Regression modelling using I-priors

Haziq Jamil Supervisors: Dr. Wicher Bergsma & Prof. Irini Moustaki

> Social Statistics (Year 1) London School of Economics & Political Science

> > 19 May 2015

PhD Presentation Event

Outline

Introduction

I-prior theory

8 Estimation methods

4 Examples of I-prior modelling

Simple linear regression 1-dimensional smoothing Multilevel modelling Longitudinal modelling

5 Further work

Structural Equation Models Models with structured error covariances Logistic models

Introduction •0000	I-prior theory 00000	Estimation methods	Examples of I-prior modelling	Further work	En

Linear regression

- Consider a set of data points $\{(y_1, x_1), \dots, (y_n, x_n)\}$.
- A model is linear if the relationship between y_i and the independent variables is linear.

•
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

•
$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$$

$$\flat \ y_i = \beta_0 x_i^{\beta_1 + 2\beta_2} + \epsilon_i \ \varkappa$$

In other words, the equations must be linear in the parameters.

-prior theory

stimation methods

Examples of I-prior modelling

Further work 000

Linear regression

• Definition (The linear regression model)

 $y_i = f(x_i) + \epsilon_i$

 $y_i \in \mathbb{R}$, real-valued observations $x_i \in \mathcal{X}$, a set of characteristics for unit i (1) $f \in \mathcal{F}$, a vector space of functions over the set \mathcal{X} $(\epsilon_1, \dots, \epsilon_n) \sim N(\mathbf{0}, \mathbf{\Psi}^{-1})$ $i = 1, \dots, n$

Note: For iid observations, $\Psi = \psi \mathbf{I}_n$. In general, $\Psi = (\psi_{ij})$.

l-prior theory 00000 Estimation methods

Examples of I-prior modelling

Further work

Linear regression



prior theory

Estimation methods

Examples of I-prior modelling

Further work

Estimation methods

How to pick the best line from the bag of stuff?

- Many ways Least squares, maximum likelihood, Bayesian...
- When dimensionality is large, may overfit. Solutions:
 - Dimension reduction
 - Random effects models
 - Regularization

...all require additional assumptions

I-priors

An I-prior on f is a distribution π on f such that its covariance matrix is the Fisher information of f. Also, assign a "best guess" on the prior mean, e.g. $f_0 = 0$.

l-prior theory 00000 Estimation methods

Examples of I-prior modelling

Further work Ei

Example: multiple regression

$$\mathbf{y} = \overbrace{\boldsymbol{\alpha} + \mathbf{X}\boldsymbol{\beta}}^{\mathrm{f}} + \boldsymbol{\epsilon}$$
$$\boldsymbol{\epsilon} \sim N(\mathbf{0}, \psi^{-1}\mathbf{I}_n)$$

We know from linear regression theory that $I[\beta] = \psi \mathbf{X}^T \mathbf{X}$. An I-prior on β is then

$$\boldsymbol{\beta} \sim N(\boldsymbol{\beta}_0, \lambda^2 \boldsymbol{\psi} \mathbf{X}^T \mathbf{X}).$$

Equivalently,

$$oldsymbol{eta} = oldsymbol{eta}_0 + \lambda \mathbf{X}^T \mathbf{w}$$

 $\mathbf{w} \sim N(\mathbf{0}, \psi \mathbf{I}_n).$

Thus, an I-prior on f is

$$\begin{split} \mathbf{f} &= \boldsymbol{\alpha} + \mathbf{X}\boldsymbol{\beta}_0 + \lambda \mathbf{X}\mathbf{X}^{\mathsf{T}}\mathbf{w} \\ & \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \psi \mathbf{I}_n). \end{split}$$

I-prior regression

I-prior theory

Functional vector spaces

Estimation methods

Examples of I-prior modelling

Further work

Inner products

Kernel methods

Reproducing kernels

Hilbert spaces

Gaussian random vectors

I-prior theory

Fisher Information

Krein spaces

Means of random functions

Feature maps

Variances of random functions

Random functions

Moore-Aronszajn Theorem

Haziq Jamil (LSE)

I-prior regression

19 May 2015 8 / 27

l-prior theory ●0000 Estimation methods

Examples of I-prior modelling

Further work En

Definitions & theorem

• Definition (Inner products)

Let \mathcal{F} be a vector space \mathbb{R} . A function $\langle \cdot, \cdot \rangle_{\mathcal{F}} : \mathcal{F} \times \mathcal{F} \to \mathbb{R}$ is said to be an inner product on \mathcal{F} if all of the following are satisfied:

- Symmetry: $\langle f,g \rangle_{\mathcal{F}} = \langle g,f \rangle_{\mathcal{F}}$
- Linearity: $\langle af_1 + bf_2, g \rangle_{\mathcal{F}} = a \langle f_1, g \rangle_{\mathcal{F}} + b \langle f_2, g \rangle_{\mathcal{F}}$
- Non-degeneracy: $\langle f, g \rangle_{\mathcal{F}} = 0 \Rightarrow f = 0$

for all $f, f_1, f_2, g \in \mathcal{F}$ and $a, b \in \mathbb{R}$. Additionally, an inner product is positive definite (negative definite) if $\langle f, f \rangle_{\mathcal{F}} \ge 0$ (≤ 0). An inner product is indefinite if it is neither positive nor negative definite.

• Definition (Hilbert space)

A positive definite inner product space which is complete, i.e. contains the limits of all Cauchy sequences.

• Definition (Krein space)

An (indefinite) inner product space which generalizes Hilbert spaces by dropping the positive definite restriction.

Haziq Jamil (LSE)

I-prior regression

I-prior theory

Estimation methods

Examples of I-prior modelling

Further work

Definitions & theorem

• Definition (Kernels)

Let \mathcal{X} be a non-empty set. A function $h : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a kernel if there exists a Hilbert space \mathcal{F} and a map $\phi : \mathcal{X} \to \mathcal{F}$ such that $\forall x, x' \in \mathcal{X}$,

$$h(x,x') = \langle \phi(x), \phi(x') \rangle.$$

• Definition (Reproducing kernels)

Let \mathcal{F} be a Hilbert/Krein space of functions over a non-empty set \mathcal{X} . A function $h: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a reproducing kernel of \mathcal{F} , and \mathcal{F} a RKHS/RKKS, if h satisfies

- $\flat \quad \forall x \in \mathcal{X}, \ h(\cdot, x) \in \mathcal{F}$
- $\forall x \in \mathcal{X}, f \in \mathcal{F}, \langle f, h(\cdot, x) \rangle_{\mathcal{F}} = f(x).$
- Kernel algorithms have many important uses in Machine Learning literature, such as pattern recognition, kernel PCA, finding distances of means in feature space, and many more.

Haziq Jamil (LSE)

I-prior regression

I-prior theory

Stimation methods

Examples of I-prior modelling

Further work

Definitions & theorem

Theorem (Gaussian I-priors) [Bergsma, 2014]
 For the linear regression model (1), let *F* be the RKKS with kernel
 h : X × X → ℝ. Then, assuming it exists, the Fisher information for
 f is given by

$$I[f](x_i, x_i') = \sum_{k=1}^n \sum_{l=1}^n \psi_{kl} h(x_i, x_k) h(x_i', x_l).$$

Let π be a Gaussian I-prior on f with prior mean f_0 and variance I[f]. Then π is called an I-prior for f, and a random vector $f \sim \pi$ has the random effect representation

$$f(x_i) = f_0(x_i) + \sum_{k=1}^n h(x_i, x_k) w_k$$

 $(w_1, \ldots, w_n) \sim N(\mathbf{0}, \mathbf{\Psi}).$

l-prior theory 000●0 Estimation methods

Examples of I-prior modelling

Further work

Back to the multiple regression example

• We saw the I-prior method applied to multiple regression:

$$f(\mathbf{x}_i) = \overbrace{\alpha + \mathbf{x}_i \beta_0}^{f_0(\mathbf{x}_i)} + \overbrace{\lambda(\mathbf{X}\mathbf{X}^T)_i \mathbf{w}}^{\sum_{k=1}^n h(\mathbf{x}_i, \mathbf{x}_k) w_k} \\ \mathbf{w} := (w_1, \dots, w_n) \sim N(\mathbf{0}, \psi \mathbf{I}_n).$$

I-prior theory

Estimation methods 0000 Examples of I-prior modelling

Further work

Back to the multiple regression example

• We saw the I-prior method applied to multiple regression:

$$f(\mathbf{x}_i) = \overbrace{\alpha + \mathbf{x}_i \beta_0}^{f_0(\mathbf{x}_i)} + \overbrace{\lambda(\mathbf{X}\mathbf{X}^{\mathsf{T}})_i \mathbf{w}}^{\sum_{k=1}^n h(\mathbf{x}_i, \mathbf{x}_k) w_k} \\ \mathbf{w} := (w_1, \dots, w_n) \sim N(\mathbf{0}, \psi \mathbf{I}_n).$$

• Choose different RKHS/RKKS \mathcal{F} and corresponding *h* to suit the type/characteristic of the **x**s in order to do regression modelling.

I-prior theory

Estimation methods

Examples of I-prior modelling

Further work

Back to the multiple regression example

• We saw the I-prior method applied to multiple regression:

$$f(\mathbf{x}_i) = \overbrace{\alpha + \mathbf{x}_i \beta_0}^{f_0(\mathbf{x}_i)} + \overbrace{\lambda(\mathbf{X}\mathbf{X}^T)_i \mathbf{w}}^{\sum_{k=1}^n h(\mathbf{x}_i, \mathbf{x}_k) w_k}$$
$$\mathbf{w} := (w_1, \dots, w_n) \sim N(\mathbf{0}, \psi \mathbf{I}_n).$$

• Choose different RKHS/RKKS \mathcal{F} and corresponding *h* to suit the type/characteristic of the **x**s in order to do regression modelling.



I-prior theory 0000● Estimation methods

Examples of I-prior modelling

Further work E

Toolbox of RKHS/RKKS

$\mathcal{X} = \{x_i\}$	Characteristic/Uses	Vector space ${\cal F}$	Kernel $h(x_i, x_k)$
Nominal	1) Categorical covariates; 2) In a multilevel setting, x_i = group no. of unit <i>i</i> .	Pearson	$\frac{\mathbb{I}[x_i = x_k]}{\substack{p_i \\ \mathcal{P}[X = x_i]}} - 1 \text{ where } p_i = $
Real	As in classical regression, x_i = real-valued covariate associated with unit <i>i</i> .	Canonical	x _i x _k
Real	As in (1-dim) smoothing, $x_i = \text{data point associated}$ with observation y_i .	Fractional Brownian Motion (FBM)	$ x_i ^{2\gamma}+ x_k ^{2\gamma}- x_i-x_k ^{2\gamma}$ with $\gamma \in (0,1)$

- We can construct new RKHS/RKKS from existing ones.
 - ► Example (ANOVA RKKS) Set of x_i = (x_{1i}, x_{2i}) of Nominal + Real characteristics. Then

$$h(x_i, x_i') = h_1(x_{1i}, x_{1i}') + h_2(x_{2i}, x_{2i}') + h_1(x_{1i}, x_{1i}')h_2(x_{2i}, x_{2i}')$$

prior theory

Estimation methods

Examples of I-prior modelling

Further work En

Parameters to be estimated

- Let's choose a prior mean of zero (or set an overall constant/intercept to be estimated).
- For the I-prior linear model

$$y_{i} = \alpha + \sum_{k=1}^{n} h_{\lambda}(x_{i}, x_{k})w_{k} + \epsilon_{i}$$

$$\epsilon_{i} \sim N(0, \psi^{-1})$$

$$w_{i} \sim N(0, \psi)$$

$$i = 1, \dots, n,$$
(2)

the parameters to be estimated are $\boldsymbol{\theta} = (\alpha, \lambda, \psi)^T$.

• λ is introduced to resolve the arbitrary scale of an RKKS/RKHS \mathcal{F} over a set \mathcal{X} . Number of λ parameters = number of kernels used, not interactions nor covariates.

Haziq Jamil (LSE)

EM algorithm

rior theory

Estimation methods •000 Examples of I-prior modelling

Further work

- For the I-prior model in (2), treat the w_is as missing.
- The distributions are easy enough to obtain:

•
$$\mathbf{y} \sim N(\boldsymbol{\alpha}, \mathbf{V}_y)$$
, where $\mathbf{V}_y := \mathbf{H}_\lambda \mathbf{\Psi} \mathbf{H}_\lambda + \mathbf{\Psi}^{-1}$
• $\mathbf{w} \sim N(\mathbf{0}, \mathbf{\Psi})$

$$\blacktriangleright \begin{pmatrix} \mathbf{y} \\ \mathbf{w} \end{pmatrix} \sim N \left(\begin{pmatrix} \boldsymbol{\alpha} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{y} & \mathbf{H}_{\lambda} \mathbf{\Psi} \\ \mathbf{\Psi} \mathbf{H}_{\lambda} & \mathbf{\Psi} \end{pmatrix} \right)$$

•
$$\mathbf{w}|\mathbf{y} \sim N\left(\Psi \mathbf{H}_{\lambda} \mathbf{V}_{y}^{-1}(\mathbf{y} - \boldsymbol{\alpha}), \mathbf{V}_{y}^{-1}\right)$$

where $\mathbf{H}_{\lambda}(i, j) = h_{\lambda}(x_{i}, x_{j})$ and $\Psi = \psi \mathbf{I}_{n}$

• E-step: Calculate
$$Q(\theta) = \mathsf{E}_{w} [\log f(\mathbf{y}, \mathbf{w}; \theta) | \mathbf{y}; \theta_{t}].$$

• M-step:
$$\boldsymbol{\theta}_{t+1} \leftarrow \operatorname{arg\,max}_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}).$$

-prior theory

Estimation methods 0000

Examples of I-prior modelling

Further work

Generalised least square estimator for α

- Write Model (2) as $\mathbf{y} = \alpha \mathbf{1} + \mathbf{H}_{\lambda} \mathbf{w} + \boldsymbol{\epsilon}$, where $\mathbf{y} \sim N(\boldsymbol{\alpha}, \mathbf{V}_{y})$.
- Assume values for λ and ψ are known, and thus too $\mathbf{V}_{y}(\lambda,\psi) = \mathbf{H}_{\lambda} \Psi \mathbf{H}_{\lambda} + \Psi^{-1}$.
- GLS estimator for α is

$$\hat{\alpha} = (\mathbf{1}^T \mathbf{V}_y^{-1} \mathbf{1})^{-1} (\mathbf{1}^T \mathbf{V}_y^{-1} \mathbf{y}).$$

• This turns out to be identical to the MLE.

-prior theory

Estimation methods

Examples of I-prior modelling

Further work

Exponential family EM algorithm

- Consider a density function belonging to the exponential family with the (canonical) form $f_{\mathbf{X}}(\mathbf{x}; \theta) = \exp[\theta \cdot \mathbf{T}(\mathbf{x}) A(\theta)]h(\mathbf{x})$.
 - The MLE is found by solving the set of equations $T(x) = A'(\theta)$.
 - It is also know that $A'(\theta) = \mathsf{E}[\mathsf{T}(\mathsf{x}); \theta]$.
- In the EM algorithm, the "full" data is $\mathbf{x} = (\mathbf{y}, \mathbf{w})$. The E-step involves calculating $Q(\boldsymbol{\theta})$, and for the exponential family, this turns out to be

$$Q(\boldsymbol{\theta}) = \mathsf{E}_{\mathsf{w}} \left[\boldsymbol{\theta} \cdot \mathsf{T}(\mathsf{y}, \mathsf{w}) - A(\boldsymbol{\theta}) + \log h(\mathsf{y}, \mathsf{w}) | \mathsf{y}; \boldsymbol{\theta}_t \right].$$

• Maximising this over heta, we arrive at the FOC

$$\begin{aligned} Q'(\boldsymbol{\theta}) &= \mathsf{E}_{\mathbf{w}} \left[\boldsymbol{\theta} \cdot \mathbf{T}(\mathbf{y}, \mathbf{w}) | \mathbf{y}; \boldsymbol{\theta}_t \right] - \mathcal{A}'(\boldsymbol{\theta}) = \mathbf{0} \\ &\Rightarrow \mathsf{E}_{\mathbf{w}} \left[\boldsymbol{\theta} \cdot \mathbf{T}(\mathbf{y}, \mathbf{w}) | \mathbf{y}; \boldsymbol{\theta}_t \right] = \mathsf{E}[\mathbf{T}(\mathbf{y}, \mathbf{w}); \boldsymbol{\theta}]. \end{aligned}$$

prior theory 0000 Estimation methods

Examples of I-prior modelling

Further work

Full Bayesian approach

- Assign prior distributions to the parameters, for example
 - $\alpha \sim N(a, b^2)$
 - $\lambda \sim U(0, c)$
 - ψ ~ Γ(d, e)
- Draw from the posterior densities $f(\theta|\mathbf{y})$ using Metropolis-Hastings algorithm. Estimates for the parameters are the posterior means.
- Easy to implement in R using JAGS (rjags or R2jags), but...



-prior theory

Estimation methods

Examples of I-prior modelling

Further work

Example: Simple linear regression

Classical model

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
$$\epsilon_i \sim N(0, \sigma)$$

I-prior model

$$y_i = \alpha + \sum_{k=1}^n h_\lambda(x_i, x_k) w_k + \epsilon_i$$
$$\epsilon_i \sim N(0, \psi^{-1})$$
$$w_i \sim N(0, \psi)$$

 h_{λ} is the Canonical kernel



Haziq Jamil (LSE)

I-prior regression

19 May 2015 19 / 27

Introduction I 00000 0

prior theory

Estimation methods

Examples of I-prior modelling

Further work

Example: 1-dimensional smoothing

Classical model

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3$$
$$\epsilon_i \sim N(0, \sigma)$$

I-prior model

$$y_i = \alpha + \sum_{k=1}^n h_{\lambda,\gamma}(x_i, x_k) w_k + \epsilon_i$$
$$\epsilon_i \sim N(0, \psi^{-1})$$
$$w_i \sim N(0, \psi)$$

 $h_{\lambda,\gamma}$ is the FBM kernel



Haziq Jamil (LSE)

I-prior regression

19 May 2015 20 / 27

prior theory

Estimation methods

Examples of I-prior modelling $\circ \circ \bullet \circ$

Further work

Example: Multilevel modelling

Classical model

$$y_{ij} = \beta_{0j} + \beta_{1j} x_{ij} + \epsilon_{ij}$$
$$\begin{pmatrix} \beta_{0j} \\ \beta_{1j} \end{pmatrix} \sim N\left(\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \begin{pmatrix} \phi_0 & \phi_{01} \\ \phi_{01} & \phi_1 \end{pmatrix} \right)$$
$$\epsilon_{ij} \sim N(0, \sigma)$$

I-prior model

$$y_i = \alpha + \sum_{k=1}^n h_\lambda(x_i, x_k) w_k + \epsilon_i$$
$$\epsilon_i \sim N(0, \psi^{-1})$$
$$w_i \sim N(0, \psi)$$

 h_{λ} is the ANOVA kernel





MSE(classical) = 0.227 MSE(I-prior) = 0.226

Haziq Jamil (LSE)

I-prior regression

19 May 2015 21 / 27

-prior theory

Estimation methods

Examples of I-prior modelling

Further work

Example: Longitudinal modelling



prior theory

Estimation methods

Examples of I-prior modelling

Further work

Further work: Structural Equation Models

• The 1-factor model

$$egin{aligned} x_{ij} &= \mu_j + \lambda_j f_i + \delta_{ij} \ f_i &\sim \mathcal{N}(0,1) \ \delta_{ij} &\sim \mathcal{N}(0, heta_j) \end{aligned}$$

• Relationship to longitudinal random intercept model:

• Set
$$\mu_j = \mu$$
, $\forall j$.

- Set $\lambda_j = 1$, , $\forall j$ and estimate variance of f_i instead.
- Set $\theta_j = \theta$, $\forall j$ We already know how to estimate this model using I-prior.
- Further work:
 - Uses of this very restricted CFA model? Rasch model?
 - Post estimation work, e.g. obtaining factor scores.
 - Can we estimate both the λ_j s and f_i simultaneously?

prior theory

Estimation methods

Examples of I-prior modelling

Further work En ○●○

Further work: Structured error covariances

- Sometimes, the responses may be correlated in a way that the model specification can't account for completely. Extend model to allow for dependence between errors, such as autocorrelations.
- Example: AR(1) covariance matrix with equal gaps between observations:

$$\Psi = \frac{\sigma^2}{1 - \phi^2} \begin{pmatrix} 1 & \phi & \phi^2 & \cdots & \phi^{n-1} \\ \phi & 1 & \phi & \cdots & \phi^{n-2} \\ \phi^2 & \phi & 1 & \cdots & \phi^{n-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi^{n-1} & \phi^{n-2} & \phi^{n-3} & \cdots & 1 \end{pmatrix}$$

• Others: Heteroskedastic errors?

-prior theory

Estimation methods

Examples of I-prior modelling

Further work ○○●

00

En

Further work: Logistic models

• Extending the I-prior methodology to GLMs, e.g. logit models:

$$y_i \sim \mathsf{Bern}(\pi_i)$$

logit $\pi_i = lpha + \sum_{k=1}^n h_\lambda(x_i, x_k) w_k$
 $w_i \sim N(0, \pi_i(1 - \pi_i))$
 $i = 1, \dots, n$

i.e. putting an I-prior on the linear predictor, and setting the Fisher information as the variance.

- Difficulties faced
 - Unable to estimate this model using JAGS due to a circular dependence of the parameters.
 - Performing ML yields a high-dimensional intractable integral. Poor results from approximation methods like Laplace and Gauss-Hermite Quadrature.

Haziq Jamil (LSE)

I-prior regression

19 May 2015 25 / 27

Introduction 00000	l-prior theory 00000	Estimation methods	Examples of I-prior modelling	Further work	Enc			
Summary								

- The I-prior methodology is a modelling technique that guards against overfitting linear models when dimensionality is large relative to sample size, with advantages such as
 - Model parsimony
 - Requires no additional assumptions
 - Simpler estimation
- Many models shown to work with using I-priors such as multiple regression, smoothing models, random effects models and growth curve models.
- Areas of research include
 - Extension to GLMs
 - Structural Equation Models
 - Models with structured error covariances

-prior theory 00000 Estimation methods 0000

Examples of I-prior modelling

Further work

End



Thank you!