

# Validációs paraméterek mintamérettől való függésének vizsgálata

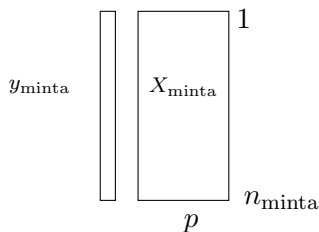
Kovács Dániel, Király Péter, Tóth Gergely

ELTE, Kémiai Intézet

2019. június 6.

- Validációs paraméterek (VPk) mintamérettől való függése
  - összevetése gyenge és jó illeszkedésű lineáris modelleken
  - $R^2$  scrambling, randomizációs változatok
  - $Q^2$  leave one out (LOO) - leave many out (LMO)
  - modellek Roy-Ojha diagramon ( $Q_{F2}^2$ )
  - VPk rangkorrelációja
- Összegzés

## Belső validáció

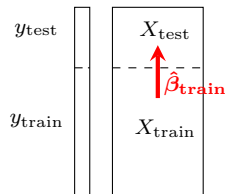


illeszkedés jósága, robusztusság

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

## Külső validáció

$$y = X\beta + e_i$$



prediktivitás

$$r_i = \left( y_{i,\text{test}} - \hat{y}_{i,\text{test}} \left( \hat{\beta}_{\text{train}} \right) \right)$$

belső VPk

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

$$R^2_{\text{adj}} = \left(1 + \frac{p-1}{n-p}\right) R^2 - \left(\frac{p-1}{n-p}\right)$$

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

$$Q^2_{\text{LOO}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_{i/i})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\text{PRESS}}{\text{TSS}}$$

$$\text{CCC} = \frac{2 \sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2 + n(\bar{y} - \bar{\hat{y}})^2} \dots$$

külső VPk

$$Q^2_{F1} = 1 - \frac{\sum_{i=1}^{n,\text{test}} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n,\text{test}} (y_i - \bar{y}_{\text{train}})^2}$$

$$Q^2_{F2} = 1 - \frac{\sum_{i=1}^{n,\text{test}} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n,\text{test}} (y_i - \bar{y}_{\text{test}})^2}$$

$$Q^2_{F3} = 1 - \frac{\sum_{i=1}^{n,\text{test}} (y_i - \hat{y}_i)^2 / n_{\text{test}}}{\sum_{i=1}^{n,\text{train}} (y_i - \bar{y}_{\text{train}})^2 / n_{\text{train}}}$$

⋮

## Adatsorok

Beton nyomószilárdság sor

Kombinált ciklusú erőmű sor

- modellezés *többváltozós lineáris regresszióval* (TLR)

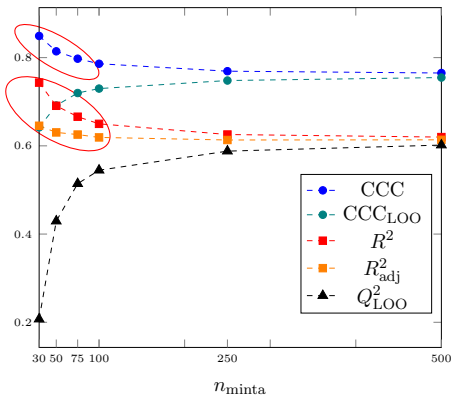
adatsor	Beton	Erőmű
# függő változó	1	1
# független változó	8	4
# megfigyelés	1030	9658
TLR alkalmasság	gyenge	jó

## Mintaméret függés - feladatok

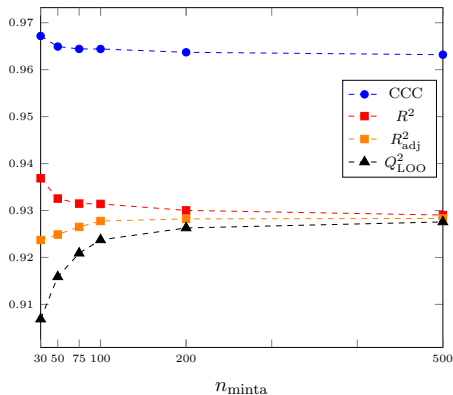
- random mintavételezés a populációból
- minta *training/test* felbontása  
külső validációhoz: 80/20
- lineáris modellezés, *validációs paraméterek* számolása

↻ fentiek ismétlése 1000×,  
 $n_{\text{minta}} = (30, 50, 75, 100, 250, 500)$ ,  
VPk átlagolása

## Beton

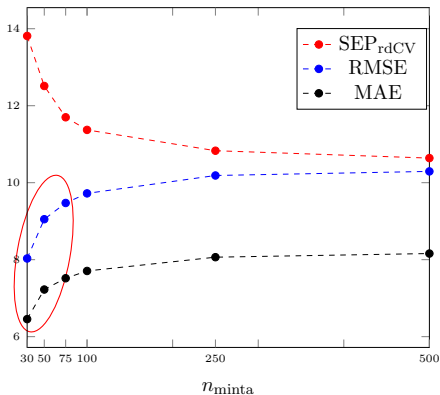


## Erőmű

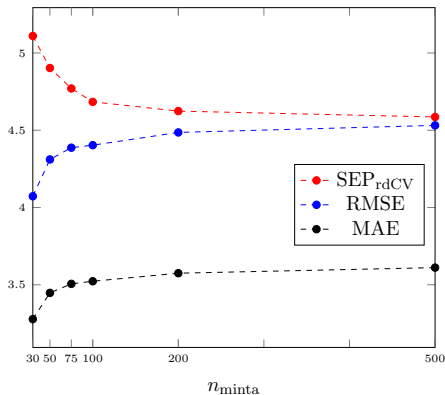


Félrevezető: Kisebb modell a jobb?

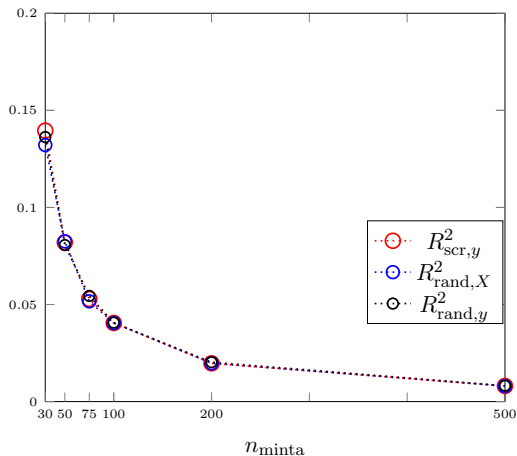
## Beton



## Erőmű



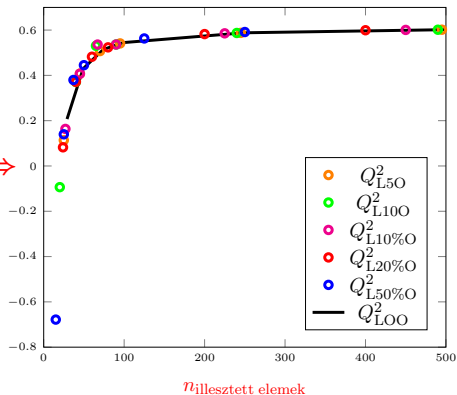
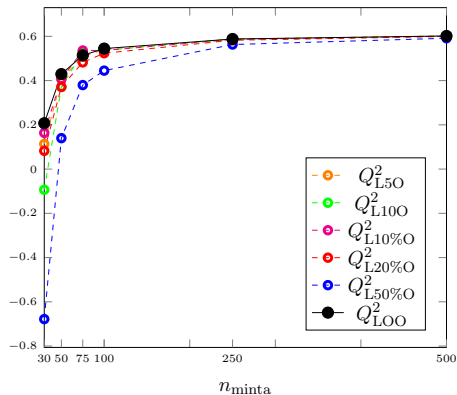
Félrevezető: Kisebb modell a jobb?



$R^2$  scrambling számolható a legkönnyebben

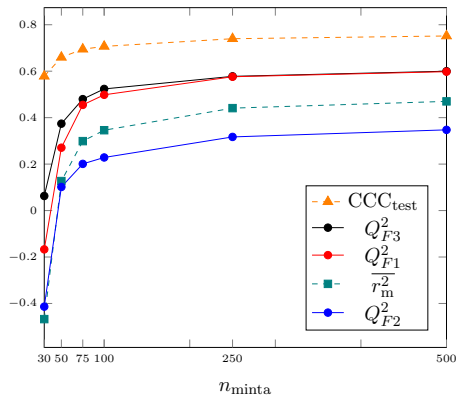


# Belső VPK - Leave many out

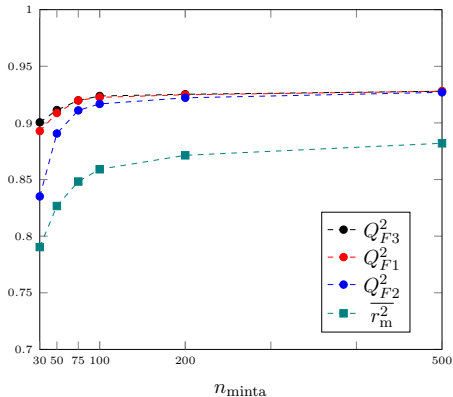


$Q^2_{LOO}$  számolása a legegyszerűbb

Beton

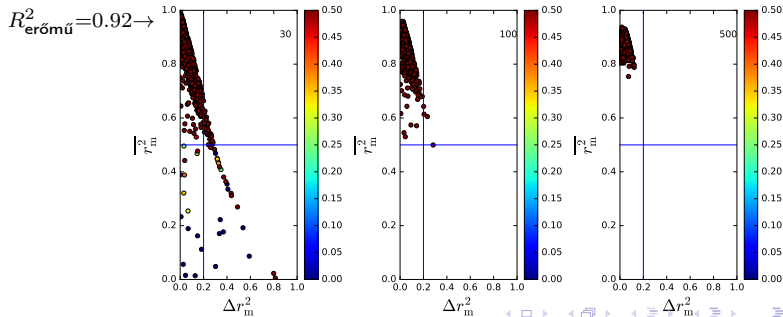
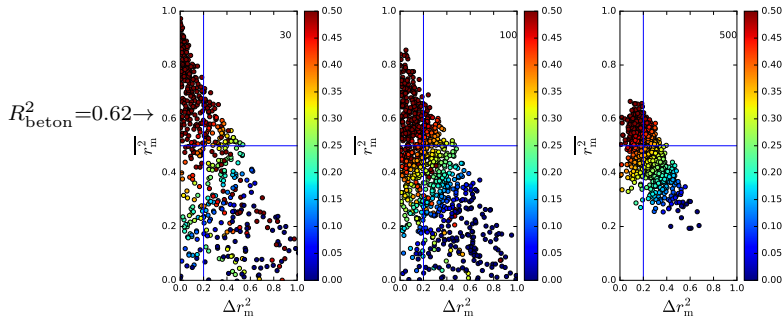


Erőmű



$Q_{F2}^2$  a legérzékenyebb

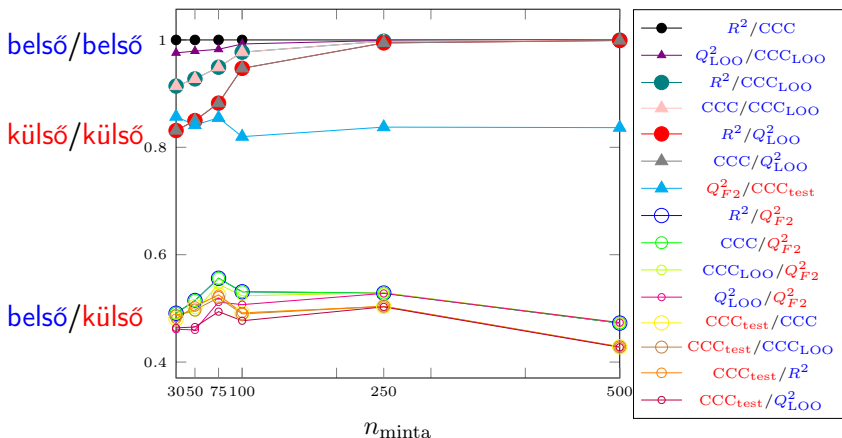
# Külső VPK - Roy-Ojha diagram



# VPk rangkorrelációja

kiválasztott VPk:  $R^2, CCC, Q_{LOO}^2, CCC_{LOO}, Q_{F2}^2, CCC_{test}$

belső
külső



- modell validáció: illesztés jósága, robusztusság & prediktivitás
- $R^2$ , CCC, RMSE, MAE mintaméret függése anomáliát mutat
- külső validáció  $Q_{F1-3}^2$  paraméterei különböző mintaméretre különböző érzékenységet mutatnak
- scrambling-gel kiválthatjuk a randomizációt
- $Q_{LOO}^2$  használható az LMO változatok helyett
- 2D Roy-Ojha diagram kategorizálja a modellek prediktivitását
- minden mintaelemszámnál eltér a külső és belső validáció információtartalma
- nemlineáris modellezés, PLS

Munkatársak: Kovács Dániel, Tóth Gergely