

## 1. VPLIV STANDARDIZACIJE PODATKOV NA REZULTATE:

Za različne metode standardizacije podatkov podane pod točko a) b) in c) boste izračunali parametre modela / enačbe po metodi LSE in PCA.

Za LSE metodo uporabite model s prostim členom:  $y = a_1 x_1 + \dots + r$

S PCA metodo izračunamo model tako, da vse meritve vključno z izhodom zložimo v matriko podatkov  $X$  ( $X = [A\_LOW \ T\_H2O \ C\_ACID \ I\_EFF]$ ). Naredimo svd ( $[P,D] = \text{svd}(F)$ ). Kjer je  $F = X'X/(N-1)$ . Izbrati morate pravi lastni vektor  $P1 = P(:,i)$ , ki nam podaja enačbo modela. S tem dobite enačbo v obliki ( $V$  = center podatkov):

$$P1' * (x-V) = 0.$$

To enačbo pretvorite v obliko  $y = a_1 x_1 + \dots + r$ .

PRIMER:

$$[p1 \ p2] * ([x1 \ y]' - [v1 \ v2]') = 0$$

$$p1 \ x1 + p2 \ y - p1 \ v1 - p2 \ v2 = 0$$

$$y = -p1/p2 \ x1 + p1/p2 \ v1 + p2/p2 \ v2$$

Za primerjavo parametrov in dobljenih modelov morate modele dobljene pod točko a) in b) pretvoriti nazaj v dejanski prostor, saj v primeru normiranih spremenljivk dobite tudi normirane parametre.

PRIMER:

$$y\_n = (y - m\_y) / n\_y \quad x\_n = (x - m\_x) / n\_x$$

$$y\_n = a\_n * x\_n + r\_n$$

$$(y - m\_y) / n\_y = a\_n * (x - m\_x) / n\_x + r\_n$$

...

$$y = (a\_n * n\_y / n\_x) * x - ( (a\_n * n\_y / n\_x) * m\_x - r\_n * n\_y - m\_y )$$

Ko to naredite imate za vsako metodo in vsako normiranje svoje parameter. Parametre lahko predstavite v razpredelnici.

Naredite primerjavo med PCA in LSE modeli glede na standardizacijo. Poglejte katera metoda je manj občutljiva na standardizacijo. Poglejte katera daje bolj podobne rezultate.

Za vsak model izračunajte napako (v originalnem prostoru):  $e = y - \text{ksi} * \text{theta}$ . Izračunajte povprečno vrednost napake in standardno deviacijo napake. Iz teh dveh vrednosti lahko ocenimo PRISTRANSKOST in KONSISTENCO modela. Dober model ni pristranski in je konsistenten. Manjša kot je standardna deviacija napake bolj je model konsistenten. Bližje kot je povprečna vrednost napake ničli, manj pristranski je model. Za primerjavo naših modelov med seboj je dovolj da primerjamo standardne deviacije napake in povprečno vrednost napake.

Če želimo na splošno oceniti kako dober je model pa uporabimo običajno NRMSE napako:

$$\text{Sqrt}(\text{Sum}((y_{\text{ocene}} - y)^2) / n) / \text{std}(y)$$

V poročilo vključite komentar katera metoda je manj občutljiva na standardizacijo podatkov in kateri model je po vašem mnenju boljši. Tisti bolj pridni boste vključili še kakšno teoretično ozadje zakaj je ena od metod bolj občutljiva na standardizacijo podatkov.

## 2. KOLINEARNOST:

Dodamo meritev, ki je odvisna od že podane meritve. Standardiziramo podatke.

Izračunamo parametre modela z LSE metodo (ni potrebno uporabiti prostega člena).

Izračunamo parametre s PCR. Sestavimo matriko X iz vhodnih podatkov. S pomočjo PCA odstranite glavno komponento z zanemarljivim vplivom. Originalne podatke transformirate v prostor glavnih komponent  $T = X * P_s$ .

Izračunate parameter z LSE metodo  $\text{THETA} = (T' T)^{-1} T' Y$

Pretvorimo parameter iz prostora glavnih komponent  $\text{THETA} = P_s \text{THETA}$ .

Tu se ustavite in pregledate dobljene parametre. Kaj boste opazili.

a) PCR in LSE parametri so čisto drugačni, iz napake težko rečemo kateri so boljši.

b) Da sta parametra pri dodani spremenljivki in tisti, od katere je odvisna zelo podobna.

c) Da je vsota teh dveh parametrov približno enaka parametru, ki nastopa pri originalnem setu podatkov (parameter za temperaturo).

d) če naredimo razporednico glavnih komponent in originalnih spremenljivk boste opazili, da imata dodana odvisna spremenljivka in originalna spremenljivka, isti vplivni koeficijent na glavne komponente.

Za izbrane glavne komponente (lastne vektorje) naredimo naslednjo razporednico:

SPREMENLJIVKA	P_1	P_2
X1	P11	P21
X2	P12	P22
X3	P13	P23
...	..	

P11 je prvi element prvega lastnega vektorja. Za lažjo predstavo lahko v tabelo vpišemo procente, ki jih dobimo kot  $\text{abs}(P_{11}) / \text{sum}(\text{abs}(P_{11}))$ . Ti nam povejo koliko prvotne spremenljivke X1 je vsebovane v glavni komponenti.