

Manažer modelování kreditních rizik

část I: pravděpodobnost selhání a ztráta ze selhání

Vojtěch FILIPEC





whoami

☢ jaderná fyzika

🏢 profesní zkušenosti:

- banky: vývoj a validace modelů pro řízení kreditního rizika, monitoring, automatizace...
- fintech, farmaceutický průmysl, FMCG: “analytik - konzultant”

🎯 data scientist



obsah dopolední části kurzu

1. komponenty očekávané ztráty: PD, LGD, EAD
2. modelování PD: logistická regrese, klasifikační stromy, náhodné lesy
3. modelování LGD: vážené průměry, lineární regrese, regresní stromy

metody vs. systém řízení rizika

Co se chcete naučit?

Proč odhadovat PD a LGD?

část 1: komponenty očekávané
ztráty



$$\mathbf{EL = PD \times LGD \times EAD}$$

[CZK]

[%]

[%]

[CZK]

EL ... očekávaná ztráta, expected loss

PD ... pravděpodobnost selhání, probability of default

LGD ... loss given default, ztráta z defaultu

EAD ... exposure at default, expozice v čase defaultu



Než se pustíme do odhadů

dostupnost a relevance historických dat

horizont predikce

kalibrace: Point-in-Time vs.
Through-the-Cycle

očekávaná vs. neočekávaná ztráta: opravné
položky a kapitál

dokumenty BIS (“Basel”): standardizovaný
nebo IRB přístup

segmentace: portfolia a dostupnost dat

Jak odhadovat PD?

část 2: o modelech, logistická
regrese, klasifikační stromy,
náhodné lesy



O modelech (statistické minimum)

spojité, ordinální, nominální proměnné

trénování a validace, overfitting, chybu
měříme na **validačním** vzorku

historická vs. budoucí data: stabilita
populace, population stability index



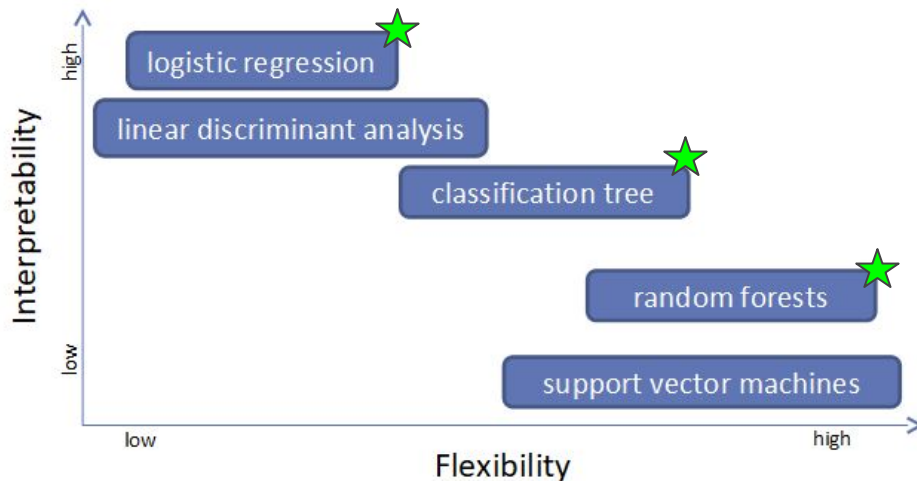
O modelech (statistické minimum)

spojité, ordinální, nominální proměnné

trénování a validace, overfitting, chybu měříme na **validačním** vzorku

historická vs. budoucí data: stabilita populace, population stability index

komplexita vs. interpretabilita





Strategie odhadu defaultu

klasifikace vs. regrese

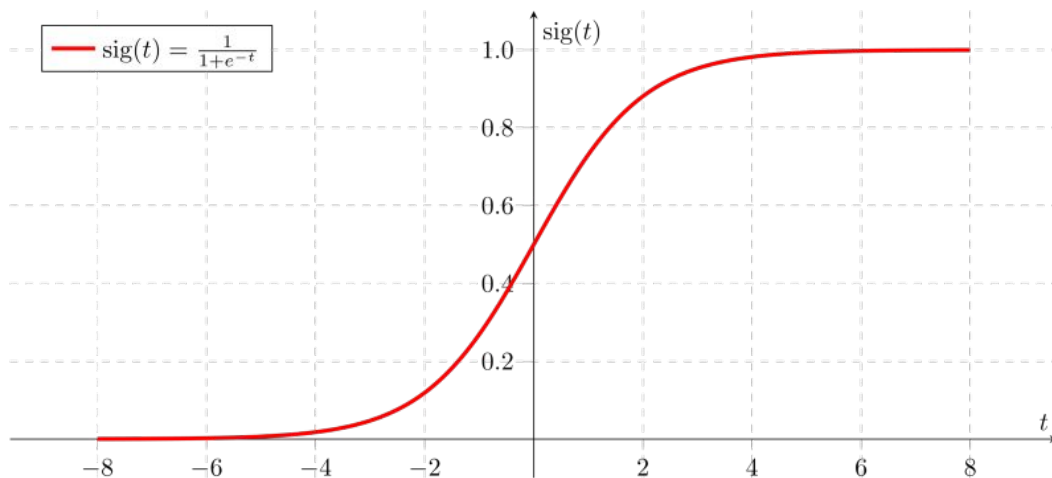
diskrétní veličina × modelování spojité pravděpodobnosti (PD)

jak sestavit trénovací vzorek pro PD:

- co nejnovější data, avšak
- dostatečně stará na to, abychom věděli, zda default nastal
- nové proměnné ve starých datech?!



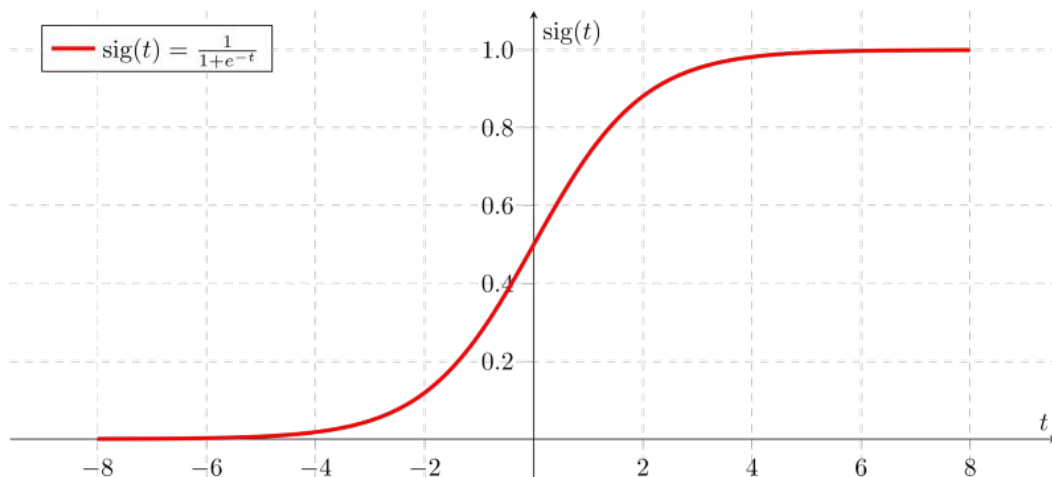
Logistická regrese



zdroj: <https://towardsdatascience.com/@iArunava>



Logistická regrese



zdroj: <https://towardsdatascience.com/@iArunava>

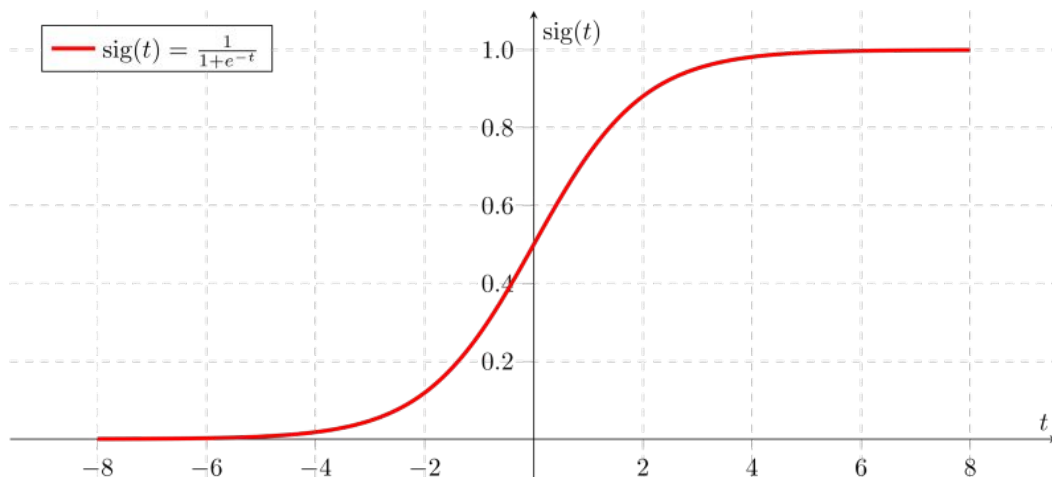
metoda odhadu: sestrojit $t = f(x_0, \dots, x_j)$

předpoklady:

- $x \propto \log\text{odds}(y(x))$... transformace
- extrémní hodnoty ... transformace
- korelace prediktorů ... výběr



Logistická regrese



zdroj: <https://towardsdatascience.com/@iArunava>

metoda odhadu: sestrojit $t = f(x_0, \dots, x_j)$

předpoklady:

- $x \propto \log\text{odds}(y(x))$... transformace
- extrémní hodnoty ... transformace
- korelace prediktorů ... výběr

výhoda: interpretabilní

nevýhody: transformace: nutné, leč obtížná validace



Logistická regrese: výstup

model: **trénovací** data:

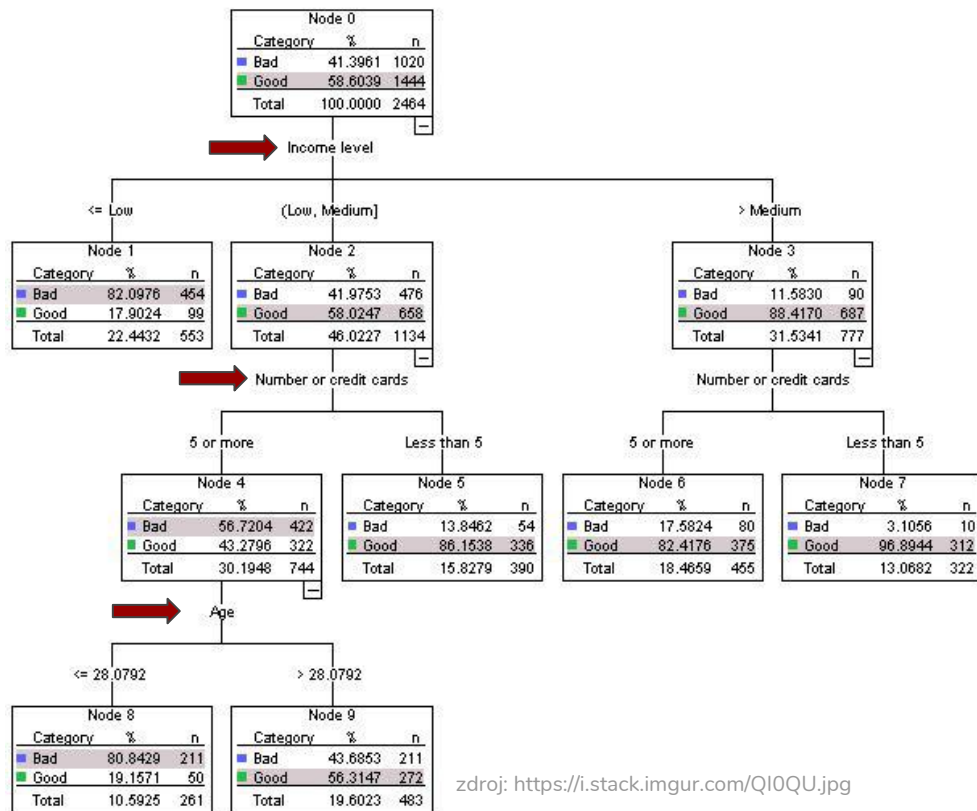
Parameter	DF	Estimate
Intercept	1	-1.99
Credit TOV to debt (WOE)	1	-1.03
Credit registry scores (WOE)	1	-0.55
Family status (WOE)	1	-0.6
Behavioral scores (WOE)	1	-0.49
Affordability (WOE)	1	-0.47
Time with bank (WOE)	1	-0.63

výkonnost: **testovací** data:

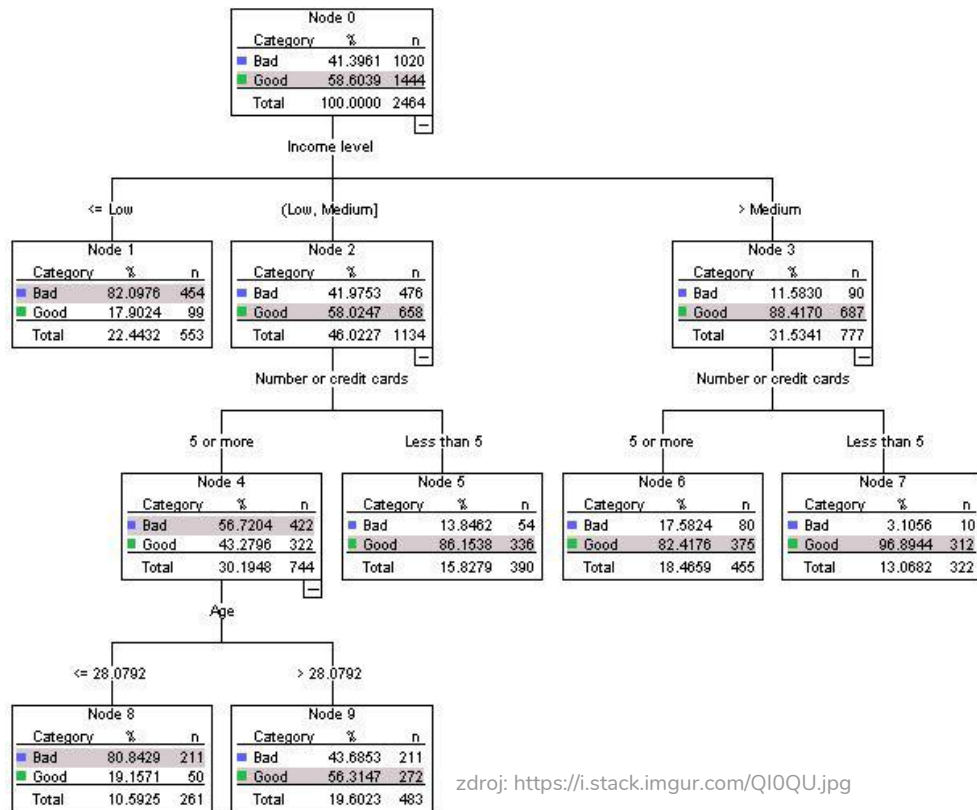
Statistics	Estimate
% Concordant	74.3
% Discordant	25.1
% Tied	0.6
Somers' D	0.49
C-statistics	0.75
Kolm-Smirn	0.28

(znaménka, DF, standardní odhady
transformace,
definice statistik)

Klasifikační strom



Klasifikační strom



zdroj: <https://i.stack.imgur.com/Q10QU.jpg>

metoda odhadu: rozdělit pozorování dle x_1, \dots, x_j do homogenních skupin

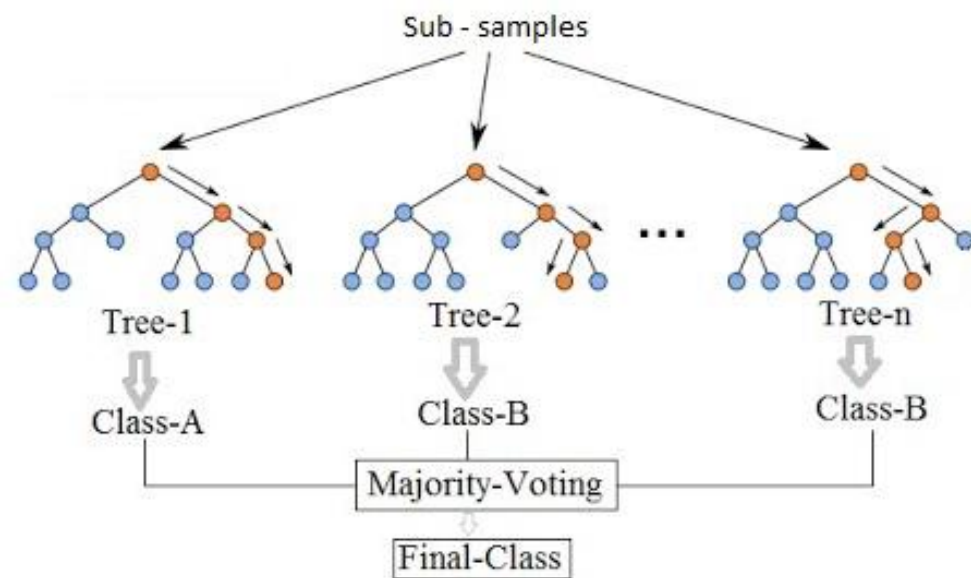
předpoklady: dostatek dat (“greedy” algoritmus)

výhoda: neparametrické, interpretabilní

nevýhody: overfitting

(jak trénovat, lineární separabilita)

Náhodný les



metoda odhadu: dvojí sampling a mnoho stromů -> průměrovat výsledek

předpoklady: žádné!

výhody: robustní, neparametrické (hyperparametry!), implicitní validace

nevýhoda: interpretace

(bagging, Variable Importance Factor)

Jak odhadovat LGD?

část 3: vážené průměry,
lineární regrese, rozhodovací
stromy



Strategie odhadu ztráty ze selhání

spojitá veličina

dva typy LGD:

- defaultní pohledávky: málo pozorování
- ne-defaultní pohledávky: málo pozorování a málo proměnných

segmentace

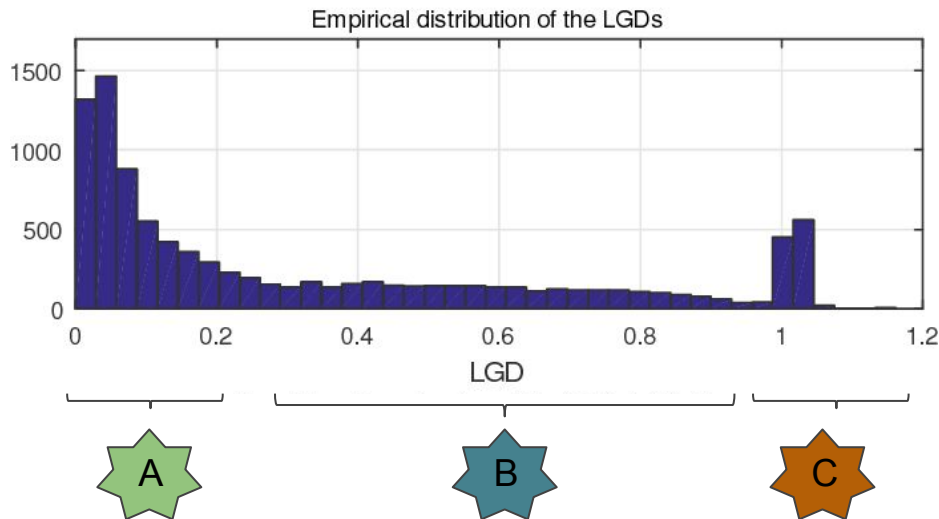
$LGD > 100 \%$?

trénovací vzorek pro odhad LGD:

- co nejnovější data, avšak
- musíme znát skutečnou ztrátu: co s neuzavřenými případy?



Expertní segmentace



Umím rozdělit pohledávky do skupin?

- ano => odhad po skupinách
- ne => odhad dohromady

Potřebuji predikovat skupinu?

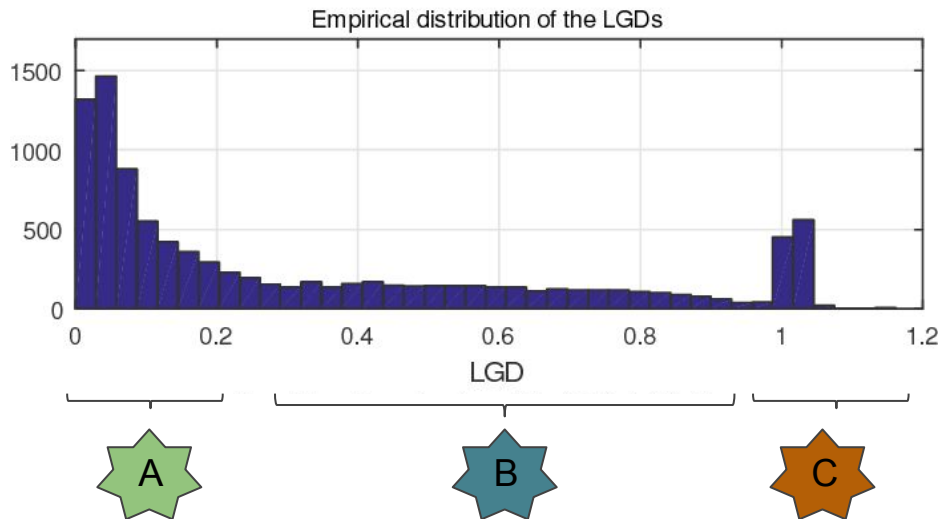
- zajištění
- soc.-dem. atributy
- kreditní toky
- přístup: vícestupňový model

zdroj:

<https://www.semanticscholar.org/paper/Loss-functions-for-LGD-model-comparison-Hurlin-Leymarie/bb392ddb155cbb4d45635789f5fb9a70ec4d3070>



Vážený průměr a regrese



metody odhadu (se segmentací nebo bez ní):

- průměr: vážení historického LGD
- lineární / beta regrese, $\text{LGD} \propto f(x_0, \dots, x_j)$

předpoklady regrese:

- statistické (lineární: homoskedasticita)
- praktické: dostatek proměnných?

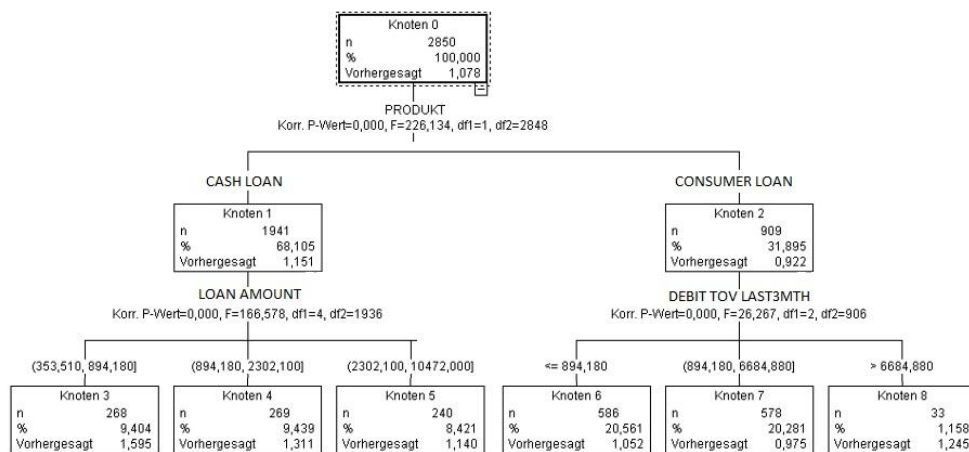
výhoda: interpretace

nevýhody: omezená síla + vliv extrémů

zdroj:

<https://www.semanticscholar.org/paper/Loss-functions-for-LGD-model-comparison-Hurlin-Leymarie/bb392ddb155cbb4d45635789f5fb9a70ec4d3070>

Regresní strom



metoda odhadu: rozdělit pozorování dle x_1, \dots, x_j do homogenních skupin

předpoklady: dostatek dat (“greedy” algoritmus)

výhoda: segmentaci vytvoří strom sám!

nevýhody:

- těžké zohlednit vliv expertů
- overfitting



Děkuji za pozornost

□ [linkedin.com/in/vojtech-filipec/](https://www.linkedin.com/in/vojtech-filipec/)

🔍 github.com/vojtech-filipec





Doporučené zdroje

- Vývoj skórkaret prakticky:
<https://support.sas.com/documentation/cdl/en/emcsgs/66024/PDF/default/emcsgs.pdf>
- Další klasifikační a regresní algoritmy:
<https://towardsdatascience.com/10-machine-learning-algorithms-you-need-to-know-77fb0055fe0>