

Deep Bayesian regression models

Hubin A.A., Storvik G.O., Frommlet F.

Department of Mathematics, University of Oslo

aliaksah@math.uio.no



UiO : Universitetet i Oslo

Yandex School of Data Analysis,
Moscow 2018

16.03.2018

- Regression models are addressed for inference and prediction in a wide range of applications;
- More and more sources of data are becoming available introducing a variety of hypothetical explanatory variables for these models to be considered;
- Model selection and averaging of different combinations of covariates in this context becomes extremely important for both good inference and prediction;

Introduction. Issues

- It is often the case that linear relations between the explanatory variables and the response are not sufficient;
- One has to avoid unreasonably deep non-linearities to avoid overfitting;
- Ideally models should remain as transparent and dense as possible, or quoting Einstein's famous *"It can scarcely be denied that the supreme goal of all theory is to make the irreducible basic elements as simple and as few as possible without having to surrender the adequate representation of a single datum of experience."*

In this work we introduce a class of deep Bayesian regression models and suggest algorithmic approaches for fitting them.

Deep Bayesian regression models

$$Y_i | \mu_i, \phi \sim f(y | \mu_i; \phi), \quad i \in \{1, \dots, n\} \quad (1)$$

$$\mu_i = h^{-1} \left(\beta_0 + \sum_{j=1}^p \gamma_j \beta_j F_j(\mathbf{x}) + \sum_{k=1}^r \gamma_{k+p} \delta_{ik} \right), \quad (2)$$

$$\delta_{\mathbf{k}} = (\delta_{1k}, \dots, \delta_{nk}) \sim N_n(\mathbf{0}, \mathbf{\Sigma}_k). \quad (3)$$

- $f(\cdot | \mu, \phi)$ is a density/distribution with expectation μ and dispersion parameter ϕ ;
- $F_j(\mathbf{x})$ are all features based on the input explanatory variables ordered w.r.t. complexity, p is the finite number of allowed features;
- $\beta_j \in \mathbb{R}, j \in \{1, \dots, p\}$ are regression coefficients of the features;
- $h(\cdot)$ is a proper link function;
- $\gamma_j \in \{0, 1\}, j \in \{p+1, \dots, q = r+p\}$ are latent indicators defining if a feature is included into the model ($\gamma_j = 1$) or not ($\gamma_j = 0$).

Hierarchy of the features

A feature $F_j(\mathbf{x})$ can be constructed recursively through:

$$F_j(\mathbf{x}) = \begin{cases} v(F_k(\mathbf{x})), & \text{for a modification;} \\ F_k(\mathbf{x}) * F_l(\mathbf{x}), & \text{for a crossover;} \\ v(\alpha^T \mathbf{F}(\mathbf{x})), & \text{for a projection;} \end{cases}$$

- $F_k(\mathbf{x})$ and $F_l(\mathbf{x})$ are previously defined features ($k, l < j$);
- $v \in \mathcal{G}$ is one of the allowed basic function from set \mathcal{G} ;
- $\mathbf{F}(\mathbf{x})$ is a sub-vector of all possible features with indexes lower than j ;
- A constraint on the complexity of feature $F_j(\mathbf{x})$ is defined by a finite number q of all possible features;
- Projections include modifications and crossovers as particular cases.

Types and meaning of functions in \mathcal{G}

- ANN: $\text{logit}(x)$, $\tanh(x)$, $\text{erf}(x)$, $\text{ReLU}(x)$;
- Polynomials: $F_k(\mathbf{x}) * F_l(\mathbf{x}) = \exp(\log(F_k(\mathbf{x})) + \log(F_l(\mathbf{x})))$;
- Logical *AND* and *OR*: $L_k \wedge L_l = L_k * L_l$ and $L_k \vee L_l = L_k + L_l - L_k * L_l$;
- CART: $I(x \geq 1)$;
- Fourier series: $\sin(x)$ and $\cos(x)$;
- Fractional polynomials: $x^{\frac{1}{n}}$;
- RNN and Lagged features: $\text{lag}^k(x)$.

The universal approximation theorem by *Hornik (1991)* is applicable if at least one of v functions is monotonous and bounded.

$$p(\gamma) \propto \mathbb{I}(|\gamma_{1:p}| \leq Q) \mathbb{I}(|\lambda_{p+1:q}| \leq R) \prod_{j=1}^p a^{\gamma_j w_j c(F_j(\mathbf{x}))} \prod_{k=p+1}^q b^{\gamma_k \omega_k c(\delta_k)}. \quad (4)$$

- $a, b \in (0, 1)$, $Q \leq p$, $R \leq r$;
- $|\gamma_{1:K}| = \sum_{j=1}^K \gamma_j$ is the number of active features in subset $\{\gamma_1, \dots, \gamma_K\}$;
- $c(F_j(\mathbf{x})) \geq 0$ is a measure of complexity for a feature $F_j(\mathbf{x})$;
- $c(\delta_k) \geq 0$ is a measure of complexity for a latent Gaussian variable δ_k ;
- w_j are weights for complexities of the corresponding features;
- ω_k are weights for the complexities of the corresponding latent Gaussian variables.

Parameter priors

Particular choices are given in the applications to follow:

$$\beta|\gamma \sim \pi_\beta(\beta), \quad (5)$$

$$\psi_k|\gamma \sim \pi_k(\psi_k), \quad (6)$$

$$\phi \sim \pi_\phi(\phi). \quad (7)$$

Prior distributions on $\beta_j|\gamma$, ϕ (if present) and $\psi_k|\gamma$ (if latent Gaussian variables are present) are usually selected in a way to efficiently compute marginal likelihoods of the models (by for example specifying conjugate priors) and should be carefully specified for the applications of interest.

Inference on the model

Let:

- $\gamma = \{\gamma_1, \dots, \gamma_q\}$ define a model itself, i.e. which components(features or latent variables) are addressed;
- $\theta_\gamma = (\beta, \psi, \phi)$ define parameters of the model γ .

Goals:

- $p(\gamma, \theta_\gamma | \mathbb{D})$ posterior distribution of parameters and models;
- $p(\gamma | \mathbb{D})$ marginal posterior probabilities of the models;
- $p(\Delta | \mathbb{D})$ marginal posterior probabilities of the quantiles of interest Δ .

But:

- $\exists 2^q$ different models in Ω_γ ;
- q is huge;
- Ω_γ is extremely difficult to specify.

Marginal likelihood approximation.

Marginal likelihood:

$$p(\mathbb{D}|\gamma) = \int_{\Theta} p(\mathbb{D}|\theta_{\gamma}, \gamma) p(\theta_{\gamma}|\gamma) d\theta_{\gamma} \quad (8)$$

can be obtained analytically or approximated by various methods:

- Laplace approximations;
- Integrated nested Laplace approximations;
- Chib's method (via Gibbs sampling);
- Chib and Jeliazkov's method (via MH algorithm).

Computation of the marginal likelihood is the major computational bottleneck currently. No efficient methods with subsampling are currently available.

Possible pipeline

- **Notice that** $p(\gamma, \theta_\gamma | \mathbb{D}) = p(\theta_\gamma | \gamma, \mathbb{D}) p(\gamma | \mathbb{D})$;
- Here $p(\mathbb{D} | \gamma)$ can be obtained by LA or similarly;
- **Notice that** $p(\gamma | \mathbb{D}) = \frac{p(\mathbb{D} | \gamma) p(\gamma)}{\sum_{\gamma' \in \Omega_\gamma} p(\mathbb{D} | \gamma') p(\gamma')}$;
- **Approximate with**

$$\widehat{p}(\gamma | \mathbb{D}) = \frac{p(\mathbb{D} | \gamma) p(\gamma)}{\sum_{\gamma' \in \mathbb{V}} p(\mathbb{D} | \gamma') p(\gamma')} \quad (9)$$

- \mathbb{V} is the **subspace** of Ω_γ to be **efficiently explored**;
- **Near modal values in terms of "MLIK \times prior" are particularly important** for construction of reasonable $\mathbb{V} \subset \Omega_\gamma$, **missing them can dramatically influence** posterior in the original space Ω_γ .

MJMCMC is efficient, but...

In Hubin and Storvik [6] we suggested efficient mode jumping proposals in the **discrete parameter spaces**. But Ω_γ must be clearly specified for MJMCMC. The later is **not feasible** in **Deep Regression Models**.

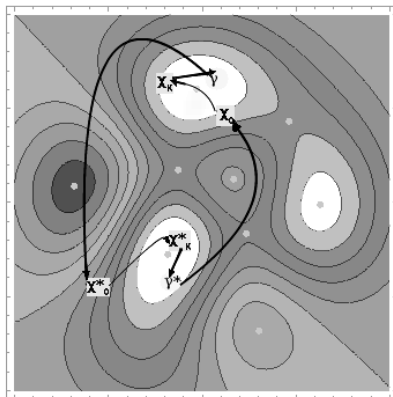


Figure: Locally optimized with randomization proposals

Genetically modified MJMCMC. Idea

- MJMCMC is embedded in the iterative setting of a genetic algorithm. In each iteration only a given set \mathcal{S} of features (of fixed size d) is considered;
- Each \mathcal{S} then induces a separate **search space** for MJMCMC or in the language of genetic algorithms \mathcal{S} is the **population**;
- \mathcal{S} dynamically evolves as a Markov chain of populations through $\{\mathcal{S}_0, \dots, \mathcal{S}_{t_{\max}}\}$ to allow MJMCMC explore different reasonable parts of the in-feasibly large total search space;
- Each $\mathcal{S}_t, t \in \{1, \dots, t_{\max}\}$ is selected from the neighborhood \mathcal{N}_{t-1} of \mathcal{S}_{t-1} . \mathcal{N}_{t-1} includes all populations feasible by performing **mutation**, **crossover**, **reduction**, **modification**, **projection** and **filtration** operations to the current \mathcal{S}_{t-1} and a special set \mathcal{F} ;
- Utilization of the approximation (9) allows us to compute marginal inclusion probabilities

$$\widehat{p}(\gamma_j = 1 | \mathbb{D}) = \sum_{\gamma \in \mathbb{V}} \mathbf{I}(\gamma_j = 1) \widehat{p}(\gamma | \mathbb{D}). \quad (10)$$

Genetically modified MJMCMC. Pipeline

- \mathcal{S}_0 is the set of p_0 input features;
- \mathcal{S}_1 is constructed by:
 - ① Running MJMCMC for a given number of iterations N_{init} on \mathcal{S}_0 ;
 - ② The first $d_1 < d$ members of population \mathcal{S}_1 are then defined by **filtration** operation, whilst $p - d_1$ filtered features from \mathcal{S}_0 are kept in \mathcal{F} ;
 - ③ The remaining $d - d_1$ members of \mathcal{S}_1 are obtained by means of the **crossover**, **mutation**, **modification**, **projection** and **reduction** operations applied to \mathcal{S}_0 ;
- All other $\mathcal{S}_t, t \in \{2, \dots, t_{max}\}$ are constructed by:
 - ① Running MJMCMC for a given number of iterations N_{expl} on \mathcal{S}_{t-1} ;
 - ② The first $d_t \leq d$ members of population \mathcal{S}_t are then defined by **filtration** operation;
 - ③ The remaining $d - d_t$ members of \mathcal{S}_t are obtained by means of the **crossover**, **mutation**, **modification**, **projection** and **reduction** operations applied to \mathcal{S}_{t-1} and \mathcal{F} .

Filtration operation. \mathcal{S}_0 case

$$\widehat{p}(\mathcal{F}_1|\mathbb{D}) \leq \dots \leq \widehat{p}(\mathcal{F}_{p_0-d_1}|\mathbb{D}) \leq \widehat{p}(F_1^1|\mathbb{D}) \leq \widehat{p}(F_2^1|\mathbb{D}) \leq \dots \leq \widehat{p}(F_{d_1}^1|\mathbb{D}) \quad (11)$$

$$\widehat{p}(F_1^1|\mathbb{D}) \geq p_s^o \quad (12)$$

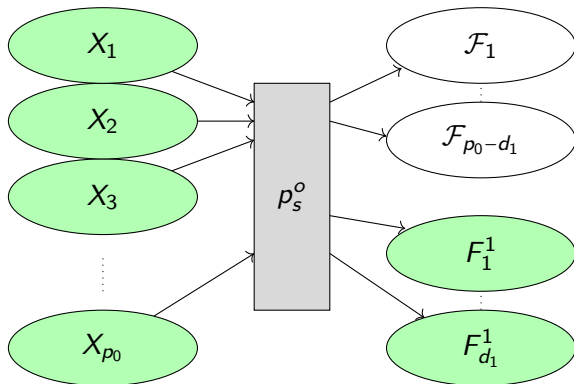


Figure: Feature filtering.

Filtration operation. $\mathcal{S}_t, t \in \{1, \dots, t_{max}\}$ case

$$\widehat{p}(D_1^{t+1}|\mathbb{D}) \leq \dots \leq \widehat{p}(D_{p-d_{t+1}}^{t+1}|\mathbb{D}) \leq \widehat{p}(F_{d_1+1}^{t+1}|\mathbb{D}) \leq \dots \leq \widehat{p}(F_{d_{t+1}}^{t+1}|\mathbb{D}) \quad (13)$$

$$\widehat{p}(F_{d_1+1}^{t+1}|\mathbb{D}) \geq p_s^t \quad (14)$$

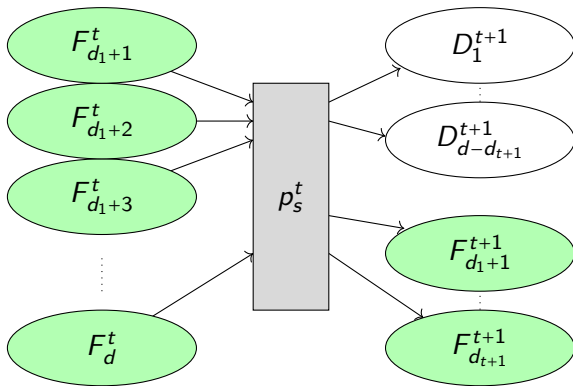


Figure: Feature filtering.

Crossover. Inbreeding of parents

Within each **mutation** or **crossover** $*$ is used for inbreeding.

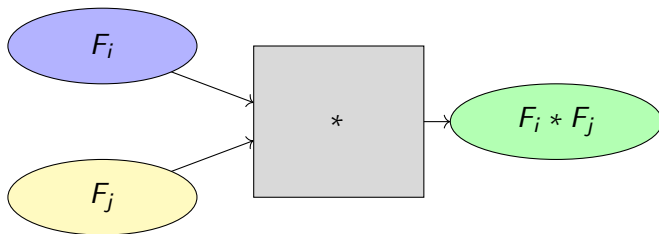


Figure: Crossover feature engineering step illustration.

Crossover operations. Parents selection

Crossovers inbreed parents from population \mathcal{S}_t and allow parents from \mathcal{F} with a smaller probability.

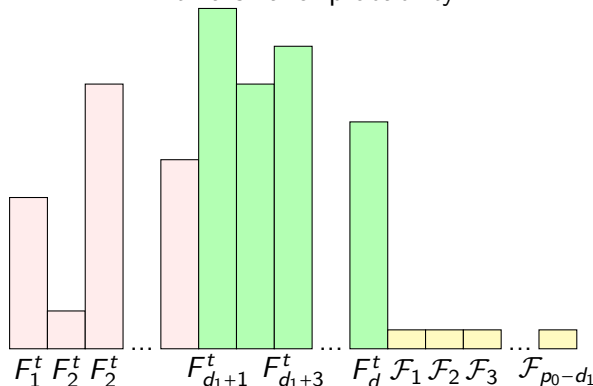


Figure: Class probabilities for selection of parents proportional to current marginal inclusion probabilities.

Mutation operations

Mutations replace features from population \mathcal{S}_t that are filtered by members of \mathcal{F} .

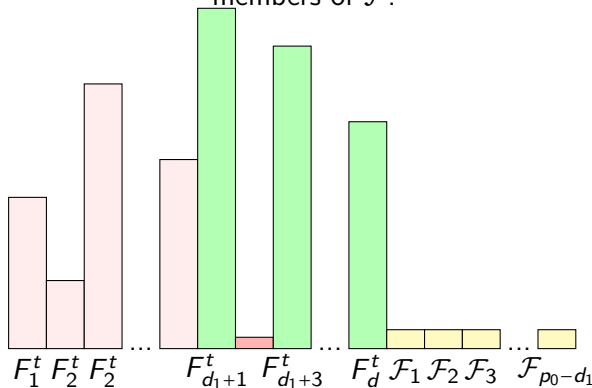


Figure: Red - a feature to be replaced

Modification operator

- Perform a functional transformation of some existing feature.

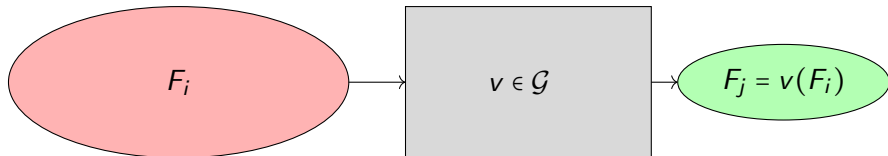


Figure: Feature modification step illustration.

Semi affine projection operator

- Perform a functional transformation of a linear combination of some existing features.

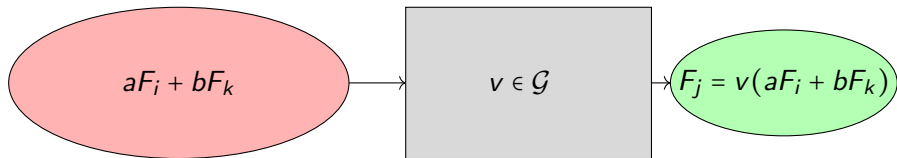


Figure: Feature projection step illustration.

Reduction operator

- Reductions are applied for the features greater than C ;
- Some parts are independently deleted with Bernoulli probability p_d ;
- The survived parts are stucked together.

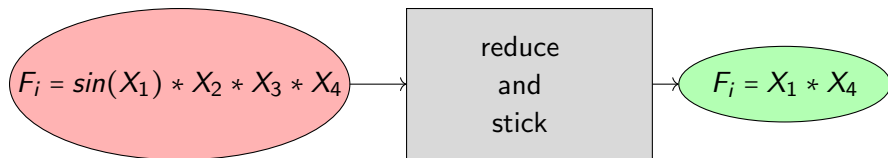


Figure: Feature reduction step illustration.

Reversible version of GMJMCMC - RGMJMCMC

- For every search space we suggest running MJMCMC for a sufficient amount of time;
- Otherwise a transition $\gamma \rightarrow \mathcal{S}' \rightarrow \gamma^1 \rightarrow \dots \rightarrow \gamma^k \rightarrow \gamma'$ is considered with a given probability kernel;
 - $q(\mathcal{S}'|\gamma)$ is the proposal for the new search space;
 - Transitions $\gamma^1 \rightarrow \dots \rightarrow \gamma^k$ are generated by local MJMCMC \mathcal{S}' ;
 - Transition $\gamma^k \rightarrow \gamma'$ is some randomization at the end of the procedure.
- Acceptance probability for such a procedure is $r_m = \min \{1, \alpha_m\}$

$$\alpha_m = \frac{\pi(\gamma')q(\gamma|\gamma'_k)}{\pi(\gamma)q(\gamma'|\gamma_k)}. \quad (15)$$

Genetically modified MJMCMC. Embarrassingly parallelized

- 1 Run B GMJMCMC chains in parallel with different seeds on separate CPUs or clusters;
- 2 Combine all unique models visited by all B chains into \mathbb{V} ;
- 3 Compute model posteriors as (9);
- 4 Compute marginal inclusion probabilities as (10) ;
- 5 Compute posteriors of other parameters of interest as

$$\widehat{p}(\Delta|\mathbb{D}) = \sum_{\gamma \in \mathbb{V}} p(\Delta|\gamma, \mathbb{D}) \widehat{p}(\gamma|\mathbb{D}) . \quad (16)$$

Simulation scenario

For this scenario we generated $N = 100$ datasets with $n = 1000$ observations and $p_0 = 50$ binary covariates. The covariates were assumed to be independent and were simulated as $X_j \sim \text{Bernoulli}(0.5)$ for $j \in \{1, \dots, 50\}$.

Gaussian responses

Gaussian observations with error variance $\sigma^2 = 1$ and individual expectations:
$$E(Y) = 1 + 1.5 X_7 + 1.5 X_8 + 6.6 X_{18} * X_{21} + 3.5 X_2 * X_9 + 9 X_{12} * X_{20} * X_{37} + 7 X_1 * X_3 * X_{27} + 7 X_4 * X_{10} * X_{17} * X_{30} + 7 X_{11} * X_{13} * X_{19} * X_{50}$$

The model

For this example as well as 2 following examples we have addressed the following DGLMM with conditionally independent Gaussian observations:

$$Y_i | \mu_i \sim N(\mu_i, \sigma^2), i \in \{1, \dots, n\}, \quad (17)$$

$$\mu_i = \gamma_0 \beta_0 + \sum_{j=1}^p \gamma_j \beta_j F_j(\mathbf{x}_i). \quad (18)$$

The following priors were addressed the first 3 examples γ :

$$p(\gamma) \propto \prod_{j=1}^p \exp(-2 \log n \gamma_j c(F_j(\mathbf{x}))) \quad (19)$$

$$p(\beta | \gamma) = |J_n^\gamma(\hat{\beta})|^{\frac{1}{2}}, \quad (20)$$

$$\pi(\sigma^2) = \sigma^{-2}, \quad (21)$$

where $c(F_j(\mathbf{x}))$ is the sum of the number additive, multiplicative and non-linear complications in feature j and $|J_n^\gamma(\hat{\beta})|$ is the determinant of the corresponding Fisher information matrix. The non-linear modifications in this example include $\cos(x)$, $\text{sigmoid}(x)$, $\tanh(x)$, $\text{atan}(x)$, $\sin(x)$, and $|x|^{\frac{1}{3}}$.

Simulation scenario. Results over 100 datasets

Table: Results for the three simulation scenario. Power for individual expression, overall power, expected number of false positives (FP), and FDP.

	GMJ	RGMJ	GMJ(logic)
X_7	1.0000	1.0000	0.9900
X_8	1.0000	1.0000	1.0000
$X_2 * X_9$	1.0000	0.9600	1.0000
$X_{18} * X_{21}$	1.0000	1.0000	0.9600
$X_1 * X_3 * X_{27}$	1.0000	1.0000	1.0000
$X_{12} * X_{20} * X_{37}$	1.0000	1.0000	0.9900
$X_4 * X_{10} * X_{17} * X_{30}$	0.9900	0.9200	0.9100
$X_{11} * X_{13} * X_{19} * X_{50}$	0.9800	0.8900	0.3800
Overall Power	0.9963	0.9712	0.9038
FP	0.5100	1.1400	1.0900
FDP	0.0601	0.1279	0.1310

Ground physical laws inference

From the **ground physical** laws:

$$m_p = \frac{4}{3}\pi \times R_p^3 \times \rho_p,$$

m_p is *PlanetaryMassJpt*, R_p^3 is *RadiusJpt*³, and ρ_p is *PlanetaryDensJpt*.
And:

$$a = \left(\frac{GP^2}{4\pi^2} (M_h + m_p) \right)^{\frac{1}{3}} \approx \left(\frac{GP^2}{4\pi^2} (M_s M_h^a) \right)^{\frac{1}{3}},$$

a is semi major axes of the ellipses of the orbits, M_h^a is *HostStarMassSlrMass*, a is *SemiMajorAxisAU*, and P is *PeriodDays*.

Generally the data ha the **following other variables**:

TypeFlag, *RadiusJpt*, *PeriodDays*, *PlanetaryMassJpt*, *Eccentricity*, *HostStarMassSlrMass*, *HostStarRadiusSlrRad*, *HostStarMetallicity*, *HostStarTemp*,
PlanetaryDensJpt denoted as x_1 - x_{10}

Planet mass inference. Results over 100 simulations

Table: Power, False Positives (FP) and FDP based on the decision rule that the posterior probability of a feature is larger than 0.25. The feature *PlanetaryRadiusJpt*³ *PlanetaryDensJpt* is counted as true positive, all other selected features as false positive.

GMJMCMC				RGMJMCMC		
Threads	Power	FP	FDP	Power	FP	FDP
16	1.00	0.00	0.00	0.97	0.06	0.058
4	0.79	0.40	0.34	0.61	0.73	0.54
1	0.43	1.21	0.74	0.32	1.67	0.84

3rd Kepler's law inference. Results over 100 simulations

Table: Comparison of Results on Example 3 for GMJMCMC and RGMJMCMC using different number of threads. F_1, F_2 and F_3 refer to the number of times the specific features $(HostStarMassSlrMass \times PeriodDays^2)^{\frac{1}{3}}$, $(HostStarRadiusSlrRad \times PeriodDays^2)^{\frac{1}{3}}$ and $(HostStarTempK \times PeriodDays^2)^{\frac{1}{3}}$ had a posterior probability larger than 0.25. Power gives the percentage of runs where at least one of these three features was detected. FP counts the number of other features and FDP is the corresponding false discovery proportion.

GMJMCMC							RGMJMCMC						
Th	F_1	F_2	F_3	Power	FP	FDP	F_1	F_2	F_3	Power	FP	FDP	
64	81	71	1	1.00	0.02	0.013	78	75	2	0.99	0.03	0.019	
32	63	58	11	0.99	0.14	0.11	55	57	9	0.95	0.12	0.09	
16	34	41	32	0.84	0.46	0.30	31	38	18	0.79	0.68	0.44	
4	15	10	16	0.38	1.05	0.62	8	14	8	0.29	1.47	0.83	
1	6	5	3	0.13	1.46	0.82	6	4	2	0.12	1.81	0.94	

What would have happened with simpler ANN

- 1 When $\mathcal{G} = \{\text{sigmoid}(x)\}$;
- 2 When $\mathcal{G} = \{\text{sigmoid}(x)\}$, $D_{\max} = 300$, and $P_c = 0$;
- 3 When $\mathcal{G} = \{\text{sigmoid}(x)\}$, $D_{\max} = 300$, and $P_c = 0$ and $p(\gamma_j) \propto 1$.

Table: 10 most frequent features detected under scenarios 1, 2 and 3

Fq	Feature	Fq	Feature	Fq	Feature
99	x_3	100	x_3	100	x_3
98	$x_3 * x_3$	72	$g_{\sigma}(-10.33+0.24x_4-8.83x_8)$	54	x_2
93	$x_3 * x_{10}$	64	x_{10}	21	$g_{\sigma}(-16.91-4.94x_2)$
4	$x_3 * x_3 * x_{10}$	62	x_2	19	x_9
1	$x_9 * x_3$	16	$g_{\sigma}(0.21+0.01x_3+0.20x_7)$	16	x_5
1	$x_9 * x_3 * x_3$	9	x_4	14	x_{10}
1	$x_{10} * x_{10} * x_3$	7	$g_{\sigma}(-13.11-7.76x_8-3.33x_2+0.40x_{10})$	10	$g_{\sigma}(6.88 \times 10^9 - 3.92x_2 + 3.44 \times 10^9 g_{\sigma}(-13.57-0.17x_4 - 2.84x_2 - 7.66x_8 + 0.54x_{10}) - 13.76 \times 10^9 g_{\sigma}(g_{\sigma}(-13.57 - 0.17x_4 - 2.84x_2 - 7.66x_8 + 0.54x_{10})))$
1	$x_7 * x_3 * x_3$	5	$g_{\sigma}(-3.36+2.83x_3+0.21x_3-3.36x_9)$	9	x_4
1	$x_6 * x_3 * x_3$	3	$g_{\sigma}(g_{\sigma}(-10.33+0.24x_4)-8.83x_8)$	8	$g_{\sigma}(-13.57-0.17x_4 - 2.84x_2 - 7.66x_8 + 0.54x_{10})$
1	$x_3 * x_3 * x_3$	3	$g_{\sigma}(0.15+0.05x_4-0.01x_3+0.15x_7)$	7	$g_{\sigma}(0.21+0.21x_3)$
0	Others	4	Others	> 300	Others

Trash features under the non regularized case

[illegible]

Figure: A snapshot of features detected under Scenario 3

Application. NEO objects classification. Problem

- **Observations:** Asteroid is a NEO (PHA) object or not (Phocaea)
- **Covariates:** 8 different covariates describing objects
- **Logistic deep Bayesian regression** addressed

$$y_i = y | \rho_i \sim \text{Binom}(1, \rho_i) \quad (22)$$

$$\rho_i = \frac{e^{\gamma_0 \beta_0 + \sum_{j=1}^p \gamma_j \beta_j F_j(\mathbf{x}_i)}}{1 + e^{\gamma_0 \beta_0 + \sum_{j=1}^q \gamma_j \beta_j F_j(\mathbf{x}_i)}} \quad (23)$$

$$p(\gamma) \propto \prod_{j=1}^p \exp(-\gamma_j 2c(F_j(\mathbf{x}))) \quad (24)$$

$$p(\beta | \gamma) = |J_n^\gamma(\hat{\beta})|^{\frac{1}{2}}, \quad (25)$$

Input variables include:

Mean anomaly, Inclination, Argument of perihelion, Longitude of the ascending node, Rms residual, Semi major axis, Eccentricity, Mean motion, Absolute magnitude.

The **allowed modifications** are: $\cos(x)$, $\sin(x)$, $\text{sigmoid}(x)$, $\tanh(x)$, $\text{atan}(x)$, and $\text{erf}(x)$.

Real data analysis. Neo asteroids classification

- Short runs of GMJMCMC with 1 thread are addressed;
- $\|\text{training set}\| = 64$, $\|\text{test set}\| = 20720$;
- Prediction is based on marginalized over all models' probabilities, namely $\hat{Y} = \mathbf{I} \{ \hat{p}(Y|\mathbb{D}) \geq 0.5 \}$, $\hat{p}(Y|\mathbb{D}) = \sum_{\gamma \in \mathbb{V}} \hat{p}(Y|\gamma, \mathbb{D}) \hat{p}(\gamma|\mathbb{D})$;

Algorithm	min.p	med.p	max.p	min.fn	med.fn	max.fn	min.fp	med.fp	max.fp
GMJMCMC	0.9949	0.9998	1.0000	0.0001	0.0002	0.0073	0.0000	0.0002	0.0032
RGMJMCMC	0.9939	0.9998	1.0000	0.0001	0.0002	0.0088	0.0000	0.0002	0.0072
LASSO	0.9991	0.9991	0.9991	0.0013	0.0013	0.0013	0.0000	0.0000	0.0000
RIDGE	0.9982	0.9982	0.9982	0.0026	0.0026	0.0026	0.0000	0.0000	0.0000
LXGBOOST	0.9980	0.9980	0.9980	0.0029	0.0029	0.0029	0.0000	0.0000	0.0000
NBAYESS	0.9963	0.9963	0.9963	0.0054	0.0054	0.0054	0.0000	0.0000	0.0000
MJMCMC	0.9943	0.9946	0.9947	0.0002	0.0002	0.0002	0.0159	0.0162	0.0172
DEEPNETS	0.8979	0.9728	0.9979	0.0018	0.0384	0.1305	0.0000	0.0000	0.0153
TXGBOOST	0.8283	0.8283	0.8283	0.0005	0.0005	0.0005	0.3488	0.3488	0.3488
RFOREST	0.6761	0.8150	0.9991	0.0003	0.1972	0.3225	0.0000	0.0162	0.3557
LR	0.6471	0.6471	0.6471	0.0471	0.0471	0.0471	0.4996	0.4996	0.4996

Application. Breast cancer classification. Problem

- **Observations:** 357 benign, 212 malignant tissues
- **Covariates:** 10 different covariates describing objects
- **Logistic deep Bayesian regression** addressed

$$y_i = y | \rho_i \sim \text{Binom}(1, \rho_i) \quad (26)$$

$$\rho_i = \frac{e^{\gamma_0 \beta_0 + \sum_{j=1}^p \gamma_j \beta_j F_j(\mathbf{x}_i)}}{1 + e^{\gamma_0 \beta_0 + \sum_{j=1}^q \gamma_j \beta_j F_j(\mathbf{x}_i)}} \quad (27)$$

$$p(\gamma) \propto \prod_{j=1}^p \exp(-\gamma_j 2c(F_j(\mathbf{x}))) \quad (28)$$

$$p(\beta | \gamma) = |J_n^\gamma(\hat{\beta})|^{\frac{1}{2}}, \quad (29)$$

Input variables include such parameters of the tissues as:

Radius, texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, fractal dimension.

The **allowed modifications** are: $\text{sigmoid}(x)$, $\mathbb{I}(x > 1)$, $\text{ReLu}(x)$, $x^{\frac{1}{3}}$, $x^{\frac{1}{5}}$.

Real data analysis. Breast cancer classification

- Short runs of GMJMCMC with 1 thread are addressed;
- $\|\text{training set}\| = 142$, $\|\text{test set}\| = 427$;
- Prediction is based on marginalized over all models' probabilities, namely $\hat{Y} = \mathbf{I} \{ \hat{p}(Y|\mathbb{D}) \geq 0.5 \}$, $\hat{p}(Y|\mathbb{D}) = \sum_{\gamma \in \mathbb{V}} \hat{p}(Y|\gamma, \mathbb{D}) \hat{p}(\gamma|\mathbb{D})$;

Algorithm	min.p	med.p	max.p	min.fn	med.fn	max.fn	min.fp	med.fp	max.fp
RIDGE	0.9742	0.9742	0.9742	0.0592	0.0592	0.0592	0.0037	0.0037	0.0037
GMJMCMC	0.9437	0.9695	0.9812	0.0479	0.0536	0.1067	0.0000	0.0148	0.0361
DEEPNETS	0.9225	0.9695	0.9789	0.0305	0.0674	0.1167	0.0000	0.0074	0.0949
RGMJMCMC	0.9554	0.9683	0.9789	0.0479	0.0536	0.0809	0.0037	0.0148	0.0361
NAIVEBAYESS	0.9671	0.9671	0.9671	0.0479	0.0479	0.0479	0.0220	0.0220	0.0220
MJMCMC	0.9624	0.9624	0.9624	0.0756	0.0756	0.0756	0.0111	0.0111	0.0111
LASSO	0.9577	0.9577	0.9577	0.0756	0.0756	0.0756	0.0184	0.0184	0.0184
LXGBOOST	0.9554	0.9554	0.9554	0.0809	0.0809	0.0809	0.0184	0.0184	0.0184
TXGBOOST	0.9484	0.9531	0.9601	0.0536	0.0647	0.0756	0.0291	0.0326	0.0361
RFOREST	0.9038	0.9343	0.9624	0.0422	0.0914	0.1675	0.0000	0.0361	0.1010
LR	0.9272	0.9272	0.9272	0.0305	0.0305	0.0305	0.0887	0.0887	0.0887

Breast cancer classification. 64 CPUs

GMJMCMC

Δ	min.p	med.p	max.p	min.fn	med.fn	max.fn	min.fp	med.fp	max.fp
0.6000	0.9718	0.9765	0.9812	0.0479	0.0592	0.0702	0.0000	0.0000	0.0074
0.5000	0.9671	0.9742	0.9812	0.0479	0.0479	0.0536	0.0000	0.0111	0.0220
0.7000	0.9648	0.9695	0.9742	0.0647	0.0756	0.0862	0.0000	0.0000	0.0000
0.4000	0.9577	0.9624	0.9671	0.0479	0.0479	0.0479	0.0220	0.0291	0.0361
0.8000	0.9554	0.9601	0.9648	0.0862	0.0966	0.1067	0.0000	0.0000	0.0000
0.3000	0.9507	0.9531	0.9601	0.0422	0.0479	0.0479	0.0361	0.0430	0.0464
0.2000	0.9413	0.9460	0.9531	0.0305	0.0422	0.0422	0.0498	0.0565	0.0632
0.9000	0.9366	0.9460	0.9531	0.1117	0.1264	0.1452	0.0000	0.0000	0.0000
0.1000	0.9272	0.9319	0.9413	0.0245	0.0305	0.0305	0.0729	0.0825	0.0918
1.0000	0.6268	0.6268	0.6268	0.5000	0.5000	0.5000	0.0000	0.0000	0.0000

RGMJMCMC

Δ	min.p	med.p	max.p	min.fn	med.fn	max.fn	min.fp	med.fp	max.fp
0.5000	0.9695	0.9765	0.9812	0.0479	0.0479	0.0536	0.0000	0.0074	0.0184
0.6000	0.9695	0.9765	0.9789	0.0536	0.0592	0.0756	0.0000	0.0000	0.0037
0.7000	0.9671	0.9695	0.9742	0.0647	0.0756	0.0809	0.0000	0.0000	0.0000
0.4000	0.9577	0.9624	0.9695	0.0479	0.0479	0.0479	0.0184	0.0291	0.0361
0.8000	0.9554	0.9601	0.9648	0.0862	0.0966	0.1067	0.0000	0.0000	0.0000
0.3000	0.9507	0.9531	0.9577	0.0422	0.0479	0.0479	0.0361	0.0430	0.0464
0.2000	0.9413	0.9460	0.9507	0.0364	0.0422	0.0479	0.0498	0.0565	0.0632
0.9000	0.9366	0.9460	0.9531	0.1117	0.1264	0.1452	0.0000	0.0000	0.0000
0.1000	0.9249	0.9319	0.9390	0.0245	0.0305	0.0305	0.0729	0.0825	0.0918
1.0000	0.6268	0.6268	0.6268	0.5000	0.5000	0.5000	0.0000	0.0000	0.0000

Application. Epigenetic study

- **Observations:** Epigenetically modified reads
- **Covariates:** 14 different covariates
- **Poisson deep Bayesian mixed regression** addressed

$$y_i \sim \text{Poisson}(\mu_i) \quad (30)$$


$$\mu_i = e^{\gamma_0 \beta_0 + \sum_{j=1}^p \gamma_j \beta_j F_j(\mathbf{x}_i) + \sum_{k=1}^r \gamma_{k+p} \delta_{ik}} \quad (31)$$

$$p(\gamma_j) \propto \exp(-2 \log n \gamma_j s_j(\mathbf{x})) \quad (32)$$

$$\beta | \gamma \sim N_{p_\gamma}(\mu_{\beta_\gamma}, \Sigma_{\beta_\gamma}). \quad (33)$$

Latent variables include: **AR(1)**, **RW(1)**, **OU**, **IG** processes

Input variables include: *3 levels corresponding to whether a location belongs to a CGH, CHH or CHG, whether a distance to the previous (C) in DNA is 1, 2, 3, 4, 5, from 6 to 20 or greater than 20, whether a location belongs to a gene from a particular group of genes of biological interest, whether expression is greater than 3000 or greater than 10000, the offset for the number of total bases per location offset(log(total.bases)).*

The **allowed modifications** are: $\cos(x)$, $\text{sigmoid}(x)$, $\tanh(x)$, $\text{atan}(x)$, 

Epigenetic study. Results

No non-linearities found!?

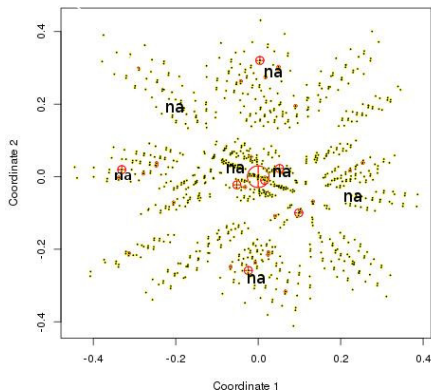
Variable	Posterior
f(model="rw1")	1.00000e+00
offset(log(total.bases))	1.00000e+00
CG	9.990546e-01
CHG	9.515581e-01

Table: Features with posterior probability above 0.25 found by GMJMCMC with 16 threads for the epigenetic data example

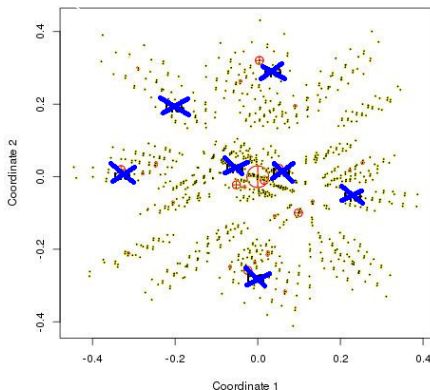
Missing data handling in predictions is easy and intuitive

- Delete models containing NA for the corresponding prediction from \mathbb{V} .
- Recalculate the posteriors.
- Get model averaged predictions.

Metric MDS



Metric MDS



Concluding remarks

- We introduced the (R)GMJMCMC algorithm for deep regression models capable of
 - estimating posterior model probabilities
 - Bayesian model averaging and selection
- *EMJMCMC* R-package is available
 - <http://aliaksah.github.io/EMJMCMC2016/>
 - flexibility in the choice of methods
 - marginal likelihoods
 - model selection criteria
 - extensive parallel computing is available
 - vectorized predictions with NA handling is incorporated
- Results showed that (R)GMJMCMC
 - performs well in terms of the search speed and quality
 - addresses a more general class of models than competitors
 - provides nice predictive and inferential performance in the applications

References



A. Hubin, G.O. Storvik, F. Frommlet

A novel algorithmic approach to Bayesian Logic Regression.

arXiv:1705.07616v1, 2017.



A. Hubin and G.O. Storvik

Efficient mode jumping MCMC for Bayesian variable selection in GLMM.

arXiv:1604.06398v3, 2016.



C. Kooperberg, and I. Ruczinski.

Identifying Interacting SNPs Using Monte Carlo Logic Regression.

Genetic Epidemiology, 28:157–170, 2005.



A. Fritsch.

A Full Bayesian Version of Logic regression for SNP Data.

PhD thesis, 2006.

The End.



Thank you.