

# 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, farma, 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

Co se chcete dovědět?

# **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 odhadu

dostupnost a relevance historických dat

horizont predikce

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

očekávaná vs. neočekávaná ztráta: opravné  
položky a kapitál

dokumenty Basel II a III: standardizovaný  
nebo IRB přístup

# **Jak odhadovat PD?**

část 2: logistická regrese,  
klasifikační stromy, náhodné  
lesy



# Strategie odhadu defaultu

klasifikace vs. regrese

diskrétní veličina × modelování spojitě pravděpodobnosti (PD)

trénovací vzorek vs. horizont odhadu PD:

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





# 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



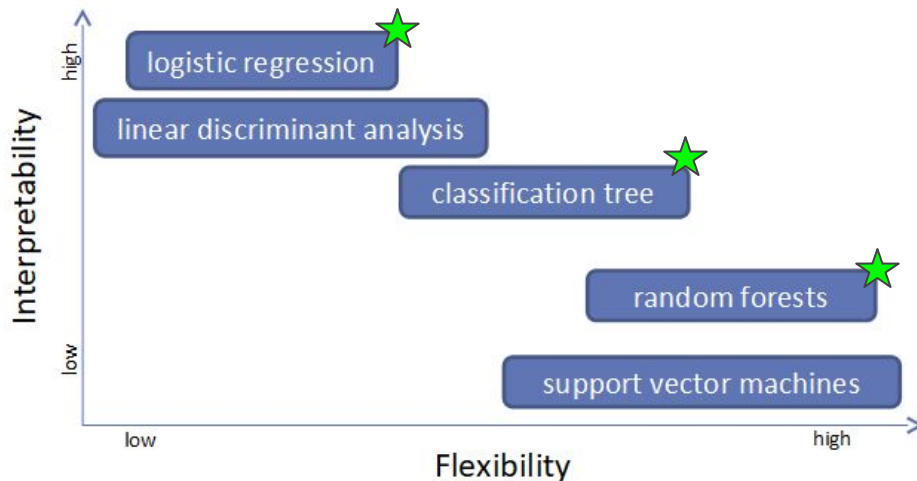
# O statistických modelech

spojité, ordinální, nominální proměnné

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

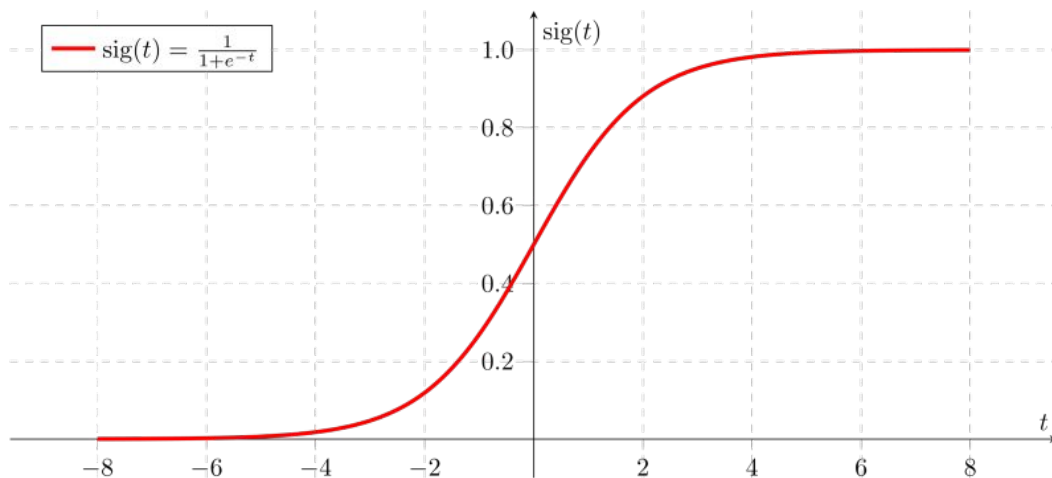
historická vs. budoucí data: stabilita populace, population stability index

komplexita vs. interpretabilita





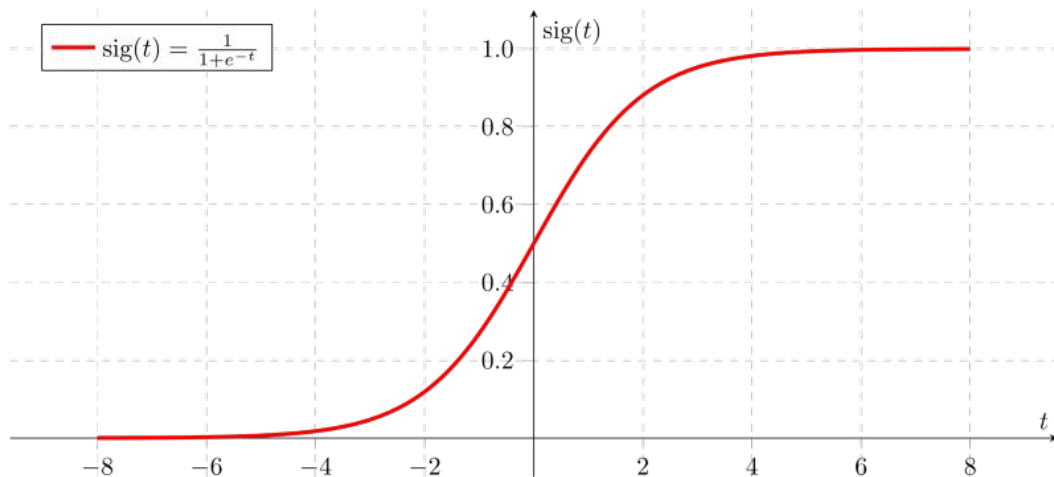
# Logistická regrese



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



# Logistická regrese



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

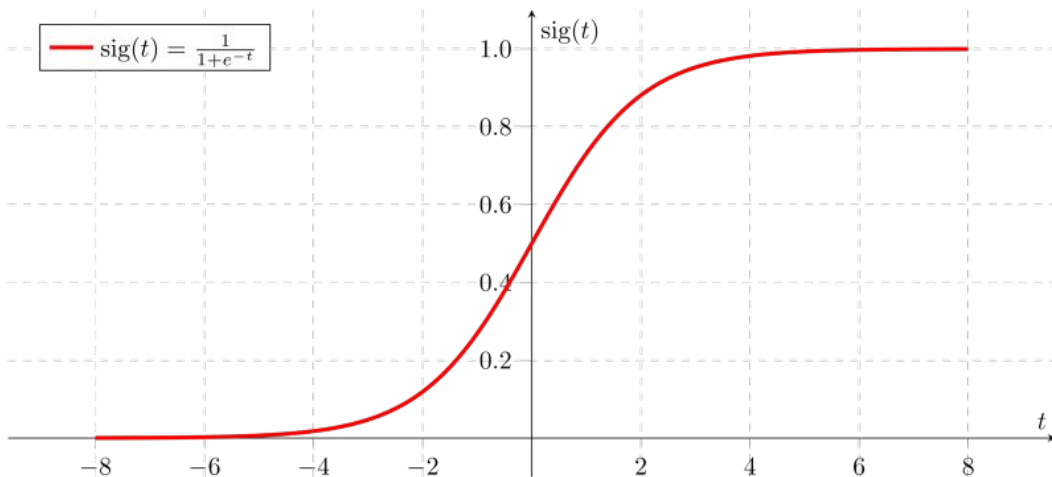
**metoda odhadu:** sestrojit  $t = f(x_0, \dots, x_j)$

**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:

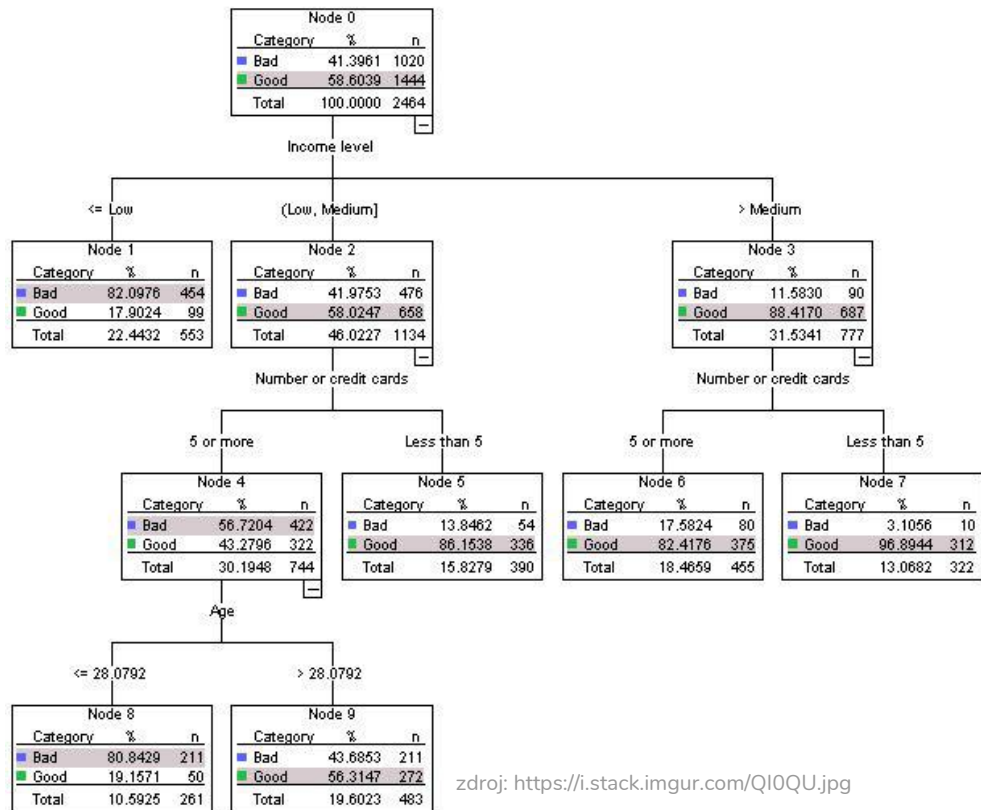
Parameter	DF	Estimate
Intercept	1	<b>-1.99</b>
Credit TOV to debt (WOE)	1	<b>-1.03</b>
Credit registry scores (WOE)	1	<b>-0.55</b>
Family status (WOE)	1	<b>-0.6</b>
Behavioral scores (WOE)	1	<b>-0.49</b>
Affordability (WOE)	1	<b>-0.47</b>
Time with bank (WOE)	1	<b>-0.63</b>

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



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

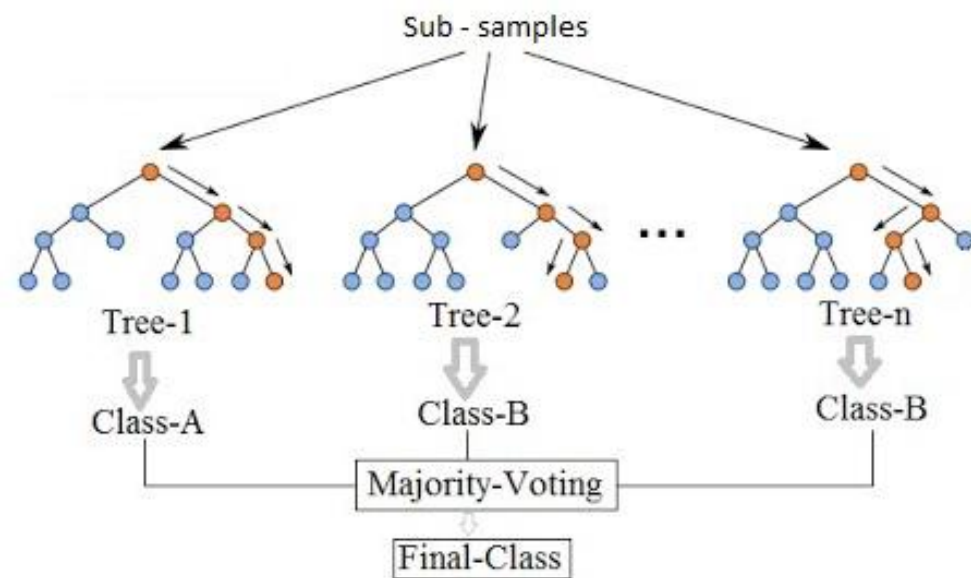
**předpoklady:** dostatek dat (“greedy” algoritmus)

**výhoda:** neparametrické, interpretabilní

**nevýhody:** overfitting

(jak trénovat, lineární separabilita)

# Náhodný les



**metoda odhadu:** mnoho stromů -> průměrovat výsledek

**předpoklady:** žádné!

**výhody:** robustní, neparametrické (hyperparametry: dvojí sampling), 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

málo pozorování a proměnných: LGD = loss **given** default = ztráta podmíněná selháním

**segmentace**

LGD > 100 %?

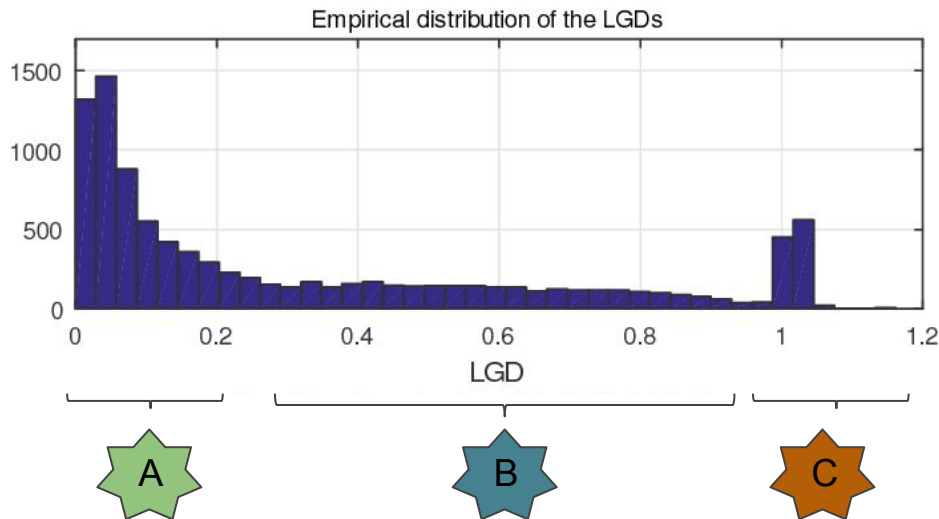
vazba na ekonomický cyklus: silnější než u PD

trénovací vzorek pro odhad LGD:

- co nejnovější data, avšak
- musíme znát skutečnou ztrátu



# Expertní segmentace



Umím rozdělit pohledávky do skupin?

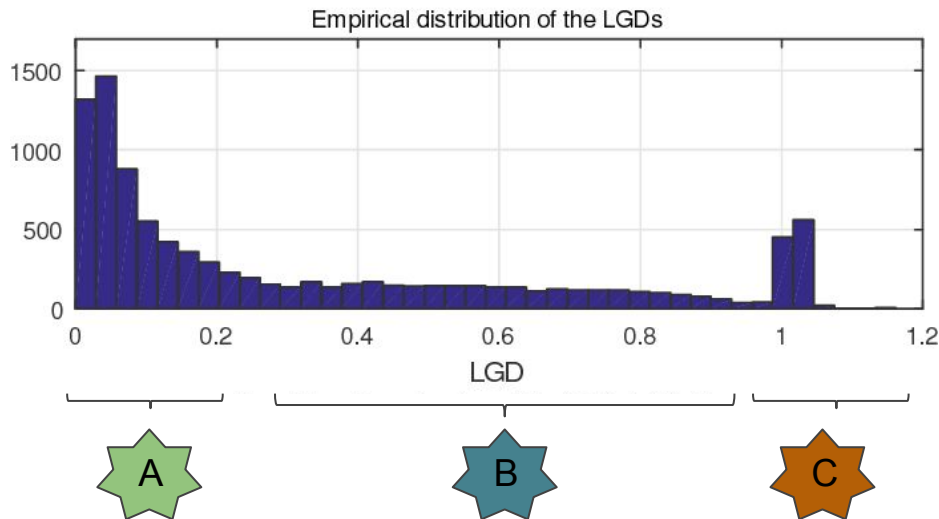
- ano => odhad po skupinách
- ne => odhad dohromady

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?

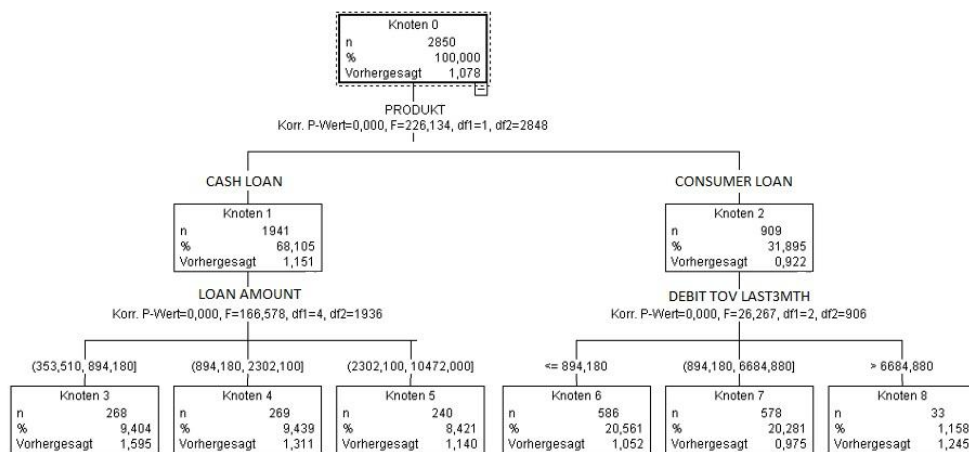
**výhoda:** interpretace

**nevýhody:** omezená síla + vliv extr. hodnot

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



## 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>



# Děkuji za pozornost

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

🔍 [github.com/vojtech-filipec](https://github.com/vojtech-filipec)

