

Sparse Signal Recovery with Exponential-Family Noise

Irina Rish

IBM T.J. Watson Research Center
Yorktown, NY

rish@us.ibm.com

Genady Grabarnik
CUNY, NY

genadyg@gmail.com

Abstract—The problem of sparse signal recovery from a relatively small number of noisy measurements has been studied extensively in the recent literature on compressed sensing. However, the focus of those studies appears to be limited to the case of linear projections disturbed by *Gaussian* noise, and the sparse signal reconstruction problem is treated as *linear* regression with l_1 -norm regularization constraint. A natural question to ask is whether one can accurately recover sparse signals under different noise assumptions. Herein, we extend the results of [13] to the more general case of *exponential-family noise* that includes Gaussian noise as a particular case, and yields l_1 -regularized *Generalized Linear Model (GLM)* regression problem. We show that, under standard restricted isometry property (RIP) assumptions on the design matrix, l_1 -minimization can provide stable recovery of a sparse signal in presence of the exponential-family noise, provided that certain sufficient conditions on the noise distribution are satisfied.

I. INTRODUCTION

Accurate and efficient recovery of sparse high-dimensional signals from low-dimensional linear measurements received much attention in the recent compressed sensing literature [4]–[7], [10], [12]. While finding the sparsest solution satisfying a set of linear constraints is NP-hard as it involves a combinatorial problem of l_0 -norm minimization, its l_1 -relaxation can be exact under certain conditions and allows for efficient optimization algorithms. Particularly interesting for real-life applications is the case of signal recovery from *noisy* measurements [11], [13], which relates to practically all modern applications of compressed sensing in image processing, sensor networks, biology, and medical imaging, just to name a few (see [16] for an extensive list of references on compressed sensing and its recent applications).

The problem of sparse signal recovery from noisy observations is typically formulated as minimization of the l_1 -norm of an unobserved signal \mathbf{x} subject to the sum-squared loss constraint $\|\mathbf{y} - A\mathbf{x}\|_{l_2} < \epsilon$ [11], [13]. From a probabilistic point of view, this problem is equivalent to loglikelihood maximization under the assumption of (1) linear Gaussian noise $P(\mathbf{y}) \sim N(\mu = A\mathbf{x}, \Sigma = I)$ (which results into sum-squared function), and the sparsity-promoting Laplace prior on the input signal, $p(\mathbf{x}) \sim e^{-\lambda\|\mathbf{x}\|_{l_1}}$, which produces the l_1 -norm. However, in many practical applications, it might be more appropriate to use non-Gaussian models of noise: for example, Bernoulli or multinomial distributions are better suited for describing such measurements as (binary) failures or multilevel performance degradations of end-to-end test transactions (“probes”) in a distributed computer systems [17], [19]; exponential distribution is better suited

for describing nonnegative measurements such as end-to-end response time in such systems [3], [9]. Non-Gaussian observations, including binary, discrete, non-negative, etc., variables, are common in various other applications such as, for example, computational biology and medical imaging: e.g., predicting the presence or absence of a certain disease given DNA microarrays, or predicting a particular “state of mind” (e.g., an emotional state of being angry, happy, anxious, etc.) or response from fMRI images [8], [14]. A common challenge in such applications is a very high number (on the order of 10,000–100,000) of variables, but a much lower number of samples (a few hundreds), and the need to reconstruct not only predictive, but also interpretable regression models that can select most relevant variables such as genes or brain voxels.

In this paper, we will consider the general class of exponential-family distributions that includes, besides Gaussian, a wide variety of other commonly used distributions, such as exponential, Bernoulli, multinomial, gamma, chi-square, beta, Weibull, Dirichlet, and Poisson, just to name a few. The corresponding regression problem of recovering the unobserved vector \mathbf{x} from the measurements \mathbf{y} contaminated by an exponential-family noise is known as the *Generalized Linear Model (GLM)* regression. GLM regression problem maximizes the exponential-family loglikelihood of the observations (target variable \mathbf{y}) with respect to the unobserved parameters (signal \mathbf{x}), which turns out to be equivalent to minimizing the corresponding *Bregman divergence* $d(\mathbf{y}, \mu(A\mathbf{x}))$, where μ is the mean parameter of the exponential-family distribution, and $\theta = A\mathbf{x}$ is the corresponding natural parameter (there is a one-to-one correspondence between those two parameters). In case of Gaussian likelihood, for example, $\mu = \theta$ and (assuming independent, unit-variance noise) the corresponding Bregman divergence is simply the squared Euclidean distance $\|\mathbf{y} - A\mathbf{x}\|_{l_2}^2$. Adding l_1 -norm constraint to GLM regression allows for an efficient method of sparse signal recovery, and is often used in statistical literature [15]. Thus, a natural question to ask is to what extent stable signal recovery results from the compressed sensing literature apply to the linear measurements corrupted by an exponential-family noise? This work provides an initial investigation of this question, deriving some conditions for stable sparse signal recovery from exponential-family observations.

We show that accurate recovery of sparse signals under the exponential-family noise assumption is possible in many cases, and derive sufficient conditions on such recovery: for a class of exponential-family distributions with bounded

$\phi''(y)$, where $\phi(y)$ is the Legendre conjugate of the log-partition function that uniquely determines the distribution, and for several individual exponential-family members that do not satisfy the above condition. Essentially, we show that, if the signal x^0 is sufficiently sparse, the measurement noise is sufficiently small (where the noise is expressed as Bregman divergence between the measurement y and the mean μ^0 of the distribution determined by the natural parameter $\theta^0 = Ax^0$), and the matrix A obeys the restricted isometry property (RIP) with appropriate RIP constant, then the solution to the sparse GLM regression problem approximates the true signal well in the l_2 sense. Moreover, we show that the results of [13] for a more general case of compressible, rather than sparse, signals can be also extended to the exponential-family noise.

II. BACKGROUND

A. Sparse Signal Recovery from Noisy Observations

We assume that $x^0 \in R^m$ is an s -sparse signal, i.e. a signal with no more than s nonzero entries, where $s \ll m$. Let A be an n by m matrix that produces a vector of linear projections $y^0 = Ax^0$, where $n \ll m$, and let y be a vector of n noisy measurements that follow some noise distribution $P(y|Ax^0)$. It is often assumed that A satisfies the so-called "restricted isometry property" (RIP) at the sparsity level S (or S -restricted isometry property), that essentially says that every subset of columns of A with cardinality less than S behaves like an almost orthonormal system. Formally, following [7]

Definition 1 (Restricted Isometry Property) Let A_T , where T subset $\{1, \dots, m\}$ denote an $n \times T$ submatrix of A that contains columns with indexes in T . The S -restricted isometry constant δ_S of A is the smallest quantity such that

$$(1 - \delta_S) \|c\|_{l_2}^2 \leq \|A_T c\|_{l_2}^2 \leq (1 + \delta_S) \|c\|_{l_2}^2 \quad (1)$$

for any all subsets T with $|T| \leq S$ and for any vector $(c_j)_{j \in T}$ defined over coordinates in T . The matrix A is said to satisfy the restricted isometry property if there exists such constant δ_S that the eq. 1 is satisfied.

It was shown (e.g., in [7]) that if

$$\delta_S + \delta_{2S} + \delta_{3S} < 1,$$

then solving the l_1 -minimization problem in eq. 2 below can recover any signal x that is S -sparse (contains no more than S non-zero entries).

Our question is: can one recover x^0 from y , given that noise is "sufficiently small" (to be defined precisely below)? This question has been answered in the compressed sensing literature for the particular case when the noise distribution is Gaussian. Indeed, [13] show that, if: (1) $\|y - Ax^0\|_{l_2} \leq \epsilon$ (small noise assumption), (2) x^0 is sufficiently sparse and the (3) matrix A obeys the restricted isometry property (RIP) with appropriate RIP constants, then the solution to the following l_1 -optimization problem

$$x^* = \arg \min_x \|x\|_{l_1} \quad \text{subject to } \|y - Ax\|_{l_2} \leq \epsilon \quad (2)$$

approximates the true signal well. More formally, Theorem 1 in [13] states:

Theorem 1 [13] Let S be such that $\delta_{3S} + 3\delta_{4S} < 2$, where δ_S is the S -restricted isometry constant of the matrix A , as defined above. Then for any signal x^0 with the support $T^0 = \{t : x^0 \neq 0\}$, where $|T^0| \leq S$ and any noise vector (perturbation) e with $\|e\|_{l_2} \leq \epsilon$, the solution x^* to the problem in eq. 2 obeys

$$\|x^* - x^0\|_{l_2} \leq C_S \cdot \epsilon, \quad (3)$$

where the constant C_S may only depend on δ_{4S} . For reasonable values of δ_{4S} , C_S is well-behaved; e.g. $C_S \approx 8.82$ for $\delta_{4S} = 1/5$ and $C_S \approx 10.47$ for $\delta_{4S} = 1/4$.

Moreover, [13] show that (1) no other recovery method "can perform fundamentally better for arbitrary perturbations of size ϵ , i.e. even if an oracle would make the actual support T^0 of x^0 available to us, making the problem well-posed, the least-squares solution \hat{x} (i.e., the maximum-likelihood solution which is optimal in the absence of any other information) would approximate the true signal x^0 with the error proportional to ϵ ".

Finally, [13] extend their result from sparse to approximately sparse vectors in the following

Theorem 2 [13] Let $x^0 \in R^m$ be an arbitrary vector, and let $x_{0,S}$ be the truncated vector corresponding to the S largest values of x^0 (in absolute value). Under the assumptions of Theorem 1, the solution x^* to the problem in eq. 2 obeys

$$\|x^* - x^0\|_{l_2} \leq C_{1,S} \cdot \epsilon + C_{2,S} \cdot \frac{\|x^0 - x_{0,S}\|_{l_1}}{\sqrt{S}}. \quad (4)$$

For reasonable values of δ_{4S} the constants above are well-behaved; e.g. $C_{1,S} \approx 12.04$ and $C_{2,S} \approx 8.77$ for $\delta_{4S} = 1/5$.

B. Exponential-family distributions and Bregman divergences

Herein, we will generalize the above results in the case of exponential-family noise. Note that $\|y - Ax\|_{l_2} \leq \epsilon$ is a constraint on the negative log-likelihood of a Gaussian variable $y \sim N(\mu, \Sigma)$ with $\mu = Ax$ and $\Sigma = I$ (assuming independent unit-variance noise), i.e.

$$-\log P(y|Ax^0) = f(y) + \frac{1}{2} \|y - Ax\|_{l_2}^2. \quad (5)$$

Gaussian distribution is a particular member of the exponential family of distributions.

Definition 2 An exponential family is a parametric family of probability distributions where the probability density has the form

$$\log p_{\psi, \theta}(\mathbf{y}) = \mathbf{x}\theta - \psi(\theta) + \log p_0(\mathbf{y}), \quad (6)$$

where θ is called the natural parameter, $\psi(\theta)$ is the (strictly convex and differentiable) cumulant function, or

the log-partition function, that uniquely determines the member distribution of the exponential family, and $p_0(\mathbf{y})$ is a non-negative function called base measure that does not depend on the parameter θ .

As shown by [2], there is a bijection between the exponential-family densities $p_{\psi,\theta}(\mathbf{y})$ and Bregman divergences $d_\phi(\mathbf{y}, \mu)$, so that each exponential-family density can be also expressed as

$$p_{\psi,\theta}(\mathbf{y}) = \exp(-d_\phi(\mathbf{y}, \mu)) f_\phi(\mathbf{y}), \quad (7)$$

where $\mu = \mu(\theta) = E_{p_{\psi,\theta}}(Y)$ is the expectation parameter corresponding to θ , ϕ is the (strictly convex and differentiable) Legendre conjugate of ψ , $f_\phi(\mathbf{y})$ is a uniquely determined function, and $d_\phi(\mathbf{y}, \mu)$ is the corresponding Bregman divergence defined as follows.

Definition 3 Given a strictly convex function $\phi : S \rightarrow \mathbb{R}$ defined on a convex set $S \subseteq \mathbb{R}$, and differentiable on the interior of S , $\text{int}(S)$ [18], the **Bregman divergence** $d_\phi : S \times \text{int}(S) \rightarrow [0, \infty)$ is defined as

$$d_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \mathbf{x} - \mathbf{y}, \nabla \phi(\mathbf{y}) \rangle, \quad (8)$$

where $\nabla \phi(\mathbf{y})$ is the gradient of ϕ .

In other words, the Bregman divergence can be thought of as the difference between the value of ϕ at point \mathbf{x} and the value of the first-order Taylor expansion of ϕ around point \mathbf{y} evaluated at point \mathbf{x} (see Figures 1 and 2, where $h(x) = \phi(y) + \langle \mathbf{x} - \mathbf{y}, \nabla \phi(\mathbf{y}) \rangle$).

Table I (derived from Tables 1 and 2 in [1]) shows particular examples of commonly used exponential-family distributions and their corresponding Bregman divergences. For example, the unit-variance Gaussian distribution leads to square loss, multivariate spherical Gaussian (diagonal covariance/independent variables) gives rise to Euclidean distance, an multivariate Gaussian with the inverse-covariance (concentration) matrix C leads to Mahalanobis distance, Bernoulli distribution corresponds to logistic loss, exponential distribution leads to Itakura-Saito distance, while a multinomial distribution corresponds to the KL-divergence.

III. OUR CONTRIBUTION

We now extend the result in Theorem 1 to the case of exponential-family noise. Let us consider the following constrained l_1 -regularization problem that generalizes the standard noisy compressed sensing problem of [13] to the following:

$$\min \|x\|_1 \quad \text{subject to} \quad \sum_i d(y_i, \mu(A_i x)) \leq \epsilon, \quad (9)$$

where $d(y_i, \mu(A_i x))$ is Bregman divergence between the noisy observation y_i and the mean parameter of the corresponding exponential-family distribution with the natural parameter $\theta_i = A_i x$. Note that this problem corresponds to l_1 -regularized *Generalized Linear Model (GLM)* regression,

that includes as a particular case the standard compressed-sensing formulation, i.e. the l_1 regularized linear regression (in that case, Bregman divergence is simply the Euclidean distance, and $\mu(A_i x) = A_i x$).

We show that, if: (1) the noise is small, (2) x^0 is sufficiently sparse and the (3) matrix A obeys the restricted isometry property (RIP) with appropriate RIP constants, then the solution to the above problem approximates the true signal well. More formally,

Theorem 3 Let S be such that $\delta_{3S} + 3\delta_{4S} < 2$, where δ_S is the S -restricted isometry constant of the matrix A , as defined above. Then for any signal x^0 with the support $T^0 = \{t : x^0 \neq 0\}$, where $|T^0| \leq S$, and for any vector $\mathbf{y} = (y_1, \dots, y_n)$ of noisy linear measurements where

- 1) the noise follows exponential-family distributions $p_{\theta_i}(y_i)$, with the natural parameter $\theta_i = (A_{i,:} x^0)$,
- 2) the noise (in terms of Bregman divergence) is sufficiently small, i.e. $\forall i, d_{\phi_i}(y_i, \mu(A_{i,:} x^0)) \leq \epsilon$, and
- 3) each function $\phi_i(\cdot)$ (i.e., the Legendre conjugate of the corresponding log-partition function, uniquely defining the Bregman divergence), satisfies the conditions imposed by at least one of the Lemmas below,

the solution x^* to the problem in eq. 9 obeys

$$\|x^* - x^0\|_{l_2} \leq C_S \cdot \delta(\epsilon), \quad (10)$$

where C_S is the constant from Theorem 1 of [13], and $\delta(\epsilon)$ is a continuous monotone increasing function of ϵ s.t. $\delta(0) = 0$ (and thus $\delta(\epsilon)$ is small when ϵ is small). A particular form of this function depends on particular members of exponential family.

Proof: Following the proof of Theorem 1 in [13], we will only have to show that the ‘‘tube constraint’’ (condition 1) still holds (the rest of the proof remains unchanged), i.e. that

$$\|Ax^* - Ax^0\|_{l_2} \leq \delta(\epsilon) \quad (11)$$

where δ is some continuous monotone increasing function of ϵ , and $\delta(0) = 0$, so its small when ϵ is small. It was a trivial consequence of the triangle inequality in case of Euclidean distance; however, triangle inequality does not hold, in general, for Bregman divergences, and thus we must provide a different proof for the tube constraint, possibly for each type of Bregman divergence (exponential-family distribution). Since

$$\|Ax^* - Ax^0\|_{l_2}^2 = \sum_{i=1}^m (A_{i,:} x^* - A_{i,:} x^0)^2 = \sum_{i=1}^m (\theta_i^* - \theta_i^0)^2,$$

we will need to show that $|\theta_i^* - \theta_i^0| < \beta(\epsilon)$, where $\beta(\epsilon)$ is a continuous monotone increasing function of ϵ s.t. $\beta(0) = 0$ (and thus $\beta(\epsilon)$ is small when ϵ is small), then in eq. 11 we get $\delta(\epsilon) = \sqrt{m \cdot \beta(\epsilon)}$. Lemma 1 provides the proof of this fact for a class of exponential-family distributions with bounded $\phi''(y)$ (where $\phi(y)$ is the Legendre conjugate of the log-partition function that uniquely determines the

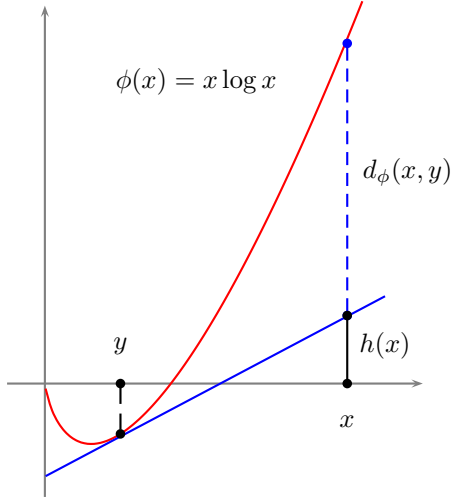


Fig. 1. Relative entropy (KL-divergence)

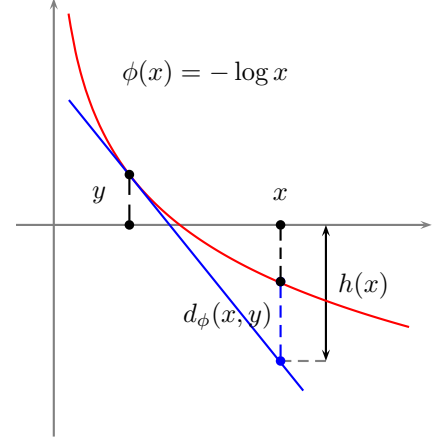


Fig. 2. Itakura-Saito distance (Burg divergence)

TABLE I. Examples of commonly-used exponential-family distributions and their corresponding Bregman divergences.

Domain	Distribution	$p_{\theta}(y)$	μ	$\phi(\mu)$	$d_{\phi}(y, \mu)$	Divergence
\mathbb{R}	1D Gaussian	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-a)^2}{2\sigma^2}}$	a	$\frac{1}{2\sigma^2}\mu^2$	$\frac{1}{2\sigma^2}(y - \mu)^2$	square loss
$\{0, 1\}$	Bernoulli	$q^y(1-q)^{1-y}$	q	$\mu \log \mu + (1-\mu) \log(1-\mu)$	$y \log(\frac{y}{\mu}) + (1-y) \log(\frac{1-y}{1-\mu})$	logistic loss
\mathbb{R}_{++}	Exponential	$\lambda e^{-\lambda y}$	$1/\lambda$	$-\log \mu - 1$	$\frac{y}{\mu} - \log(\frac{y}{\mu}) - 1$	Itakura-Saito distance
n-simplex	nD Multinomial	$\frac{1}{\prod_{j=1}^n y_j!} \prod_{j=1}^n q_j^{y_j}$	$[Nq_j]_{j=1}^{n-1}$	$\sum_{j=1}^n \mu_j \log(\frac{\mu_j}{N})$	$\sum_{j=1}^n y_j \log(\frac{y_j}{\mu_j})$	KL-divergence
\mathbb{R}^n	nD Sph. Gaussian	$\frac{1}{\sqrt{(2\pi\sigma^2)^n}} e^{-\frac{\ x-a\ _2^2}{2\sigma^2}}$	\mathbf{a}	$\frac{1}{2\sigma^2}\ \mu\ _2^2$	$\frac{1}{2\sigma^2}\ y - \mu\ _2^2$	squared Euclidean distance
\mathbb{R}^n	nD Gaussian	$\frac{\sqrt{\det(C)}}{\sqrt{(2\pi)^n}} e^{-\frac{(y-a)^T C (y-a)}{2}}$	\mathbf{a}	$\frac{\mu^T C \mu}{2}$	$\frac{(y-\mu)^T C (y-\mu)}{2}$	Mahalanobis distance ¹

distribution). However, for several members of the exponential family (e.g., Bernoulli distribution) this condition is not satisfied, and those cases must be handled individually. Thus, we provide separate proofs for several different members of the exponential family in Lemmas 2.1, 2.2 and 2.3, and obtain particular expressions for $\beta(\epsilon)$ in each case. Note that for simplicity sake, we only consider univariate exponential-family distributions, corresponding to the case of independent noise for each measurement y_i , which was effectively assumed in standard problem formulation that used Euclidean distance corresponding to a spherical Gaussian distribution, i.e. a vector of independent Gaussian variables. However, Lemma 1 below can be extended from scalar to vector case, i.e. to multivariate exponential-family distributions that do not necessarily imply independent noise. Lemma 2.3 will provide a specific case of such distribution - a multivariate Gaussian with concentration matrix C .

The ‘‘cone constraint’’ part of the proof in [13] remains intact; it is easy to see that it does not depend on the particular constraint in the l_1 -minimization problem 9, and only makes use of the sparsity of x^0 and l_1 -optimality of x^* .

Thus, we can simply substitute $\|Ah\|_{l_2}$ by $\delta(\epsilon)$ in eq. 13 on page 8 in the proof of Theorem 1 of [13], or, equivalently, replace 2ϵ (that was shown to bound $\|Ah\|_{l_2}$) by $\delta(\epsilon)$ in the eq. 14. ■

Just like for the sparse signal case (Theorem 1 in [13]), the only change we have to make in the proof of the Theorem 2 (general case of approximable, rather than sparse, signals), when generalizing it from Euclidean distance to Bregman divergence in eq. 9, is the tube constraint. Thus, once we showed it for the Theorem 3 above, the generalization to approximable signals follows automatically:

Theorem 4 *Let $x^0 \in \mathbb{R}^m$ be an arbitrary vector, and let $x_{0,S}$ be the truncated vector corresponding to the S largest values of x^0 (in absolute value). Under the assumptions of Theorem 3, the solution x^* to the problem in eq. 9 obeys*

$$\|x^* - x^0\|_{l_2} \leq C_{1,S} \cdot \delta(\epsilon) + C_{2,S} \cdot \frac{\|x^0 - x_{0,S}\|_{l_1}}{\sqrt{S}}. \quad (12)$$

where $C_{1,S}$ and $C_{2,S}$ are the constants from Theorem 2 of [13], and $\delta(\epsilon)$ is a continuous monotone increasing function

of ϵ s.t. $\delta(0) = 0$ (and thus $\delta(\epsilon)$ is small when ϵ is small). A particular form of this function depends on particular members of exponential family.

The following lemma states the sufficient conditions for the ‘‘tube constraint’’ in eq. 11 to hold in general case of arbitrary exponential-family noise, provided that $\phi''(y)$ exists and is bounded on the appropriate intervals.

Lemma 1 *Let y denote a random variable following an exponential-family distribution $p_\theta(y)$, with the natural parameter θ , and the corresponding mean parameters $\mu(\theta)$. Let $d_\phi(y, \mu(\theta))$ denote the Bregman divergence associated with this distribution. If*

- 1) $d_\phi(y, \mu^0(\theta^0)) \leq \epsilon$ (small noise),
- 2) $d_\phi(y, \mu^*(\theta^*)) \leq \epsilon$ (constraint in GLM problem eq. 9), and
- 3) $\phi''(y)$ exists and is bounded on $[y_{min}, y_{max}]$, where $y_{min} = \min\{y, \mu^0, \mu^*\}$ and $y_{max} = \max\{y, \mu^0, \mu^*\}$,

then

$$|\theta^* - \theta^0| \leq \beta(\epsilon) = \sqrt{\epsilon} \cdot \frac{2\sqrt{2} \max_{\hat{\mu} \in [\mu^*, \mu^0]} |\phi''(\hat{\mu})|}{\sqrt{\min_{\hat{y} \in [y_{min}, y_{max}]} \phi''(\hat{y})}} \quad (13)$$

Proof: We prove the lemma in two steps: first, we show that $|\mu^*(\theta^*) - \mu^0(\theta^0)|$ is small if ϵ is small, and then infer $|\theta^* - \theta^0|$ is small.

- 1) By definition in eq. 8, Bregman divergence is the non-linear tail of the Taylor expansion of $\phi(y)$ at point μ , i.e., the *Lagrange remainder* of the linear approximation:

$$d_\phi(y, \mu) = \phi''(\hat{y})(y - \mu)^2/2, \quad \hat{y} \in [y_1; y_2],$$

$$\text{where } y_1 = \min\{y, \mu\}, \quad y_2 = \max\{y, \mu\}.$$

Let $y_1^0 = \min\{y, \mu^0\}$, $y_2^0 = \max\{y, \mu^0\}$ and $y_1^* = \min\{y, \mu^*\}$, $y_2^* = \max\{y, \mu^*\}$. Using the conditions $0 \leq d_\phi(y, \mu^0) \leq \epsilon$ and $0 \leq d_\phi(y, \mu^*) \leq \epsilon$, and observing that

$$\min_{\hat{y} \in [y_{min}; y_{max}]} \phi''(\hat{y}) \leq \min_{\hat{y} \in [y_1^0; y_2^0]} \phi''(\hat{y})$$

$$\text{and } \min_{\hat{y} \in [y_{min}; y_{max}]} \phi''(\hat{y}) \leq \min_{\hat{y} \in [y_1^*; y_2^*]} \phi''(\hat{y}),$$

we get

$$\phi''(\hat{y})(y - \mu^0)^2/2 \leq \epsilon \Leftrightarrow (y - \mu^0)^2 \leq \frac{2\epsilon}{\phi''(\hat{y})} \Leftrightarrow$$

$$\Leftrightarrow |y - \mu^0| \leq \frac{\sqrt{2\epsilon}}{\sqrt{\min_{\hat{y} \in [y_1^0; y_2^0]} \phi''(\hat{y})}} \leq$$

$$\leq \frac{\sqrt{2\epsilon}}{\sqrt{\min_{\hat{y} \in [y_{min}; y_{max}]} \phi''(\hat{y})}}$$

$$\text{and, similarly, } |y - \mu^*| \leq \frac{\sqrt{2\epsilon}}{\sqrt{\min_{\hat{y} \in [y_1^*; y_2^*]} \phi''(\hat{y})}} \leq$$

$$\leq \frac{\sqrt{2\epsilon}}{\sqrt{\min_{\hat{y} \in [y_{min}; y_{max}]} \phi''(\hat{y})}},$$

from which, using the triangle inequality, we conclude

$$|\mu^* - \mu^0| \leq |y - \mu^*| + |y - \mu^0| \leq$$

$$\leq \frac{2\sqrt{2\epsilon}}{\sqrt{\min_{\hat{y} \in [y_{min}; y_{max}]} \phi''(\hat{y})}} \quad (14)$$

Note that $\phi''(\hat{y})$ under the square root is always positive since ϕ is strictly convex.

- 2) The mean and the natural parameters of an exponential-family distribution relate to each other as follows: $\theta(\mu) = \phi'(\mu)$ (respectively, $\theta(\mu) = \nabla\phi(\mu)$ for vector μ), where $\phi'(\mu)$ is called the *link function*. Therefore, we can write

$$|\theta^* - \theta^0| = |\phi'(\mu^*) - \phi'(\mu^0)| = |\phi''(\hat{\mu})(\mu^* - \mu^0)|,$$

$$\text{where } \hat{\mu} \in [\mu^*, \mu^0],$$

and thus, using the above result in eq. 14, we get

$$|\theta^* - \theta^0| \leq \beta(\epsilon) = \sqrt{\epsilon} \cdot \frac{2\sqrt{2} \max_{\hat{\mu} \in [\mu^*, \mu^0]} |\phi''(\hat{\mu})|}{\sqrt{\min_{\hat{y} \in [y_{min}; y_{max}]} \phi''(\hat{y})}}$$

which concludes the proof. ■

The condition (3) in the above lemma requires that $\phi''(y)$ exists and is bounded on the intervals between y and both μ^0 and μ^* . However, even when this condition is not satisfied, as it happens for the logistic loss, where $\phi''(y) = \frac{1}{y(1-y)}$ is unbounded at 0 and 1, and for several other Bregman divergences shown in Table 1, we may still be able to prove similar results using specific properties of each $\phi(y)$, as shown by the following lemmas.

Lemma 2.1 (Bernoulli noise / Logistic loss) *Let the conditions (1) and (2) of Lemma 1 be satisfied, and let $\phi(y) = y \log y + (1-y) \log(1-y)$, which corresponds to the logistic-loss Bregman divergence and Bernoulli distribution $p(y) = \mu^y(1-\mu)^{1-y}$, where the mean parameter $\mu = P(y = 1)$. We assume that $0 < \mu^* < 1$, and $0 < \mu^0 < 1$. Then*

$$|\theta^0 - \theta^*| \leq \beta(\epsilon) = 4\epsilon.$$

Proof: Using the definition of the logistic-loss Bregman divergence from Table 1, and the conditions (1) and (2) of Lemma 1, we can write:

$$d_\phi(y, \mu^0) = y \log\left(\frac{y}{\mu^0}\right) + (1-y) \log\left(\frac{1-y}{1-\mu^0}\right) \leq \epsilon,$$

$$d_\phi(y, \mu^*) = y \log\left(\frac{y}{\mu^*}\right) + (1-y) \log\left(\frac{1-y}{1-\mu^*}\right) \leq \epsilon, \quad (15)$$

which implies

$$|d_\phi(y, \mu^0) - d_\phi(y, \mu^*)| \leq 2\epsilon, \quad (16)$$

and, after substituting the expressions 15 into eq. 16, and simplifying, we get

$$|y \log\left(\frac{\mu^0}{\mu^*}\right) + (1-y) \log\left(\frac{1-\mu^0}{1-\mu^*}\right)| \leq 2\epsilon. \quad (17)$$

The above must be satisfied for each $y \in \{0, 1\}$ (the domain of Bernoulli distribution). Thus, we get:

$$(1) \left| \log\left(\frac{1-\mu^0}{1-\mu^*}\right) \right| \leq 2\epsilon \text{ if } y = 0, \text{ and}$$

$$(2) \left| \log\left(\frac{\mu^0}{\mu^*}\right) \right| \leq 2\epsilon \text{ if } y = 1,$$

or, equivalently

$$(1) e^{-2\epsilon} \leq \frac{1-\mu^0}{1-\mu^*} \leq e^{2\epsilon} \text{ if } y = 0, \text{ and}$$

$$(2) e^{-2\epsilon} \leq \frac{\mu^0}{\mu^*} \leq e^{2\epsilon} \text{ if } y = 1.$$

Let us first consider the case of $y = 0$; subtracting 1 from the corresponding inequalities yields

$$e^{-2\epsilon} - 1 \leq \frac{\mu^* - \mu^0}{1 - \mu^*} \leq e^{2\epsilon} - 1 \Leftrightarrow$$

$$\Leftrightarrow (1 - \mu^*)(e^{-2\epsilon} - 1) \leq \mu^* - \mu^0 \leq (1 - \mu^*)(e^{2\epsilon} - 1).$$

By the mean value theorem, $e^x - 1 = e^x - e^0 = \frac{d(e^x)}{dx}|_{\hat{x}} \cdot (x - 0) = e^{\hat{x}}x$, for some $\hat{x} \in [0, x]$ if $x > 0$, or for some $\hat{x} \in [x, 0]$ if $x < 0$. Thus, $e^{-2\epsilon} - 1 = -e^{\hat{x}} \cdot 2\epsilon$, for some $\hat{x} \in [-2\epsilon, 0]$, and since e^x is a continuous monotone increasing function, $e^{\hat{x}} \leq 1$ and thus $e^{-2\epsilon} - 1 \geq -2\epsilon$. Similarly, $e^{2\epsilon} - 1 = e^{\hat{x}} \cdot 2\epsilon$, for some $\hat{x} \in [0, 2\epsilon]$, and since $e^{\hat{x}} \leq e^{2\epsilon}$, we get $e^{2\epsilon} - 1 \leq 2\epsilon \cdot e^{2\epsilon}$. Thus,

$$\begin{aligned} -2\epsilon(1 - \mu^*) \leq \mu^* - \mu^0 \leq 2\epsilon e^{2\epsilon}(1 - \mu^*) &\Rightarrow \\ \Rightarrow |\mu^* - \mu^0| \leq 2\epsilon \cdot e^{2\epsilon}. \end{aligned} \quad (19)$$

Similarly, in case of $y = 1$, we get

$$e^{-2\epsilon} - 1 \leq \frac{\mu^0 - \mu^*}{\mu^*} \leq e^{2\epsilon} - 1.$$

and can apply same derivation as above, and get same result for $|\mu^* - \mu^0|$ as in eq. 19. Finally, since $\theta(\mu) = \phi'(\mu) = \log\left(\frac{\mu}{1-\mu}\right)$, we get

$$\begin{aligned} |\theta^0 - \theta^*| &= \left| \log\left(\frac{\mu^0}{1-\mu^0}\right) - \log\left(\frac{\mu^*}{1-\mu^*}\right) \right| = \\ &= \left| \log\left(\frac{\mu^0}{\mu^*}\right) - \log\left(\frac{1-\mu^0}{1-\mu^*}\right) \right|. \end{aligned}$$

From the eq. 18 we get $|\log\left(\frac{\mu^0}{\mu^*}\right)| \leq 2\epsilon$ and $|\log\left(\frac{1-\mu^0}{1-\mu^*}\right)| \leq 2\epsilon$, which implies

$$|\theta^0 - \theta^*| = \left| \log\left(\frac{\mu^0}{\mu^*}\right) - \log\left(\frac{1-\mu^0}{1-\mu^*}\right) \right| \leq 4\epsilon. \quad \blacksquare$$

Lemma 2.2 (Exponential noise/ Itakura-Saito distance)
Let the conditions (1) and (2) of Lemma 1 be satisfied, and

let $\phi(y) = -\log \mu - 1$, which corresponds to the Itakura-Saito distance $d_\phi(y, \mu) = \frac{y}{\mu} - \log\left(\frac{y}{\mu}\right) - 1$ and exponential distribution $p(y) = \lambda e^{-\lambda y}$, where the mean parameter $\mu = 1/\lambda$. We will also assume that the mean parameter is always separated from zero, i.e. $\exists c_\mu > 0$ such that $\mu \geq c_\mu$. Then

$$|\theta^* - \theta^0| \leq \beta(\epsilon) = \frac{\sqrt{6}\epsilon}{c_\mu}.$$

(18) *Proof:* To establish the result of the lemma we start with inequality $|u - \log u - 1| \leq \epsilon$, where u is $\frac{y}{\mu}$. Replacing u by $z = u - 1$, $z > -1$ gives us $|z - \log(1+z)| \leq \epsilon$. Without loss of generality, let us assume that $\epsilon \leq \frac{1}{18}$. Then the Taylor decomposition of function $z - \log(1+z)$ at the point $z = 0$

$$z - \log(1+z) = \frac{z^2}{2} - \frac{z^3}{3} + \frac{\theta^4}{4}, \text{ for } \theta \in [0, z] \text{ or } [z, 0]$$

implies that

$$\epsilon \geq z - \log(1+z) \geq \frac{z^2}{2} - \frac{z^3}{3} \text{ (since } \frac{\theta^4}{4} \geq 0 \text{)}.$$

This, in turns, implies that $z \leq \frac{1}{3}$ and $\frac{z^2}{2} - \frac{z^3}{3} \geq \frac{z^2}{6}$ for $0 \leq z \leq \frac{1}{3}$.

Hence

$$z - \log(1+z) \geq \frac{z^2}{2} \text{ for } -\frac{1}{3} \leq z \leq 0, \quad (20)$$

$$z - \log(1+z) \geq \frac{z^2}{6} \text{ for } 0 \leq z \leq \frac{1}{3}. \quad (21)$$

Combining together both estimates we get $|z| \leq \sqrt{6}\epsilon$, or

$$|y - \mu| \leq \sqrt{6}\epsilon \cdot \mu,$$

and

$$|\mu^0 - \mu^*| \leq \sqrt{6}\epsilon \cdot \max\{\mu^0, \mu^*\}.$$

Then

$$|\theta^* - \theta^0| = \left| \frac{1}{\mu^0} - \frac{1}{\mu^*} \right| = \left| \frac{\mu^* - \mu^0}{\mu^* \mu^0} \right| \leq \frac{\sqrt{6}\epsilon}{\min\{\mu^*, \mu^0\}} \leq \frac{\sqrt{6}\epsilon}{c_\mu},$$

since by the assumption of the lemma $\min\{\mu^*, \mu^0\} \geq c_\mu$. \blacksquare

We now consider multivariate exponential-family distributions; the next lemma handles the general case of a multivariate Gaussian distribution (not necessarily spherical one that had a diagonal covariance matrix and corresponded to the standard Euclidean distance (see Table 1).

Lemma 2.3 (Non-i.i.d. Multivariate Gaussian noise / Mahalanobis distance) Let $\phi(\mathbf{y}) = \mathbf{y}^T C \mathbf{y}$, which corresponds to the general multivariate Gaussian with concentration matrix C , and Mahalanobis distance $d_\phi(\mathbf{y}, \mu) = \frac{1}{2}(\mathbf{y} - \mu)^T C (\mathbf{y} - \mu)$. If $d_\phi(\mathbf{y}, \mu^0) \leq \epsilon$ and $d_\phi(\mathbf{y}, \mu^*) \leq \epsilon$, then

$$\|\theta^0 - \theta^*\| \leq \sqrt{2\epsilon} \|C^{-1}\|^{1/2} \cdot \|C\|,$$

where $\|C\|$ is the operator norm.

Proof: Since C is (symmetric) positive definite, it can be written as $C = L^T L$ where L defines a linear operator on \mathbf{y} space, and thus

$$\begin{aligned} \epsilon/2 &\geq (\mathbf{y} - \mu)^T C (\mathbf{y} - \mu) = (L(\mathbf{y} - \mu))^T (L(\mathbf{y} - \mu)) = \\ &= \|L(\mathbf{y} - \mu)\|^2. \end{aligned}$$

Also, it is easy to show that $\|C^{-1}\|I \leq C \leq \|C\|I$ (where $\|B\|$ denote the operator norm of B), and that

$$\begin{aligned} \epsilon/2 &\geq \|L(\mathbf{y} - \mu)\|^2 \geq \|L^{-1}\|^{-2} \|\mathbf{y} - \mu\|^2 \Rightarrow \\ &\Rightarrow \|\mathbf{y} - \mu\| \leq \sqrt{\frac{\epsilon}{2}} \|L^{-1}\|. \end{aligned}$$

Then, using triangle inequality, we get

$$\|\mu^* - \mu^0\| \leq \|\mathbf{y} - \mu^0\| + \|\mathbf{y} - \mu^*\| \leq \sqrt{2\epsilon} \|L^{-1}\|.$$

Finally, since $\theta(\mu) = \nabla \phi(\mu) = C\mu$, we get

$$\begin{aligned} \|\theta^0 - \theta^*\| &= \|C\mu^0 - C\mu^*\| \leq \|C\| \cdot \|\mu^0 - \mu^*\| = \\ &= \|C\| \cdot \sqrt{2\epsilon} \|L^{-1}\| \cdot \|C\|. \end{aligned}$$

Note that $\|L^{-1}\| = \|C^{-1}\|^{1/2}$, which concludes the proof. ■

IV. SUMMARY

In this paper, we extend the results of [13] to the more general case of *exponential-family noise* that includes Gaussian noise as a particular case, and yields l_1 -regularized *Generalized Linear Model (GLM)* regression problem. We show that, under standard restricted isometry property (RIP) assumptions on the design matrix, l_1 -minimization can provide a stable recovery of a sparse signal under exponential-family noise assumptions, provided that the noise is sufficiently small and the distribution satisfies certain (sufficient) conditions, such as bounded second derivative of the Legendre conjugate $\phi(y)$ of the log-partition function that uniquely determines the distribution. We also provide distribution-specific proofs for several members of exponential family that do not satisfy the above conditions. Moreover, we show that the results of [13] for a more general case of compressible (rather than sparse) signals can be extended to the exponential-family noise in a similar way.

Clearly, this is work in progress, since so far we were able to demonstrated recoverability of a sparse signal only for a subset of distributions in the exponential family that satisfying certain conditions. Directions for future work include exploring the recoverability conditions for other members of the exponential family, as well as investigating alternative recoverability criteria such as, for example, support recovery (versus recovery in the l_2 -norm sense), which is often a more relevant measure of success in practical applications of sparse regression and compressed sensing.

REFERENCES

- [1] A. Banerjee and S. Merugu and I. S. Dhillon and J. Ghosh. Clustering with Bregman Divergences. *Journal of Machine Learning Research*, 6:1705–1749, October 2005.
- [2] A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh. Clustering with Bregman divergences. In *Proceedings of the Fourth SIAM International Conference on Data Mining*, pages 234–245, April 2004.
- [3] A. Beygelzimer, J. Kephart, and I. Rish. Evaluation of optimization methods for network bottleneck diagnosis. In *In Proc. of ICAC-07*, 2007.
- [4] E. Candes. Compressive sampling. In *Int. Congress of Mathematics*, volume 3, pages 1433–1452, 2006.
- [5] E. Candes and J. Romberg. Quantitative robust uncertainty principles and optimally sparse decompositions. *Foundations of Comput. Math.*, 6(2):227–254, April 2006.
- [6] E. Candes, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. on Information Theory*, 52(2):489–509, February 2006.
- [7] E. Candes and T. Tao. Decoding by linear programming. *IEEE Trans. on Information Theory*, (51(12)):4203–4215, December 2005.
- [8] M.K. Carroll, G.A.Cecchi, I. Rish, R. Garg, and A.R. Rao. Prediction and interpretation of distributed neural activity with sparse models. *Neuroimage*, (44(1)):112–22, 2009.
- [9] G. Chandalia and I. Rish. Blind Source Separation Approach to Performance Diagnosis and Dependency Discovery. In *In Proceedings of IMC-2007*, 2007.
- [10] D. Donoho. Compressed sensing. *IEEE Trans. on Information Theory*, 52(4):1289–1306, April 2006.
- [11] D. Donoho. For most large underdetermined systems of linear equations, the minimal ell_1 norm near-solution approximates the sparsest near-solution. *Communications on Pure and Applied Mathematics*, 59(7):907–934, July 2006.
- [12] D. Donoho. For most large underdetermined systems of linear equations, the minimal ell_1 norm solution is also the sparsest solution. *Communications on Pure and Applied Mathematics*, (59(6)):797–829, June 2006.
- [13] E. Candes and J. Romberg and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics*, 59(8):1207–1223, August 2006.
- [14] T.M. Mitchell, R. Hutchinson, R.S. Niculescu, F. Pereira, X. Wang, M. Just, and S. Newman. Learning to decode cognitive states from brain images. *Machine Learning*, 57:145–175, 2004.
- [15] Mee-Young Park and Trevor Hastie. An L_1 Regularization-path Algorithm for Generalized Linear Models. *JRSSB*, 69(4):659–677, 2007.
- [16] Compressive Sensing Resources. <http://dsp.rice.edu/cs>.
- [17] I. Rish, M. Brodie, S. Ma, N. Odintsova, A. Beygelzimer, G. Grabarnik, and K. Hernandez. Adaptive diagnosis in distributed systems. *IEEE Transactions on Neural Networks (special issue on Adaptive Learning Systems in Communication Networks)*, 16(5):1088–1109, 2005.
- [18] R.T. Rockafeller. *Convex Analysis*. Princeton University Press, 1970.
- [19] A. Zheng, I. Rish, and A. Beygelzimer. Efficient Test Selection in Active Diagnosis via Entropy Approximation. In *Proceedings of UAI-05*, 2005.