Supplement: On Prior Distributions and Approximate Inference for Structured Variables

Oluwasanmi Koyejo Psychology Dept., Stanford sanmi@stanford.edu

Joydeep Ghosh ECE Dept., UT Austin ghosh@ece.utexas.edu Rajiv Khanna ECE Dept., UT Austin rajivak@utexas.edu

Russell A. Poldrack Psychology Dept., Stanford poldrack@stanford.edu

This supplement contains the proofs of theorems and additional exposition not included in the manuscript.

Preliminaries

We begin additional definitions that will be useful for the proofs.

Let X be either a countable set, or a complete separable metric space equipped with the standard Borel σ -algebra of measurable sets, and let Q and P be probability measures on X. The *relative entropy*, also known as the *Kullback-Leibler divergence* (KL divergence) of Q with respect to P is given by:

$$\mathrm{KL}(\boldsymbol{Q} \| \boldsymbol{P}) = \int_{\mathsf{X}} \frac{d\boldsymbol{Q}}{d\boldsymbol{P}}(x) \log \frac{d\boldsymbol{Q}}{d\boldsymbol{P}}(x) d\boldsymbol{P}(x),$$

where $\frac{dQ}{dP}$ is the Radon-Nikodym derivative, existence of which requires that Q is absolutely continuous with respect to P. Let $\mathcal{P} \ni p$ denote the set of probability densities on X.

We make the following assumption:

Assumption 1. All distributions P are absolutely continuous with respect to the dominating measure ν so there exists a density $p \in \mathcal{P}$ that satisfies $dP = pd\nu$.

and the condition:

Condition 2. The sample space X is a subset of Euclidean space with ν given by the Lebesgue measure. Alternatively, X is a countable set with ν given by the counting measure.

Let \mathbb{E} be the *expectation operator*, which we denote using densities as $\mathbb{E}_p[f] = \int_X p(x)f(x)dx$ to simplify notation. We suppress the dependence on the random variable when the expectation and the relative entropy are clear from context. We define a *characteristic function* of the set $A \subset X$, denoted by ϕ_A as the function $\phi_A : X \mapsto \mathfrak{R}_+$ that satisfies $\phi_A(\mathbf{x}) > 0$ for $x \notin A$, and $\phi_A(\mathbf{x}) = 0$ otherwise e.g. $\delta_{X\setminus A}$ is a characteristic function of A. Note that this definition differs slightly from the standard definition in convex analysis, where the characteristic function evaluates to ∞ outside the set.

Priors for structured variables

Our first result is that the information projection a domain restricted density set is given by the restriction of the base measure to the support set. To begin, we show that the support constraint is equivalent to a particular expectation constraint.

Lemma. Let \mathcal{F}_A be the domain restricted density set of A, ϕ_A be its characteristic function, and $\mathcal{G} = \{q \in \mathcal{P} \mid \mathbb{E}_q \left[\phi_A(x) \right] = 0 \} \subset \mathcal{P}.$ Then $\mathcal{F}_A = \mathcal{G}.$

Proof. $[\mathcal{F}_A \subset \mathcal{G}]$: Let $q \in \mathcal{F}_A$, then $\mathbb{E}_q[\phi_A] = 0$, thus $q \in \mathcal{G}$. $[\mathcal{G} \subset \mathcal{F}_A]$: Let $q \in \mathcal{G}$, the nonnegativity of ϕ_A implies that for each $x \in X$, either q = 0 or $\phi_A = 0$.

The following Lemma from Altun and Smola [1] characterizes relative entropy minimization subject to norm ball expectation constraints. For simplicity, the theorem is modified to address the special case of the result where a solution exists and the infimum is attained.

Lemma (Altun and Smola [1]).

$$\min_{q \in \mathcal{P}} \operatorname{KL}(q \| p) \text{ s.t. } \mathbb{E}_{q} \left[\boldsymbol{\beta} \right] = \mathbf{b}$$
$$= \max_{\boldsymbol{\lambda}} \left\langle \boldsymbol{\lambda}, \mathbf{b} \right\rangle - \log \int_{\mathsf{X}} p(x) e^{\left\langle \boldsymbol{\lambda}, \boldsymbol{\beta}(x) \right\rangle} dx + e^{-1}$$

and the unique solution is given by $q_*(x) = p(x)e^{\langle \lambda_*, \beta(z) \rangle - G(\lambda_*)}$ where λ_* is the dual solution and $G(\boldsymbol{\lambda}_*)$ ensures normalization.

We can now show the first main result, investigating the relationship between restriction of densities and information projection subject to domain constraints.

Theorem 3. Under Condition 2, the information projection of the density p to the constraint set \mathcal{F}_{A} , if it exists, is the restriction of p to the domain A.

Proof. The information projection of p to \mathcal{F}_{A} is equivalent to:

$$\min_{q \in \mathcal{P}} \operatorname{KL}(q \| p) \text{ s.t. } \mathbb{E}_q \left[\phi_{\mathsf{A}} \right] = 0. \tag{1}$$

The solution is given by $q_*(z) = p(x)e^{\langle \lambda_*, \phi_A \rangle - G(\lambda_*)}$, where:

$$\lambda_* = rg\max_{\lambda} \langle \lambda, 0 \rangle - \log \int_{\mathbb{X}} p(x) e^{\langle \lambda, \phi_{\mathsf{A}}(x) \rangle} dx.$$

Clearly $\lambda_* = -\infty$, thus $e^{\langle \lambda_*, \phi_A(x) \rangle} \rightarrow \delta_A(x)$ via standard limit arguments, so q_* $p(x)\delta_{\mathsf{A}}(x)/\int_{\mathsf{A}}p(x)dx.$

Corollary 4. Consider the product space $X = W \times Y$. Let domain constraint be given by $W \times \{\hat{y}\}$ for some $\hat{y} \in \mathsf{Y}$. Under Condition 2, the information projection of p to $\mathcal{F}_{\mathsf{w} \times \{\hat{y}\}}$ is given by $p(w|y = \hat{y})\delta_{\hat{y}}$.

Proof. Follows directly from Theorem 3.

Theorem 5. Let $\pi : [n] \mapsto [n]$ be a permutation function and $\{C_{\pi(i)} \mid C_{\pi(i)} \subset X\}$ represent a sequence of sets with non empty intersection $\mathsf{B} = \bigcap \mathsf{C}_i \neq \emptyset$. Given a base density p, let $q_0 = p$, and define the sequence of information projections:

$$q_i = \operatorname*{arg\,min}_{q \in \mathcal{F}_{\mathsf{C}_{\pi(i)}}} \mathrm{KL}(q \| q_{i-1}),$$

Under Condition 2, $q_* = q_N$ is independent of π . Further $q_* = \min_{q \in \mathcal{F}_B} \operatorname{KL}(q||p)$.

Proof. Consider the case when n = 2, then $\pi : [1,2] \mapsto \{[1,2],[2,1]\}$. From Theorem 3, we have that $q_2(x) \propto p(x)\delta_{c_1}(x)\delta_{c_2}(x) = p(x)\delta_{B}(x)$ independent of π . The proof is extended to n > 2 by induction.

We propose approximate inference via the following rule:

$$p_{\mathsf{S}_{*},\hat{y}} = \underset{q \in \mathcal{D} \times \mathcal{F}_{\{\hat{y}\}}}{\operatorname{arg\,min}} \operatorname{KL}(q(w,y) \| p_{\mathsf{A}}(w,y)) = \underset{\mathsf{S}}{\operatorname{arg\,min}} \left[\underset{q \in \mathcal{F}_{\mathsf{S} \times \{\hat{y}\}}}{\operatorname{min}} \operatorname{KL}(q(w,y) \| p_{\mathsf{A}}(w,y)) \right].$$
(2)

Г

The solution is given by $p_{s_*,\hat{y}}(w,y) = p_{s_*}(w|\hat{y})\delta_{\hat{y}}$.

Corollary 6. Let $p_{s_*,\hat{y}}(w, y)$ be the solution of (2), then the posterior distribution $p_{s_*}(w|\hat{y})$ is given by:

$$p_{\mathsf{S}_*}(w|\hat{y}) = \underset{q \in \mathcal{D}}{\arg\min} \operatorname{KL}(q(w) \| p_{\mathsf{A}}(w|\hat{y})) = \underset{q \in \mathcal{D}}{\arg\min} \operatorname{KL}(q(w) \| p(w|\hat{y})).$$
(3)

Proof. The proof follows from repeated application of Theorem 5. The first equality follows from solving (2) in two steps. The first step involves the information projection of $p_A(w, y)$ to $\mathcal{F}_{A \times \{\hat{y}\}}$. The second step is the information projection of the resulting solution $p_A(w|\hat{y})\delta_{\hat{y}}$ to the set $D \times \mathcal{F}_{\{\hat{y}\}}$. The solution is equivalent to the direct information projection of $p_A(w, y)$ to $D \times \mathcal{F}_{\{\hat{y}\}}$ as solved by (2) by Theorem 5.

The second equality follows from the projection of p(w, y) to $\mathcal{F}_{W \times \{\hat{y}\}}$, followed by the information projection of the resulting solution $p(w|\hat{y})\delta_{\hat{y}}$ to the set $\mathsf{D} \times \mathcal{F}_{\{\hat{y}\}}$. Since $\mathsf{D} \subset \mathsf{A}$, the equality holds by Theorem 5.

Corollary 7. Let $p_{s_*,\hat{y}}(w, y)$ be the solution of (2). Define the density $p_{s_*}(w)$ as:

$$p_{\mathsf{S}_*}(w) = \underset{q \in \mathcal{F}_{\mathsf{S}_*}}{\arg\min} \operatorname{KL}(q(w) \| p_\mathsf{A}(w)) = \underset{q \in \mathcal{F}_{\mathsf{S}_*}}{\arg\min} \operatorname{KL}(q(w) \| p(w)). \tag{4}$$

then $p_{s_*}(w)$ is the prior distribution corresponding to the Bayesian posterior $p_{s_*}(w|\hat{y})$

Proof. Let $p_{s_*}(w)$ be the projection of $p_A(w)$ to the set \mathcal{F}_{s_*} , then by Corollary 4, the posterior associated with the observation of \hat{y} is given by $p_{s_*}(w|\hat{y})$.

The second equality follows from the projection of p(w) to \mathcal{F}_{s_*} . As $S_* \subset A$, the equality holds by Theorem 5.

Priors for structured variables

The following theorem explores the submodularity of subset selection using relative entropy. The presented theorem is a special case of the result of Madiman and Tetali [2] with the application of Assumption 1. For simplicity, the theorem is modified to address the special case.

Theorem (Madiman and Tetali [2]). Let $q \in \mathcal{P}$ and $p \in \mathcal{P}$ be probability densities on \tilde{X}^d , and let q_s and p_s be their marginals on $\tilde{X}^{|s|}$, such that the set function $F(s) : \wp(d) \mapsto [0, \infty]$, $F(s) = -\mathrm{KL}(q_s || p_s)$ does not take the value $-\infty$ for any $s \in \wp(d)$, then F(s) is submodular.

In the following, we show that subset selection loss function J(s) is monotone submodular.

Theorem 8. Let $J : \wp(d) \mapsto \Re$, $J(s) = \log p(\mathbf{x}_{s'} = \mathbf{c}_{s'})$, and define $\tilde{J}(s) = J(s) - J(\emptyset)$, then $\tilde{J}(s)$ is normalized and monotone submodular.

Proof. Normalized: By definition $0 \leq \tilde{J}(s) \leq -J(\emptyset)$.

Monotone: Let $c \subset s$, then:

$$p(oldsymbol{x}_{\mathsf{M} ackslash \mathsf{C}}) = p(oldsymbol{x}_{\mathsf{M} ackslash \mathsf{S}}, oldsymbol{x}_{\mathsf{M} ackslash \mathsf{S}}) \leq p(oldsymbol{x}_{\mathsf{M} ackslash \mathsf{S}})$$

Submodular: Consider $F(s) = \log p(x_s = c_s)$. Recall that p is bounded, so F(s) is bounded above and below.

Recall the following identity: $\operatorname{KL}(\delta_{\{a\}} \| p) = \mathbb{E}_{\delta_{\{a\}}} [\log \delta_{\{a\}}] - \mathbb{E}_{\delta_{\{a\}}} [\log p] = -\log p(x = a)$, which follows by standard limit arguments for $\mathbb{E}_{\delta_{\{a\}}} [\log g(x)] \to 0$ with an appropriate $g(x) \to \delta_{\{a\}}$. Thus, we may define $F(\mathsf{s}) = \operatorname{KL}(\delta_{\mathsf{s}} \| p)$.

Applying the theorem of Madiman and Tetali [2], it follows that F(s) is submodular. Finally, we note that if F(s) is submodular, so is its reflection $J(s) = F([d] \setminus s)$.

While maximization of submodular functions is generally NP-hard, a simple greedy forward selection heuristic (Algorithm 1), has been shown to perform almost as well as the optimal in practice, and is known to have strong theoretical guarantees.

Algorithm 1 Greedy selection, $\max J(s)$ s.t. |s| = k

Input: $k, s = \emptyset$ while |s| < k do foreach $i \in [d] \setminus s, f_i = J(s \cup i) - J(s)$ $s = s \cup \{ \arg \max f_i \}$ end while Return: s.

Theorem 9 (Nemhauser et al. [3]). In the case of any normalized, monotonic submodular function F, the set s_* obtained by the greedy algorithm achieves at least a constant fraction $\left(1 - \frac{1}{e}\right)$ of the objective value obtained by the optimal solution i.e. $F(s_*) = \left(1 - \frac{1}{e}\right) \max_{|s| \le k} F(s)$.

Corollary 10. Let J(s) be defined as in Theorem 8 and suppose the base density is product form i.e. $p(\boldsymbol{x}) = \prod_{i=1}^{d} p(x_i)$, then J(s) is linear.

Proof. Define $\mathbf{1}_{s} \in \mathfrak{R}^{d}$ as the vector $(\mathbf{1}_{s})_{i} = 1$ if $i \in s$ and zero otherwise, and define the vector $\mathbf{h} \in \mathfrak{R}^{d}$ taking values $\mathbf{h}_{i} = p(x_{i} = c_{i})$. When $p(\mathbf{x}) = \prod_{i=1}^{d} p(x_{i})$, we have that $J(s) = \log p(\mathbf{x}_{s'} = \mathbf{c}_{s'}) = \sum_{i \in s'} \log p(x_{i}) = \langle \mathbf{1}_{s'}, \mathbf{h} \rangle$, a linear function.

Gaussian linear regression with sparse structure

Consider a generative model for *n* samples given by a linear model combined with Gaussian noise $y_i|_{w,z_i} = w^{\top} z_i + \epsilon$, where the response $y_i \in \mathfrak{R}$, the feature vector $z_i \in \mathfrak{R}^d$, and the weight vector $w \in \mathfrak{R}^d$. The weights are drawn from the zero mean Gaussian distribution $w \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$, where $\mathbf{C} \in \mathfrak{R}^{d \times d}$ is the prior covariance matrix, and its inverse $\mathbf{D} = \mathbf{C}^{-1} \in \mathfrak{R}^{d \times d}$ is the corresponding prior precision matrix. The noise is drawn from a univariate Gaussian $\epsilon \in \mathfrak{R}, \epsilon \sim \mathcal{N}(0, \sigma^2)$. We set $\lambda = \frac{1}{\sigma^2}$. Let $y \in \mathfrak{R}^n$ represent the responses collected into a vector with $y(i) = y_i$, and $\mathbf{Z} \in \mathfrak{R}^{n \times d}$ represent the features in matrix form, so $\mathbf{Z}(i, :) = z_i^{\top}$.

As the prior and the likelihood are Gaussian, the unconstrained posterior distribution P(w|y) is Gaussian, represented as $\mathcal{N}(\mu, \Sigma)$, where $\Sigma \in \mathfrak{R}^{d \times d}$ is the posterior covariance matrix, and its inverse $\Lambda = \mathbf{S}^{-1} \in \mathfrak{R}^{d \times d}$ is the corresponding precision matrix. The posterior precision is given by $\Lambda = \mathbf{D} + \lambda \mathbf{Z}^{\top} \mathbf{Z}$. The posterior mean $\mu \in \mathfrak{R}^d$ is given by $\mu = \lambda \Lambda \mathbf{Z}^{\top} \mathbf{y}$. Recall that $\mu_s \in \mathfrak{R}^{|s|}$ is the subvector given by $\mu_s = {\mu(i) \mid i \in s}$. For matrices, define $\Sigma_{s,c} \in \mathfrak{R}^{|s| \times |c|}$ as the submatrix ${\Sigma(i, j) \mid i \in s, j \in c}$. We also define the linear projection matrix $\mathbf{P}_s \in \mathfrak{R}^{d \times |s|}$, $\mathbf{P}_s : \mathfrak{R}^{|s|} \mapsto \mathfrak{R}^d$ by imputing zeros as missing entries.

We set the default value $c_i = 0 \forall i \in [d]$. For any fixed s, the information projection is given by the restriction of P(w|y) as the conditional distribution:

$$\boldsymbol{P}(\boldsymbol{w}_{\mathsf{S}}|\boldsymbol{w}_{\mathsf{S}'}=\boldsymbol{0},\mathbf{y}) = \mathcal{N}\left(\boldsymbol{m}_{\mathsf{S}|\mathsf{S}'},\boldsymbol{\Sigma}_{\mathsf{S}|\mathsf{S}'}\right),\tag{5}$$

where $m_{s|s'} \in \Re^{|s|}$, given by $m_{s|s'} = \mu_s + \Lambda_{s,s'}^{-1} \Lambda_{s,s'} \mu_s$. Approximate inference is applied using the convex sparse subsets, equivalent to the submodular optimization. The resulting cost function (up to constants) is given by:

$$J(\mathbf{s}) = \boldsymbol{\mu}_{\mathbf{s}'}^{\top} \boldsymbol{\Sigma}_{\mathbf{s}',\mathbf{s}'}^{-1} \boldsymbol{\mu}_{\mathbf{s}'} - \log |\boldsymbol{\Sigma}_{\mathbf{s}',\mathbf{s}'}| + a_1 = (\boldsymbol{\mu} - \mathbf{P}_{\mathbf{s}} \boldsymbol{m}_{\mathbf{s}|\mathbf{s}'})^{\top} \boldsymbol{\Lambda} (\boldsymbol{\mu} - \mathbf{P}_{\mathbf{s}} \boldsymbol{m}_{\mathbf{s}|\mathbf{s}'}) - \log |\boldsymbol{\Lambda}_{\mathbf{s},\mathbf{s}}|, \quad (6)$$

where a_1 is an additive constant. The second equality may be computed directly, but is most easily derived by directly solving the information projection (3) for a fixed s. Interpreting the resulting form, we find that the subsets are selected in order to minimize the distance between the conditional mean and marginal mean vector and the determinant term adds regularity depending on the coupling between the variables.

Additional Experiments

Functional Neuroimaging Data

To evaluate all the models, we applied additional dimensionality reduction to the fMRI data resulting in 10,000 voxels. The spatially constrained Ward hierarchical clustering approach of Michel et al. [4] was used for all dimensionality reduction. The predictive R^2 are: *Sparse-G* (0.0167), *Lasso* (-0.299), *Ridge* (-0.542), *ARD* (-0.155), *SpikeSlabFull* (-0.250), *SpikeSlab0.5* (-0.270) and *SpikeSlabKL* (-0.183). As expected, all models (apart from *ARD*) had worse results upon dimensionality reduction. Again, the projection approach improved the performance of the *Spike and Slab* models.

References

- [1] Y. Altun and A.J. Smola. Unifying divergence minimization and statistical inference via convex duality. In *COLT*, 2006.
- [2] M. Madiman and P. Tetali. Information inequalities for joint distributions, with interpretations and applications. *IEEE Trans. IT*, 56(6):2699–2713, 2010.
- [3] G.L. Nemhauser, L.A. Wolsey, and M.L. Fisher. An analysis of approximations for maximizing submodular set functions. *Mathematical Programming*, 14(1):265–294, 1978.
- [4] V. Michel, A. Gramfort, G. Varoquaux, E. Eger, C. Keribin, and B. Thirion. A supervised clustering approach for fMRI-based inference of brain states. *Pattern Recognition*, 45(6):2041– 2049, 2012.





20

(b) R^2 as a function of n:k ratio

Figure 1: Full size simulated data performance in terms of support recovery (AUC)





Figure 2: Full size simulated data performance in terms of regression performance (\mathbb{R}^2)