

# Metody pro kombinování modelů a klasifikátorů

## Diplomová práce

### Cíle práce

- Zlepšení kvality výstupu regresních a klasifikačních modelů za pomoci ensembleů.
- Automatizace data mining procesu.

**Autor**

**Jan Černý**

cernynj@seznam.cz

České vysoké učení technické v Praze

Fakulta elektrotechnická

Katedra počítačů

## Základní modely

### Třídy problémů

#### 1) Regresní

- Pro každý vstupní vektor je výstupní hodnota jedno reálné číslo, kterému se model snaží co nejvíce přiblížit.

#### 1) Klasifikační

- Výstupem je zařazení vstupního vektoru do jedné třídy z konečné množiny tříd.

- Obecně základní model může být cokoliv co nám poskytne na daná vstupní data nějaký výstup.  
- V našem softwaru (FAKE-GAME) jsou základní modely reprezentovány jednoduchými funkcemi (exponenciální, polynomiální, sinusová, gaussovská...), ale lze využít libovonně složitější modely i modely z jiných aplikací a vylepšit jejich výstupy ensemble algoritmy.



### Princip ensemble algoritmů

- Kombinace modelů takovým způsobem, aby výsledek byl lepší než nejlepší z modelů.  
- Využívá toho, že každý model dělá různé chyby.

### Implementované algoritmy

- Bagging
- Boosting
- Stacking
- Cascade Generalization
- Cascading
- Delegating

### Area Specialization (AS)

- Nově navržený algoritmus využívající model v oblasti kde je nejlepší.  
- Využití modelů relevantních pouze v některé oblasti dat.

### Divide ensemble (DIV)

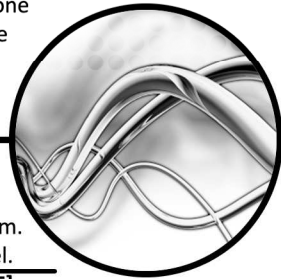
- Nově navržený algoritmus, který přiděluje každému modelu pouze tu oblast dat na které se optimálně naučí a zvládne ji aproximovat.

## Ensemble

### Experimenty

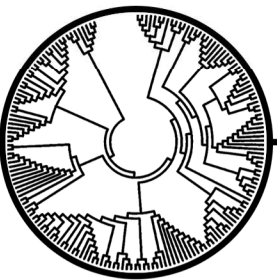
- Ens. = sloupec s nejlepším základním ensemble algoritmem.  
- Zakl. = nejlepší základní model.

data	modely [RMSE]			
	zakl.	ens.	AS	DIV
block	6.6	4.02	<b>1.21</b>	3.2
bump	6.85	6.7	<b>1.39</b>	4.8
doppler	6.32	5.24	<b>1.24</b>	3.21
hevysine	5.1	2.19	0.52	<b>0.4</b>



## Hierarchický ensemble

- Zobecnění principu ensembleů takovým způsobem, že jako základní modely můžeme využít i ensemble algoritmy.  
- Vznik stromových struktur s nekonečným stavovým prostorem => nutnost neúplného prohledávání (evoluce).  
- Experimenty jak s ensembley tak s hierarchickými ensembley ukazují na velikou datovou závislost všech algoritmů pro tvorbu modelů.  
- Neexistuje tedy jediný algoritmus (nebo jejich kombinace), který by byl nejlepší na všechna data.  
=> Potvrzení správnosti našeho přístupu sestavování optimálního modelu výběrem kombinací algoritmů pomocí evoluce.



- Algoritmus pro evoluci stromových struktur a optimalizaci proměnných.  
- Algoritmus je problémově nezávislý a využívá funkcí problémově závislého kontextu.

### Evoluční algoritmus

- Problémově nezávislý (pracuje s obecnými stromy).  
- Provádí evoluční operace:

- mutace uzlu,
- přidání uzlu,
- mutace proměnné.

### Kontext

- Problémově závislá část algoritmu.  
- Načítání a předzpracování dat.  
- Výpočet fitness.

### Experimenty

data	modely [RMSE]				data	modely [RMSE]			
	zakl.	ens.	evol.	zakl.		ens.	evol.		
block	6.6	1.21	<b>0.0015</b>	buildings 2	0.525	0.517	<b>0.46</b>		
bump	6.85	1.39	<b>0.678</b>	buildings 3	0.63	0.597	<b>0.52</b>		
doppler	6.32	1.24	<b>0.65</b>	flare 1	0.094	0.095	<b>0.085</b>		
hevysine	5.1	0.4	<b>0.0116</b>	flare 2	0.033	0.032	<b>0.018</b>		
bosthouse	3.31	3.33	<b>3.23</b>	flare 3	0.026	0.024	<b>0.013</b>		
buildings 1	121.2	115.4	<b>112.4</b>	mandarin	0.103	0.097	<b>0.081</b>		



## Evoluční algoritmus

Experimenty ukazují, že přístup popsáný v mé práci dokáže zlepšit kvalitu výstupu základních modelů o jednotky až desítky procent a zároveň automatizovat proces tvorby modelů.