

Efficient Bayesian Multivariate Surface Regression



Feng Li

feng.li@cufe.edu.cn

**School of Statistics and Mathematics
Central University of Finance and Economics**

Outline of the talk

- 1 Introduction to flexible regression models
- 2 The multivariate surface model
- 3 Application to firm leverage data
- 4 Extensions and future work

Flexible regression models

↪ Introduction

- Flexible models of the regression function $E(y|x)$ has been an active research field for decades.
- Attention has shifted from kernel regression methods to spline-based models.
- Splines are regression models with flexible mean functions.
- Example: a simple spline regression with only one explanatory variable with truncated linear basis function can be like this

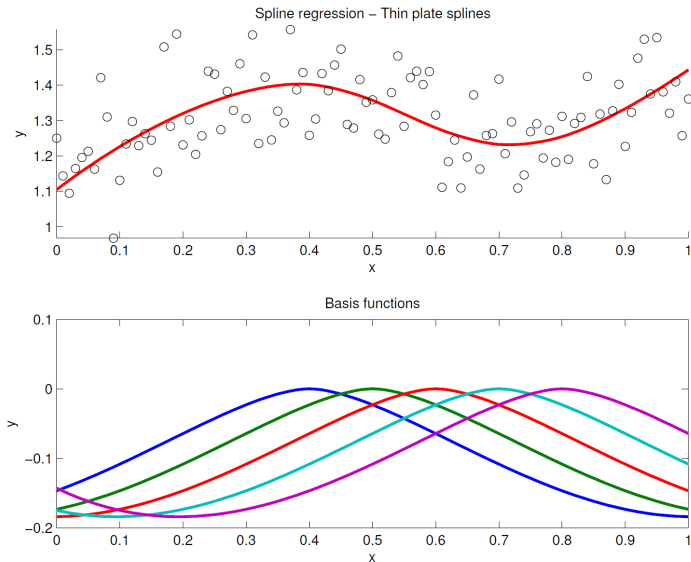
$$y = \alpha_0 + \alpha_1 x + \beta_1(x - \xi_1)_+ + \dots + \beta_q(x - \xi_q)_+ + \varepsilon$$

where

- $(x - \xi_i)_+$ are called the basis functions,
- ξ_i are called knots (the location of the basis function).

Flexible regression models

→ Spline example (single covariate with thinplate bases)



Flexible regression models

→ Spline regression with multiple covariates

- Additive spline model

- Each knot ξ_j (scalar) is connected with only one covariate

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_q x_q + \left[\sum_{j_1=1}^{m_1} \beta_{j_1} f(x_1, \xi_{j_1}) + \dots + \sum_{j_q=1}^{m_q} \beta_{j_q} f(x_q, \xi_{j_q}) \right] + \varepsilon$$

- Good and simple if you know there is no interactions in the data a priori.

- Surface spline model

- Each knot ξ_j (vector) is connected with more than one covariate

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_q x_q + \left[\sum_{j=1}^m \beta_j g(x_1, \dots, x_q, \xi_j) \right] + \varepsilon$$

- A popular choice of $g(x_1, \dots, x_q, \xi_j)$ can be e.g. the multi-dimensional thinplate spline

$$g(x_1, \dots, x_q, \xi_j) = \|\mathbf{x} - \xi_j\|^2 \ln \|\mathbf{x} - \xi_j\|$$

- Can handle the interactions but the model complexity increase dramatically with the interactive knots.

The challenges

- How many knots are needed?
 - Too few knots lead to a bad approximation; too many knots yield overfitting.
- Where to place those knots?
 - Equal spacing for the additive model,
 - which is obviously not efficient with the surface model.
- Common approaches to the two problems:
 - place enough many knots and use variable selection to pick up useful ones.
 - ★ not truly flexible
 - use reversible jump MCMC to move among the model spaces with different numbers of knots
 - ★ very sensitive to the prior and not computational efficient
 - clustering the covariates to select knots
 - ★ does not use the information from the responses
- How to choose between additive spline and surface spline?
 - NA

The multivariate surface model

↪ The model

- The multivariate surface model consists of three different components, *linear*, *surface* and *additive* as

$$\mathbf{Y} = \mathbf{X}_o \mathbf{B}_o + \mathbf{X}_s(\xi_s) \mathbf{B}_s + \mathbf{X}_a(\xi_a) \mathbf{B}_a + \mathbf{E}.$$

- We treat the knots ξ_i as unknown parameters and let them move freely.
 - A model with a minimal number of free knots outperforms model with lots of fixed knots.
- For notational convenience, we sometimes write model in compact form

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E},$$

where $\mathbf{X} = [\mathbf{X}_o, \mathbf{X}_s, \mathbf{X}_a]$ and $\mathbf{B} = [\mathbf{B}_o', \mathbf{B}_s', \mathbf{B}_a']'$ and $\mathbf{E} \sim \mathbf{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$

The multivariate surface model

↪ The prior

- Conditional on the knots, the prior for \mathbf{B} and Σ are set as

$$\text{vec}\mathbf{B}_i | \Sigma, \lambda_i \sim \mathbf{N}_q \left[\mu_i, \Lambda_i^{1/2} \Sigma \Lambda_i^{1/2} \otimes \mathbf{P}_i^{-1} \right], \quad i \in \{o, s, a\},$$
$$\Sigma \sim \text{IW}[n_0 \mathbf{S}_0, n_0],$$

- $\Lambda_i = \text{diag}(\lambda_i)$ are called the shrinkage parameters, which is used for overcome overfitting through the prior.
 - If $\mathbf{P}_i = \mathbf{I}$, can prevent singularity problem, like the ridge regression estimate.
 - If $\mathbf{P}_i = \mathbf{X}_i' \mathbf{X}_i$: use the covariates information, also a compressed version of least squares estimate when λ_i is large.
- The shrinkage parameters are estimated in MCMC
 - A small λ_i shrinks the variance of the conditional posterior for \mathbf{B}_i
 - It is another approach to selection important variables (knots) and components.
- We allow to mixed use the two types priors ($\mathbf{P}_i = \mathbf{I}$, $\mathbf{P}_i = \mathbf{X}_i' \mathbf{X}_i$) in different components in order to take the both the advantages of them.

The multivariate surface model

↪ The Bayesian posterior

- The posterior distribution is conveniently decomposed as

$$p(\mathbf{B}, \boldsymbol{\Sigma}, \boldsymbol{\xi}, \boldsymbol{\lambda} | \mathbf{Y}, \mathbf{X}) = p(\mathbf{B} | \boldsymbol{\Sigma}, \boldsymbol{\xi}, \boldsymbol{\lambda}, \mathbf{Y}, \mathbf{X}) p(\boldsymbol{\Sigma} | \boldsymbol{\xi}, \boldsymbol{\lambda}, \mathbf{Y}, \mathbf{X}) p(\boldsymbol{\xi}, \boldsymbol{\lambda} | \mathbf{Y}, \mathbf{X}).$$

- Hence $p(\mathbf{B} | \boldsymbol{\Sigma}, \boldsymbol{\xi}, \boldsymbol{\lambda}, \mathbf{Y}, \mathbf{X})$ follows the multivariate normal distribution according to the conjugacy;
- When $p = 1$, $p(\boldsymbol{\Sigma} | \boldsymbol{\xi}, \boldsymbol{\lambda}, \mathbf{Y}, \mathbf{X})$ follows the inverse Wishart distribution

$$\text{IW} \left[n_0 + n, \left\{ n_0 \mathbf{S}_0 + n \tilde{\mathbf{S}} + \sum_{i \in \{0, s, a\}} \boldsymbol{\Lambda}_i^{-1/2} (\tilde{\mathbf{B}}_i - \mathbf{M}_i)' \mathbf{P}_i (\tilde{\mathbf{B}}_i - \mathbf{M}_i) \boldsymbol{\Lambda}_i^{-1/2} \right\} \right]$$

- When $p \geq 2$, no closed form of $p(\boldsymbol{\Sigma} | \boldsymbol{\xi}, \boldsymbol{\lambda}, \mathbf{Y}, \mathbf{X})$, the above result is a very accurate approximation. Then the marginal posterior of $\boldsymbol{\Sigma}$, $\boldsymbol{\xi}$ and $\boldsymbol{\lambda}$ is

$$p(\boldsymbol{\Sigma}, \boldsymbol{\xi}, \boldsymbol{\lambda} | \mathbf{Y}, \mathbf{X}) = c \times p(\boldsymbol{\xi}, \boldsymbol{\lambda}) \times |\boldsymbol{\Sigma}_{\boldsymbol{\beta}}|^{-1/2} |\boldsymbol{\Sigma}|^{-(n+n_0+p+1)/2} |\boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}}|^{-1/2} \\ \times \exp \left\{ -\frac{1}{2} \left[\text{tr} \boldsymbol{\Sigma}^{-1} (n_0 \mathbf{S}_0 + n \tilde{\mathbf{S}}) + (\tilde{\boldsymbol{\beta}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}}^{-1} (\tilde{\boldsymbol{\beta}} - \boldsymbol{\mu}) \right] \right\}$$

The MCMC algorithm

↪ Metropolis-Hastings within Gibbs

- The coefficients (\mathbf{B}) are directly sampled from normal distribution.
- We update covariance ($\mathbf{\Sigma}$), all knots (ξ) and shrinkages (λ) jointly by using Metropolis-Hastings within Gibbs.
- The proposal density for $\mathbf{\Sigma}$ is the inverse Wishart density on previous slide.
- The proposal density for ξ and λ is a multivariate t -density with $\nu > 2$ df,

$$\theta_p | \theta_c \sim \text{MVT} \left[\hat{\theta}, - \left(\frac{\partial^2 \ln p(\theta | \mathbf{Y})}{\partial \theta \partial \theta'} \right)^{-1} \bigg|_{\theta = \hat{\theta}}, \nu \right],$$

where $\hat{\theta}$ is obtained by R steps ($R \leq 3$) Newton's iterations during the proposal with analytical gradients for matrices.

- The analytical gradients are very complicated and we have implemented it in an efficient way (**the key!**).

Application to firm leverage data

→ The data

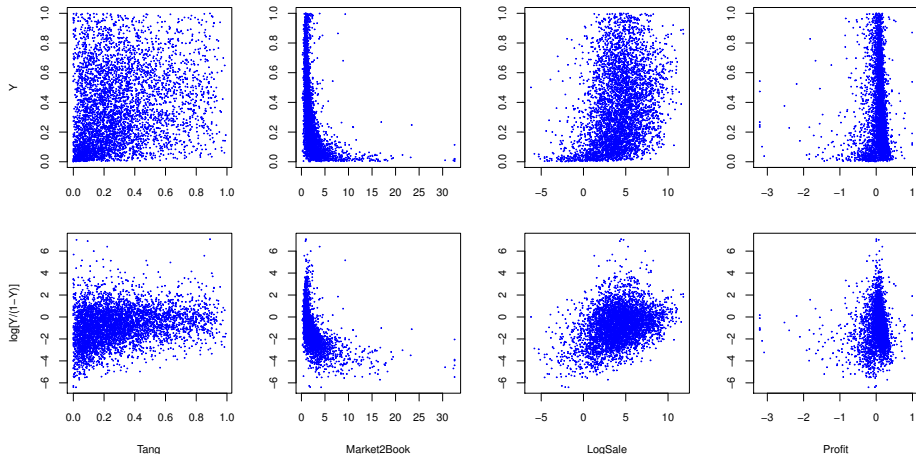
leverage (Y): total debt/(total debt+book value of equity), 4405 observations;

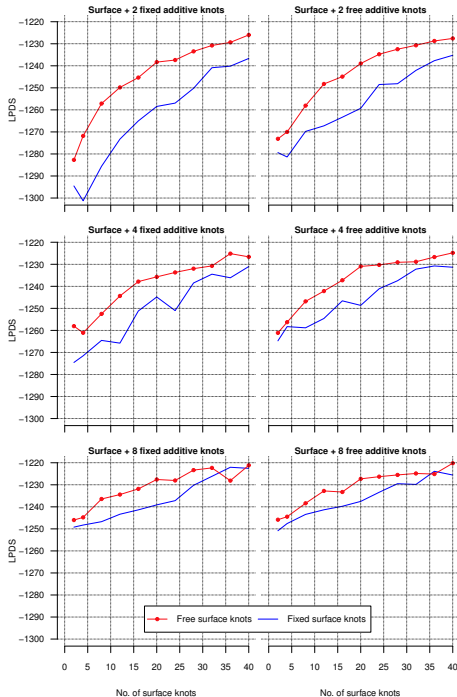
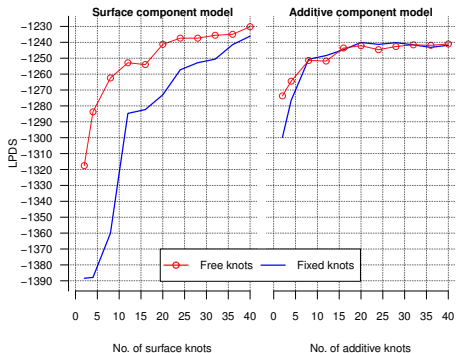
tang: tangible assets/book value of total assets;

market2book: (book value of total assets - book value of equity + market value of equity) / book value of total assets;

logSales: logarithm of sales;

profit: (earnings before interest, taxes, depreciation, and amortization) / book value of total assets.





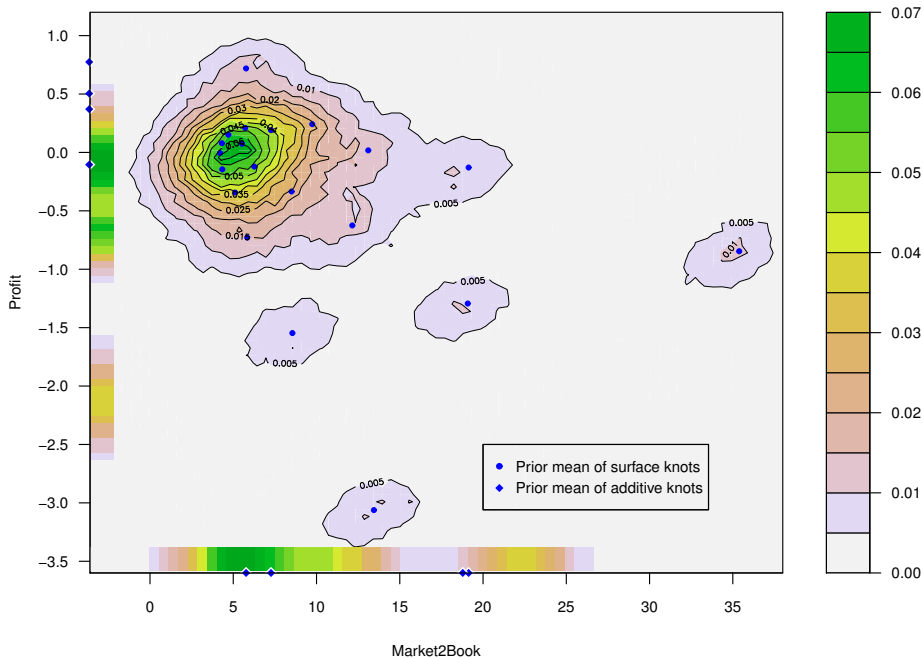
↑ Models with only surface or additive components
 → Model with both additive and surface components.

LPDS Log predictive density score which is defined as

$$\begin{aligned} \text{LPDS} &= \frac{1}{D} \sum_{d=1}^D \ln p(\tilde{Y}_d | \tilde{Y}_{-d}, \mathbf{X}) \\ &= \int \prod_{i \in \tau_d} p(\mathbf{y}_i | \boldsymbol{\theta}, \mathbf{x}_i) p(\boldsymbol{\theta} | \tilde{Y}_{-d}) d\boldsymbol{\theta}, \end{aligned}$$

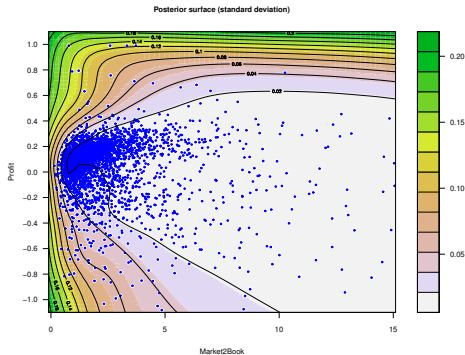
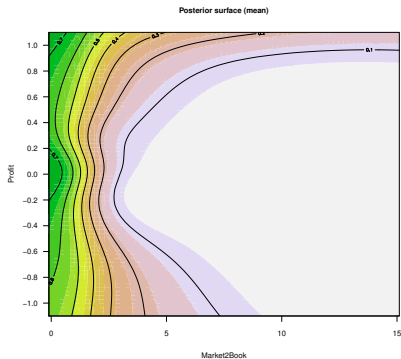
and $D = 5$ in the cross-validation.

Posterior locations of knots



Application to firm leverage data

➤ Posterior mean surface(left) and standard deviation(right)



Extensions and future work

- The model and the methods we used are very general.
- It is easy to generalize the model to GLM framework.
- Variable selection is possible for knots.
- Dirichlet precess prior can be plugged into the model when heteroscedasticity is the problem.
- And the copula...

Thank you!