozorovanie), ktoré je možné porovnať. Z porovnania týchto dvojíc získame numerickú hodnotu, ktorá nám hovorí o veľkosti chyby predpovede daného modelu pre vybrané predpovedané časy.

Chybu predpovede  $e_i$  pre  $i$ -tu dvojicu  $(y_i, \hat{y}_i)$  definujeme takto:

$$e_i = (y_i - \hat{y}_i)$$

Kde  $y_i$  je predpoveď a  $\hat{y}_i$  je pozorovanie. Takýmto spôsobom z  $n$  párov získame  $n$  chýb, ktoré agregujeme pomocou rôznych štatistických metód, ktoré sú bežne používané pri verifikácií predpovedí, ako sa spomína v [Nur03], [FJB12] a [Cas09]. Výsledkom agregácie je numerická hodnota, ktorá sa nazýva *skóre* predpovede.

### 2.3.1 Stredná chyba predpovede

Budeme ju označovať ako *MFE* z anglického *Mean Forecast Error*, ale v literatúre je možné ju nájsť ako *ME* [Nur03], teda *stredná chyba* alebo ako *Linear Bias* [Cas09], [FJB12]. Vzorec pre výpočet MFE vyzerá nasledovne:

$$MFE = \frac{1}{n} \sum_{i=0}^n e_i$$

MFE je možné vypočítať aj ako rozdiel priemerov predpovedí a pozorovaní.

$$MFE = \bar{y} - \bar{\hat{y}}$$

MFE vyjadruje priemerný smer chyby. To znamená, že pozitívny výsledok indikuje *over-forecast*, teda nadhodnotenú predpoveď a negatívny výsledok *under-forecast*, teda podhodnotenú predpoveď. Avšak MFE **nevyjadruje veľkosť** chyby v tomto smere, keďže kladné a záporné chyby sa navzájom môžu zrušiť. Napríklad máme množinu chýb  $E = \{2, -5\}$ , tak MFE pre  $E$  je -1.5, ale priemerná veľkosť chyby je 3.5.

### 2.3.2 Stredná absolútna chyba

Budeme ju označovať ako *MAE* z anglického *Mean Absolute Error*. Vzorec pre výpočet MAE vyzerá nasledovne:

$$MAE = \frac{1}{n} \sum_{i=0}^n |e_i|$$

Narozdiel od MFE **neurčuje smer chyby**, ale vyjadruje veľkosť chyby. Z týchto dôvodov je v praxi odporúčané zobrazovať MFE a MAE súčasne [Nur03].

### 2.3.3 Stredná kvadratická chyba

Budeme ju označovať ako *RMSE* z anglického *Root Mean Square Error*. Vzorec pre výpočet RMSE vyzerá nasledovne:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=0}^n e_i^2}$$

Z povahy vzorca pre RMSE je jasné, že rovnako ako MAE, ani RMSE neurčuje smer chyby, pretože nadobúda vždy iba kladné hodnoty. Ďalšou vlastnosťou RMSE je, že nadobúda hodnoty vždy väčšie alebo rovné ako MAE, pričom výsledok RMSE je citlivý na veľké hodnoty chýb.

V praxi sa zvykne používať aj *MSE* (*Mean Square Error*):

$$MSE = \frac{1}{n} \sum_{i=0}^n e_i^2$$

Má podobné vlastnosti ako RMSE s jediným rozdielom, že RMSE meria veľkosť chyby zachovávajúc jednotky danej veličiny (napr. °C), zatiaľ čo MSE jednotky nezachováva [Nur03]. Preto sme si pre náš účel zvolili RMSE, ktoré je jednoduchšie zobrazíť spolu s MFE a MAE v jednom grafe, keďže sa zachováva konzistentnosť jednotiek veličín.

### 2.3.4 Všeobecná kumulovaná chyba

V našom systéme sme navrhli všeobecný vzorec na výpočet kumulovaného skóre, ktorým možno vyjadriť ľubovoľnú zo spomenutých štatistických metód. Takéto vyjadrenie umožňuje

nie len všeobecnosť, ale aj jednoduché rozšírenie systému o ďalšie metódy a to nie len programátorom, ale aj samotným užívateľom systému.

Všeobecný vzorec na výpočet *skóre* pre danú predpoveď vyzerá takto:

$$Score = \Phi\left(\sum_{i=0}^n \varepsilon(e_i)\right)$$

Kde  $\Phi$  je ľubovoľná funkcia z  $\mathbb{R}$  do  $\mathbb{R}$ , teda  $\Phi : \mathbb{R} \rightarrow \mathbb{R}$  a podobne funkcia  $\varepsilon : \mathbb{R} \rightarrow \mathbb{R}$ . Spomenuté metódy môžeme teda skonštruovať zadaním správneho  $\Phi$  a  $\varepsilon$ .

Napríklad pre *MFE*:

$$\begin{aligned}\Phi(x) &= \frac{x}{n} \\ \varepsilon(e) &= e\end{aligned}$$

Pre *MAE*:

$$\begin{aligned}\Phi(x) &= \frac{x}{n} \\ \varepsilon(e) &= |e|\end{aligned}$$

Pre *RMSE*:

$$\begin{aligned}\Phi(x) &= \sqrt{\frac{x}{n}} \\ \varepsilon(e) &= e^2\end{aligned}$$

Pre *MSE*:

$$\begin{aligned}\Phi(x) &= \frac{x}{n} \\ \varepsilon(e) &= e^2\end{aligned}$$

Ako sme spomenuli, je možné rozšírenie o ďalšie metódy a to napríklad o Brownov a Triggov *signál chybných predikcií*, ktorý budeme označovať ako *TS* z anglického *Tracking Signal*. Tieto metódy sme vyššie nespomenuli, keďže sa v meteorologickej praxi nepoužívajú. Uvádzame ich však ako možné rozšírenie, keďže sú tieto metódy bežne používané pri verifikácii iných predpovedných modeloch, ako sú tie meteorologické.

### 2.3.5 Medián absolútnych chýb

Budeme ju označovať ako  $MAD$  z anglického *Median Absolute Deviation*. Vzorec pre výpočet  $MAD$  vyzerá nasledovne:

$$MAD = median(|e|) = |\tilde{e}|$$

Nech je daná usporiadaná postupnosť  $Y_1, \dots, Y_N$ , tak potom *median* náhodnej premennej  $x$  je definovaný rovnako ako v [Wei14]:

$$median(x) = \tilde{x} \equiv \begin{cases} Y_{(N+1)/2} & \text{ak } N \bmod 2 = 0 \\ \frac{1}{2}(Y_{(N+1)/2} + Y_{(N+1)/2+1}) & \text{ak } N \bmod 2 = 1 \end{cases}$$

Z daného vzorca môžeme vidieť podobné vlastnosti ako má MAE, avšak  $MAD$  je robustnejší a extrémne chyby nemajú na skóre žiaden efekt.

# Kapitola 3

## Predchádzajúce riešenia

### 3.1 Verifikačný softvér

Verifikácia predpovedných modelov počasia je úloha dokonale stvorená pre automatizáciu. Z tohto dôvodu meteorológovia začali využívať dostupný štatistický softvér a neskôr boli taktiež vyvíjané špecializované nástroje určené pre verifikáciu. Môžeme teda rozdeliť verifikačný softvér do dvoch základných kategórií a to *štatistický* a *špecializovaný*, ktorý je zväčša podporovaný rôznymi národnými a medzinárodnými organizáciami.

#### 3.1.1 Štatistický softvér

Spoločnou črtou:

- obmedzená funkcionálna
- obmedzená vizualizácia
- slabé / žiadne GUI
- vyžaduje znalosť špecifického programovacieho jazyka
- ...

#### Tabuľkový softvér

Napriek tomu, že je tabuľkový softvér na výpočet štatistík zamietnutý komunitou vedcov a štatistikov ako nevhodný a neprofesionálny, tak je využívaný, a to pomerne často, aj vo vedeckých kruhoch. Výhodou je, že novému užívateľovi umožňuje okamžite vidieť všetky

kroky v základných procedúrach verifikácie a teda je výborný pre výučbové účely. [Poc11] Najznámejší kus softvéru z pomedzi komerčných produktov je *Microsoft Excel* [Mic15] a z voľne dostupných je jeho opensource náprotivok *Open Office Calculate* [Ope15]. Oba programy zahrňujú základné štatistické funkcie ako napríklad stredná kvadratická chyba (*MSE*) pre spojené predpovede (pozri odsek 2.3.3) a taktiež umožňujú generovanie jednoduchých grafov na základe tabuľkových dát. Tabuľkový softvér neposkytuje priamo funkcionality na výpočet ďalších sofistikovanejších verifikačných štatistík, avšak umožňuje ich implementáciu pomocou makro programovania v špecifickom jazyku. Pre Microsoft Excel je to *Microsoft Visual Basic for Applications*(VBA) [Mic13] a pre Open Office Calculate zasa *OpenOffice.org Basic* [Ope13]. Oba jazyky patria do rodiny *Basic* jazykov, takže majú mnoho podobných prvkov.

## MATLAB

*MATLAB* je interaktívne prostredie s vlastným programovacím jazykom, ktorý je využívaný miliónmi inžinierov a vedcov po celom svete [TM15] a tým nevynímajúc meteorológov a ďalších odborníkov pracujúcich v atmosférickom výskume. Zvyčajne sa *MATLAB* využíva na výskum a prototypovanie nových metód a procedúr [Poc11], pretože umožňuje rýchlu a jednoduchú implementáciu, keďže jeho súčasťou je mnoho matematických knižníc a je prispôsobený na prácu s maticami dát. Výhodou *MATLABU* je, že umožňuje tvorbu GUI a taktiež poskytuje kreslenie rôznorodých grafov a diagramov. Mali by sme však podotknúť, že podobne ako väčšina štatistického softvéru, aj *MATLAB* je komerčný produkt. Jeho cena za jednu licenciu je \$2,650 (k roku 2015), čo je pomerne vysoká suma, ak vezmeme do úvahy za akým účelom chceme tento softvér využívať a ako dobre je naň prispôsobený.

## R

Často používaným a pomerne mocným nástrojom je *open source* skriptovací jazyk *R* [www.r-project.org]. V posledných desaťročiach sa stal dominantným jazykom v oblasti štatistického výskumu. Napriek tomu, že ide o voľne stiahnuteľný softvér, tak jeho základný balík obsahuje všetky funkcie, ktoré obsahujú aj platené produkty. *R*-ko však nezostáva len pri tom, pretože v dobe písania tejto práce (marec 2015) bolo dostupných vyše 6400



užívateľských balíkov s rôznorodou funkcionalitou. Pre nás je dôležité, že medzi týmito balíkmi sa objavil aj balík určený na verifikáciu s názvom *verification* [NCAR, 2010]. Tento balík obsahuje základné funkcie verifikácie na výpočet štatistík pre spojité, kategorické ale i pravdepodobnostné predpovede.

Jazyk R neslúži iba na rôznorodé štatistické výpočty, ale poskytuje aj veľmi dobre parametrizovateľnú vizualizáciu. V balíkoch jazyka sa nachádzajú funkcie pre čiarové diagramy, krabicové diagramy, bodové grafy a mnohé iné komplexnejšie vizualizácie, ale tiež funkcie na zobrazenie základných vizuálnych prvkov, ktorými možno vytvoriť úplne novú osobitnú vizualizáciu.

## Statistical Analysis Software (SAS)

*Statistical Analysis Software* [?], skratene SAS, je opäť štatistický programovací jazyk aj so svojim vývojovým prostredím. V oblasti bioštatistiky a farmakológie je veľmi uznávaným a často používaným jazykom. Keďže je veľká podobnosť v používaných metódach medzi verifikáciou predpovedí a spomínanými odvetvami [Poc11], SAS poskytuje funkcionalitu použiteľnú aj pre verifikáciu. Okrem iného SAS ponúka základné, ale aj niektoré pokročilejšie nástroje na vizualizáciu dát. Opäť však musíme podotknúť, že ide o komerčný produkt, ktorého cena licencie je pomerne vysoká.

## Interactive Data Language (IDL)

IDL, teda *Interactive Data Language* je opäť jeden z matematických programovacích jazykov, ktoré patria medzi menej používané v komunite atmosferického výskumu. Napriek tomu niektorí výskumníci medzi, ktorými je aj *Beth Ebert* z *Centre for Australian Weather and Climate Research* (CAWCR) uverejnili na svojich webstránkach kód obsahujúci metódy verifikácie napísané v IDL:

- Metódy pre verifikáciu pravdepodobnosti zrážok (<http://www.cawcr.gov.au/projects/verification/POP3/POP3.html>)
- Priestorové metódy (<http://www.cawcr.gov.au/staff/eee/#Interests>)

IDL na používanie požaduje taktiež získanie platenej licencie, čo obmedzuje počet užívateľov a rovnako aj zdieľanie kódu medzi, ktorý by si mohol ktokoľvek spustiť.

### 3.1.2 Špecializovaný softvér

NCAR Command Language (NCL)

Model Evaluation Tools (MET)

Ensemble Verification System (EVS)

## 3.2 Techniky vizualizácie vo verifikácii

### 3.2.1 Bodový graf

Najjednoduchším spôsobom ako analyzovať vzťah dvoch náhodných premenných je *bodový graf*, ktorý je inak nazývaný aj *korelačný diagram* a známy je tiež pod svojim anglickým pomenovaním *scatter plot*. Bodový graf je vhodný na štúdium kolerácie dvoch premenných a taktiež je výborný pri odhaľovaní takzvaných *outlier*-v, teda hodnôt, ktoré sa nejakým spôsobom výrazne odlišujú od tých ostatných. Taktiež nám nepriamo podáva správu o distribúcii hodnôt, čo však pri ich veľkom počte môže byť skreslené, keďže sa body začínú postupne prekrývať, a tak nemožno určiť v akej oblasti je viacej, či menej bodov.

#### Konštrukcia bodového grafu

Bodový graf na svoju konštrukciu využíva karteziánsku sústavu súradníc. Dve náhodné premenné, ktoré chceme porovnať vizualizujeme tak, že spravíme zobrazenie jednej premennej na  $x$ -ovú súradnicovú os a zobrazenie druhej premennej na  $y$ -ovú os. Následne v danom bode nakreslíme určený symbol, čo zvyčajne býva čierna bodka alebo krúžok. Na obrázku XYZ vidíme príklad výsledného bodového grafu, ktorý vznikne takýmto postupom.

Bodový graf ponúka mnoho spôsobov, ako pridať ďalší rozmer informácie do vizualizácie. Môžeme zvoliť odlišnú súradnicovú sústavu, čím veľmi jednoducho vytvoríme

napríklad 3D bodový graf [ TODO ref ]. Ďalšou možnosťou rozšírenia je zmena vykresľovaného symbolu, ktorému môžeme nastavovať rôzne parametre, do ktorých možno zakódovať nové informácie. Zvyčajne sú týmito parametrami farba[ TODO ref ], alfa-transparencia[ TODO ref ], veľkosť [VWvH<sup>+</sup>07], tvar [GSS, ESPUEPPAS] a podobne. Tieto modifikácie si v komunite vizuálnej analýzy vyslúžili aj vlastné názvy a teoretické zázemie, avšak v našej práci sa týmto druhom diagramov nebudeme venovať.

### Kantil-kvantil graf

Dôležitou a často využívanou variáciou pre bodový graf je *kvantil-kvantil graf*, skrátene *Q-Q graf* (Q z anglického *quantil*).

### Úloha bodového grafu vo verifikácii

#### 3.2.2 Histogram

#### 3.2.3 Krabicový diagram

Krabicový diagram je v anglickej literatúre zvyčajne nazývaný *box plot* alebo na niektorých miestach označovaný tiež ako *box and whisker<sup>1</sup> plot*. Odkedy bol prvýkrát publikovaný v roku 1977 [Tuk77], uplynulo už takmer 40 rokov a dnes ho považujeme už za štandardnú techniku ako vizualizovať distribúciu hodnôt kompaktným spôsobom. Na svoju reprezentáciu využíva súbor 5 čísel (tzv. *5-number summary*) [Pot06], ktoré charakterizujú distribúciu dát robustným spôsobom. Tým, že zredukujeme zvyčajne veľkú dátovú množinu na týchto pár hodnôt ušetríme nielen vzácny vizuálny priestor [WS12], ale taktiež námahu analytika, ktorý sa snaží preskúmať iba niektoré vybrané charakteristiky.

### Konštrukcia krabicového diagramu

Na zostavenie krabicového diagramu potrebujeme týchto 5 hodnôt: medián, horný a dolný kvartil, maximum, minimum. (pozri obrázok 3.1) Prvé tri hodnoty sú takzvané kvartily

---

<sup>1</sup>Slovo *whisker* znamená po slovensky fúz, čo naznačuje, že čiary, ktoré spájajú horný a dolný kvartál s hraničnými hodnotami pripomínajú fúzy.

(Q1, Q2, Q3), ktoré rozdeľujú súbor dát na 4 rovnako veľké časti a ďalšie dve sú extrémne hodnoty, ktoré ohraničujú celú dátovú množinu.

Zovšeobecnenie kvartilov sú kvantily ktoré rozdeľujú množinu na  $n$  rovnakých častí a preto môžeme o kvartile hovoriť ako o 4-quantile. Medián hodnôt je 2-quantil a teda rozdeľuje množinu na 2 rovnaké časti a je definovaný rovnako ako v časti 2.3.5. Ďalej horný (Q1) a dolný (Q3) kvartil získame ako medián hodnôt pod a nad hodnotou Q2, pričom hodnotu Q2 nezahrňame do výpočtov.

Na obrázku 3.1 vidíme, že krabica v grafe určuje pozície horného a dolného kvartilu, zatiaľ čo vnútro krabice znázorňuje takzvané *IQR*. Táto skratka označuje *interquartile range*, čo sa dá preložiť ako *medzikvartilový rozsah*. IQR definujeme ako rozdiel kvartilov Q3 a Q1:

$$IQR = Q3 - Q1$$

IQR nám hovorí o vzdialenosti týchto dvoch kvartilov, preto nám môže byť tento vzorec na pohľad podozrivý, keďže sa javí, že IQR by mohlo nadobúdať aj záporné hodnoty. My však vieme z definície Q3 a Q1, že  $Q3 > Q1$  a ich rozdiel je teda vždy nezáporný (Hovoríme o rozdieli Q3 od Q1, tak ako je definované IQR).

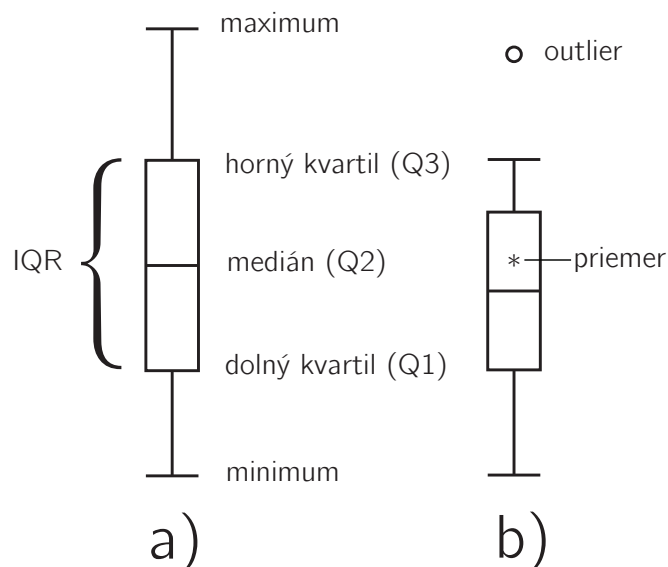
Malú obmenu pôvodného návrhu krabicového diagramu od Tukeyho, vidíme na obrázku 3.1 b), kde malé bodky znázorňujú hodnoty nazývané *outlier*, teda hodnoty ležiace ďaleko od hlavného dátového tela, a hviezdička v strede diagramu určuje priemer hodnôt. Môžeme si všimnúť, že konce čiar vychádzajúcich z boxu nemôžu byť extrémne celej množiny dát, ale sú iba extrémami vypočítaných z dát bez *outlier*-ov.

Otázkou zostáva ako určiť, ktorá hodnota je *outlier* a ktorá nie je. Na zodpovedanie tejto otázky sa využíva už spomínaný rozsah IQR. Pomocou neho sa definujú hranice *inner fences* ( $f_1, f_2$ ) a *outer fences* ( $F_1, F_2$ ), za ktorými hovoríme už o *outlier*-och alebo o ďalekých *outlier*-och [SOA04]. Definované sú nasledovne:

$$f_1 = Q1 - c \times IQR$$

$$f_2 = Q3 + c \times IQR$$

$$F_1 = Q1 - C \times IQR$$



Obr. 3.1: Pôvodný návrh krabicového diagramu, ako bol prezentovaný v práci *Exploratory Data Analysis*(1977) [Tuk77]

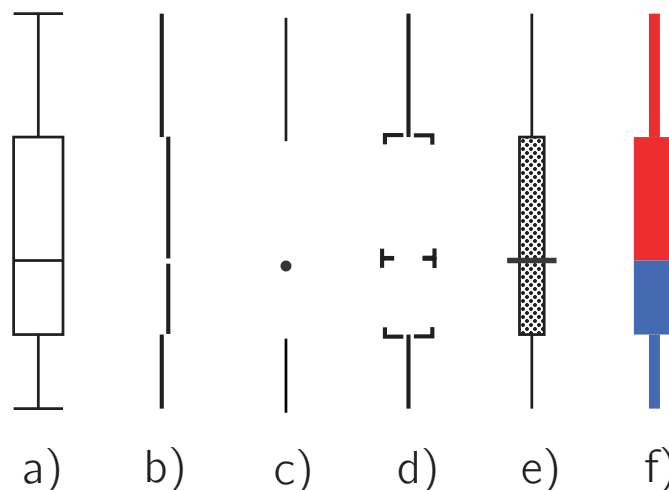
$$F_2 = Q3 + C \times IQR$$

Konštanty  $c$  a  $C$  sú v niektorých zdrojoch definované rôzne. Najčastejšie sa však vyskytujú hodnoty  $c = 1.5$  a  $C = 3$ , tak ako ich určil pôvodný autor krabicového diagramu [Tuk77].

### Ďalšie variácie krabicového diagramu

Popularita krabicového diagramu nevyhnutne viedla k jeho vývoji a modifikáciám. Môžeme hovoriť o dvoch druhoch modifikácií. Jednak *syntaktickej* (vizuálnej), kedy sa zachováajú všetky vlastnosti a informácie ako v pôvodnom diagrame, len sa menia vizuálne prvky grafu. A modifikácii *sémantickej* pridaním ďalšej popisnej informácie do grafu, čo má na záver vplyv aj na jeho vizuálnu stránku.

**N**a obrázkoch 3.2 b-f) môžeme vidieť niektoré vizuálne variácie krabicového diagramu. Vznik prvých troch motivovala snaha maximalizovať takzvaný *data-ink* [Tuf83], teda množstvo atramentu alebo počet pixlov, ktoré zodpovedajú nejakým dátam. Autori sa teda snažili čo najviac znížiť počet vizuálnych prvkov a ponechať len tie, ktoré skutočne nesú nejakú informáciu.



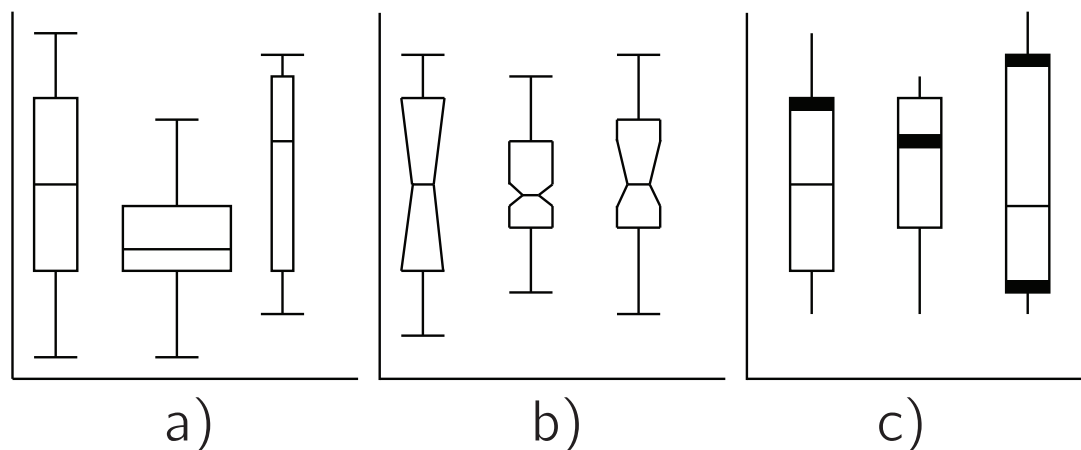
Obr. 3.2: a) Klasický krabicový diagram b-f) Vizuálne variácie krabicového diagramu b-c) 2 variácie pre kvartilový graf [Tuf83] c) Skrátený krabicový diagram [PKR07] e) Range-bar chart [Spe52] f) Farebná variácia [Car94]

Na obrázku 3.2 b) a c) vidíme dve z viacerých riešení navrhnutých v knihe *The Visual Display of Quantitative* [Tuf83], ktoré autor nazýva *kvartilový graf*. Preceptuálne štúdie [SB91] však ukázali, že tieto variácie sú výrazne menej presné ako originálny návrh.

Návrh d) ukazuje skrátený krabicový diagram [PKR07], ktorý sa taktiež snažil zredukovať množstvo okupovaného vizuálneho priestoru. Rozdielom je však to, že jeho účelom nie je existovať samostatne, ale ako súčasť *summary plot*-u, ktorý zahŕňa histogram a ďalšie glify znázorňujúce informácie ako priemer, štandardná odchylka alebo koeficient asymetrie.

Obrázok 3.2 e) znázorňuje predchodcu krabicového diagramu *range* graf alebo tiež nazývaný *range-bar* graf, ktorého autorkou je Mary Eleanor Spear [Spe52].

Ako posledný príklad vizuálnej modifikácie uvádzame pridanie farieb do krabicového diagramu. Táto farebná variácia uchováva tvar diagramu, avšak časť nad mediánom je zafarbená inou farbou ako hodnoty pod. Autori článku odporúčajú červenú a modrú farbu, tak ako na obrázku 3.2 f). Cieľom tohto prístupu bolo chápať krabicový diagram ako jednu percepčnú jednotku a nahradiť 5 symbolov jedným, čím sa mala znížiť námaha pozorovateľa pri analýze a taktiež uľahčiť porovnávanie viacerých diagramov navzájom.



Obr. 3.3: Krabicové diagramy s pridanou informáciou a) Krabicový diagram s variabilnou šírkou [MTL78] b) Vrúbkovaný krabicový diagram [MTL78] c) Krabicový diagram s informáciou o šikmosti dát [CM05]

**K**rabicový diagram umožňuje svojim vzhľadom zakódovanie ďalšej informácie do grafu. Na obrázku 3.3 vidíme aspoň niektoré najčastejšie úpravy...

Len rok po oficiálnom publikovaní krabicového diagramu vznikol článok [MTL78], ktorý zhrňa jeho tri najčastejšie používané modifikácie, z ktorých prvé dve môžeme vidieť na obrázkoch 3.3a) a 3.3b) a tretí je ich kombináciou. V prvom prípade sa využíva šírka boxu na zakódovanie veľkosti množiny, ktorú skrýva za sebou diagram. Takýto graf sa nazýva *Krabicový diagram s variabilnou šírkou*. V druhom prípade ide o takzvaný *Vrúbkovaný krabicový diagram*. V tomto grafe sú pridané *vrúbky*, ktoré zhruba naznačujú ako výrazné sú rozdiely v miere spoľahlivosti rôznych dátových množín.

V niektorých prípadoch krabicový diagram zakrýva skutočný tvar dát, teda jeho šikmosť alebo modalitu, keďže jeho vzhľad nabáda k tomu, aby si užívateľ myslel, že sú dát zamerané na stred a unimodálne. V práci s názvom *Can the Box Plot be Improved?* [CM05] autor uvádza príklad kedy skutočne rôznorodé dáta generujú rovnaký krabicový diagram. Tento problém rieši elegantným a čistým spôsobom pridaním hrubej čiary na základe koeficientu asymetrickosti  $\gamma$ . Na obrázku 3.3c) môžeme vidieť zľava asymetrické dáta, na stred zarovnané dáta a bimodálne dáta.

Krabicový diagram umožňuje mnoho ďalších rozšírení napríklad pridaním informácie

o hustote dát (histogramový krabicový diagram, vázový diagram [Yoa88] , huslový diagram [HN98]), rozšírením pre viacrozmerné dáta (vrecový graf [RRT99], 2D krabicový diagram [Ton05] ) alebo zobrazením hodnôt v inej súradnicovej sústave (napríklad polárnej - vejárový graf [Fis10]). Opis týchto techník je však nad rámec tejto práce, preto ho tu ani nebudeme uvádzať.

### Úloha krabicového diagramu vo verifikácii

Ako sme spomenuli v úvode tejto sekcie, vo všeobecnosti je úlohou krabicového diagramu zobrazíť distribúciu dát v kompaktnom tvare a teda slúži na rýchle porovnanie distribúcií viacerých skupín dát. Pri verifikácii spojitej predpovede sa používa na viacero účelov a my tu spomenieme len niekoľko z nich.

V prvom rade ide o porovnanie distribúcie predpovedí s distribúciou pozorovaní za istý časový interval. V takomto prípade máme vedľa seba iba dva krabicové diagramy, ktoré navzájom porovnávame. Použitie krabicového diagramu v takomto prípade, kedy jeden graf pozostáva iba z niekoľkých (2 až 4) krabicových diagramov považujeme za zbytočné. Pri takomto počte nie je potreba na redukcii vizuálnych prvkov a existujú lepšie techniky na vizualizáciu distribúcie, ktoré sprostredkujúviac informácie a teda môže byť analýza efektívnejšia.

Ďalším použitím je porovnanie distribúcie chýb predpovedí, či už pre rôzne merania, konkrétne predpovedané časy, rôzne predpovedné modely a podobne. Tu považujeme použitie krabicového diagramu za opodstatnené, keďže ide zväčša o porovnávanie väčšieho množstva distribúcií, a tak je jeho jednoduchosť, čitateľnosť, kompaktnosť a iné jeho vlastnosti potrebné.

#### 3.2.4 Time series plot



# Kapitola 4

## Návrh riešenia

### 4.1 Návrh systému

### 4.2 Návrh vizualizácie

#### 4.2.1 Charakteristika dát

#### 4.2.2 Špecifikácia požiadaviek na vizualizáciu

#### 4.2.3 Návrh rozloženia prvkov vizualizácie

#### 4.2.4 Návrh vizualizácie štatistík verifikácie

#### 4.2.5 Návrh vizualizácie distribúcie chýb

Pri verifikácii predpovede spojitej premennej sme použili štatistické metódy spomenuté v sekcii 2.3, ktorých výsledok sme následne vizualizovali. Pôvodné dáta však zostali skryté za použitým matematickým modelom, a tak sme stratili informáciu o distribúcii chyby. Pri verifikácii sa štandardne používajú dve metódy na priamu, či nepriamu vizualizáciu a analýzu distribúcie, ktoré sme opísali v sekcii 3.2. Týmito metódami sú bodový graf (pozri podsekciiu 3.2.1) a krabicový diagram (pozri podsekciiu 3.2.3).

Pri návrhu vizualizácie sme vyskúšali niekoľko vizualizačných techník a zvažili ich silné a slabé stránky.

## Graf hustoty

Jedným z viacerých spôsobov, ako pomerne presne určiť distribúciu chýb je pomocou *grafu hustoty*. Ten sa skonštruuje jednoducho z *funkcie hustoty*, ktorú získame *odhadom hustoty* z dát.

Prvý pohľad na dáta by nám vravel, že ide o dvojrozmerné dáta a teda je potrebné použiť odhad hustoty dvoch premenných. Takýto postup by samozrejme bol možný, ale doviedol by nás k chybnnej vizualizácii a tak aj k mylnej predstave o dátach. Dôvodom je to, že máme záujem o analýzu distribúcie chýb pre každú hodinu predpovede zvlášť, čo znamená, že chceme zistiť distribúciu iba v jednom smere.

Pre vytvorenie grafu hustoty, v prvom rade je potrebné vybrať správny spôsob odhadu hustoty. Jedným z bežne používaným štandardným spôsobom je odhad hustoty pomocou jadra, po anglicky známy ako *kernel density estimation* (KDE) [ref Rosenblatt 56, parzen 62].

Nech máme  $n$  hodnôt  $x_i, 1 \leq i \leq n$ , z ktorých chceme určiť odhad hustoty, potom estimátor hustoty  $\hat{f}_h(x)$ , ktorý aproximuje *funkciu hustoty pravdepodobnosti* (PDF)  $f$ , sa vypočíta takto:

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

kde  $h$  je šírka jadra a  $K(x)$  je funkcia jadra (skrátene iba jadro), ktorá by mala spĺňať nasledovné vlastnosti:

$$K(x) \geq 0$$
$$\int K(x)dx = 1$$

tieto vlastnosti hovoria, že  $K(x)$  je na celom definičnom obore nezáporná a jej integrál je rovný 1, teda sa jedná o normalizovanú funkciu. Bolo preštudovaných mnoho jadier, ako napríklad uniformné, tri-angulárne, Epanechnikovo [29], kvadratické, Gaussové, kosínusové a veľa ďalších. Najbežnejšie a zrejme aj najpraktickejšie [18] je Gaussovo (normálne) jadro, ktoré sme použili aj my:

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$$

Voľba jadra však nemá na výsledok až taký vplyv, ako voľba šírky jadra  $h$ . My sme použili

výpočet šírky jadra na základe dátovej množiny, ktorý aproximuje optimálnu šírku jadra [Scott 1992; Bowman and Azzalini 1997]<sup>1</sup>. Všeobecne pre  $d$  dimenzionálne dáta je vzorec nasledovný:

$$h = \sigma \left( \frac{4}{(d+2)n} \right)^{\frac{1}{d+4}}$$

kde  $\sigma$  je smerodajná odchýlka vypočítaná z daných dát a  $n$  je veľkosť dátovej množiny. V našom prípade je  $d = 1$ , a tak sa nám vzorec zjednoduší na

$$h = \sigma \left( \frac{4}{3n} \right)^{\frac{1}{5}}$$

Aby sme zjednodušili výpočet, tak sme si konštanty vypočítali predom a zaokrúhlili na 2 desatinné miesta, čo považujeme za dostačujúce. Výsledný vzorec, ktorý sa nakoniec objavil v aplikácii je takýto:

$$h = 1.06 \times \sigma \times n^{-\frac{1}{5}}$$

Takýmto spôsobom sme si pre každú hodinu predpovede určili samostatnú funkciu  $\hat{f}_h(x)$ , ktorú môžeme vizualizovať. Zvyčajne sa funkcie hustoty vizualizujú ako bežné funkcie, teda pomocou čiarového diagramu, tak ako na obrázku ???. V našom prípade by bol tento prístup nepraktický, keďže máme veľké množstvo funkcií, tak jednak by bola takáto vizualizácia nepraktická pri porovnávaní distribúcií a taktiež na to nemáme potrebný vizuálny priestor.

Opäť sme teda zvolili štandardné riešenie, ako ušetriť vzácny priestor a to tak, že hodnoty, ktoré by boli zobrazené na  $y$ -ovej os zobrazíme na zvolenú farebnú škálu. Vďaka tomu by teoreticky mohol mať graf hustoty pre jeden čas predpovede šírku 1 pixel bez straty akejkoľvek informácie.

Vieme, že pre ľudí nie je také jednoduché pozorovať malé rozdiely medzi dvoma farbami. Testovaním sme zistili, že pri takejto vizualizácii, že tento fakt spôsobuje problémy aj pri našej vizualizácii, kedy sa pre pozorovateľa strácajú výkyvy hodnôt

## Pruhový kvantilový diagram

Z časti 3.2.1 sme už dobre oboznámený s pojmom kvantil. Klasický kvantilový diagram [ref] zobrazuje kvantil hodnôt pre jednu dimenziu. Ak si vezmeme naše dáta, kde pre

---

<sup>1</sup>Optimálna šírka jadra je taká, ktorá minimalizuje *strednú integrovanú kvadratickú chybu*.

každý čas je niekoľko chýb predpovedí, tak kvantilový diagram skonštruujeme tak, že pre každý čas vypočítame kvantil, ktorý zobrazíme ako bod alebo ako súčasť lomenej čiary v grafe. Prirodzeným rozšírením je zobrazovať nielen jeden kvantil, ale mnoho kvantilov súčasne. Zvyčajne sú to tieto kvantily  $Q_{0.02}, Q_{0.98}, Q_{0.25}, Q_{0.75}, Q_{0.5}$ . V našej práci sme využili toto rozšírenie na lepšie zobrazenie distribúcie a navrhli sme takzvaný *pruhový kvantilový diagram*.

Jeden *pruh* v grafe definujeme pomocou dvojíc hodnôt v čase - spodným a jeho protiľahlým kvantilom. Spodným kvantilom je kvantil  $Q_\alpha$  a k nemu protiľahlý je kvantil  $Q_{1-\alpha}$ , kde  $0 \leq \alpha \leq \frac{1}{2}$ . Vidíme teda, že pruh ohraničuje hodnoty v okolí stredu usporiadanej množiny dát. Špeciálnym prípadom pruhu je pre  $\alpha = \frac{1}{2}$ , vtedy spodný aj horný kvantil je  $Q_{0.5}$ , čo je vlastne medián.

Pri návrhu vizualizácie sme sa snažili, aby mohol mať diagram variabilný počet pruhov a taktiež, aby rozstup medzi pruhmi bol pravidelný. Pri riešení tohto problému, sme sa inšpirovali krabicovým diagramom, kde sa hodnoty delia mediánom na dve časti, ktoré sa ďalej taktiež delia ich mediánom. Ide teda o rekurzívne delenie usporiadanej množiny na polovicu do hĺbky 2. Túto myšlienku sme rozšírili na ľubovoľnú hĺbku delenia  $d$ . Potom  $i$ -ty pruh  $\mathcal{P}_d(i)$  pre hĺbku  $d$  definujeme takto:

$$p_d(i) = (Q_\alpha, Q_{1-\alpha}), \alpha = i \times (0.5)^d$$

$$\mathcal{P}_d(i) = \{(t, p_d(i)) : t \in I\}$$

a množina všetkých pruhov grafu pre hĺbku delenia  $d$  je definovaná takto:

$$\{\mathcal{P}_d(i) : 0 \leq i \leq 2^{d-1}, i \in \mathbb{N}\}$$

Aby sme sa vyhli rekruzii, vypočítali sme si krok medzi susednými kvantilmi pri hĺbke  $d$ , ktorý je  $(0.5)^d$ , a jednotlivé pruhy sme generovali s týmto krokom. Pri rekruzívnom delení sa vygeneruje  $2^d$  kvantilov a z nich je možné vyrobiť  $2^{d-1}$  pruhov, preto sme index  $i$  obmedzili na  $i \leq 2^{d-1}$ . Z tohto vidíme, že počet pruhov grafu s rastúcou hĺbkou rastie exponenciálne, preto odporúčame, aby  $d$  bolo maximálne 4, kedy sa nám množina rozdelí na 16 častí 15 hexadecimmi a tak vznikne 8 pruhov.

Na obrázku ?? vidíme, že takto definovaný pruh sa potom vizualizuje, ako plocha medzi krivkami, ktoré tvoria dvojice hodnôt patriace danému pruhu, s výnimkou špeciálneho prípadu  $Q_{0.5}$ , ktorý vizualizujeme ako krivku.

## Funkčný krabicový diagram

Pre pochopenie dát je dôležité, aby sme sa vedeli pozrieť na hodnoty v ich kontexte. Všetky predošlé techniky uvažovali o chybe ako o samostatnej hodnote pre určitý predpovedný čas, avšak chyby sa nenachádzajú len v kontexte predpovedného času, ale aj v kontexte konkrétnej predpovede. Preto môžeme uvažovať o predpovediach ako o funkciách  $x_i(t)$ , kde  $i \in \{1..n\}$  je poradie predpovede a  $t \in I$  je čas predpovede, kde  $I$  je časový interval predpovede z  $\mathbb{R}$  (v našom prípade sa jednalo o dvojdnňový, teda 48 hodinový predpoveď).

Takýmto spôsobom sme sa dostali do novej situácie, kedy nechceme vizualizovať distribúciu jednotlivých chýb, ale celých predpovedí, ktoré chápeme ako funkcie. Na riešenie tohto problému existuje niekoľko spôsobov, z ktorých sme si zvolili *funkčný krabicový diagram* [SG11], keďže myšlienково vychádza z klasického krabicového diagramu, ktorý je jednak na túto situáciu vhodný ale je tiež medzi užívateľmi dobre známy a zaužívaný.

Ako sme spomenuli v časti 3.2.3, klasický krabicový diagram potrebuje na svoju konštrukciu 5 hodnôt: 3 kvartily a 2 extrémny. Aby sme tieto hodnoty našli pre funkcie, musíme ich vedieť porovnať a povedať, ktorá je “väčšia” alebo “menšia”. Autori funkčného krabicového diagramu riešia problém s využitím takzvanej pásmovej hĺbky (*band depth*) [LPR09]. Grafom  $G$  funkcie  $x$  je množina bodov  $G = \{(t, x(t)) : t \in I\}$ . Pásmo  $\mathcal{B}$  (*band*) v  $\mathbb{R}^2$  ohraničené krivkami  $x_{i_1}, x_{i_2}, \dots, x_{i_k}$ , kde  $k \geq 2$  je definované takto:

$$\mathcal{B}(x_{i_1}, x_{i_2}, \dots, x_{i_k}) = \{(t, y) : t \in I, \min_{r=1..k} x_{i_r}(t) \leq y \leq \max_{r=1..k} x_{i_r}(t)\}$$

Pásmo  $\mathcal{B}$  je teda množina všetkých bodov existujúcich medzi extrémami všetkých kriviek, ktoré doň vstupujú ako parameter. Pomocou týchto dvoch funkcií môžeme definovať pomocnú funkciu  $BD_n^{(j)}(x)$  pre krivku  $x$ , ktorá vyzerá takto:

$$BD_n^{(j)}(x) = \binom{n}{j}^{-1} \sum_{1 \leq i_1 < i_2 < \dots < i_j \leq n} \mathcal{I}\{G(x) \subset \mathcal{B}(x_{i_1}, \dots, x_{i_j})\}, j \geq 2$$

kde  $j$  je počet kriviek definujúce pásmo  $\mathcal{B}$ ,  $n$  je celkový počet kriviek a  $\mathcal{J}$  je takáto funkcia:

$$\mathcal{J}(x) = \begin{cases} 1 & \text{ak platí } x \\ 0 & \text{ak neplatí } x \end{cases}$$

Pomocná funkcia  $BD$  pre krivku  $x$  definuje pomer všetkých pásem zložených z  $j$  kriviek, v ktorých sa graf  $G(x)$  nachádza, ku všetkým možným  $j$ -ticiam kriviek vybraným z  $n$ . Samotná funkcia pásmovej hĺbky  $\mathcal{BD}$  pre krivku  $x$  je definovaná takto:

$$\mathcal{BD}_{n,J}(x) = \sum_{j=2}^J BD_n^{(j)}(x), J \geq 2$$

Hĺbka pásma  $\mathcal{BD}$  je teda suma všetkých  $BD$  pre počet kriviek 2 až  $J$ .

Autor článku definujúci pojem pásmová hĺbka navrhol taktiež flexibilnejšiu verziu s použitím pomocnej funkcie  $MBD$  (*modified band depth*) [LPR09]. V pravom rade je potrebné zdefinovať si funkciu  $A$ , ktorá určí všetky časové body, kedy sa krivka  $x$  nachádza v pásme  $B$ .

$$A(x, B) = \{t \in I : (t, x(t)) \in G(x) \wedge (t, x(t)) \in B\}$$

V spomínanom článku autori využívajú alternatívnu definíciu funkcie  $A$ , do ktorej vstupuje  $j + 1$  kriviek. Jej význam zostáva rovnaký ako pri našej definícii, avšak náš prístup považujeme za jednoduchší a zrozumiteľnejší. S využitím *Lebesgueovej miery*  $\lambda$  autori ďalej definujú funkciu  $\lambda_r$ , ktorá nám dáva "pomer času, ktorý krivka strávi v pásme":

$$\lambda_r(A) = \frac{\lambda(A)}{\lambda(I)}$$

Nová pomocná funkcia  $MBD$  je definovaná nasledovne:

$$MBD_n^{(j)}(x) = \binom{n}{j}^{-1} \sum_{1 \leq i_1 < i_2 < \dots < i_j \leq n} \lambda_r\{A(x, \mathcal{B}(x_{i_1}, \dots, x_{i_j}))\}, j \geq 2$$

Ak platí, že  $G(x) \subset \mathcal{B}(x_{i_1}, \dots, x_{i_j})$ , tak funkcia  $MBD$  sa degeneruje na  $BD$  [SG11].

V našej aplikácii sme sa rozhodli, že budeme pásmo definovať pomocou iba dvoch kriviek, čo nám vzorec výrazne zjednodušilo. Taktiež to implikovalo fakt, že pri výpočte  $MBD$  nie je potrebné  $\binom{n}{j}^{-1}$ , keďže berieme pásma zložené vždy z rovnakého počtu kriviek. Pre naše účely sme si taktiež zjednodušili funkciu  $\lambda_r$  na  $\lambda$ , keďže nepotrebujeme vlastnosť

tejto funkcie, ktorá dosahovala to, že  $MBD$  pre špeciálny prípad degeneruje na  $BD$ . Po týchto úpravách výsledný vzorec pre naše  $\mathcal{BD}'$  vyzerá nasledovne:

$$\mathcal{BD}'_n(x) = \sum_{1 \leq i_1 < i_2 \leq n} \lambda\{A(x, \mathcal{B}(x_{i_1}, x_{i_2}))\}$$

Teraz, keď sme úspešne definovali mieru, podľa ktorej môžeme usporiadať funkcie resp. ich krivky, je veľmi ľahké skonštruovať funkčný krabicový diagram.

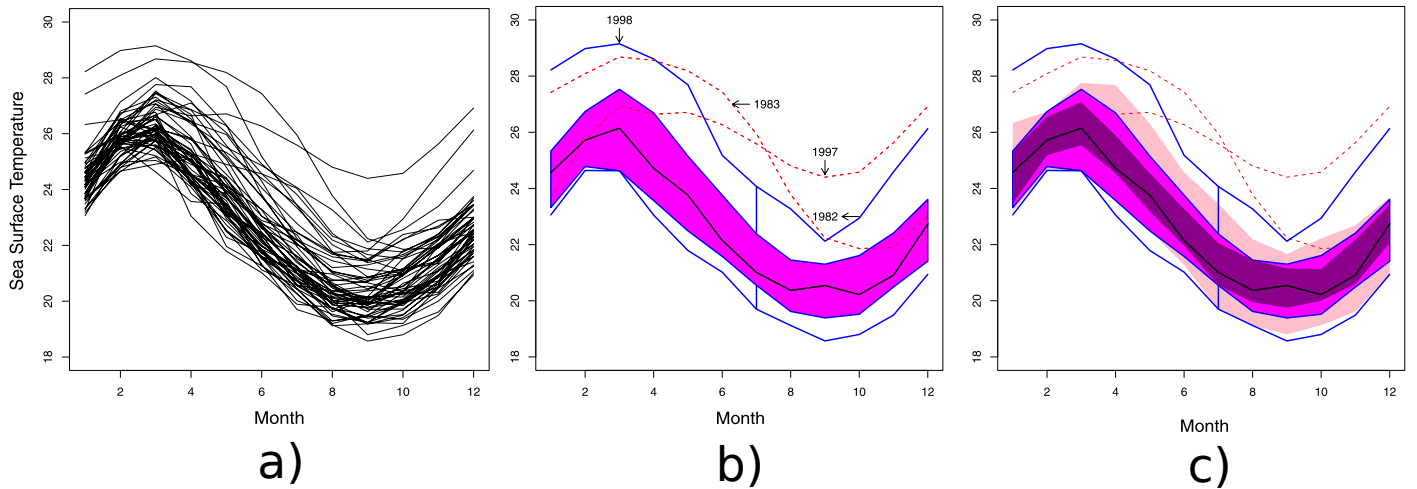
Nech sú naše funkcie predpovedí  $x_1, \dots, x_n$  usporiadané zostupne podľa  $\mathcal{BD}'$ . Potom krivka pre funkciu  $x_1$  má najvyššiu pásmovú hĺbku a predstavuje strednú hodnotu pre množinu funkcií, teda niečo podobné ako medián hodnôt pri krabicovom diagrame. Pri konštrukcii funkcionálneho diagramu nám taktiež pomôže koncept centrálného regiónu [LPS99]. Centrálny región  $C$  pre 50% kriviek je pásmo vytvorené z 50% najhlbších kriviek, teda:

$$C_{0.5} = B(x_1, x_2, \dots, x_{\lceil n/2 \rceil})$$

Vidíme, že Centrálny región  $C_{0.5}$  zaobahuje 50% najhlbších kriviek, a teda sa jedná o akúsi analógiu pre medzikvartálový rozsah (IQR) v klasickom krabicovom diagrame (pozri sekciu 3.2.3), ktorý ohraničoval 50% centrálnych dát. Zobrazením hraničných bodov tohto regiónu získame obálku myšlienkovito totožnú boxu v krabicovom diagrame (pozri obrázok 4.1b). Táto myšlienka sa dá použiť ďalej a môžeme, tak ako na obrázku 4.1c), zobraziť taktiež 25%-ný a 75%-ný centrálny región  $C_{0.25}$ ,  $C_{0.75}$ , avšak kvôli nižšej čitateľnosti grafu sme túto alternatívu nepoužili.

V prípade, že nepotrebujeme identifikovať *outlier*-ov, tak extrémálne hodnoty je už veľmi jednoduché získať, pretože ich tvorí pásmo zložené zo všetkých kriviek  $B(x_1, \dots, x_n)$ . V opačnom prípade musíme najprv identifikovať *outlier*-ov, ktorých potom vylúčime z výpočtov. Opäť sa využíva myšlienka z klasického krabicového diagramu, kedy sa *outlier* určil pomocou hodnoty  $c \times IQR$ , kde  $c$  bolo zvyčajne 1.5. Hranice sa teda získajú naškálovaním centrálného regiónu so škálovacím faktorom 1.5 a všetky krivky, ktoré sa v tomto regióne nenachádzajú, budú považované za *outlier*-ov. Test na *outlier*-a vyzerá teda takto:

$$y_i(t) = 1.5 \times x_i(t)$$



Obr. 4.1: Obrázky z článku *Functional Boxplots* [SG11] a) Funkcie meraní teploty hladiny mora b) Funkčný krabicový diagram c) Rozšírený Funkčný krabicový diagram o centrálné regióny  $C_{0.25}$  a  $C_{0.75}$

$$isOutlier(x) = [G(x) \not\subseteq B(y_1, \dots, y_{\lceil n/2 \rceil})]$$

Na obrázku 4.1b) môžeme vidieť červené prerušované čiary, ktorými sú *outlier*-i znázornené.

V našej práci sme do diagramu pridali ešte jednu krivku pre ľubovoľnú štatistiku vypočítanú z chýb ako napríklad MFE, MAE alebo RMSE (pozri sekciu 2.3). Takto môžeme spraviť porovnanie distribúcie chýb s vypočítanou štatistikou.

## Porovnanie metód

Tu bude pekný obrázok a obkeci.

Graf hustoty je presnejši, ale menej citateľny.

Pruhový kvantilový diagram nahradzuje krabicový diagram.

Funkcionalny krabicovy riesi distribuciu funkcii.



### 4.2.6 Návrh farebnej palety

- rainbow a preco zle
- monochrome cez 3 farby
- monochrome 2 farby ostro rozdelené
- ekvalizácia histogramom / boxplotom
- návrh farby pre density, púhový, functional boxplot

# Kapitola 5

## Implementácia

### 5.1 Použité technológie

#### 5.1.1 Java

#### 5.1.2 JavaScript

#### 5.1.3 d3.js

## Kapitola 6

### Výsledky

# Literatúra

- [Car94] Daniel B. Carr. A Colorful Variation On Box Plots. *Statistical Computing & Statistical Graphics Newsletter*, 5(3):19–23, December 1994.
- [Cas09] Barbara Casati. *Verification of continuous predictands*. Joint Working Group on Forecast Verification Research (JWGFVR), Jún 2009.
- [CM05] Chamnein Choonpradub and Don McNeil. Can the box plot be improved? *Songklanakarin Journal of Science and Technology*, 27(3):649–657, 2005.
- [CWS<sup>+</sup>08] B. Casati, L. J. Wilson, D. B. Stephenson, P. Nurmi, A. Ghelli, M. Pocerlich, U. Damrath, E. E. Ebert, B. G. Brown, and S. Mason. Forecast verification: current status and future directions. *Meteorological Applications*, 15(1):3–18, 2008.
- [Fis10] Wolfram Fischer. *Neue Grafiken zur Datenvisualisierung*. Z I M – Zentrum für Informatik und wirtschaftliche Medizin., 2010.
- [FJB12] Tressa L. Fowler, Tara L. Jensen, and Barbara G. Brown. *Introduction to Forecast Verification*. 2012.
- [Gol] Professor Brian Golding. Weather forecasting part 1. <http://www.rmets.org/weather-and-climate/weather/weather-forecasting>. [Prístupné online: 6.12.2014].
- [HN98] Jerry L. Hintze and Ray D. Nelson. Violin plots: A box plot-density trace synergism. *The American Statistician*, 52(2):181–184, Máj 1998.

- [LE11] Dr. Arlene Laing and Dr. Jenni-Louise Evans. *Introduction to Tropical Meteorology 2nd Edition*. UCAR, Október 2011.
- [LPR09] Sara López-Pintado and Juan Romo. On the concept of depth for functional data. *Journal of the American Statistical Association*, 104(486):718–734, 2009.
- [LPS99] Regina Y. Liu, Jesse M. Parelius, and Kesar Singh. Multivariate analysis by data depth: descriptive statistics, graphics and inference, (with discussion and a rejoinder by Liu and Singh). *Ann. Statist.*, 27(3):783–858, 06 1999.
- [Lyn07] Peter Lynch. The origins of computer weather prediction and climate modeling. *Journal of Computational Physics* 227 (2008) 3431–3444, Febrúar 2007.
- [Mic13] Microsoft. ??? visual basic for applications. <https://>, Január 2013. [Prístupné online: 19.1. 2013].
- [Mic15] Microsoft. Microsoft Excel. <https://products.office.com/en-us/excel>, 2015. [Prístupné online: 18.3.2015].
- [MTL78] Robert McGill, John W. Tukey, and Wayne A. Larsen. Variations of box plots. *The American Statistician*, 32(1):12–16, Febrúar 1978.
- [NCE14] NCEP. Inventory of Data Products on the NOAA Servers. <http://www.nco.ncep.noaa.gov/pmb/products/>, November 2014. [Prístupné online: 10.11.2014].
- [Nur03] Pertti Nurmi. *Recommendations on the verification of local weather forecasts*. European Centre for Medium Range Weather Forecasts, Decmeber 2003.

- [Ope13] Apache OpenOffice. Openoffice.org basic programming guide. [https://wiki.openoffice.org/wiki/Documentation/BASIC\\_Guide](https://wiki.openoffice.org/wiki/Documentation/BASIC_Guide), Január 2013. [Přístupné online: 19.1. 2013].
- [Ope15] Apache OpenOffice. OpenOffice.org Calculate. <https://www.openoffice.org/product/calc.html>, 2015. [Přístupné online: 18.3.2015].
- [PKR07] Kristin Potter, Joe Kniss, and Richard Riesenfeld. Visual summary statistics. Technical Report UUCS-07-004, University of Utah, 2007.
- [Poc11] Matthew Pocernich. *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, chapter Appendix - Verification Software, pages 232–240. John Wiley & Sons, Ltd., 2nd edition, December 2011.
- [Pot06] Kristin Potter. Methods for presenting statistical information: The box plot. In Hans Hagen, Andreas Kerren, and Peter Dannenmann, editors, *Visualization of Large and Unstructured Data Sets*, volume S-4 of *GI-Edition Lecture Notes in Informatics (LNI)*, pages 97–106. 2006.
- [RRT99] Peter J. Rousseeuw, Ida Ruts, and John W. Tukey. The bagplot: A bivariate boxplot. *The American Statistician*, 53(4):382–287, November 1999.
- [SB91] William Stock and John Behrens. Box, Line, and Midgap Plots: Effects of Display Characteristics on the Accuracy and Bias of Estimates of Whisker Length. *Journal of Educational Statistics*, 16(1):1–20, 1991.
- [SG11] Ying Sun and Marc G. Genton. Functional Boxplots. *Journal of Computational and Graphical Statistics*, 20(2):316–334, Jún 2011.
- [SOA04] Neil C. Schwertman, Margaret Ann Owens, and Robiah Adnan. A simple more general boxplot method for identifying outliers. *Computational Statistics & Data Analysis*, 47(1):165–174, 2004.

- [Spe52] Mary Eleanor Spear. *Charting Statistics*. McGraw-Hill Book Company, Inc., 1952.
- [TM15] Inc. The MathWorks. Matlab the language of technical computing. <http://www.mathworks.com/products/matlab/>, Marec 2015. [Prístupné online: 18.3. 2015].
- [Ton05] Phattrawan Tongkumchum. Two-dimensional box plot. *Songklanakarin Journal of Science and Technology*, 27(4):859–866, 2005.
- [Tuf83] Edward R. Tufte. *The Visual Display of Quantitative Information*. Graphics Press, Cheshire, Connecticut, 1983.
- [Tuk77] John W. Tukey. *Exploratory Data Analysis*. Addison-Wesley, 1977.
- [Vas98] Tim Vasquez. *Observer Handbook*. International Weather Watchers, 1995, 1998.
- [VWvH<sup>+</sup>07] Fernanda B. Viegas, Martin Wattenberg, Frank van Ham, Jesse Kriss, and Matt McKeon. Manyeyes: A site for visualization at internet scale. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1121–1128, November 2007.
- [WCSDG<sup>+</sup>08] Joseph B. Klemp William C. Skamarock, Jimy Dudhia, David O. Gill, Dale M. Barker, Michael G. Duda, Xiang-Yu Huang, Wei Wang, and Jordan G. Powers. *A Description of the Advanced Research WRF Version 3*. National Center for Atmospheric Research, Jún 2008.
- [Wei14] Eric W. Weisstein. Statistical Median. <http://mathworld.wolfram.com/StatisticalMedian.html>, 2014.
- [WMO94] WMO. *A GUIDE TO THE CODE FORM FM 92-IX Ext. GRIB Edition 1*. WMO, Máj 1994.
- [WMO03] WMO. *Introduction to GRIB Edition1 and GRIB Edition 2*. WMO, Jún 2003.

- [WS12] Hadley Wickham and Lisa Stryjewski. 40 years of boxplots. Technical report, had.co.nz, 2012.
- [Yoa88] Yoav Benjamini. Opening the box of a boxplot. *The American Statistician*, 42(4):257–262, November 1988.