

Přednáška 3: Zobecněný lineární regresní model, regularizace volbou apriorní

Ondřej Tichý

Bayesovské metody v machine learningu (BML)

March 2, 2018

Obsah přednášky

- ▶ Zobecněný lineární model a k čemu je to dobré
- ▶ Nejmenší čtverce a v čem je problém
- ▶ Regularizace
- ▶ Formulace pomocí optimalizace
 - ▶ Tichonovova regularizace
 - ▶ LASSO
 - ▶ elastic net
- ▶ Bayesovská formulace problému
 - ▶ bayesovská regrese
 - ▶ ridge regression
 - ▶ sparse regression

Zobecněný lineární regresní model

Opakování

Předpoklad vysvětlení každé naměřené hodnoty:

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + e_i. \quad (1)$$

My budeme využívat maticový zápis:

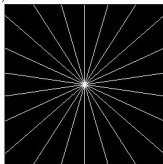
$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & & x_{2n} \\ \vdots & & \ddots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{pmatrix}}_{\boldsymbol{\beta}} + \mathbf{e}. \quad (2)$$

Motivace

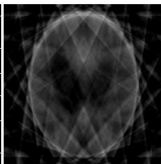
K čemu je to dobré?

- ▶ kdekoliv kde vzniká špatně podmíněná soustava lineárních rovnic
- ▶ rekonstrukce obrazu v tomografii:

- Sparse projections: 11 radial lines



available portion of the spectrum
(11 radial lines)

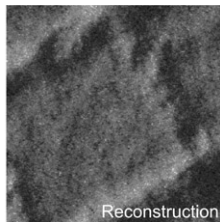
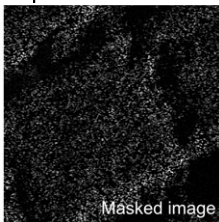
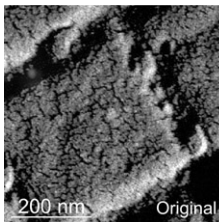


Back-projection estimate



Estimate after convergence
(exact reconstruction)

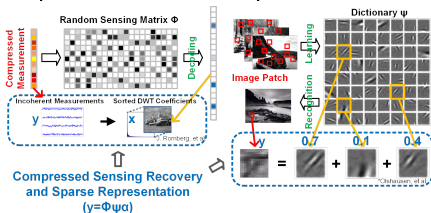
- ▶ elektronová mikroskopie:



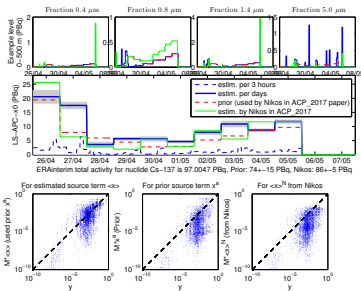
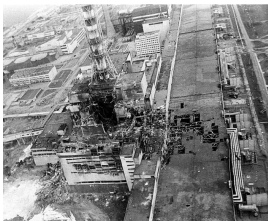
Motivace

K čemu je to dobré?

► rozpoznání a řídká reprezentace



► určení průběhu úniku látek do prostředí



Zobecněný lineární regresní model

Nejmenší čtverce (OLS)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad (3)$$

- ▶ minimalizujeme součet čtverců (kvadrátů) odchylek:

$$\sum_j e_j^2 = \mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad (4)$$

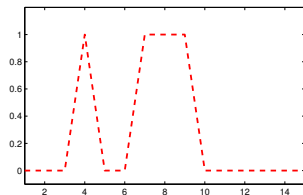
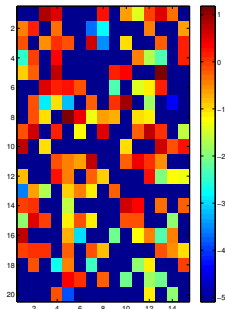
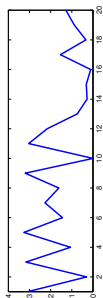
- ▶ derivace podle $\boldsymbol{\beta}$:

$$2\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} - 2\mathbf{X}^T \mathbf{y} = \mathbf{0} \quad (5)$$

$$\hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (6)$$

Zobecněný lineární regresní model

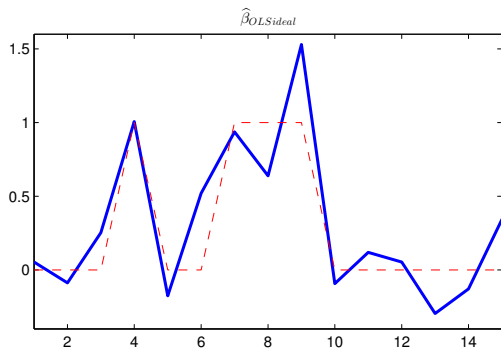
Cvičný příklad - data



$$y = X\beta + \text{šum}$$

Zobecněný lineární regresní model

Řešení OLS



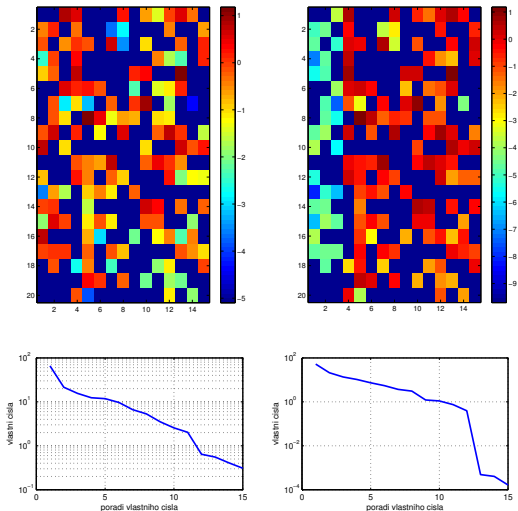
Potud je vše jednoduché a krásné, proč se s tím nespokojíme?

Inverze $(X^T X)^{-1}$.

Zobecněný lineární regresní model

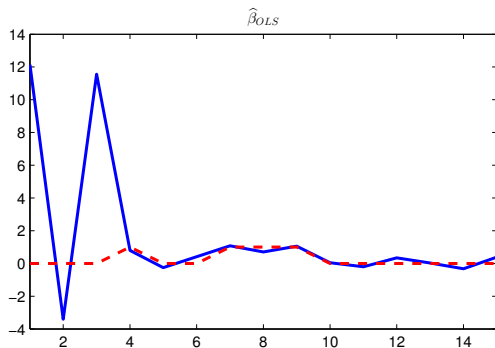
Cvičný příklad

Co když si problém trochu ztížíme (přiblížíme realitě)?



Zobecněný lineární regresní model

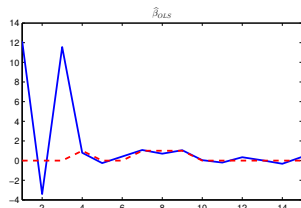
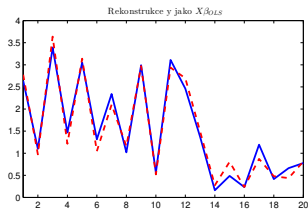
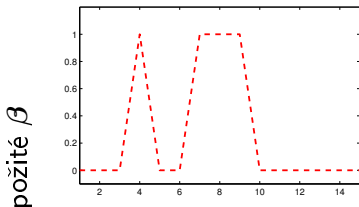
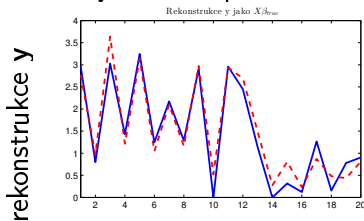
Řešení OLS



Regularizace

Principy

- ▶ chceme vložit dodatečnou informaci do problému,
- ▶ snaha vyhnout se přefitování:



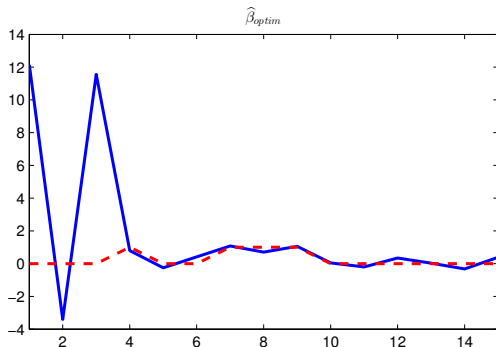
- ▶ dodatečná informace příliš slabá x silná,

Optimalizační přístup

Formulace OLS jako optimalizační úlohy

$$\hat{\beta}_{\text{optim}} = \arg \min_{\beta} \|\mathbf{y} - X\beta\|_2^2 \quad (7)$$

- ▶ $\|\cdot\|_2$ je Eukleidovská norma,
- ▶ funkce $\|\mathbf{y} - X\beta\|_2^2$ je v β (naštěstí) konvexní,
- ▶ vede na identické řešení jako OLS



Optimalizace s Tichonovovou regularizací

Formulace

- ▶ též “Ridge regression”.
- ▶ obecný zápis:

$$\hat{\beta}_{Tikhonov} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \|\Gamma\beta\|_2^2, \quad (8)$$

- ▶ ovšem velmi časté použití pro (α skalár, I jednotková matice):

$$\Gamma = \alpha I, \quad (9)$$

- ▶ a pozorný student snadno nahlédne řešení ve tvaru:

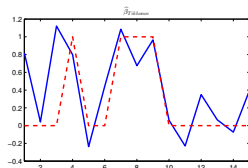
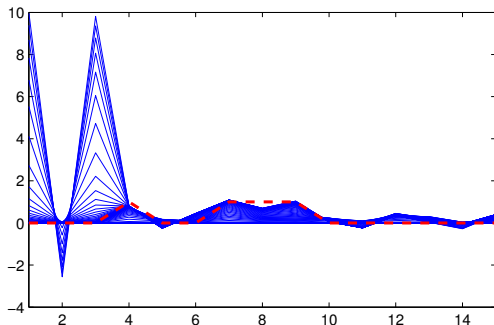
$$\hat{\beta}_{Tikhonov} = \left(\mathbf{X}^T \mathbf{X} + \alpha^2 I \right)^{-1} \mathbf{X}^T \mathbf{y}. \quad (10)$$

- ▶ co to vyřešilo za problém a v čem je potenciální nedostatek?

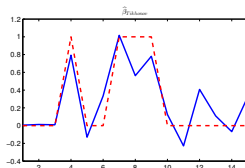
Optimalizace s Tichonovovou regularizací

Řešení

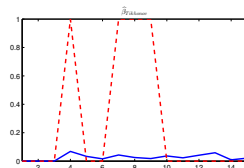
Pro Tichonovův faktor $\alpha = 10^{-3}$ až $+3$:



$\alpha = 10^{-2}$



$\alpha = 10^0$

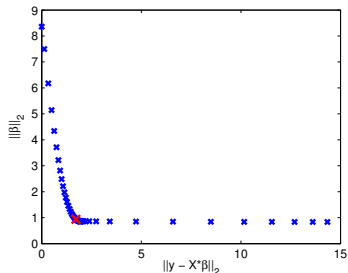


$\alpha = 10^{+2}$

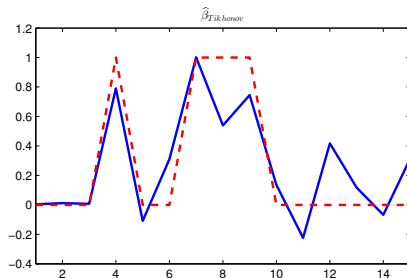
Optimalizace s Tichonovovou regularizací

Volba Tichonovova faktoru

- ▶ např. L-křivka



pro $\alpha = 1.2589$



- ▶ cross-validace
- ▶ ovšem velmi často metoda “kouknu a vidím”...

Optimalizace s LASSO regularizací

Formulace

- ▶ LASSO (least absolute shrinkage and selection operator),
- ▶ provádí zároveň regularizaci a selekci proměnných,
- ▶ zápis problému jako

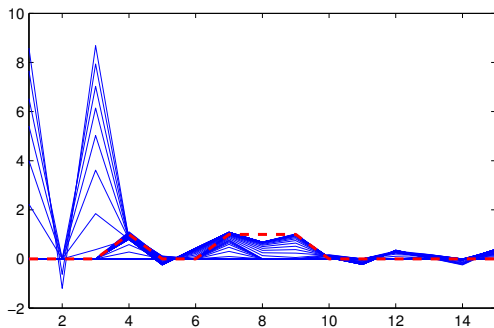
$$\hat{\beta}_{LASSO} = \arg \min_{\beta} \|\mathbf{y} - X\beta\|_2^2 + \alpha \|\beta\|_1, \quad (11)$$

- ▶ kde $\|\mathbf{a}\|_1 = \sum_j |a_j|$.
- ▶ princip...

Optimalizace s LASSO regularizací

Řešení

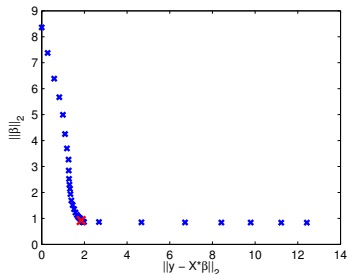
Pro parametr $\alpha = 10^{-3}$ až $+3$:



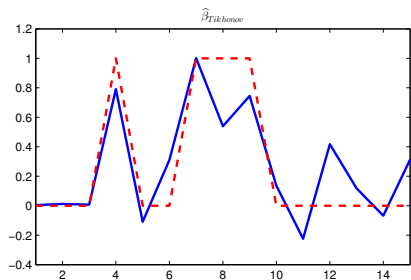
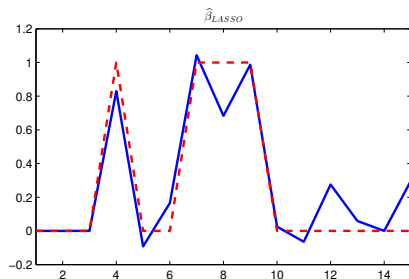
Optimalizace s LASSO regularizací

Řešení

► opět např. L-křivka



pro $\alpha = 0.2512$



Optimalizace a další možnosti

Elastic net

- ▶ proč řešení není ideální?
- ▶ nic nám nebrání v kombinaci:

$$\hat{\beta} = \arg \min_{\beta} ||\mathbf{y} - \mathbf{X}\beta||_2^2 + \alpha_1 ||\beta||_2^2 + \alpha_2 ||\beta||_1, \quad (12)$$

Bayesovský přístup

Snadno a rychle...

- ▶ zvolíme model dat,
- ▶ zvolíme apriorní rozdělení parametrů modelu,
- ▶ spočítáme odhady posteriorních rozdělení parametrů.

Metedologická vsuvka: Variační Bayes (VB)

Variační Bayesova aproximace

- ▶ máme model dat $f(\mathbf{y}|\theta)$,
- ▶ máme q podmíněně nezávislých proměnných, tzn. pro jejich arpiorna platí

$$f(\theta_1, \theta_2, \dots, \theta_q) = \prod_{i=1}^q f(\theta_i), \quad (13)$$

- ▶ pak aposteriorní rozdělení $\tilde{f}(\theta_i|\mathbf{y})$ pro i tý parametr získáme VB aproximací jako

$$\tilde{f}(\theta_i|\mathbf{y}) \propto \exp \left[E_{\tilde{f}(\theta_{\setminus i}|\mathbf{y})} (\ln f(\theta_1, \theta_2, \dots, \theta_q, \mathbf{y})) \right], \quad \forall i. \quad (14)$$

VB: příklad skalární dekompozice

Skalární dekompozice

- ▶ skalární model:

$$d = ax + e, \quad e \sim \mathcal{N}(0, r_e), \quad (15)$$

- ▶ tzn. model dat

$$f(d|a, x) = \mathcal{N}(ax, r_e), \quad (16)$$

- ▶ volba apriorních rozdělení pro parametry a a x :

$$f(a) = \mathcal{N}(0, r_a), \quad (17)$$

$$f(x) = \mathcal{N}(0, r_x). \quad (18)$$

- ▶ r_e, r_a, r_x jsou nyní konstanty (+ data d), \hat{a} a \hat{x} chceme odhadnout.

VB: příklad skalární dekompozice

Skalární dekompozice

- ▶ vyjádříme si logaritmus sdružené pravděpodobnosti $\ln f(\theta_1, \dots, \theta_q, \mathbf{y})$

$$\ln f(a, x, d) \propto -\frac{1}{2} \left(\frac{(ax - d)^2}{r_e} + \frac{a^2}{r_a} + \frac{x^2}{r_x} \right), \quad (19)$$

- ▶ a prostým dosazením do

$$\tilde{f}(\theta_i | \mathbf{y}) \propto \exp \left[\mathbb{E}_{\tilde{f}(\theta_{\setminus i} | \mathbf{y})} (\ln f(\theta_1, \theta_2, \dots, \theta_q, \mathbf{y})) \right] \text{ dostáváme}$$

$$\tilde{f}(a|d) \propto \exp \left[-\frac{1}{2} a^2 \left(\hat{x}^2 r_e^{-1} + r_a^{-1} \right) - a \left(d \hat{x} r_e^{-1} \right) \right], \quad (20)$$

$$\tilde{f}(x|d) \propto \exp \left[-\frac{1}{2} x^2 \left(\hat{a}^2 r_e^{-1} + r_x^{-1} \right) - x \left(d \hat{a} r_e^{-1} \right) \right], \quad (21)$$

- ▶ a všimneme si, že to je opět tvar normálního rozdělení:

$$\tilde{f}(a|d) = \mathcal{N}_a(\mu_a, \sigma_a), \quad \tilde{f}(x|d) = \mathcal{N}_x(\mu_x, \sigma_x). \quad (22)$$

VB: příklad skalární dekompozice

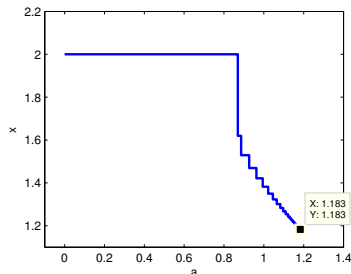
Skalární dekompozice

- ▶ dopočítáme tzv. tvarovací parametry $\mu_a, \sigma_a, \mu_x, \sigma_x$:

$$\sigma_a = \left(\hat{x}^2 r_e^{-1} + r_a^{-1} \right)^{-1}, \quad \mu_a = \sigma_a d \hat{x} r_e^{-1}, \quad (23)$$

$$\sigma_x = \left(\hat{a}^2 r_e^{-1} + r_x^{-1} \right)^{-1}, \quad \mu_x = \sigma_x d \hat{a} r_e^{-1}. \quad (24)$$

- ▶ a to je (skoro) vše...



Bayesovská formulace lineární regrese

Bayesovský pohled

- ▶ připomenutí: máme rovnici

$$\mathbf{y} = X\boldsymbol{\beta} + \mathbf{e}, \quad (25)$$

- ▶ předpokládáme rozdělení chyb ve tvaru

$$e_i \sim \mathcal{N}(0, \omega^{-1}), \quad \forall i, \quad (26)$$

- ▶ tím dostáváme přirozeně model dat jako

$$f(\mathbf{y}|\boldsymbol{\beta}, \omega) = \mathcal{N}(X\boldsymbol{\beta}, \omega^{-1}I_p), \quad (27)$$

- ▶ zajímá nás: $\boldsymbol{\beta}, \omega$.

Bayesovská formulace lineární regrese

Apriorno pro parametry

- ▶ bayesovský pohled: β a ω jsou pro nás neznámé parametry a budeme je odhadovat,
- ▶ volba apriorních rozdělání pro β a ω ,
- ▶ my si pro začátek vystačíme s nejjednodušší možností:

$$f(\beta) = \mathcal{N}(\mathbf{0}, I_n) = \underbrace{(2\pi)^{-\frac{n}{2}}}_{coef} \underbrace{|I_n|^{-\frac{1}{2}}}_{coef} \exp\left(-\frac{1}{2}\beta' I_n^{-1} \beta\right), \quad (28)$$

$$f(\omega) = \mathcal{G}(c_0, d_0) = \frac{c_0^{d_0}}{\underbrace{\Gamma(c_0)}_{coef}} \omega^{c_0-1} \exp(-d_0 \omega). \quad (29)$$

- ▶ naším cílem je určit $\hat{\beta} = E[\beta]$ a $\hat{\omega} = E[\omega]$.

Bayesovská formulace lineární regrese

VB aproximace

- ▶ posteriorně ve tvaru $\tilde{f}(\beta) = \mathcal{N}(\mu_\beta, \Sigma_\beta)$ a $\tilde{f}(\omega) = \mathcal{G}(c, d)$ s tvarovacími parametry:

$$\Sigma_\beta = \left(\hat{\omega} X^T X + I_n \right)^{-1} \quad (30)$$

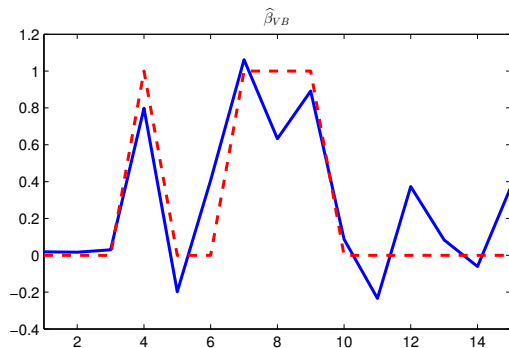
$$\mu_\beta = \Sigma_\beta \hat{\omega} X^T \mathbf{y} = \left(\hat{\omega} X^T X + I_n \right)^{-1} \hat{\omega} X^T \mathbf{y} \quad (31)$$

$$c = c_0 + \frac{p}{2} \quad (32)$$

$$d = d_0 + \frac{1}{2} \left(\mathbf{y} \mathbf{y}^T - 2 \mathbf{y}^T X \hat{\beta} + \beta^T \widehat{X^T X} \beta \right) \quad (33)$$

Bayesovská formulace lineární regrese

Řešení



Bayesovská ridge regression

Formulace

- ▶ model dat a šumu zachováme,

$$f(\mathbf{y}|\boldsymbol{\beta}, \omega) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \omega^{-1}I_p), \quad (34)$$

$$f(\omega) = \mathcal{G}(c_0, d_0). \quad (35)$$

- ▶ do modelu $\boldsymbol{\beta}$ zavedeme hyperparametr $v \in \mathbf{R}$ následovně:

$$f(\boldsymbol{\beta}) = \mathcal{N}(\mathbf{0}, I_n) \quad \rightarrow \quad f(\boldsymbol{\beta}|v) = \mathcal{N}(\mathbf{0}, v^{-1}I_n) \quad (36)$$

- ▶ v se stává parametrem modelu, zvolíme jeho apriorní rozdělení

$$f(v) = \mathcal{G}(a_0, b_0). \quad (37)$$

Bayesovská ridge regression

VB aproximace

- ▶ aposteriorně $\tilde{f}(\omega)$ zůstane stejné jako v předchozím případě.
- ▶ aposteriorně $\tilde{f}(\beta|v) = \mathcal{N}(\mu_\beta, \Sigma_\beta)$ a $\tilde{f}(v) = \mathcal{G}(a, b)$ má (po odvození) následující tvarovací parametry:

$$\Sigma_\beta = \left(\hat{\omega} X^T X + \hat{v} I_n \right)^{-1} \quad (38)$$

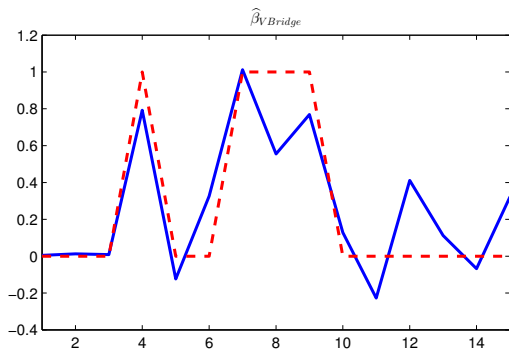
$$\mu_\beta = \Sigma_\beta \hat{\omega} X^T \mathbf{y} = \left(\hat{\omega} X^T X + \hat{v} I_n \right)^{-1} \hat{\omega} X^T \mathbf{y} \quad (39)$$

$$a = a_0 + \frac{n}{2} \quad (40)$$

$$b = b_0 + \frac{1}{2} \text{trace} \left(\widehat{\beta \beta'} \right) \quad (41)$$

Bayesovská ridge regression

Řešení



- přílišná jednoduchost modelu

Bayesovská sparse regression

Formulace

- ▶ též “Sparse bayesian learning” nebo “Relevance vector machine”.
- ▶ rozptyl prvků vektoru β nebudeme modelovat pouze jedním parametrem, ale přiřadíme parametr každému prvku β :

$$f(\beta|\mathbf{v}) = \mathcal{N}\left(\mathbf{0}, \text{diag}(\mathbf{v})^{-1}\right) = \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} v_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & v_n \end{pmatrix}^{-1}\right), \quad (42)$$

- ▶ model pro prvky v_j volíme obdobně jako u ridge regression:

$$f(v_j) = \mathcal{G}(a_0, b_0), \quad \forall j. \quad (43)$$

Bayesovská sparse regression

VB aproximace

- ▶ a posteriori $\tilde{f}(\beta|v) = \mathcal{N}(\mu_\beta, \Sigma_\beta)$ a $\tilde{f}(v_j) = \mathcal{G}(a_j, b_j)$, $\forall j$, má (po odvození) následující tvarovací parametry:

$$\Sigma_\beta = \left(\hat{\omega} X^T X + \text{diag}(\hat{\mathbf{v}}) \right)^{-1} \quad (44)$$

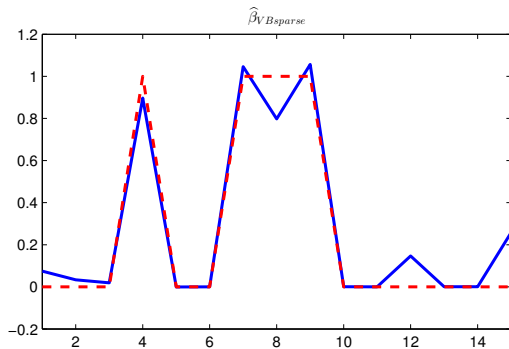
$$\mu_\beta = \Sigma_\beta \hat{\omega} X^T \mathbf{y} = \left(\hat{\omega} X^T X + \text{diag}(\hat{\mathbf{v}}) \right)^{-1} \hat{\omega} X^T \mathbf{y} \quad (45)$$

$$a_j = a_0 + \frac{1}{2} \quad (46)$$

$$b_j = b_0 + \frac{1}{2} \left(\widehat{\beta\beta'} \right)_{jj} \quad (47)$$

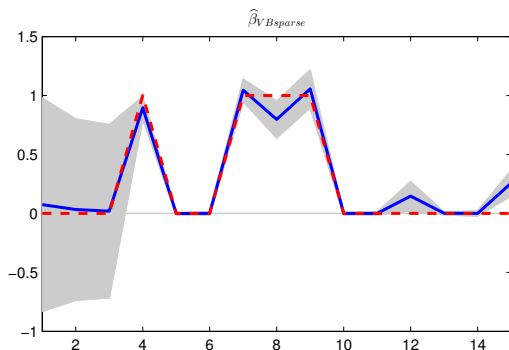
Bayesovská sparse regression

Řešení



Výhody bayesovského přístupu

- ▶ kvantifikace neurčitosti



- ▶ potlačení “ladicích” knoflíků, i když...
- ▶ nepřeborné modelovací možnosti

Dodatek

Co jsme neřešili a možná bychom měli...

- ▶ Pracovali jsme pouze s nagenеровaným příkladem bez fyzikální podstaty
- ▶ Co jsme vynechali?
 - ▶ pozitivita
 - ▶ vztah jednotlivých prvků β
 - ▶ informace o měřeních y
 - ▶ ...

Co nás čeká přístě

Aplikace zobecněných lineárních modelů

