To-do list					
(Coi	nte	nts		
4	Reg	ressio	n modelling using I-priors	1	
	4.1	Variou	us regression models	. 1	
		4.1.1	Multiple linear regression	. 1	
		4.1.2	Multilevel linear modelling	. 4	
		4.1.3	Longitudinal modelling	. 4	
		4.1.4	Smoothing models	. 4	
		4.1.5	Models with functional covariates	. 4	
	4.2	Estim	ation	. 4	
		4.2.1	The intercept \dots	. 4	
		4.2.2	Direct optimisation	. 4	
		4.2.3	Expectation-maximisation algorithm	. 5	
		4.2.4	Markov chain Monte Carlo methods	. 6	
		4.2.5	Computational considerations	. 6	
		4.2.6	Comparison of estimation methods	. 8	
	4.3	Exam	ples	. 8	
	4.4	Concl	usion	. 8	
Bi	ibliog	graphy	•	11	
				1.0	
L1	st of	Figure	es	12	
Γ.i	st of	Table	S	13	

List of Theorems	14	
List of Definitions	15	
List of Symbols	16	

Haziq Jamil

Department of Statistics

London School of Economics and Political Science

March 11, 2018

Chapter 4

Regression modelling using I-priors

In the previous chapter, we defined an I-prior for the normal regression model (1.1) subject to (1.2) and f belonging to a reproducing kernel Hilbert or Krein space of functions. We also saw how new function spaces can be constructed via the polynomial and ANOVA RKKS. In this chapter, we shall describe various regression models, and connect them to an appropriate RKKS, so that an I-prior may be defined on it. Methods for estimating I-prior models will also be described. Finally, several examples of I-prior modelling are presented.

4.1 Various regression models

4.1.1 Multiple linear regression

Let $\mathcal{X} \equiv \mathbb{R}^p$ equipped with the regular dot product, and \mathcal{F}_{λ} the scaled canonical RKHS of functions over \mathcal{X} with kernel $h_{\lambda}(\mathbf{x}, \mathbf{x}') = \lambda \mathbf{x}^{\top} \mathbf{x}'$. Then, an I-prior on f implies that

$$f(\mathbf{x}_i) = \sum_{j=1}^n \lambda \mathbf{x}_i^{\top} \mathbf{x}_j w_j$$
$$= \sum_{j=1}^n \lambda \sum_{k=1}^p x_{ik} x_{jk} w_j$$
$$= \beta_1 x_{i1} + \dots + \beta_p x_{ip},$$

where each $\beta_k = \lambda \sum_{j=1}^n x_{jk} w_j$. This implies a multivariate normal prior distribution for the regression coefficients

$$\boldsymbol{\beta} := (\beta_1, \dots, \beta_p) \sim \mathrm{N}_p(\mathbf{0}, \lambda^2 \mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X}),$$

where **X** is the $n \times p$ design matrix for the covariates, excluding the column of ones at the beginning typically reserved for the intercept. The covariance matrix for β is recognised as the scaled Fisher information matrix for the regression coefficients.

The I-prior for β resembles the objective g-prior (Zellner, 1986) for regression coefficients,

$$\boldsymbol{\beta} \sim \mathrm{N}_p \left(\mathbf{0}, g(\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X})^{-1} \right),$$

although they are quite different objects. The g-prior for β has the *inverse* (scaled) Fisher information matrix as its covariance matrix. This, in itself, has a much different and arguably counterintuitive meaning: large amounts of Fisher information about β corresponds to a small prior variance, and hence less deviation away from the prior mean of zero in estimating β . The choice of the hyperparameter g has been the subject of much debate, with choices ranging from fixing g = n (corresponding to the concept of unit Fisher information), to fully Bayesian and empirical Bayesian methods of estimating g from the data.

On the other hand, we note that the g-prior has an I-prior interpretation when argues as follows. Assume that the regression function f lies in the continual dual space of \mathbb{R}^p equipped with the inner product $\langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}} = \mathbf{x}^{\top} (\mathbf{X}^{\top} \mathbf{\Psi} \mathbf{X})^{-1} \mathbf{x}$. With this inner product and from (3.3), the Fisher information for β is

$$\mathcal{I}_g(\boldsymbol{\beta}) = \sum_{i=1}^n \sum_{j=1}^n \psi_{ij} (\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X})^{-1} \mathbf{x}_i \otimes (\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X})^{-1} \mathbf{x}_j$$
$$= (\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X})^{-1} (\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X}) (\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X})^{-1}$$
$$= (\mathbf{X}^\top \boldsymbol{\Psi} \mathbf{X})^{-1},$$

and this rather than the usual $\mathbf{X}^{\top} \mathbf{\Psi} \mathbf{X}$ as the prior covariance matrix for $\boldsymbol{\beta}$ means that the I-prior is in fact the standard g-prior.

The metric induced by the inner product is actually the *Mahalanobis distance*, a scale-invariant natural distance if the covariates are measured on different scales. To expand on this idea, circle back to the regression function and write it as $f(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\beta} \rangle_{\mathcal{X}}$. In

usual least squares regression, the choice of inner product is irrelevant, so the usual dot product is commonly used (however, as we have seen above, the choice of inner product determines the form of the Fisher information for β). In particular, suppose that all the x_{ik} 's, k = 1, ..., p for each unit i = 1, ..., n are measured on the same scale; for instance, these could be measurements in centimetres. In this case, the dot product is reasonable, because $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = \sum_{k=1}^p x_{ik} x_{jk}$ and the inner product has a coherent unit, namely the squared unit of the x_{ik} 's. However, if they were a mix of various scaled measurements, then obviously the inner product's unit is incoherent—one would be resorted to adding measurements in different units, for example, cm² and kg² and so on. In such a case, a unitless inner product is appropriate, like the Mahalonobis inner product, which technically rescales the x_{ik} 's to unity. In summary, if the covariates are all measured on the same scale, then the I-prior is appropriate, and if not, the g-prior is appropriate.

A different approach for covariate measurements of differing scales, without resorting to g-priors, is the ANOVA approach. By considering only the main effects, one decomposes the regression function into

$$f(\mathbf{x}_i) = \beta_0 + f_1(x_{i1}) + \dots + f_p(x_{ip})$$

for which $(f - \beta_0) \in \mathcal{F}_{\lambda} \equiv \mathcal{F}_{\lambda_1} \oplus \cdots \oplus \mathcal{F}_{\lambda_p}$, and $\mathcal{F}_{\lambda_k}, k = 1, \ldots, p$ are unidimensional centred canonical RKKSs with kernels $h_{\lambda_k}(x_{ik}, x_{jk}) = \lambda_k(x_{ik} - \bar{x}_k)(x_{jk} - \bar{x}_k)$, where $\bar{x}_k = \sum_{i=1}^n x_{ik}/n$. In effect, we now have p scale parameters, one for each of the RKKSs associated with the p covariates. Denote $\tilde{x}_{ik} = x_{ik} - \bar{x}_k$, the centred covariates. The RKKS \mathcal{F}_{λ} therefore has the kernel

$$h(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^p \lambda_k \tilde{x}_{ik} \tilde{x}_{jk},$$

and hence each regression coefficient can now be written as $\beta_k = \sum_{j=1}^n \lambda_k \tilde{x}_{jk} w_j$. Thus, the corresponding I-prior for β is

$$\boldsymbol{\beta} \sim N_p(\mathbf{0}, \lambda^2 \tilde{\mathbf{X}}^{\top} \boldsymbol{\Lambda} \boldsymbol{\Psi} \boldsymbol{\Lambda} \tilde{\mathbf{X}}),$$

with $\mathbf{\Lambda} = \operatorname{diag}(\lambda_1, \dots, \lambda_p)$. Note that the overall effect β_0 can be treated in a number of ways, the simplest of which is as an intercept to be estimated. This approach is disadvantageous when p is large, in which case there would be numerous scale parameters to estimate.

2. Can I just

4.1.2 Multilevel linear modelling

4.1.3 Longitudinal modelling

4.1.4 Smoothing models

4.1.5 Models with functional covariates

4.2 Estimation

4.2.1 The intercept

Given the regression model (??) subject to an I-prior (1.5), the marginal likelihood of the intercept α (after integrating out the I-prior) can be maximised with respect to α , which yields the sample mean for y as the ML estimate for intercept.

4.2.2 Direct optimisation

The kernel parameter η and the error precision ψ (which we collectively refer to as the model hyperparameters of the covariance kernel θ) can be estimated in several ways. One of these is direct optimisation of the marginal log-likelihood—the most common method in the Gaussian process literature.

$$\log L(\theta) = \log \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f}$$
$$= -\frac{n}{2}\log 2\pi - \frac{1}{2}\log |\mathbf{\Sigma}_{\theta}| - \frac{1}{2}\mathbf{y}^{\top}\mathbf{\Sigma}_{\theta}^{-1}\mathbf{y}$$

where $\Sigma_{\theta} = \psi \mathbf{H}_{\eta}^2 + \psi^{-1} \mathbf{I}_n$. This is typically done using conjugate gradients with a Cholesky decomposition on the covariance kernel to maintain stability, but the **iprior** package opts for an eigendecomposition of the kernel matrix (Gram matrix) $\mathbf{H}_{\eta} = \mathbf{V} \cdot \operatorname{diag}(u_1, \ldots, u_n) \cdot \mathbf{V}^{\top}$ instead. Since \mathbf{H}_{η} is a symmetrix matrix, we have that $\mathbf{V}\mathbf{V}^{\top} = \mathbf{I}_n$, and thus

$$\Sigma_{\theta} = \mathbf{V} \cdot \operatorname{diag}(\psi u_1^2 + \psi^{-1}, \dots, \psi u_n^2 + \psi^{-1}) \cdot \mathbf{V}^{\top}$$

for which the inverse and log-determinant is easily obtainable. This method is relatively robust to numerical instabilities and is better at ensuring positive definiteness of the covariance kernel. The eigendecomposition is performed using the **Eigen C++** template

library and linked to **iprior** using **Rcpp** (Eddelbuettel and Francois, 2011). The hyperparameters are transformed by the **iprior** package so that an unrestricted optimisation using the quasi-Newton L-BFGS algorithm provided by optim() in R. Note that minimisation is done on the deviance scale, i.e., minus twice the log-likelihood. The direct optimisation method can be prone to local optima, in which case repeating the optimisation at different starting points and choosing the one which yields the highest likelihood is one way around this.

4.2.3 Expectation-maximisation algorithm

Alternatively, the expectation-maximisation (EM) algorithm may be used to estimate the hyperparameters, in which case the I-prior formulation in (??) is convenient. Substituting this into (??) we get something that resembles a random effects model. By treating the w_i as "missing", the tth iteration of the E-step entails computing

$$Q(\theta) = E \left[\log p(\mathbf{y}, \mathbf{w} | \theta) | \mathbf{y}, \theta^{(t)} \right]. \tag{4.1}$$

As a consequence of the properties of the normal distribution, the required joint and posterior distributions $p(\mathbf{y}, \mathbf{w})$ and $p(\mathbf{w}|\mathbf{y})$ are easily obtained. The M-step then maximises the Q function above, which boils down to solving the first order conditions

$$\frac{\partial Q}{\partial n} = -\frac{1}{2} \operatorname{tr} \left(\frac{\partial \mathbf{\Sigma}_{\theta}}{\partial n} \tilde{\mathbf{W}}^{(t)} \right) + \psi \cdot \mathbf{y}^{\top} \frac{\partial \mathbf{H}_{\eta}}{\partial n} \tilde{\mathbf{w}}^{(t)}$$
(4.2)

$$\frac{\partial Q}{\partial \psi} = -\frac{1}{2} \mathbf{y}^{\mathsf{T}} \mathbf{y} - \operatorname{tr} \left(\frac{\partial \mathbf{\Sigma}_{\theta}}{\partial \psi} \tilde{\mathbf{W}}^{(t)} \right) + \mathbf{y}^{\mathsf{T}} \mathbf{H}_{\eta} \tilde{\mathbf{w}}^{(t)}$$
(4.3)

equated to zero. Here, $\tilde{\mathbf{w}}$ and \mathbf{W} are the first and second posterior moments of \mathbf{w} . The solution to (4.3) can be found in closed-form, but not necessarily for (4.2). In cases where closed-form solutions exist, then it is just a matter of iterating the update equations until a suitable convergence criterion is met (e.g. no more sizeable increase in successive log-likelihood values). In cases where closed-form solutions do not exist for θ , the Q function is again optimised with respect to θ using the L-BFGS algorithm.

The EM algorithm is more stable than direct maximization, and is especially suitable if there are many scale parameters. However, it is typically slow to converge. The **iprior** package provides a method to automatically switch to the direct optimisation method after running several EM iterations. This then combines the stability of the EM with the speed of direct optimisation.

4.2.4 Markov chain Monte Carlo methods

For completeness, it should be mentioned that a full Bayesian treatment of the model is possible, with additional priors on the hyperparameters set. Markov chain Monte Carlo (MCMC) methods can then be employed to sample from the posteriors of the hyperparameters, with point estimates obtained using the posterior mean or mode, for instance. Additionally, the posterior distribution encapsulates the uncertainty about the parameter, for which inference can be made. Posterior sampling can be done using Gibbsbased methods in WinBUGS (Lunn et al., 2000) or JAGS (Plummer, 2003), and both have interfaces to R via R2WinBUGS (Sturtz et al., 2005) and runjags (Denwood, 2016) respectively. Hamiltonian Monte Carlo (HMC) sampling is also a possibility, and the Stan project (Carpenter et al., 2017) together with the package rstan (Stan Development Team, 2016) makes this possible in R. All of these MCMC packages require the user to code the model individually, and we are not aware of the existence of MCMC-based packages which are able to estimate GPR models. This makes it inconvenient for GPR and I-prior models, because in addition to the model itself, the kernel functions need to be coded as well and ensuring computational efficiency would be a difficult task. Note that this full Bayesian method is not implemented in **iprior**, but described here for completeness.

4.2.5 Computational considerations

Computational complexity for estimating I-prior models (and in fact, for GPR in general) is dominated by the inversion of the $n \times n$ matrix $\Sigma_{\theta} = \psi \mathbf{H}_{\eta}^2 + \psi^{-1} \mathbf{I}_n$, which scales as $O(n^3)$ in time. As mentioned earlier, the **iprior** package inverts this by way of the eigendecomposition of \mathbf{H}_{η} , but this operation is also $O(n^3)$. For the direct optimisation method, this matrix inversion is called when computing the log-likelihood, and thus must be computed at each Newton step. For the EM algorithm, this matrix inversion appears when calculating $\tilde{\mathbf{w}}$, the posterior mean of the I-prior random effects. Furthermore, storage requirements for I-priors models are similar to that of GPR models, which is $O(n^2)$.

The machine learning literature is rich in ways to resolve this issue, as summarised by Quiñonero-Candela and Rasmussen, 2005. One such method is to use low-rank matrix approximations. Let \mathbf{Q} be a matrix with rank q < n, and that $\mathbf{Q}\mathbf{Q}^{\top}$ can be used to

approximate the kernel matrix \mathbf{H}_{η} . Then

$$(\psi \mathbf{H}_{\eta}^2 + \psi^{-1} \mathbf{I}_n)^{-1} \approx \psi \left[\mathbf{I}_n - \mathbf{Q} \left((\psi^2 \mathbf{Q}^{\top} \mathbf{Q})^{-1} + \mathbf{Q}^{\top} \mathbf{Q} \right)^{-1} \mathbf{Q}^{\top} \right],$$

obtained via the Woodbury matrix identity, is a potentially much cheaper operation which scales $O(nq^2)$ — $O(q^3)$ to do the inversion, and O(nq) to do the multiplication (because typically the inverse is premultiplied to a vector). When the kernel matrix itself is sufficiently low ranked (for instance, when using the linear kernel for a low-dimensional covariate) then the above method is exact. However, other interesting kernels such as the fractional Brownian motion (fBm) kernel or the squared exponential kernel results in kernel matrices which are full rank.

Another method of approximating the kernel matrix, and the method implemented by our package, is the Nyström method (Williams and Seeger, 2001). The theory has its roots in approximating eigenfunctions, but this has since been adopted to speed up kernel machines. The main idea is to obtain an (approximation to the true) eigendecomposition of \mathbf{H}_{η} based on a small subset $m \ll n$ of the data points. Reorder the rows and columns and partition the kernel matrix as

$$\mathbf{H}_{\eta} = egin{pmatrix} \mathbf{A}_{m imes m} & \mathbf{B}_{m imes (n-m)} \ \mathbf{B}_{m imes (n-m)} & \mathbf{C}_{(n-m) imes (n-m)} \end{pmatrix}.$$

The Nyström method provides an approximation to the lower right block \mathbf{C} by manipulating the eigenvectors and eigenvalues of \mathbf{A} , an $m \times m$ matrix, together with the matrix \mathbf{B} to give

$$\mathbf{H}_{\eta} pprox egin{pmatrix} \mathbf{V}_m \ \mathbf{B}^{ op} \mathbf{V}_m \mathbf{U}_m^{-1} \end{pmatrix} \mathbf{U}_m egin{pmatrix} \mathbf{V}_m^{ op} & \mathbf{U}_m^{-1} \mathbf{V}_m^{ op} \mathbf{B} \end{pmatrix}$$

where \mathbf{U}_m is the diagonal matrix containing the m eigenvalues of \mathbf{A} , and \mathbf{V}_m is the corresponding matrix of eigenvectors. An orthogonal version of this approximation is of interest, which has been studied by Fowlkes et al., 2001, which allows us to easily calculate the inverse of Σ_{θ} . Estimating I-prior models using the Nyström method takes $O(nm^2)$ times and O(nm) storage.

4.2.6 Comparison of estimation methods

4.3 Examples

4.4 Conclusion

The steps for I-prior modelling are basically three-fold:

- 1. Select an appropriate function space; equivalently, the kernels for which a specific effect is desired on the covariates. Several modelling examples are described in Section 4.1.
- 2. Estimate the hyperparameters (these included the RKHS scale parameter(s), error precision, and any other kernel parameters such as the Hurst index of fBm) of the I-prior model and obtain the posterior regression function.
- 3. Post-estimation procedures include
 - Posterior predictive checks;
 - Model comparison via log-likelihood ratio tests/empirical Bayes factors; and
 - Prediction of new data point.

The main sticking point with the estimation procedure is the involvement of the $n \times n$ kernel matrix, for which its inverse is needed. This requires $O(n^2)$ storage and $O(n^3)$ computational time. The Nyström method of approximating the kernel matrix reduces complexity to O(nm) storage and approximately $O(nm^2)$, and is highly advantageous if $m \ll n$. The computational issue faced by I-priors are mirrored in Gaussian process regression, so the methods to overcome these computational challenges in GPR can be explored further. However, most efficient computational solutions exploit the nature of the SE kernel structure, which is the most common kernel used in GPR.

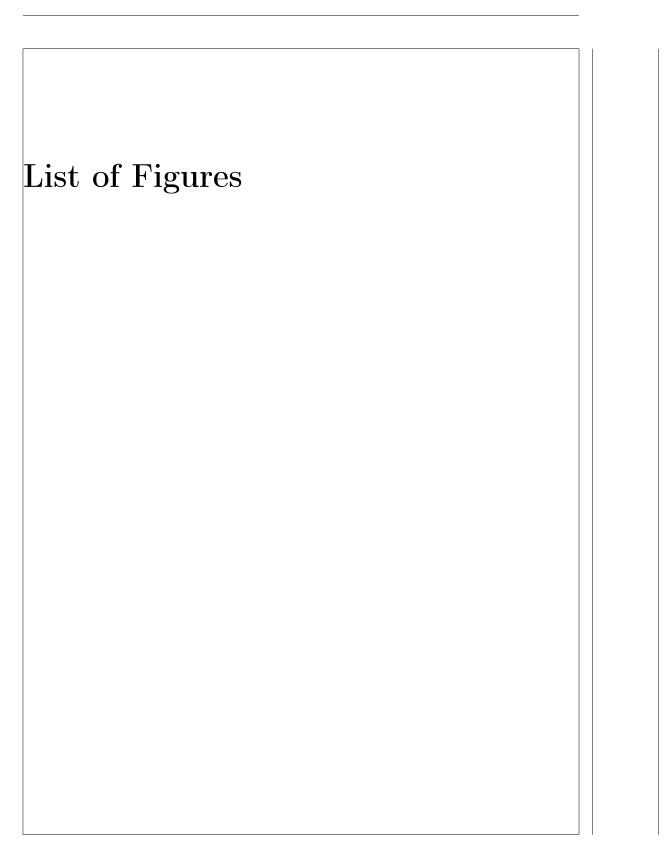
One promising avenue to achieve efficient computation for I-prior models is by using variational methods. A sparse variational approximation (typically by using inducing points) or stochastic variational inference can greatly reduce computational storage and speed requirements. A recent paper by Cheng and Boots (2017) suggested a variational algorithm with linear complexity for GPR-type models.

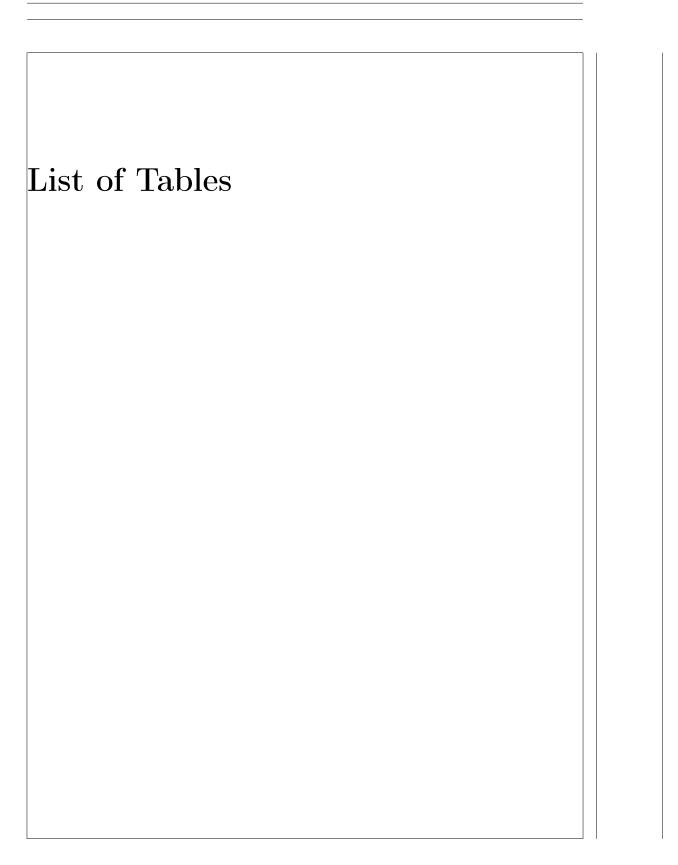
Omitted	

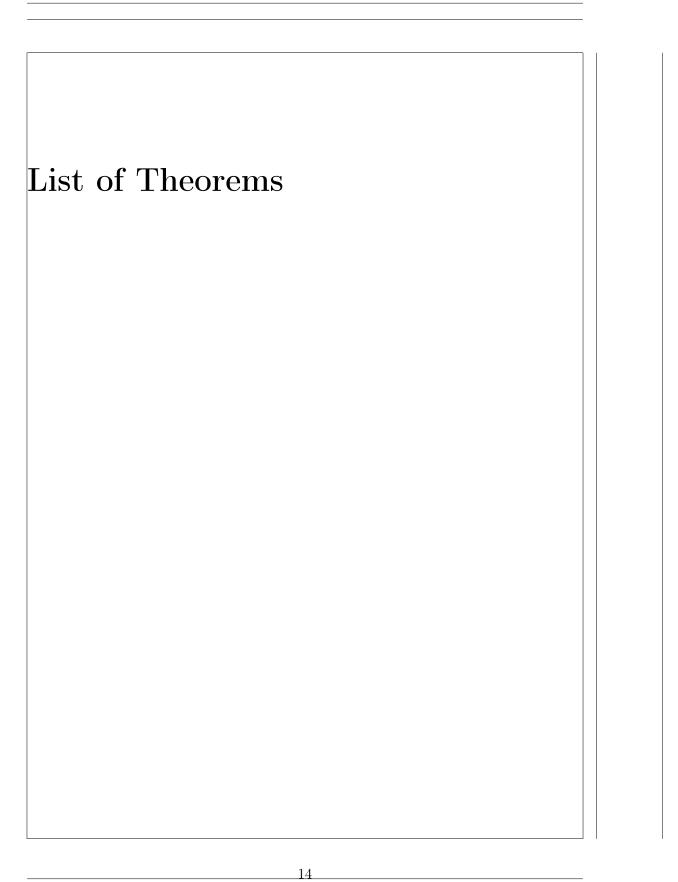
Bibliography

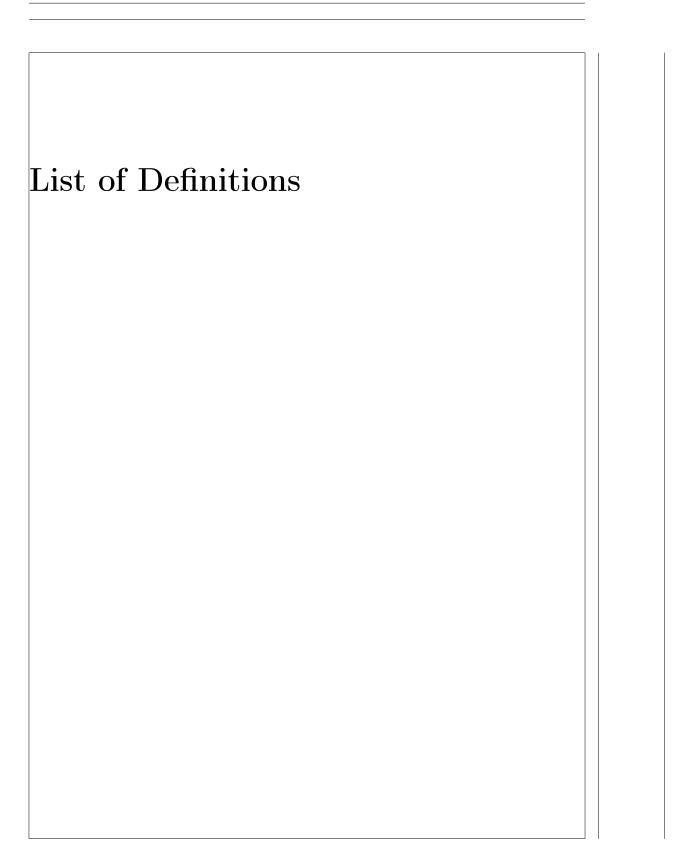
- Carpenter, B., A. Gelman, M. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell (2017). "Stan: A Probabilistic Programming Language". In: Journal of Statistical Software, Articles 76.1, pp. 1–32. DOI: 10.18637/jss.v076.i01.
- Cheng, C.-A. and B. Boots (2017). "Variational Inference for Gaussian Process Models with Linear Complexity". In: Advances in Neural Information Processing Systems, pp. 5190–5200.
- Denwood, M. (2016). "**runjags**: An R Package Providing Interface Utilities, Model Templates, Parallel Computing Methods and Additional Distributions for MCMC Models in JAGS". In: *Journal of Statistical Software* 71.9, pp. 1–25. DOI: 10.18637/jss.v071.i09.
- Eddelbuettel, D. and R. Francois (2011). "Rcpp: Seamless R and C++ Integration". In: Journal of Statistical Software 40.8, pp. 1–18. DOI: 10.18637/jss.v040.i08.
- Fowlkes, C., S. Belongie, and J. Malik (2001). "Efficient Spatiotemporal Grouping Using the Nyström Method". In: Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2001). Vol. 1, pp. 231–238. DOI: 10.1109/CVPR.2001.990481.
- Lunn, D. J., A. Thomas, N. Best, and D. Spiegelhalter (Oct. 2000). "WinBUGS A Bayesian modelling framework: Concepts, structure, and extensibility". In: Statistics and Computing 10.4, pp. 325–337. DOI: 10.1023/A:1008929526011.
- Plummer, M. (2003). "JAGS: A Program for Analysis of Bayesian Graphical Models Using Gibbs Sampling". In: *Proceedings of the 3rd International Workshop on Distributed Statistical Computing*. Vol. 124. Vienna, Austria, p. 125.

Quiñonero-Candela, J. and C. E. Rasmussen (Dec. 2005). "A Unifying View of Sparse Approximate Gaussian Process Regression". In: Journal of Machine Learning Research 6, pp. 1939-1959. Stan Development Team (2016). RStan: The R Interface to Stan. R package version 2.14.1. URL: http://mc-stan.org/. Sturtz, S., U. Ligges, and A. Gelman (2005). "R2WinBUGS: A Package for Running WinBUGS from R". In: Journal of Statistical Software 12.3, pp. 1–16. DOI: 10.18637/ jss.v012.i03. Williams, C. K. I. and M. Seeger (2001). "Using the Nyström Method to Speed Up Kernel Machines". In: Advances in Neural Information Processing Systems 13. The MIT Press, pp. 682–688. Zellner, A. (1986). "On assessing prior distributions and Bayesian regression analysis with g-prior distributions". In: Bayesian inference and decision techniques.









List of Symbols

 $N_p(\mu, \Sigma)$ p-dimensional multivariate normal distribution with mean vector μ and covariance Σ .

 \sim Is distributed as.

 $\delta_{xx'}$ The Kronecker delta: $\delta_{xx'} = 1$ if x = x', and 0 otherwise.

 \otimes The tensor product.

ndex	
10071	
alysis of variance, see ANOVA	reproducing kernel Hilbert space, <i>see</i> RKHS
ctional Brownian motion, see fBm	squared exponential, $see~\mathrm{SE}$