

PŘEDPOVĚĎ PRAVDĚPODOBNOСТИ VÝSKYTU METEOROLOGICKÝCH JEVŮ

Martin Dubrovský

Ústav fyziky atmosféry ČSAV, Hradec Králové¹

1 Úvod

V rámci grantového projektu ČSAV je na našem pracovišti vyvíjen PC systém MESOMAP, jehož součástí je velmi krátkodobá předpověď ($\Delta t \sim 12$ h) pravděpodobnosti výskytu 'významných' meteorologických jevů. Tvar prognostické funkce pro vybraný jev,

$$(1.1) \quad f(\cdot) = \Pr(Y=1|X=\cdot)$$

kde X je r -rozměrný prediktor (vysvětlující proměnná) a Y je binární prediktand (vysvětlovaná proměnná), je hledán pomocí analýzy učebního datového souboru, $T = \{[x_i, y_i]; i=1, \dots, N\}$. Práce na vývoji prognostických procedur se v zásadě ubírá dvěma směry:

- ▶ výběr vhodných prediktorů nesoucích co největší informaci o prediktandu,
- ▶ výběr vhodného typu prognostické funkce a způsobu nalezení jejích parametrů.

Předmětem tohoto příspěvku jsou dvě metody odhadu pravděpodobnostní prognostické funkce (1.1) a porovnání jejich úspěšnosti při aplikaci na předpověď pravděpodobnosti výskytu bouřky. První metodou je konstrukce prognostické funkce ve tvaru binárního rozhodovacího stromu (BRS). Zde prezentovaná verze konstrukčního algoritmu vychází z práce ŘEZÁČOVÉ a MOTLA (1990). Algoritmus, který mezitím doznal některých změn, je podrobně popsán v připravované publikaci (DUBROVSKÝ, 1997). Druhou metodou je konstrukce prognostické funkce metodou k -nejbližších sousedů (KNS). Vodítkem při konkretizaci algoritmu mi byl vztah (2.1) v práci ANTOCHA (1988) a vlastní fantazie.

V této práci se předpokládá, že jednotlivé složky prediktorového vektoru $X = (X_1, X_2, \dots, X_r)$ jsou měřitelné či (v případě rozhodovacího stromu) dichotomické veličiny.

2 Binární rozhodovací strom

Na Obr. 1 je jednoduchá verze binárního rozhodovacího stromu, který vstupnímu prediktorovému vektoru x přiřazuje pravděpodobnost výskytu jevu reprezentovaného binárním prediktandem. Každý průchozí uzel (bifurkaci) stromu lze chápat jako kořen podstromu a zapsat: $t_k = \{b_k(\cdot), t_{k<}, t_{k>}\}$, kde $b_k(\cdot)$ je bifurkační funkce a $t_{k<}, t_{k>}$ jsou dva dceřinné podstromy. Koncové uzly stromu (listy) nesou informaci o odhadu podmíněné pravděpodobnosti výskytu jevu, \hat{P}_k . Stavba BRS začíná nalezením bifurkační funkce kořenového uzlu, t_1 , s využitím datového souboru $T_1 = T$, a pokračuje rekurzivní bifurkací dceřinných podstromů dokud všechny větve stromu nejsou ukončeny listy.

Algoritmus konstrukce podstromu t_k z učebního souboru T_k :

1.krok: Nalezení bifurkační funkce $b_k(\cdot)$, jež co nejúčinněji (ve smyslu zavedené míry diskriminace) odděluje prvky

T_k s výskytem jevu ($Y=1$) od prvků bez výskytu jevu ($Y=0$).

2.krok: Test významnosti diskriminace.

3.krok (a): Je-li diskriminace významná, uzel t_k se stává bifurkací s bifurkační funkcí $b_k(\cdot)$, množina T_k je rozdělena

na $T_{k<} = \{[x_i, y_i] \in T_k \mid b_k(x_i) < 0\}$ a $T_{k>} = \{[x_i, y_i] \in T_k \mid b_k(x_i) \geq 0\}$. Z učebních souborů $T_{k<}$ a $T_{k>}$ jsou zkonstruovány dceřinné podstromy $t_{k<}$ a $t_{k>}$.

¹ Husova 456, 50008 Hradec Králové

(b) Je-li diskriminace nevýznamná, uzel t_k se stává listem a je mu přiřazena prognostická pravděpodobnost výskytu jevu:

$$(2.1) \quad \hat{p}_k = n_k / N_k$$

kde N_k a n_k je počet prvků v T_k celkem, resp. s $Y=1$.

V tomto obecném schématu konstrukčního algoritmu je třeba specifikovat:

- ▶ Míru kvality prognostické funkce a s ní související míru kvality bifurkační funkce
- ▶ Metodu štěpení stromu (množinu bifurkačních funkcí)
- ▶ Kritérium pro uzavření větve (test významnosti diskriminace)

2.1 Míra kvality prognostické funkce

Je-li ztráta způsobená nepřesnou předpovědí popsána kvadratickou ztrátovou funkcí, $l(y, f(x)) = (y_i - f(x_i))^2$, kvalitu prognostické funkce $f(\cdot)$ lze vyjádřit střední kvadratickou chybou (MSE):

$$(2.2) \quad e_f^2 = E[(y - f(x))^2]$$

Názornějším vyjádřením kvality prognostické funkce může být relativní skóre úspěšnosti:

$$(2.3) \quad RV(f) = 1 - e_f^2 / e_0^2$$

kde e_0^2 je střední kvadratická chyba průměrné klimatické předpovědi [$f_0(x) = n(T)/N(T) = \text{konst}(x)$] a RV se nazývá redukce variance.

Kvalitu bifurkační funkce lze definovat jako pokles MSE vzhledem k bifurkaci podle $b(\cdot)$:

$$(2.4) \quad m(k|b) = p_k e_k^2 - (p_{k<} e_{k<}^2 + p_{k>} e_{k>}^2)$$

kde p_k ($\cdot = k, k<, k>$) je pravděpodobnost, že vektor X 'padne' do uzlu t_k , a e_k^2 je MSE daného uzlu:

$$(2.5) \quad e_k^2 = P_{\cdot} (1 - \hat{P}_{\cdot})^2 + (1 - P_{\cdot}) \hat{P}_{\cdot}^2$$

kde P_{\cdot} je podmíněná pravděpodobnost výskytu jevu pro daný uzel.

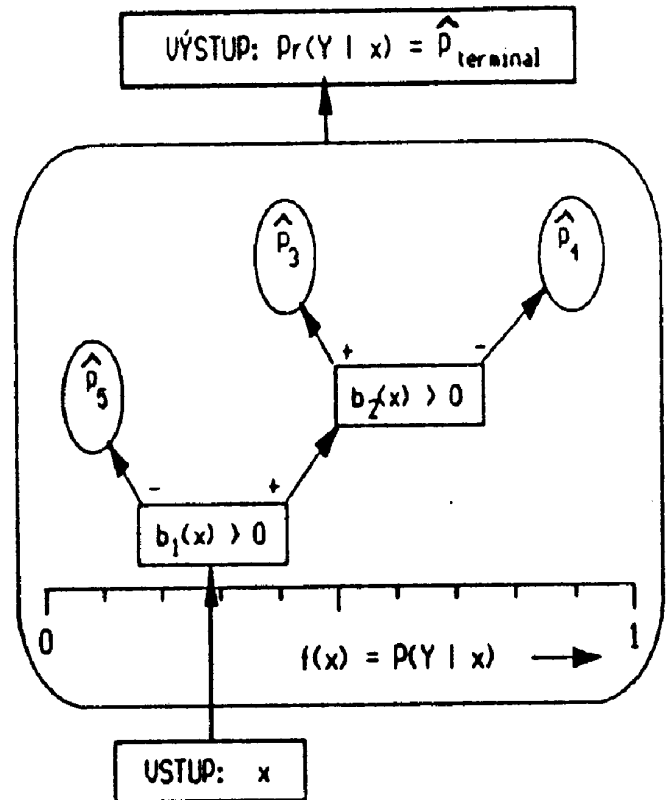
Při odhadu výše uvedených veličin z učebního souboru lze vycházet z odhadu MSE pro daný uzel:

$$(2.6) \quad \hat{e}_k^2 = [n_{\cdot} (1 - \hat{P}_{\cdot})^2 + (N_{\cdot} - n_{\cdot}) \hat{P}_{\cdot}^2] / N_{\cdot}$$

Odhady příslušných veličin budou dále značeny střísčkou. Je-li příslušná veličina odhadnuta na testovacím souboru, nezávislém na vývojovém souboru (to je ten, na kterém byla nalezena bifurkační funkce a byly odhadnuty prognostické pravděpodobnosti podle 2.1), bude odlišena čarou (\hat{e}^2 , $R\hat{P}$, \hat{m} , N' , n').

2.2 Štěpení stromu

Z množiny 'povolených' bifurkačních funkcí je v každém uzlu stromu hledána ta, která zajišťuje co největší míru diskriminace.



Obr. 1 Schéma binárního rozhodovacího stromu pro předpověď pravděpodobnosti výskytu jevu. b_k ($k=1,2$) jsou bifurkační funkce průchozích uzlů (bifurkací), \hat{P}_k ($k=3,4,5$) jsou podmíněné pravděpodobnosti výskytu jevu příslušné koncovým uzlům (listům) rozhodovacího stromu.

2.2.1 Typ bifurkační funkce

Standardní množinou bifurkačních funkcí je množina funkcí typu

$$(2.7) \quad b(x) = x_{ix} - x^*$$

kde ix je index diskriminačního prediktoru (ix -tá složka prediktorového vektoru) a x^* je jeho kritická hodnota. V každém uzlu je bifurkační funkce volena tak, aby:

$$(2.8) \quad b_k = \arg \max_{i=1, \dots, r} [\hat{m}(k | \bar{b}_i)] \quad \text{kde } \bar{b}_i = x_i - \bar{x}_i^*$$

Kritické hodnoty \bar{x}_i^* jsou nalezeny jednou ze dvou metod popsanych v §2.2.2.

Je-li prediktorový vektor dobře vychovaný (v ideálním případě s normálním rozdělením) a ne příliš rozměrný lze hledat bifurkační funkci ve tvaru

$$(2.9) \quad b_k(x) = \hat{w}_k x - x_k^*$$

kde \hat{w}_k je diskriminační vektor (výpočet viz. např. Pudil et al. (1991), rovnice 7.2-5). Kritická hodnota x_k^* je vypočtena jako kritická hodnota jednorozměrného prediktoru $z = \hat{w}_k x$ stejným způsobem jako v případě bifurkační funkce jedné proměnné. Aplikace bifurkační funkce typu (2.9) umožňuje rychlejší diskriminaci prediktorového prostoru a ve svém důsledku i vyšší kvalitu odhadu prognostické funkce. Nevýhody vzhledem k bifurkační funkci typu (2.7) jsou: (1) nutný předpoklad o 'vychovanosti' prediktorů, (2) neumožňuje redukci dimenze prediktorového vektoru.

2.2.2 Kritická hodnota bifurkační funkce

Nechť X_j je diskriminační prediktor (j -tá složka prediktorového vektoru). Jeho kritická hodnota, x_j^* , je vypočtena podle algoritmu:

Soubor T_k je náhodně rozdělen na dva stejně velké (případně lišící se o jeden prvek) podmnožiny T_k^a and T_k^b . Na každé z těchto podmnožin je nalezena kritická hodnota x_a^* , x_b^* . Toto se třikrát opakuje a výslednou kritickou hodnotou je aritmetický průměr z takto získaných šesti hodnot. Kritická hodnota prediktoru X_j na libovolné množině T_k^* ($\cdot = a, b$) je určena jednou ze dvou metod:

momentová metoda:

$$(2.10) \quad x^* = [\bar{x}_{j_0} s_{j_1} + \bar{x}_{j_1} s_{j_0}] / [s_{j_1} + s_{j_0}]; \quad (\cdot = a, b)$$

kde \bar{x}_{j_i} a s_{j_i} ($i=0,1$) je průměr a směrodatná odchylka diskriminačního prediktoru X_j vypočtená z těch prvků T_k^* pro něž je $y=i$.

maximální metoda je založena na nalezení hodnoty prediktoru, pro níž míra diskriminace nabývá maxima:

$$(2.11) \quad x^* = \arg \max_{x^* \in R^1} [\hat{m}(k | b(x) = x_j - x^*)]; \quad (\cdot = a, b)$$

Účinnost momentové i maximální metody v závislosti na velikosti učebního souboru, apriorní pravděpodobnosti výskytu jevu a tvaru distribuční funkce diskriminačního prediktoru byla testována jednak metodou Monte Carlo na náhodně generovaných souborech, jednak na reálných datech (viz. DUBROVSKÝ, 1997). Zde jen uvedu, že v případě dobře vychovaného diskriminačního prediktoru, a dobře (resp. špatně) separovatelného souboru ((ve souboru diskriminační míry zavedené vztahem (2.4), s použitím bifurkační funkce typu (2.7)) je výhodnější použití momentové (resp. maximální) metody. V případě špatně vychovaných prediktorů je maximální metoda výhodnější vzhledem ke své robustnosti.

2.3 Kritérium pro uzavření větve

Je-li kvalita prognostické funkce odvozena od kvadratické ztrátové funkce, existuje nebezpečí, že štěpením stromu

dojde k poklesu místo k růstu RV. Proto je třeba před každým štěpením stromu testovat významnost diskriminace. Testovací statistikou může být míra diskriminace vypočtená buď na vývojovém nebo nezávislém testovacím souboru, který je štěpen paralelně s vývojovým souborem. V případě statistiky vypočtené na vývojovém souboru je kritériem pro akceptování bifurkace:

$$(2.12) \quad \hat{m}(k|b_k) \geq m^*(N_k, n_k, r)$$

přičemž síla testu klesá s rostoucím počtem prediktorů (r), z nichž je diskriminační prediktor vybrán. V případě statistiky vypočtené na nezávislém datovém souboru je síla testu nezávislá na r . Vzhledem k tomu, že výběrová míra diskriminace vypočtená na nezávislém testovacím souboru je nestranným odhadem její skutečné hodnoty, zdá se, že vhodným kritériem minimalizujícím pravděpodobnost špatného rozhodnutí (přijetí bifurkace když $m(k|b_k) < 0$ nebo zamítnutí bifurkace když $m(k|b_k) > 0$) je splnění podmínky:

$$(2.13) \quad \hat{m}'(k|b_k) \geq 0$$

Důsledek zanedbání testu významnosti diskriminace

je demonstrován na Obr. 2, kde je znázorněn vývoj redukce variance v závislosti na postupném, nekontrolovaném štěpení stromu. Z obrázku je patrná jednak rostoucí divergence větví $R\hat{V}$ a $R\hat{V}'$ s klesajícím N , jednak existence optimální velikosti stromu, která roste (stejně jako maximum hodnoty $R\hat{V}'$) s rostoucí velikostí vývojového souboru, přičemž překročení optimální velikosti stromu je doprovázeno poklesem kvality stromu ('Bugaboo' efekt).

3 Metoda k nejbližších sousedů

Odhad prognostické funkce metodou k nejbližších sousedů je váženým aritmetickým průměrem z těch pozorování, pro něž odpovídající x_j je mezi k 'nejbližšími' sousedy bodu, v němž odhadujeme (Antoch, 1988):

$$(3.1) \quad f(x) = \sum_{j=1, \dots, k} q_{ix_j} y_{ix_j} / \sum_{j=1, \dots, k} q_{ix_j}$$

kde q_{ix_j} je váha ix_j -tého vektoru, ix_j je index j -tého nejbližšího souseda vektoru x vybraného z učebního souboru $\{(x_i, y_i); i=1, \dots, N\}$ tak, že $\delta(x, x_{ix_{j-1}}) \leq \delta(x, x_{ix_j}) \leq \delta(x, x_{ix_{j+1}})$. Definice vzdálenosti $\delta(x, x_{ix_j})$, definice váhy q_{ix_j} a počet sousedů k je předmětem optimalizace, která probíhá opět s využitím metody Monte Carlo. Nejchoulostivějším prvkem se zdá být definice vzdálenosti. Jako nejvýhodnější se ukazují dva typy vzdáleností:

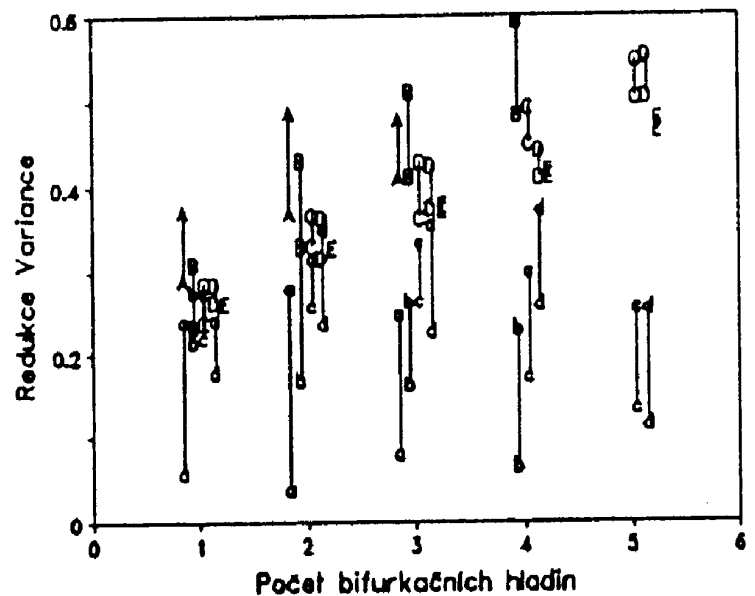
(a) Mahalanobisova vzdálenost:

$$(3.2) \quad \delta(x, x_j) = (x - x_j)^t \Sigma^{-1} (x - x_j)$$

kde Σ je kovarianční matice vektoru X . Pokus o zohlednění vzájemné vazby jednotlivých složek vektoru X s prediktandem vedl k modifikaci předešlé definice:

$$(3.3) \quad \delta(x, x_j) = (x_k - x_{jk}) c_k \Sigma_{kl}^{-1} c_l (x_l - x_{jl}) ; (k, l \text{ jsou sčítací indexy})$$

kde c_l ($l = k, l$) je váha vyjadřující vazbu mezi prediktandem a l -tou složkou vektoru X . Testování různě definovaných



Obr. 2 Závislost redukce variance na výšce stromu (počet bifurkačních hladin). Malými [velkými] písmeny jsou vyznačeny intervaly $\langle \text{avg}(R\hat{V}) \pm s(R\hat{V}) \rangle$ [$\langle \text{avg}(R\hat{V}') \pm s(R\hat{V}') \rangle$], kde $\text{avg}(\cdot)$, $s(\cdot)$ jsou průměr a směrodatná odchylka z deseti simulací. $[N, N'] = [100, 500]$ (intervaly A–A, a–a); $[200, 500]$ (B, b); $[500, 500]$ (C, c); $[700, 300]$ (D, d); $[1000, -]$ (E, –). Prediktand = $C_{<12-24>}$.

vah c však zatím nevedlo k jednoznačnému zlepšení.

(b) Vzdálenost projekcí dvou bodů na nejlepší diskriminační přímku:

$$(3.4) \quad \delta(x_i, x_j) = |\hat{w} \cdot (x_i - x_j)|$$

kde \hat{w} je diskriminační vektor určen stejným způsobem jako pro (2.9).

Pro volbu váhy q se ukázaly jako přibližně stejně dobré: (a) $q_{ix_j} = 1$; (b) $q_{ix_j} = [1/\delta(x_i, x_j)]^{1/2}$; (c) $q_{ix_j} = (1/j)^{1/2}$

Ke zrobusnění metody lze složky vektoru x_j transformovat dvěma způsoby: (a) $x_{ij} \rightarrow R_j(x_j)$; (b) $x_{ij} \rightarrow F^{-1}[R_j(x_j)/(N+1)]$, kde $R_j(x_j)$ je pořadí vektoru x_j v posloupnosti tvořené prvky učebního souboru uspořádanými podle j -té složky vektoru, N je počet individuí učebního souboru a F je distribuční funkce normálního rozdělení.

4 Aplikace na reálná data

Učební datový soubor sestává z 1058 pozorování čtyřrozměrného prediktorového vektoru $[X = (ADED2, KMOD, SICP, FI)]$, kde jednotlivé složky vektoru jsou indexy instability odvozené z aerologické sondáže Praha - Libuš a čtyř různě definovaných binárních prediktandů (A, B, C_{<12-18>}, C_{<12-24>}) charakterizujících výskyt bouřky na daném území, v daném časovém intervalu. Pro každý ze čtyř prediktandů byla hledána prognostická funkce (1.1) jednak ve tvaru binárního stromu s bifurkační funkcí typu buď (2.7) nebo (2.9), jednak metodou k nejbližších sousedů se vzdáleností definovanou vztahem buď (3.2) nebo (3.4). Prediktorem je buď kompletní vektor $X = (ADED2, KMOD, SICP, FI)$, nebo samotný Faustův index (FI) — index s nejvyšší individuální prediktivní schopností. Mírou úspěšnosti je redukce variance vypočtená na nezávislém datovém souboru, $R \hat{V}'$. Pro každou kombinaci vstupních podmínek bylo provedeno 10 nezávislých simulací, z nichž byla spočtena průměrná hodnota a směrodatná odchylka $R \hat{V}'$. Variabilita hodnoty $R \hat{V}'$ je dána náhodným způsobem rozdělení vstupních dat na vývojový a testovací soubor, v případě konstrukce stromu má vliv též neurčitost při výpočtu diskriminační hodnoty.

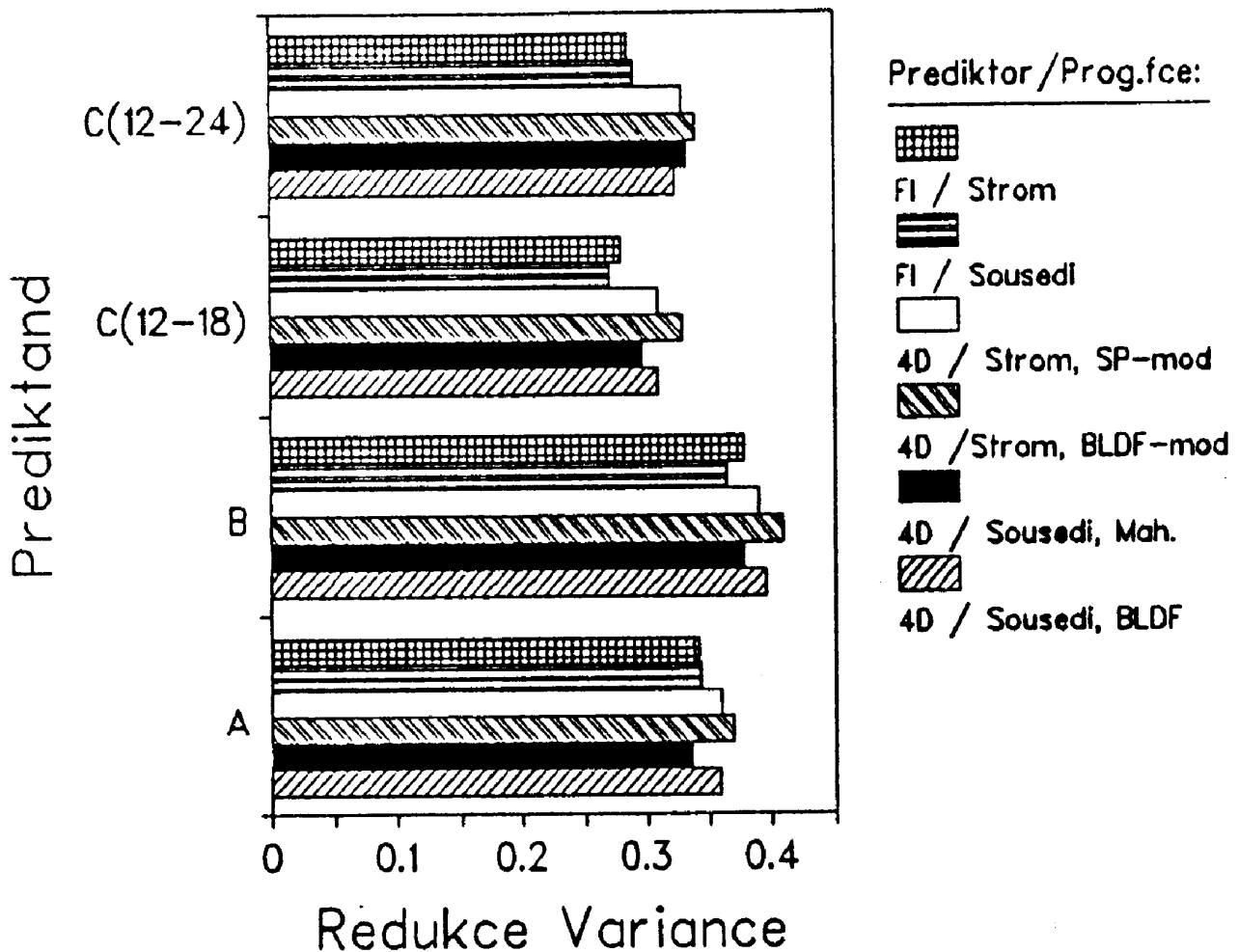
Tabulka. Výsledky 10 konstrukcí binárního rozhodovacího stromu. $R \hat{V}'$ je redukce variance vypočtená na testovacím souboru ($N' = N = 529$), L je počet listů stromu. $avg(\cdot)$ a $s(\cdot)$ je aritmetický průměr a směrodatná odchylka vypočtené z 10 hodnot.

odhad diskrim. hodnoty: prediktor/bif.fce ¹⁾ :		momentová metoda			maximální metoda		
		FI/sp	4D/sp	4D/bldf	FI/sp	4D/sp	4D/bldf
Y= A	avg($R \hat{V}'$)=	0.34	0.36	0.37	0.35	0.35	0.38
	s($R \hat{V}'$)=	0.03	0.04	0.03	0.03	0.03	0.02
	avg(L)=	12	12	9	11	10	8
	s(L)=	3	3	2	2	3	2
Y= B	avg($R \hat{V}'$)=	0.38	0.39	0.41	0.38	0.38	0.40
	s($R \hat{V}'$)=	0.03	0.05	0.04	0.02	0.02	0.03
	avg(L)=	9	9	9	10	8	8
	s(L)=	2	3	2	2	2	3
Y= C<12-18>	avg($R \hat{V}'$)=	0.28	0.31	0.33	0.29	0.30	0.31
	s($R \hat{V}'$)=	0.04	0.04	0.03	0.03	0.05	0.03
	avg(L)=	8	9	7	9	9	7
	s(L)=	3	2	2	2	2	2
Y= C<12-24>	avg($R \hat{V}'$)=	0.29	0.33	0.34	0.31	0.34	0.33
	s($R \hat{V}'$)=	0.03	0.04	0.04	0.03	0.04	0.04
	avg(L)=	7	9	7	8	8	8
	s(L)=	2	3	2	2	1	3

1) FI/sp: $x = FI$, bifurkační funkce typu (2.7);
4D/sp: $X = (ADED2, KMOD, SICP, FI)$, bifurkační funkce typu (2.7); 4D/bldf: $x = (ADED2, KMOD, SICP, FI)$, bifurkační funkce typu (2.9).

V Tabulce jsou výsledky simulací stavby stromu, z nichž vyplývá:

- úspěšnost stromů využívajících vícerozměrný prediktor je vyšší než při použití jediného indexu (FI). Zlepšení je ovšem malé, což lze přičíst na vrub značné vzájemné korelovanosti jednotlivých indexů instability.
- použití bifurkačních funkcí ve tvaru nejlepší lineární diskriminační funkce (2.9) dává nevýrazně (vzhledem k variabilitě



Obr. 3 Skóre úspěšnosti [Redukce Variance = $\text{avg}(R \hat{V}')$, kde $R \hat{V}'$ je redukce variance vypočtená na testovacím souboru, $\text{avg}(\cdot)$ značí průměr z 10 simulací] Binárního rozhodovacího stromu (Strom) a metody k nejbližších sousedů (Sousedí) aplikované na 4-rozměrný prediktor 4D = (ADED2, KMOD, SICP, FI) a samotný Faustův index (FI). Strom je staven buď v sp-módu (bifurkační funkce typu (2.7)) nebo *blďf*-módu (bifurkační funkce typu (2.9)), odhad diskriminační hodnoty je momentovou metodou. Odhad metodou kNS: $k=50$, $q_i = (1/\delta_i)^{\frac{1}{2}}$, vzdálenost δ_i je buď podle (3.2) (Mah.) nebo podle vztahu (3.4) (*blďf*).

hodnot $R \hat{V}'$) vyšší skóre úspěšnosti než při použití bifurkačních funkcí jedné proměnné (2.7). Toto lze opět vysvětlit značnou vzájemnou závislostí jednotlivých prediktorů. Za povšimnutí stojí nižší počet bifurkací u stromů využívajících bifurkační funkce typu *blďf*, což je v souladu s očekávaným rychlejším vyčerpáním informace obsažené ve vícerozměrném prediktoru.

Na Obr. 3 je porovnána úspěšnost odhadu prognostické funkce metodou stromu a k nejbližších sousedů. Z obrázku lze vyčíst:

- Úspěšnost obou metod aplikovaných na 4-rozměrný prediktor je vyšší než při použití nejsilnějšího jednorozměrného prediktoru, nicméně toto zvýšení je vzhledem ke značné vzájemné závislosti mezi jednotlivými indexy nevýrazné.
- Úspěšnost obou metod (BRS vs. kNS) se v rámci chyby neliší
- Využitím nejlepší lineární diskriminační funkce (*blďf*) lze u obou metod (bifurkační funkce typu (2.9) při konstrukci stromu; definice vzdálenosti vztahem (3.4) u metody kNS) sledovat určité zlepšení, což však může být důsledek toho, že použité prediktory jsou poměrně hodné (s trochou dobré vůle s normálním rozdělením).

5 Závěr

Kvalita prognostické funkce nalezené metodou analýzy učebního souboru je dána jednak kvalitou prediktorů (tedy mírou informace o prediktandu nesené prediktorem), jednak schopností použité statistické metody extrahovat z učebního souboru co nejoptimálnější prognostickou funkci. Vzhledem k tomu, že úspěšnost obou statistických metod testovaných v této práci vychází přibližně stejně, pro konečné rozhodnutí o výběru metody mohou hrát roli následující 'druhotná' kritéria:

- ▶ **rychlost:** konstrukce prognostické funkce ve tvaru stromu zabírá určitý časový interval, použití hotového stromu na konkrétní prediktorový vektor je však velice rychlé. Naproti tomu, v případě odhadu metodou k nejbližším sousedům je při každé aplikaci prognostické funkce vyhledání sousedů časově náročné.
- ▶ **náročnost na paměť počítače:** Odhad prognostické funkce metodou sousedů vyžaduje při každé aplikaci přístup k celému datovému souboru, kdežto prognostickou funkci ve tvaru stromu lze v paměti počítače uschovat mnohem úsporněji.
- ▶ **sebezdokonalování:** chceme-li, aby se prognostická funkce postupně zdokonaľovala (tak jak se rozrůstá učební soubor), v případě kNS je toho dosaženo prostým přidáním nového pozorování k učebnímu souboru. V případě BRS je třeba při rozšíření učebního souboru konstruovat nový strom.
- ▶ **metoda kNS se lépe vyrovnává s rostoucí dimenzí:** v práci připravované pro publikaci (Dubrovský, 199?) je pro konstrukci stromu použit 50-ti rozměrný prediktorový vektor, s čímž se algoritmus konstrukce stromu s využitím bifurkačních funkcí typu (2.7) s přehledem vypořádal. Navíc algoritmus konstrukce BRS v sobě zahrnuje mechanismus pro redukci dimenze prediktorového vektoru. Pro kNS nastávají při rostoucí dimenzi prediktoru značné potíže vzhledem k nutnosti pracovat s kovarianční maticí - v případě obou použitých definic vzdáleností.
- ▶ **metoda kNS se lépe vyrovnává s kategoriálními prediktory**

Závěrem bych chtěl dodat, že kombinace prediktorů použitá v této práci byla vybrána tak trochu subjektivně na základě vyhodnocení frekvence výskytu prediktorů v dolních uzlech stromů zkonstruovaných z učebního souboru zahrnujícího pozorované hodnoty 50 různých prediktorů. V další fázi práce se počítá jednak se zahrnutím nových druhů prediktorů do prediktorového vektoru, jednak s dalším zdokonalováním v této práci uvedených statistických metod, případně vyzkoušení jiných metod.

Literatura

- ANTOCH J., 1988: Klasifikace a regresní stromy. Sborník *ROBUST 88*, JČSMF, str.0-6.
- DUBROVSKÝ M., 199?: The binary decision tree: the growing algorithm and application to thunderstorm forecasting. Nabídnuo k publikaci v *Studia geophysica et geodaetica*.
- PUDIL P., NOVOVIČOVÁ J. and BLÁHA S., 1991: Statistical approach to pattern recognition (Theory and practical solution by means of PREDITAS system). Supplement to the Journal *Kybernetika* 27, 78pp.
- ŘEZÁČOVÁ and MOTL, 1990: The use of the simple 1D steady-state convective cloud model in the decision tree for determining the probability of thunderstorm occurrence. *Studia geoph. et geod.* 34, 147-166.