

On Separability of Loss Functions, and Revisiting Discriminative Vs Generative Models

Authored by:

Pradeep K. Ravikumar
Adarsh Prasad
Alexandru Niculescu-Mizil

Abstract

We revisit the classical analysis of generative vs discriminative models for general exponential families, and high-dimensional settings. Towards this, we develop novel technical machinery, including a notion of separability of general loss functions, which allow us to provide a general framework to obtain l_1 convergence rates for general M-estimators. We use this machinery to analyze l_1 and l_2 convergence rates of generative and discriminative models, and provide insights into their nuanced behaviors in high-dimensions. Our results are also applicable to differential parameter estimation, where the quantity of interest is the difference between generative model parameters.

1 Paper Body

Consider the classical conditional generative model setting, where we have a binary random response $Y \in \{0, 1\}$, and a random covariate vector $X \in \mathbb{R}^p$, such that $X \perp (Y = i) \mid P_i$ for $i \in \{0, 1\}$. Assuming that we know $P(Y)$ and $\{P_i\}_{i=0}^1$, we can use the Bayes rule to predict the response Y given covariates X . This is said to be the generative model approach to classification. Alternatively, consider the conditional distribution $P(Y \mid X)$ as specified by the Bayes rule, also called the discriminative model corresponding to the generative model specified above. Learning this conditional model directly is said to be the discriminative model approach to classification. In a classical paper [8], the authors provided theoretical justification for the common wisdom regarding generative and discriminative models: when the generative model assumptions hold, the generative model estimators initially converge faster as a function of the number of samples, but have the same asymptotic error rate as discriminative models. And when the generative model assumptions do not hold, the discriminative model estimators eventually overtake the generative model estimators. Their analysis however was for the specific generative-discriminative

model pair of Naive Bayes, and logistic regression models, and moreover, was not under a high-dimensional sampling regime, when the number of samples could even be smaller than the number of parameters. In this paper, we aim to extend their analysis to these more general settings. Doing so however required some novel technical and conceptual developments. To motivate the machinery we develop, consider why the Naive Bayes model estimator might initially converge faster. The Naive Bayes model makes the conditional independence assumption that $P(X-Y) = \prod_{s=1}^p P(X_s-Y)$, so that the parameters of each of the conditional distributions $P(X_s-Y)$ for $s \in \{1, \dots, p\}$ could be estimated independently. The corresponding log-likelihood loss function is thus fully ‘separable’ into multiple components. The logistic regression log-likelihood on the other hand is seemingly much less ‘separable’, and in particular, it does not split into multiple components each of which can be estimated independently. In general, we do not expect the loss functions underlying statistical estimators to be fully separable into multiple components, so that we need a more flexible notion of separability, where different losses could be shown to be separable to differing degrees. In a very related note, though it might seem unrelated at first, the analysis of ‘1 convergence rates of 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

statistical estimators considerably lags that of say ‘2 rates (see for instance, the unified framework of [7], which is suited to ‘2 rates but is highly sub-optimal for ‘1 rates). In part, the analysis of ‘1 rates is harder because it implicitly requires analysis at the level of individual coordinates of the parameter vector. While this is thus harder than an ‘2 error analysis, intuitively this would be much easier if the loss function were to split into independent components involving individual coordinates. While general loss functions might not be so ‘fully separable’, they might perhaps satisfy a softer notion of separability motivated above. In a contribution that would be of independent interest, we develop precisely such a softer notion of separability for general loss functions. We then use this notion of separability to derive ‘1 convergence rates for general M-estimators.

Given this machinery, we are then able to contrast generative and discriminative models. We focus on the case where the generative models are specified by exponential family distributions, so that the corresponding discriminative models are logistic regression models with the generative model sufficient statistics as feature functions. To compare the convergence rates of the two models, we focus on the difference of the two generative model parameters, since this difference is also precisely the model parameter for the discriminative model counterpart of the generative model, via an application of the Bayes rule. Moreover, as Li et al. [3] and others show, the ‘2 convergence rates of the difference of the two parameters is what drives the classification error rates of both generative as well as discriminative model classifiers. Incidentally, such a difference of generative model parameters has also attracted interest outside the context of classification, where it is called differential parameter learning [1, 14, 6]. We thus analyze the ‘1 as well as ‘2 rates for both the generative and discriminative models, focusing on this parameter difference. As we show, unlike the case

of Naive Bayes and logistic regression in low-dimensions as studied in [8], this general highdimensional setting is more nuanced, and in particular depends on the separability of the generative models. As we show, under some conditions on the models, generative and discriminative models not only have potentially different ‘1 rates, but also differing ‘burn in’ periods in terms of the minimum number of samples required in order for the convergence rates to hold. The choice of a generative vs discriminative model, namely that with a better sample complexity, thus depends on their corresponding separabilities. As a minor note, we also show how generative model M -estimators are not directly suitable in high-dimensions, and provide a simple methodological fix in order to obtain better ‘2 rates. We instantiate our results with two running examples of isotropic and non-isotropic Gaussian generative models, and also corroborate our theory with instructive simulations.

2

Background and Setup.

We consider the problem of differential parameter estimation under the following generative model. Let $Y \in \{0, 1\}$ denote a binary response variable, and let $X = (X_1, \dots, X_p) \in \mathbb{R}^p$ be the covariates. For simplicity, we assume $P[Y = 1] = P[Y = 0] = 1/2$. We assume that conditioned on the response variable, the covariates belong to an exponential family, $X|Y \sim P_{Y^*}(\eta)$, where: $P_{Y^*}(\eta) \propto \exp(\eta^T X) A(\eta)$.

(1)

Here, η is the vector of the true canonical parameters, $A(\eta)$ is the log-partition function and X is the sufficient statistic. We assume access to two sets of samples $X_0 = \{x_i\}_{i=1}^{n_0}$ and $X_1 = \{x_i\}_{i=1}^{n_1}$.

(1)

Given these samples, as noted in the introduction, we are particularly interested in estimating the differential parameter $\eta_{\text{diff}} := \eta_1 - \eta_0$, since this is also the model parameter corresponding to the discriminative model, as we show below. In high dimensional sampling settings, we additionally assume that η_{diff} is at most s -sparse, i.e. $\|\eta_{\text{diff}}\|_0 \leq s$.

We will be using the following two exponential family generative models as running examples: isotropic and non-isotropic multivariate Gaussian models. Isotropic Gaussians (IG) Let $X = (X_1, \dots, X_p) \in \mathbb{R}^p$ be an isotropic gaussian random variable; its density can be written as: $P(x) = \frac{1}{(2\pi)^{p/2}} \exp(-\frac{1}{2} x^T x)$.

Gaussian MRF (GMRF). Let $X = (X_1, \dots, X_p) \in \mathbb{R}^p$ denote a zero-mean gaussian random vector; its density is fully-parametrized as by the inverse covariance or concentration matrix $\Sigma^{-1} = \Theta$

and can be written as: $P(x) = \frac{1}{(2\pi)^{p/2} \det(\Theta)^{1/2}} \exp(-\frac{1}{2} x^T \Theta x)$

(2)

$$\begin{aligned} & \frac{1}{2} \mathbf{x}^T \mathbf{x} \\ & (3) \end{aligned}$$

Let $d = \max_j \|\mathbf{x}_j\|_0$ is the maximum number non-zeros in a row of \mathbf{X} . Let $\|\mathbf{X}\|_1 = \sum_{j=1}^p \|\mathbf{x}_j\|_1$, where $\|\mathbf{M}\|_1$ is the '1'/'1 operator norm given by $\|\mathbf{M}\|_1 = \max_{j=1,2,\dots,p} \sum_{k=1}^n |\mathbf{M}_{jk}|$.

Generative Model Estimation. Here, we proceed by estimating the two parameters $\{\theta_i\}_{i=1}^d$ individually. Letting $\hat{\theta}_1$ and $\hat{\theta}_0$ be the corresponding estimators, we can then estimate the difference of the parameters as $\hat{\theta}_{\text{diff}} = \hat{\theta}_1 - \hat{\theta}_0$. The most popular class of estimators for the individual parameters is based on Maximum likelihood Estimation (MLE), where we maximize the likelihood of the given data. For isotropic gaussians, the negative log-likelihood function can be written as: $-\log p(\mathbf{X}|\theta) = \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i + \frac{n}{2} \log |\Sigma|$, where $\Sigma = \sigma^2 \mathbf{I}$. In the case of GGMS the negative log-likelihood function can be written as: $-\log p(\mathbf{X}|\theta) = \frac{1}{2} \sum_{i=1}^n \mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i + \frac{n}{2} \log |\Sigma|$, where $\Sigma = \sigma^2 \mathbf{I}$. DD EE b LnGGM $(\theta) = -\log p(\mathbf{X}|\theta)$, (5) $\mathbf{P} \mathbf{P} \mathbf{b} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$ is the sample covariance matrix and $\text{tr}(\mathbf{U} \mathbf{V}) = \sum_{i=1}^n \mathbf{U}_{ii} \mathbf{V}_{ii}$ where \mathbf{U}, \mathbf{V} are symmetric matrices. In high-dimensional sampling regimes ($n \ll p$), regularized MLEs, for instance with '1'-regularization under the assumption of sparse model parameters, have been widely used [11, 10, 2]. **Discriminative Model Estimation.** Using Bayes rule, we have that: $P[X=1|Y=1]P[Y=1] = P[X=1|Y=0]P[Y=0] + P[X=0|Y=1]P[Y=1]$ (6) $\frac{1}{1 + \exp(-(\mathbf{h}^T \mathbf{x} + c))}$, where $c = A(\theta_0) - A(\theta_1)$. The conditional distribution is simply a logistic regression model, with the generative model sufficient statistics as the features, and with optimal parameters being precisely θ_{diff} of the generative model parameters. The corresponding negative log-likelihood function can be written as $n \log p(\mathbf{X}|\theta) = \sum_{i=1}^n \log p(\mathbf{x}_i|\theta)$ (7) $n \sum_{i=1}^n \log p(\mathbf{x}_i|\theta) =$

where $\log(t) = \log(1 + \exp(t))$. In high dimensional sampling regimes, under the assumption that the model parameters are sparse, we would use the '1'-penalized version $\hat{\theta}_{\text{diff}}$ of the MLE (7) to estimate θ_{diff} .

Outline. We proceed by studying the more general problem of '1' error for parameter estimation for any loss function $L(\theta)$. Specifically, consider the general M-estimation problem, where we are given n i.i.d samples Z_1, \dots, Z_n , $Z_i \in \mathcal{Z}$ from some distribution P , and we are interested in estimating some parameter θ^* of the distribution P . Let $\ell : \mathcal{R}^p \times \mathcal{Z} \rightarrow \mathbb{R}$ be a twice differentiable and convex function which assigns a loss $\ell(\theta; z)$ to any parameter $\theta \in \mathcal{R}^p$, for a given z observation z . Also assume that the loss is Fisher consistent so that $\theta^* \in \arg\min_{\theta} L(\theta)$ where $\text{def } L(\theta) = \mathbb{E} \ell(\theta; Z)$ is the population loss. We are then interested in analyzing the M-estimators $\hat{\theta}_n$ that minimize P_n the empirical loss i.e. $\hat{\theta}_n \in \arg\min_{\theta} L_n(\theta)$, or regularized versions thereof, where $L_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; Z_i)$.

We introduce a notion of the separability of a loss function, and show how more separable losses require fewer samples to establish convergence for $\hat{\theta}_n$. We then instantiate our separability results from this general setting for both generative and discriminative models. We calculate the number of samples

required for generative and discriminative approaches to estimate the differential ℓ_2 parameter θ_{diff} , for consistent convergence rates with respect to ℓ_1 and ℓ_2 norm. We also discuss the consequences of these results for high dimensional classification for Gaussian Generative models. 3

3

Separability

Let $R(\theta; \theta_0) = \mathbb{E} \ell_n(\theta + \theta_0) - \mathbb{E} \ell_n(\theta_0) - \nabla \ell_n(\theta_0)^\top \theta_0$ be the error in the first order approximation of the gradient at θ_0 . Let $B_1(r) = \{\theta \mid \|\theta\|_1 \leq r\}$ be an ℓ_1 ball of radius r . We begin by analyzing the low dimensional case, and then extend it to high dimensions. 3.1

Low Dimensional Sampling Regimes

In low dimensional sampling regimes, we assume that the number of samples $n \gg p$. In this setting, we make the standard assumption that the empirical loss function $\ell_n(\theta)$ is strongly convex. Let $\theta^* = \arg\min_{\theta} \ell_n(\theta)$ denote the unique minimizer of the empirical loss function. We begin by defining a notion of separability for any such empirical loss function ℓ_n . Definition 1. ℓ_n is (γ, κ) locally separable around θ^* if the remainder term $R(\theta; \theta^*)$ satisfies: $\|R(\theta; \theta^*)\|_1 \leq \gamma \|\theta - \theta^*\|_1^\kappa$

1

γ

$\kappa \geq 1/8$

$2 B_1(\kappa)$

This definition might seem a bit abstract, but for some general intuition, γ indicates the region where it is separable, κ indicates the conditioning of the loss, while it is that quantifies the degree of separability: the larger it is, the more separable the loss function. Next, we provide some additional intuition on how a loss function's separability is connected to (γ, κ) . Using the mean-value theorem, we can write $\|R(\theta; \theta^*)\|_1 \leq \mathbb{E} \ell_n(\theta^* + t(\theta - \theta^*)) - \mathbb{E} \ell_n(\theta^*) - \nabla \ell_n(\theta^*)^\top (\theta - \theta^*)$ for some $t \in (0, 1)$. This can be further simplified as $\|R(\theta; \theta^*)\|_1 \leq \mathbb{E} \ell_n(\theta^* + t(\theta - \theta^*)) - \mathbb{E} \ell_n(\theta^*) - \nabla \ell_n(\theta^*)^\top (\theta - \theta^*) \leq \mathbb{E} \ell_n(\theta^* + t(\theta - \theta^*)) - \mathbb{E} \ell_n(\theta^*)$. Hence, γ and $1/\kappa$ measure the smoothness of Hessian (w.r.t. the ℓ_1/ℓ_1 matrix norm) in the neighborhood of θ^* , with γ being the smoothness exponent, and $1/\kappa$ being the smoothness constant. Note that the Hessian of the loss function $\mathbb{E} \ell_n(\theta)$ is a random matrix and can vary from being a diagonal matrix for a fully-separable loss function to a dense matrix for a heavily-coupled loss function. Moreover, from standard concentration arguments, the ℓ_1/ℓ_1 matrix norm for a diagonal ("separable") subgaussian random matrix has at most logarithmic dimension dependence, but for a dense ("non-separable") random matrix, the ℓ_1/ℓ_1 matrix norm could possibly scale linearly in the dimension. Thus, the scaling of ℓ_1/ℓ_1 matrix norm gives us an indication how "separable" the matrix is. This intuition is captured by (γ, κ) , which we further elaborate in future sections by explicitly deriving (γ, κ) for different loss functions and use them to derive ℓ_2 and ℓ_1 convergence rates. Theorem 1. Let ℓ_n be a strongly convex loss function which is (γ, κ) locally separable function $\theta^* \in B_1(\kappa)$ around θ^* . Then, if $\mathbb{E} \ell_n(\theta^*) \leq \min\{\gamma^2, \gamma\}$

where $\gamma = \mathbb{E} \ell_n(\theta^*)$

1 1

.
 \hat{b}
 $\hat{\beta}$

1
 $\hat{\beta} - 2\hat{\beta} - \sqrt{r \ln(p)} - 1$

Proof. (Proof Sketch). The proof begins by constructing a suitable continuous function F , for which $\hat{b} = \hat{\beta}$ is the unique fixed point. Next, we show that $F(B_1(r)) \subset B_1(r)$ for $r = 2\sqrt{r \ln(p)} - 1$. Since F is continuous and the 1 -ball is convex and compact, the contraction property coupled with Brouwer's fixed point theorem [9], shows that there exists some fixed point of F , such that $\hat{\beta} - 1 \leq \hat{\beta} - \sqrt{r \ln(p)} - 1$. By uniqueness of the fixed point, we then establish our result. See Figure 1 for a geometric description and Section A for more details 3.2

High Dimensional Sampling Regimes

In high dimensional sampling regimes ($n \ll p$), estimation of model parameters is typically an under-determined problem. It is thus necessary to impose additional assumptions on the true model parameter β . We will focus on the popular assumption of sparsity, which entails that the number of non-zero coefficients of β is small, so that $\|\beta\|_0 \ll p$. For this setting, we will be focusing in particular on ℓ_1 -regularized empirical loss minimization:

Follows from the concentration of subgaussian maxima [12]

4
 $F(\hat{b})$
 $F(\hat{b}) = \hat{b}$
 $F(B_1(2\sqrt{r \ln(p)} - 1))$
 $B_1(2\sqrt{r \ln(p)} - 1)$

Figure 1: Under the conditions of Theorem 1, $F(\cdot) = \sqrt{r} \ln(p) + 1 (R(\cdot; \hat{\beta}) + \sqrt{r \ln(p)})$ is contractive over $B_1(2\sqrt{r \ln(p)} - 1)$ and has $\hat{b} = \hat{\beta}$ as its unique fixed point. Using these two observations, we can conclude that $\hat{b} \leq 2\sqrt{r \ln(p)} - 1$.

1
 \hat{b}
 n
 $= \argmin_n \ln(\cdot) + \cdot$
 n
(8)
 $\sqrt{r} - 1$

Let $S = \{i : \beta_i \neq 0\}$ be the support set of the true parameter and $M(S) = \{v : v_S = 0\}$ be the corresponding subspace. Note that under a high-dimensional sampling regime, we can no longer assume that the empirical loss $\ln(\cdot)$ is strongly convex. Accordingly, we make the following set of assumptions:
• Assumption 1 (A1): Positive Definite Restricted Hessian. $\sqrt{r} \ln(p) \leq \min_i \lambda_i$
• Assumption 2 (A2): Irrepresentability. There exists some $\alpha \in (0, 1]$ such that $\sqrt{r} \ln(p) \leq \sqrt{r} \ln(p) + \alpha \sqrt{r} \ln(p)$

1

?

)) are such

Proof. (Proof Sketch). The proof invokes the primal-dual witness argument [13] which when combined with Assumption 1-3, gives $\hat{b}_n \in M(S)$ and that \hat{b}_n is the unique solution of the restricted problem. The rest of the proof proceeds similar to Theorem 1, by constructing a suitable function $F : \mathbb{R}^S \rightarrow \mathbb{R}^S$ for which $b = \hat{b}_n$ is the unique fixed point, and showing that F is contractive over $B_1(r; \cdot)$ for $r = 2\epsilon \left(\frac{1}{\epsilon} \ln \left(\frac{1}{\epsilon} \right) + 1 \right)$. See Section B for more details. Discussion. Theorems 1 and 2 provide a general recipe to estimate the number of samples required by any loss $\ell(\cdot, z)$ to establish ℓ_1 convergence. The first step is to calculate the separability constants (γ, κ, η) for the corresponding empirical loss function L_n . Next, since the loss ℓ is Fisher consistent, $\ell(\hat{b}_n, z) = 0$, the upper bound on $\frac{1}{n} \sum_{i=1}^n \ell(\hat{b}_n, z_i)$ can be shown to hold by analyzing the so that $rL(\hat{b}_n)$ concentration of $rL_n(\hat{b}_n)$ around its mean. We emphasize that we do not impose any restrictions on the values of (γ, κ, η) . In particular, these can scale with the number of samples n ; our results hold so long as the number of samples n satisfy the conditions of the theorem. As a rule of thumb, the smaller that either or get for any given loss ℓ , the larger the required number of samples.

4

ℓ_1 -rates for Generative and Discriminative Model Estimation

In this section we study the ℓ_1 rates for differential parameter estimation for the discriminative and generative approaches. We do so by calculating the separability of discriminative and generative loss functions, and then instantiate our previously derived results. 4.1

Discriminative Estimation

As discussed before, the discriminative approach uses ℓ_1 -regularized logistic regression with the sufficient statistic as features to estimate the differential parameter. In addition to A1-A3, we Pn 2 assume column normalization of the sufficient statistics, i.e. $\sum_{i=1}^n \left(\sum_{j=1}^p (x_{ij})^2 \right) = n$. Let $n = \max_i \sum_{j=1}^p (x_{ij})^2$, $\gamma_n = \max_i \sum_{j=1}^p (x_{ij})^2$. Firstly, we characterize the separability of the logistic loss. Lemma 1. The logistic regression negative log-likelihood L_n^{Logistic} from (7) is $2, s \geq 1$ restricted local separable around \hat{b}_n . Combining Lemma 1 with Theorem 2, we get the following corollary. Corollary 3. (Logistic Regression) Consider the model in (1), then there exist q universal positive constants C_1, C_2 and C_3 such that for n differential estimate \hat{b}_{diff} , satisfies 4.2

$\text{support}(\hat{b}_{\text{diff}}) \leq$

$C_1 \frac{1}{\gamma_n} \frac{1}{\epsilon^2} \frac{1}{s^2}$

$\frac{1}{\epsilon} \text{support}(\hat{b}_{\text{diff}})$

$2 \frac{4}{n} \frac{1}{\gamma_n}$

and

Generative Estimation

$\log p$ and

\hat{b}_{diff}

$\frac{1}{\epsilon} \hat{b}_{\text{diff}}$

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_i) \\
& = C_2 \\
& \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{x}_i) \\
& = C_3 \\
& \log p(\mathbf{x}) \\
& \text{the discriminative} \\
& \log p(\mathbf{x})
\end{aligned}$$

We characterize the separability of Generative Exponential Families. The negative log-likelihood function can be written as: $L_n(\theta) = \sum_{i=1}^n \log p(\mathbf{x}_i | \theta)$, where $n = \sum_{i=1}^n 1(\mathbf{x}_i)$. In this setting, the remainder term is independent of the data and can be written as $R(\theta) = r_A(\theta) + r_B(\theta) - r_C(\theta)$ and $r_L(\theta) = E[\sum_{i=1}^n \log p(\mathbf{x}_i | \theta)]$. Hence, $-r_L(\theta) - 1$ is a measure of how well the sufficient statistics concentrate around their mean. Next, we show the separability of our running examples Isotropic Gaussians and Gaussian Graphical Models. Lemma 2. The isotropic Gaussian negative log-likelihood L_nIG from (4) is $(\frac{1}{2}, 1, 1)$ locally separable around θ^* . Lemma 3. The Gaussian MRF negative log-likelihood L_nGGM from (5) is $(2, 3d, 3d, 1)$ restricted locally separable around θ^* .

Comparing Lemmas 1, 2 and 3, we see that the separability of the discriminative model loss depends weakly on the feature functions. On the other hand, the separability for the generative model loss depends critically on the underlying sufficient statistics. This has consequences for their differing sample complexities for differential parameter estimation, as we show next. Corollary 4. (Isotropic Gaussians) Consider the model in (2). Then there exist universal constants C_1, C_2, C_3 such that if the number of samples scale as $n \geq C_1 \log p$, then with probability at least $1 - 1/p^{C_2}$, the generative estimate of the differential parameter θ_{diff} satisfies $\|\hat{\theta}_{diff} - \theta_{diff}\| \leq C_3 \cdot \frac{1}{n}$.

Comparing Corollary 3 and Corollary 4, we see that for isotropic gaussians, both the discriminative and generative approach achieve the same $1/n$ convergence rates, but at different sample complexities. Specifically, the sample complexity for the generative method depends only logarithmically on the dimension p , and is independent of the differential sparsity s , while the sample complexity of the discriminative method depends on the differential sparsity s . Therefore in this case, the generative method is strictly better than its discriminative counterpart, assuming that the generative model assumptions hold. Corollary 5. (Gaussian MRF) Consider the model in (3), and suppose that the scaled covariates $X_k / \sqrt{\lambda_k}$ are subgaussian with parameter 2 . Then there exist universal positive constants C_2, C_3, C_4 such that if the number of samples for the two generative models scale as $n_i \geq C_2 \cdot \frac{1}{\lambda_i} \log p$, for $i \in \{0, 1\}$, then with probability at least $1 - 1/p^{C_3}$, the generative estimate of the differential parameter, $\hat{\theta}_{diff} = \hat{\theta}_1 - \hat{\theta}_0$, satisfies $\|\hat{\theta}_{diff} - \theta_{diff}\| \leq C_4 \cdot \frac{1}{n_i}$ for $i \in \{0, 1\}$.

Comparing Corollary 3 and Corollary 5, we see that for Gaussian Graphical Models, both the discriminative and generative approach achieve the same $1/n$ convergence rates, but at different sample complexities. Specifically, the sample

complexity for the generative method depends only on row-wise sparsity of the individual models $d_{2??}$, and is independent of sparsity s of the differential ℓ_1 parameter θ_{diff} . In contrast, the sample complexity of the discriminative method depends only on the sparsity of the differential parameter, and is independent of the structural complexities of the individual model parameters. This suggests that in high dimensions, even when the generative model assumptions hold, generative methods might perform poorly if the underlying model is highly non-separable (e.g. $d = \theta(p)$), which is in contrast to the conventional wisdom in low dimensions.

Related Work. Note that results similar to Corollaries 3 and 5 have been previously reported in [11, 5] separately. Under the same set of assumptions as ours, Li et al. [5] provide a unified analysis for support recovery and ℓ_1 -bounds for ℓ_1 -regularized M-estimators. While they obtain the same rates as ours, their required sample complexities are much higher, since they do not exploit the separability of the underlying loss function. As one example, in the case of GMRFs, their results require the number of samples to scale as $n \gtrsim k^2 \log p$, where k is the total number of edges in the graph, which is sub-optimal, and in particular does not match the GMRF-specific analysis of [11]. On the other hand, our unified analysis is tighter, and in particular, does match the results of [11].

5

ℓ_2 -rates for Generative and Discriminative Model Estimation

In this section we study the ℓ_2 rates for differential parameter estimation for the discriminative and generative approaches. 5.1

Discriminative Approach

The bounds for the discriminative approach are relatively straightforward. Corollary 3 gives bounds $b \propto \text{support}(\theta_{\text{diff}})$. Since the true model parameter is on the ℓ_1 error and establishes that $\text{support}(\theta_{\text{diff}})$ is s -sparse, $\|\theta_{\text{diff}}\|_1 \leq 0$, the ℓ_2 error can be simply bounded as $s \leq k \leq b \leq k \leq 1$. 7

5.2

Generative Approach

In the previous section, we saw that the generative approach is able to exploit the inherent separability of the underlying model, and thus is able to get ℓ_1 rates for differential parameter estimation at a much lower sample complexity. Unfortunately, it does not have support consistency. Hence a naive generative estimator will have an ℓ_2 error scaling with $p \log n$, which in high dimensions, would make it unappealing. However, one can exploit the sparsity of θ_{diff} and get better rates of convergence in ℓ_2 -norm by simply soft-thresholding the generative estimate. Moreover, soft-thresholding also leads to support consistency. Definition 3. We denote the soft-thresholding operator ST ST

$$\begin{aligned} \text{ST}(\theta) &= \arg\min_w \sum_{i=1}^n |w_i - \theta_i| \\ \text{ST}(\theta), \text{ defined as:} \\ &= \end{aligned}$$

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} + \frac{1}{2} \frac{\| \mathbf{y}_i - \mathbf{w} \|_2}{\| \mathbf{y}_i \|_2} \right) .$$

Lemma 4. Suppose $\mathbf{y} = \mathbf{y}_0 + \mathbf{y}_1$ for some s -sparse \mathbf{y}_0 . Then there exists a universal constant C_1 such that for $n \geq \frac{1}{\epsilon^2} \log \frac{1}{\epsilon}$, $p \geq \frac{1}{\epsilon^2} \log \frac{1}{\epsilon}$, $\mathbf{y}_0 \in \mathbb{R}^n$, $\mathbf{y}_1 \in \mathbb{R}^n$, $\mathbf{y}_0 \perp \mathbf{y}_1$ and $\mathbf{y}_0 \in \mathbb{R}^n$, $\mathbf{y}_1 \in \mathbb{R}^n$, $\mathbf{y}_0 \perp \mathbf{y}_1$ (10) Note that this is a completely deterministic result and has no sample complexity requirement. Motivated by this, we introduce a thresholded generative estimator that has two stages: (a) compute $\hat{\mathbf{y}}_{\text{diff}}$ using the generative model estimates, and (b) soft-threshold the generative estimate with $n = \frac{1}{\epsilon^2} \log \frac{1}{\epsilon}$. An elementary application of Lemma 4 can then be shown to provide ‘2 error 1

ϵ bounds for $\hat{\mathbf{y}}_{\text{diff}}$ given its ‘1 error bounds, and that the true parameter \mathbf{y}_{diff} is s -sparse. We instantiate these ‘2 -bounds via corollaries for our running examples of Isotropic Gaussians, and Gaussian MRFs.

Lemma 5. (Isotropic Gaussians) Consider the model in (2). Then there exist universal constants C_1, C_2, C_3 such that if the number of samples scale as $n \geq C_1 \log p$, then with probability $\geq 1 - \frac{1}{p^{C_2}}$, the soft-thresholded generative estimate of the differential parameter $\mathbf{ST}_n(\hat{\mathbf{y}}_{\text{diff}})$, with q the soft-thresholding parameter set as $n = c \log p$ for some constant c , satisfies: $\| \mathbf{r} - \mathbf{y}_{\text{diff}} \|_2 \leq C_3 \cdot \frac{1}{n^2}$

Lemma 6. MRF) Consider the model in Equation 3, and suppose that the covarip (Gaussian \mathbf{y} are subgaussian with parameter $\frac{1}{2}$. Then there exist universal positive constants C_2, C_3, C_4 such that if the number of samples for the two generative models scale as $n_i \geq C_2 \log \frac{1}{\epsilon}$, for $i \in \{0, 1\}$, for $i \in \{0, 1\}$, then with probability at least $1 - \frac{1}{p^{C_3}}$, $\| \mathbf{r} - \mathbf{y}_{\text{diff}} \|_2 \leq C_4 \cdot \frac{1}{n^2}$, with the soft-thresholded generative estimate of the differential parameter, $\mathbf{ST}_n(\hat{\mathbf{y}}_{\text{diff}})$, with q thresholding parameter set as $n = c \log p$ for some constant c , satisfies: $\| \mathbf{r} - \mathbf{y}_{\text{diff}} \|_2 \leq C_4 \cdot \frac{1}{n^2}$ Comparing Lemmas 5 and 6 to Section 5.1, we can see that the additional soft-thresholding step allows the generative methods to achieve the same ‘2 -error rates as the discriminative methods, but at different sample complexities. The sample complexities of the generative estimates depend on the separabilities of the individual models, and is independent of the differential sparsity s , where as the sample complexity of the discriminative estimate depends only on the differential sparsity s .

6

Experiments: High Dimensional Classification

In this section, we corroborate our theoretical results on ‘2 -error rates for generative and discriminative model estimators, via their consequences for high dimensional classification. We focus on the case of isotropic Gaussian generative models $\mathbf{X} \sim \mathcal{N}(\mathbf{Y}, \mathbf{I}_p)$, where $\mathbf{Y} \in \mathbb{R}^p$, $\mathbf{Y} \in \mathbb{R}^p$ are unknown 8

0-1 Error for $s=4, p=512, d=1$

0.5

0.4

0-1 Error

0.3 0.25

0.35 0.3 0.25
 0.2
 50
 100
 150
 200
 250
 300
 350
 400
 0.15
 0.3
 0.2
 0
 50
 100
 150
 200
 n
 (a)
 0.35
 0.25
 0.2
 0
 Gen-Thresh Logistic
 0.45
 0.4
 0-1 Error
 0-1 Error
 0.4
 0-1 Error for $s=64, p=512, d=1$
 0.5 Gen-Thresh Logistic
 0.45
 0.35
 0.15
 0-1 Error for $s=16, p=512, d=1$
 0.5 Gen-Thresh Logistic
 0.45
 250
 300
 350
 400
 0.15
 0
 50
 100

n
 (b)
 $s = 4, p = 512$
 150
 200
 250
 300
 350
 400

n
 (c)
 $s = 16, p = 512$
 Figure 2: Effect of sparsity s on excess 0
 $s = 64, p = 512$
 1 error.

and β_1, β_0 is s -sparse. Here, we are interested in a classifier $C : \mathbb{R}^p \rightarrow \{0, 1\}$ that achieves low classification error $\mathbb{E}_{X,Y} [1 \{C(X) \neq Y\}]$. Under this setting, it can be shown that the Bayes classifier, that achieves the lowest possible classification error, is given by the linear discriminant $\beta^T \beta^T$ classifier $C^*(x) = 1 \text{ if } x^T w^* + b^* \geq 0$, where $w^* = (\beta_1, \beta_0)$ and $b^* = 0.0211$. Thus, the coefficient w^* of the linear discriminant is precisely the differential parameter, which can be estimated via both generative and discriminative approaches as detailed in the previous section. Moreover, the classification error can also be related to the ℓ_2 error of the estimates. Under some mild assumptions, Tibshirani et al. [3] showed that for any linear classifier $C(x) = 1 \text{ if } x^T w + b \geq 0$, the excess classification error can be bounded as:

$$\mathbb{E}_{X,Y} [1 \{C(X) \neq Y\}] \leq C_1 \left(\frac{\|w - w^*\|_2^2}{\|w^*\|_2^2} + \frac{\|b - b^*\|_2^2}{\|b^*\|_2^2} \right),$$

for some constant $C_1 \geq 0$, and where $E(C) = \mathbb{E}_{X,Y} [1 \{C(X) \neq Y\}]$ is the excess 0-1 error. In other words, the excess classification error is bounded by a constant times the ℓ_2 error of the differential parameter estimate. Methods. In this setting, as discussed in previous sections, the discriminative model is simply a logistic regression model with linear features (6), so that the discriminative estimate of the differential parameter w as well as the constant bias term b can be simply obtained via ℓ_1 -regularized logistic regression. For the generative estimate, we use our two stage estimator from Section 5, which proceeds by estimating β_0, β_1 using the empirical means, and then estimating the differential parameter by soft-thresholding the difference of the generative model parameter estimates $w_{bT} = S_T(n)(\beta_1 - \beta_0)$

where $\beta_{bT} =$
 n
 $1/2$

$$= C_1$$

$\log p/n$

for some constant C_1 . The corresponding estimate for b^* is given by

$$\frac{1}{n} \sum_{i=1}^n b_i, \quad b_1 + \dots + b_n.$$

Experimental Setup. For our experimental setup, we consider isotropic Gaussian models with $1/s$ means $\theta_0 = 1/p$, $\theta_1 = 1/p + 1/s$, and vary the sparsity level s . For both methods, $0 \leq s \leq p$ we set the regularization parameter λ as $\lambda = \log(p)/n$. We report the excess classification error for the two approaches, averaged over 20 trials, in Figure 2. Results. As can be seen from Figure 2, our two-staged thresholded generative estimator is always better than the discriminative estimator, across different sparsity levels s . Moreover, the sample complexity or “burn-in” period of the discriminative classifier strongly depends on the sparsity level, which makes it unsuitable when the true parameter is not highly sparse. For our two-staged generative estimator, we see that the sparsity s has no effect on the “burn-in” period of the classifier. These observations validate our theoretical results from Section 5. 2

See Appendix J for cross-validated plots.

9

Acknowledgements A.P. and P.R. acknowledge the support of ARO via W911NF-12-1-0390 and NSF via IIS-1149803, IIS-1447574, DMS-1264033, and NIH via R01 GM117594-01 as part of the Joint DMS/NIGMS Initiative to Support Research at the Interface of the Biological and Mathematical Sciences.

2 References

- [1] Alberto de la Fuente. From “differential expression” to “differential network-identification of dysfunctional regulatory networks in diseases. Trends in genetics, 26(7):326–333, 2010.
- [2] Christophe Giraud. Introduction to high-dimensional statistics, volume 138. CRC Press, 2014.
- [3] Tianyang Li, Adarsh Prasad, and Pradeep K Ravikumar. Fast classification rates for high-dimensional gaussian generative models. In Advances in Neural Information Processing Systems, pages 1054–1062, 2015.
- [4] Tianyang Li, Xinyang Yi, Constantine Carmanis, and Pradeep Ravikumar. Minimax gaussian classification & clustering. In Artificial Intelligence and Statistics, pages 1–9, 2017.
- [5] Yen-Huan Li, Jonathan Scarlett, Pradeep Ravikumar, and Volkan Cevher. Sparsistency of 1-regularized m-estimators. In AISTATS, 2015.
- [6] Song Liu, John A Quinn, Michael U Gutmann, Taiji Suzuki, and Masashi Sugiyama. Direct learning of sparse changes in markov networks by density ratio estimation. Neural computation, 26(6):1169–1197, 2014.
- [7] Sahand Negahban, Bin Yu, Martin J Wainwright, and Pradeep K Ravikumar. A unified framework for highdimensional analysis of m-estimators with decomposable regularizers. In Advances in Neural Information Processing Systems, pages 1348–1356, 2009.
- [8] Andrew Y Ng and Michael I Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. Advances in neural information processing systems, 2:841–848, 2002.
- [9] James M Ortega and Werner

C Rheinboldt. Iterative solution of nonlinear equations in several variables. SIAM, 2000. [10] Pradeep Ravikumar, Martin J Wainwright, John D Lafferty, et al. High-dimensional ising model selection using ℓ_1 -regularized logistic regression. *The Annals of Statistics*, 38(3):1287?1319, 2010. [11] Pradeep Ravikumar, Martin J Wainwright, Garvesh Raskutti, Bin Yu, et al. High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electronic Journal of Statistics*, 5: 935?980, 2011. [12] JM Wainwright. High-dimensional statistics: A non-asymptotic viewpoint. preparation. University of California, Berkeley, 2015. [13] Martin J Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using-constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183?2202, 2009. [14] Sihai Dave Zhao, T Tony Cai, and Hongzhe Li. Direct estimation of differential networks. *Biometrika*, page asu009, 2014.