

Non-Manipulable Machine Learning: The Incentive Compatibility of Lasso

MEHMET CANER*

KFIR ELIAZ†

January 5, 2021

Abstract

We consider situations where a user feeds her attributes to a machine learning method that tries to predict her best option based on a random sample of other users. The predictor is incentive-compatible if the user has no incentive to misreport her covariates. Focusing on the popular Lasso estimation technique, we borrow tools from high-dimensional statistics to characterize sufficient conditions that ensure that Lasso is incentive compatible in large samples. In particular, we show that incentive compatibility is achieved if the tuning parameter is kept above some threshold. We present simulations that illustrate how this can be done in practice.

1 Introduction

Rapid advances in machine learning methods for analyzing big data have given rise to automated systems that employ these methods to predict the best fitting outcomes for users based on their personal characteristics. For example, many online platforms try to predict which content - a song, a video, a post, or an article - is the best fit for each user. Medical providers have also begun using machine learning techniques to automate check-ups and test appointments for patients based on their medical history. Typically, these automated systems use data from past users to estimate a model that relates the best fit for a user (such as the most preferred content or the appropriate medical test) to her characteristics. These estimates are then applied to a new user's characteristics, which she discloses either actively or passively via her past online behavior (which may be reflected in his cookies or collected by his browser). Given the growing interaction of users with such automated systems, it is only natural to ask whether a user should truthfully disclose her characteristics?

*North Carolina State University, Nelson Hall, Department of Economics, NC 27695. Email: mcaner@ncsu.edu. We thank Simon Fraser University Department of Economics for the virtual seminar, and Anders Kock and Ran Spiegler for comments. We also thank Columbia University, Department of Economics for the hospitality where this research is initiated, when both authors were visitors in 2018-2019.

†School of Economics, Tel-Aviv University and Eccles School of Business, the University of Utah. Email: kfire@tauex.tau.ac.il.

If the information the user discloses is also used to exploit her (say, by providing it to third parties for advertising or price discrimination), then the user has an obvious reason not to reveal her private information. The question is whether special features of some popular machine learning methods introduce an incentive to misreport one’s personal characteristics even when this information will be used *solely* for predicting her best outcome?¹ This question is of crucial importance: If individuals submit false reports to systems that rely on these reports for estimation and predictions, then the conclusions drawn from such estimates and predictions will be wrong and may lead to quite undesirable outcomes (e.g., think of an automated medical platform that schedules tests for patients based on false reports on attributes such as smoking, drinking and physical exercise).

To address the above question, we consider a stylized environment where each user i ’s ideal option is a linear function f of her privately observed attributes $X_i = (X_{i,1}, \dots, X_{i,p})'$ such that $f(X_i) = X_i'\beta_0$. A user may not know the values of the coefficients β_0 , in which case she would have some (possibly degenerate) prior beliefs over them. A “statistician”, who represents some automated prediction platform has a sample of the attributes of n users and *noisy* observations on their ideal options. For instance, suppose $f(X_i)$ is the optimal dosage of some medication when taken immediately at the onset of symptoms, conditional on the patient’s medical history X_i , but the statistician observes the dosage that was given after some delay. Similarly, $f(X_i)$ may be the mix of news and reality shows that a user with attributes X_i actually watches, but the statistician observes only self reports by a user who may have forgotten exactly what he watched.

The statistician uses her sample to estimate the function f by computing an estimate $\hat{\beta}$ of the true coefficients β_0 . The statistician wishes to apply these estimates to predict the ideal option of a new user, $n + 1$, whose true attributes X_{n+1} are not observed by the statistician. This new user must decide what vector of attributes \tilde{X}_{n+1} (which may *differ* from the truth) to report to the statistician. In making this decision, the new user takes into account her beliefs about the statistician’s sample (the new user only knows the distribution from which the sample is drawn, but she does not observe its realization), and her beliefs about the true parameters β_0 . The statistician then plugs the new user’s reported attributes into the estimated function and gives the user the option $\tilde{X}'_{n+1}\hat{\beta}$, which is the statistician’s estimate of the user’s ideal option based her report. The new user’s expected loss from a report \tilde{X}_{n+1} is given by the mean square error between her expectation of the ideal option $X'_{n+1}\beta_0$ and her assigned option $\tilde{X}'_{n+1}\hat{\beta}$. The statistician’s estimator is (ex-ante) *incentive-compatible*, if in expectation, the new user has *no* incentive to deviate from truthful reporting for *any* prior belief on β_0 : i.e., if for every possible value of β_0 , the expected value of $(X'_{n+1}\beta_0 - \tilde{X}'_{n+1}\hat{\beta})^2$ is minimized at the truth $\tilde{X}_{n+1} = X_{n+1}$, where the expectation is taken with respect to the statistician’s sample and the possible realizations of the user’s attributes.

Intuition suggests that an individual cannot benefit from lying to a procedure that is meant to predict the

¹In a recent interview of Brian Christian, the author of *The Alignment Problem*, he notes that “computers may one day be able not only to learn our behavior but also intuit our values - figure out from our actions what it is we’re trying to optimize. ... What if an algorithm intuits the ‘wrong’ values, based on its best read of who we currently are but not of who we aspire to be? Do we really want our computers inferring our values from browser histories? See Shaywitz (2020) for this interview.

best outcome for her. To counter this intuition, Eliaz and Spiegler (2019, 2020) use the above framework to illustrate that a user may have a strict incentive to lie about her attributes when the prediction is based on a linear regression that penalizes non-zero estimated coefficients. However, these papers focus on particular examples in which attributes are *binary*, the statistician has the *same* (fixed) finite number of observations on each possible combination of attribute values, and the penalty parameter is *fixed* and does *not* adjust to the sample size. Hence, these papers leave open the following important question: For a general environment, are there conditions ensuring that a penalized regression model is incentive compatible in large samples?

Answering this question can potentially allow platforms, like those discussed above, to use machine-learning methods to predict users’ most preferred options without worrying that their data is “contaminated” by non-truthful users. Put bluntly, estimates and predictions made by methods that are *not* incentive-compatible are possibly unreliable since they may be based on false data.

This paper addresses the above open question by focusing on the most popular form of penalized regressions - the *Lasso* estimator.² Borrowing tools from high-dimensional statistics, we establish sufficient conditions for incentive compatibility of the Lasso estimator in large samples. We show that to achieve incentive compatibility, the tuning parameter must be *large* enough (i.e., it must remain above some threshold as sample size increases) so as to avoid overfitting, which is the main reason why a user may want to lie. This potential to lie implies that the standard way of choosing small enough tuning parameters to ensure consistency may lead to incentive compatibility violations. We provide simulation results that illustrate how the tuning parameter can be chosen in practice to ensure incentive compatibility. Incentive compatibility may therefore be viewed as an additional important property that should be imposed on estimators on top of consistency and unbiasedness. We also offer a new technical contribution by extending the oracle inequalities of Jankova and van de Geer (2018) to non-sub Gaussian data and providing a different proof.³

The motivation to focus on the Lasso estimator stems from the fact that this estimator is the benchmark among all high dimensional statistical estimators that predict large scale models when the number of regressors exceeds the sample size. Following its original proposal by Tibshirani (1996), econometricians and statisticians have used Lasso-based estimators to push the boundaries of economics and finance. One of the most critical issues facing these Lasso type estimators is post-inference after estimation and model selection, which require uniformly valid confidence intervals. In a seminal series of papers, Belloni et al. (2012) and Belloni et al. (2014) solved these issues by introducing the idea of “partialling out” the regressors. A different but complementary approach, via debiasing-desparsifying is proposed by van de Geer et al. (2014). Caner and Kock (2018) extended the debiasing of van de Geer et al. (2014) to heteroskedastic-non-sub-Gaussian data with strong oracle optimality property, thereby proposing a high dimensional estimator that is robust to heteroskedasticity, and with uniformly valid confidence intervals. Lasso-based debiasing are

²Our results can be extended to apply to the debiased lasso estimator, but this involves a different proof technique, and hence, is beyond the scope of the current paper.

³Even though we prove the main theorems with the bounded signal to noise ratio as in Jankova and van de Geer (2018), we can relax this ratio constraint as shown in the Appendix.

used in panel data models (see, e.g., Chernozhukov et al. (2018), Kock (2016), Kock and Tang (2019)) and for addressing quantile treatment effects and text analysis (see, e.g., Chiang and Sasaki (2019) and Chiang (2020)).

The concern that statistical procedures such as estimation, forecasting and classification are vulnerable to manipulation, has been the subject of some recent papers in the computer science literature. In contrast to us, this literature assumes there is an explicit conflict of interest between the statistician and the data providers - either because the latter are concerned about their privacy, they have to incur a cost to provide a precise report, or they have a different objective than the statistician. These papers analyze the Nash equilibria of a game where users submit private values that are used for estimation/classification, and propose incentive schemes that induce truthful reporting. Some notable works in this literature include Cai et al. (2015), Cummings et al. (2015), Dekel et al. (2010), Gao et al. (2015), Hardt et al. (2016), Meir et al. (2012) and Perte and Perote-Pena (2004). *None* of these papers consider penalized regression methods, and *none* of them characterize conditions guaranteeing incentive compatibility of regression techniques when the statistician and users have *aligned interests* (as is the case in our model).

The remainder of the paper is organized as follows. Section 2 considers the model and assumptions. Section 3 provides new oracle inequalities. Section 4 shows under what conditions lasso is incentive compatible, and Section 5 provides a simulation. Appendix A and B provide proofs of the results when $p > n$, and $p \leq n$, respectively.

2 The model

We will use the following notational conventions. For any vector $\nu \in R^d$, let $\|\nu\|_1, \|\nu\|_2, \|\nu\|_\infty$ denote its l_1, l_2, l_∞ norm respectively, and $\|\nu\|_0$ be the l_0 norm, which means the total number of nonzero entries. For a set $S \subseteq \{1, 2, \dots, d\}$, let $|S| = s$ be the cardinality of the set. Let ν_S be the modified ν such that we put 0 when the index does not belong to S (i.e., say $S = \{1, 2, 6\}$ for a 10×1 vector ν , this means that ν is modified such that now all elements are zero except elements 1, 2, 6). Let $\|A\|_{l_1}$ be the maximum absolute column-sum norm of a matrix of dimensions $m \times l$, i.e., $\|A\|_{l_1} = \max_{1 \leq k \leq l} \sum_{i=1}^m |A_{ik}|$ which is also called induced l_1 norm of A .

Our environment consists of users who are characterized by a set of p personal characteristics. For instance, in the context of medical decision making, a characteristic can represent a risk factor (obesity, smoking, etc.). For each user i , these characteristics are modeled as p explanatory variables, $X_{i,1}, \dots, X_{i,p}$, drawn from some distribution over a subset of \mathbb{R}^p . These attributes determine the ideal option for a user according to the function

$$f(X_{i,1}, \dots, X_{i,p}) = \sum_{k=1}^p X_{i,k} \beta_{0,k}$$

This function applies to all users, who differ only in the values of their characteristics. The realized values of $(X_{i,1}, \dots, X_{i,p})$ are privately observed by user i . A user may or may not know the value of the coefficients

$(\beta_{0,1}, \dots, \beta_{0,p})$. In the latter case, she has some (possibly degenerate) prior beliefs over their values.

A *statistician* (representing the automated prediction systems described in the introduction) has *private* access to a sample of n observations. Each observation $i = 1, \dots, n$ consists of the true attributes $X_i = (X_{i,1}, \dots, X_{i,p})$ of user i and a noisy signal y_i of that user’s ideal option,

$$y_i = X_i' \beta_0 + u_i, \quad (1)$$

where u_i is random noise that is drawn *i.i.d* from some distribution with zero mean.⁴

The X_i ’s are also *i.i.d.* across i , and exogenous, and will be discussed in detail in Assumption 1 in the next subsection. β_0 is a $p \times 1$ vector, representing the true parameters in f . We let $S_0 = \{j : \beta_{0,j} \neq 0\}$ denote the set of relevant regressors with s_0 being the cardinality of the set S_0 . (i.e., s_0 of the elements of β_0 are nonzero, and the rest are zero). s_0 is a nondecreasing function of n . These facts are known to an “oracle” but not to the statistician (and possibly not to a user).

Using her (privately observed) sample, the statistician estimate the function f , or equivalently, she estimates the coefficients $\beta_{0,1}, \dots, \beta_{0,p}$. When $p > n$, the least squares estimator is infeasible due to singularity of the empirical Gram matrix. Hence, the statistician uses Lasso, the penalized regression procedure that assigns costs to including explanatory variables in the regression. Specifically, she solves the following minimization problem

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{\sum_{i=1}^n (y_i - X_i' \beta)^2}{n} + 2\lambda_n \|\beta\|_1, \quad (2)$$

where $\lambda_n > 0$ is the penalty (also called tuning parameter) that decreases with the number of observations at the rate of $\lambda_n = O(\sqrt{\ln p/n})$ (an explicit expression for the sequence λ_n is given in equation (A.14) in Appendix A).⁵

Given her estimates $\hat{\beta}$, the statistician must take an action $a \in \mathbb{R}$ on behalf of a *new* user, $j = n+1$. This action is just the statistician’s prediction of the ideal option of that user. The new user’s payoff from action a is $-(a - f(X_{n+1}))^2$, where $f(X_{n+1})$ is the true ideal option associated with her personal attributes X_{n+1} . The distribution of X_{n+1} may be different from that of (X_1, \dots, X_n) , and we do not impose any restriction on its correlation with the sample distribution.

Since the statistician does not observe X_{n+1} , in order to make her prediction of $f(X_{n+1})$, she asks the $n+1$ user to report a $p \times 1$ vector, \tilde{X}_{n+1} , which is interpreted as that user’s attributes. The statistician then plugs \tilde{X}_{n+1} into her estimated model and chooses the action $a = \tilde{X}_{n+1}' \hat{\beta}$. When the $n+1$ user decides what attribute values to report, she takes into account that she does not observe the statistician’s sample, and hence, does not know the values of the estimated coefficients $\hat{\beta}$. She only knows the distribution from which the statistician’s sample is drawn, and that given her sample, the statistician chooses $\hat{\beta}$ according to (2). Given this, the user chooses the report \tilde{X}_{n+1} that minimize her expected loss $E_{\beta_0, \hat{\beta}} (\tilde{X}_{n+1}' \hat{\beta} - X_{n+1}' \beta_0)^2$,

⁴Access to such observations is a necessary condition for any platform that tries to learn about users (say, Netflix, Spotify). In the introduction, we gave a couple of examples for such data, which may be obtained from a third party, or from marketing surveys.

⁵We established this rate in Lemma A.2 in Appendix A.

where the expectation is taken with respect to the user’s prior beliefs about the true parameters β_0 , and his beliefs about the estimate $\hat{\beta}$. Hence, the new user may decide to lie and report $\tilde{X}_{n+1} \neq X_{n+1}$. In particular, she may decide to “opt out” and submit a vector of zeros.⁶ Our objective is to understand under what conditions it is in the user’s best interest to be truthful regardless of her prior beliefs on β_0 .

2.1 Incentive compatibility

We say that an estimator is *ex-ante incentive compatible* if for any belief over the true model parameters, the user’s expected payoff from truth-telling is at least as high as her expected payoff from any misreport, where the expectation is taken with respect to the user’s realized covariates, and with respect to the statistician’s sample.

Definition 1. *An estimator is (ex-ante) **incentive-compatible** if for every \tilde{X}'_{n+1} and every β_0 ,*

$$E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 \geq E[X'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2. \quad (3)$$

where the expectation E is taken with respect to the possible realizations of X_{n+1} and the possible realizations of the statistician’s sample.

An alternative notion of incentive-compatibility would be defined ex-post with respect to the realization of the user’s covariates, such that inequality (3) would be required to hold for every realization of X_{n+1} , and the expectation operator would only be with respect to the statistician’s sample. The sufficient condition for ex-ante incentive-compatibility of the Lasso estimator, which we establish in Section 4, also guarantees ex-post incentive-compatibility. Furthermore, the proof of ex-post incentive compatibility follows from our proof of ex-ante incentive-compatibility. In light of this, we shall focus on the ex-ante notion henceforth.

Incentive compatibility means that the user is unable to perform better by misreporting her personal characteristics, *regardless* of her beliefs over the true model’s parameters in mean squared sense. How should we interpret this requirement, given that we do not necessarily want to think of the user as being sophisticated enough to think in these terms? One interpretation is that lack of incentive compatibility is merely a *normative* statement about the user’s welfare - namely, given our model of how the statistician takes actions on the user’s behalf, it would be advisable for her to misrepresent her personal characteristics. Furthermore, there are opportunities for new firms to enter and offer the user paid advice for how to manipulate the procedure - in analogy to the industry of “search engine optimization”. Incentive compatibility theoretically eliminates the need for such an industry. In the context of the online content provision story, some misreporting strategies take the form of “deleting cookies”. This deviation is straightforward to implement, and the user can check if it makes her better off in the long run.

Note that incentive-compatibility is not a property that can be tested statistically. To see this, suppose each user is characterized by only a single covariate that is uniformly distributed on $\{0, 1\}$. If users are

⁶In the case in which the individual’s attributes are collected “passively” from her browsing history, then reporting a vector of zero attributes can be interpreted as the act of deleting cookies.

truthful, then one would expect a 50-50 distribution of 0's and 1's in the population. However, if each user lies about his covariate, then one would also observe a 50-50 distribution of 0's and 1's.

Recall that the statistician's sample contains the *true* attributes of n users. The idea is that the data on these users is obtained through a different process than the way the statistician obtains the data from the $n + 1$ user. For instance, as mentioned earlier, this data may be obtained from a marketing survey where there is no incentive to lie. Alternatively, one may interpret our incentive compatibility requirement as a requirement that truth-telling is a *Nash equilibrium* among all participants - such that given that everyone else is telling the truth, no user has an incentive to lie.

2.2 The statistician's data

In this subsection, we introduce a number of restrictions on the statistician's data. To describe these restrictions, we shall make use of the following notation. Define an l_0 ball $\mathcal{B}_{l_0}(s_0) = \{\|\beta_0\|_{l_0} \leq s_0\}$. Denote $\Sigma := EX_i X_i'$ for $i = 1, 2, \dots, n$ and let $\hat{\Sigma} := X'X/n$ be the sample counterpart. Our first requirement extends the sub-Gaussian data assumption used in statistics:

Assumption 1. (i). $E(u_i|X_i) = 0$, X_i, u_i are identical and independent across $i = 1, \dots, n$, and $\max_{1 \leq j \leq p} E|X_{ij}|^4, E|u_i|^l, l = \max(2k, 4)$ for all $k \geq 1$ are uniformly bounded from above (across n). (ii). The minimal eigenvalue of Σ is bounded away from zero uniformly in n .

Our second set of restrictions applies to the first and second moments. These will guarantee the consistency of the Lasso estimator, but will not ensure incentive compatibility (sufficient conditions for incentive compatibility will be introduced in Section 4). We start by defining the maximal value of certain cross products, which will be related to the behavior of moments in high dimensions in our next assumption.

$$M_1 := \max_{1 \leq i \leq n} \max_{1 \leq j \leq p} |X_{ij} u_i|,$$

$$M_2 := \max_{1 \leq i \leq n} \max_{1 \leq j \leq p} \max_{1 \leq l \leq p} |X_{il} X_{ij} - EX_{il} X_{ij}|.$$

Note that M_1 is the maximal covariance between the regressors and errors in a high dimensional context. Roughly speaking, when this covariance is small, it captures exogeneity of the regressors in the sample. M_2 is the maximal variance of the regressors in the sample. With large p and n , these covariance and variance terms can grow arbitrarily large; hence, we need a condition that restricts the growth rate of their moments. Because we are allowing for heteroskedastic data and unbounded regressors, we need to consider the growth rate of *higher-order* moments. ⁷

Assumption 2. (i).

$$\frac{\sqrt{\ln p}}{\sqrt{n}} [\max((EM_1^2)^{1/2}, (EM_2^2)^{1/2})] \rightarrow 0.$$

⁷Alternatively, we could strengthen Assumption 2 using boundedness of individual moments of X, u .

- (ii). $s_0(\frac{\ln p}{n})^{1/2} \rightarrow 0$.
- (iii). $\|\beta_0\|_2 = O(1)$.

Assumption 2(i) and 2(ii) are standard in high dimensional econometrics. 2(i) is used in Chernozhukov et al. (2017) allowing them to apply a concentration inequality, and 2(ii) is a standard sparsity condition. Assumption 2(iii) ensures that the signal to noise ratio is bounded (see p.2343 of Jankova and van de Geer (2018)). To see this clearly, set $\sigma_u^2 := \text{var}(u_i)$, which is the variance of the errors, and $\sigma_u^2 \geq c > 0$, where c is a generic positive constant. Hence,

$$\frac{\text{var}(y_i)}{\text{var}(u_i)} = \frac{\beta_0' \Sigma \beta_0}{\sigma_u^2} + 1,$$

under $E(u_i|X_i) = 0$ in Assumption 1 and $\Sigma := EX_i X_i'$. But

$$\frac{\beta_0' \Sigma \beta_0}{\sigma_u^2} + 1 \geq \frac{\|\beta_0\|_2^2 \phi_{\min}(\Sigma)}{\sigma_u^2} + 1.$$

where $\phi_{\min}(\Sigma) \geq c > 0$ is the minimum eigenvalue of Σ and is positive by Assumption 1. Hence, if Assumption 2(iii) holds, then the signal to noise ratio satisfies $\text{var}(y_i)/\text{var}(u_i) \geq C_0 + 1 > 0$, with C_0 being a positive constant, and defined as $C_0 := \frac{\|\beta_0\|_2^2 \phi_{\min}(\Sigma)}{\sigma_u^2}$. The empirical implication of this is that only a fixed number of nonzero coefficients can be constants, and the other nonzero coefficients have to be local to zero. To see this implication, note that

$$\|\beta_0\|_2 = \sqrt{\sum_{j=1}^p \beta_{0,j}^2} = \sqrt{\sum_{j \in S_0} \beta_{0,j}^2} = O(1).$$

But this last point can be achieved, in the case of s_0 growing with n , with

$$\sqrt{\sum_{j \in S_0} \beta_{0,j}^2} = \sqrt{\sum_{j \in F_1} \beta_{0,j}^2 + \sum_{j \in S_0 - F_1} \beta_{0,j}^2} \leq \sqrt{f_1 C^2 + (s_0 - f_1) \frac{C^2}{s_0 - f_1}} = O(1),$$

where $F_1 := \{j : |\beta_{0,j}| = C\}$ with $|F_1| = f_1$ being a fixed number, C is a generic positive constant and $F_2 := \{j : |\beta_{0,j}| = \frac{C}{\sqrt{s_0 - f_1}}\}$ with $|F_2| = s_0 - f_1$. For ease of exposition, we set all coefficients in F_1 and F_2 to be the same constants $C, C/\sqrt{s_0 - f_1}$ respectively. F_2 contains indices of all local to zero coefficients. This can easily be generalized without affecting our results.

In Appendix B we take a more flexible approach compared with Assumption 2(iii). There, we assume that $\|\beta_0\|_2 = O(\sqrt{s_0})$. In this case, all nonzero coefficients can be large (i.e., none of them are local to zero, as in set F_2 above). In other words, there is no index set F_2 as above, but all nonzero coefficients (their indices) are in the set F_1 above.

As p and n grow large, the total number of nonzero coefficients s_0 (also known as the *sparsity index*) can grow arbitrarily large. To guarantee consistency and unbiasedness, it is typically assumed that the product of the sparsity index and the tuning parameter should go to zero. However, this standard condition does not guarantee the incentive compatibility of the Lasso estimator as can be seen in the proof of Theorem 3 below.

3 New oracle inequalities

Oracle inequalities in high dimensional statistics are upper bounds on prediction and estimation errors. We require moment bounds on the Lasso estimator's error in l_1 norm for our main result. By taking the sample size to be large, we can show that the upper bound on the mean of higher-order moments of Lasso estimation errors tend to zero. We then use this asymptotic result to establish the incentive compatibility of the Lasso estimator in large samples. To illustrate this, we note that from the proof of Theorem 3 in Appendix A.2.4, the incentive compatibility constraint is tied to the following expression

$$E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 - E[X'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 = E[\hat{\beta}'(\tilde{X}_{n+1} - X_{n+1})(\tilde{X}_{n+1} - X_{n+1})'\hat{\beta}] \quad (4)$$

$$+ E[\hat{\beta}'(\tilde{X}_{n+1} - X_{n+1})X'_{n+1}(\hat{\beta} - \beta_0)] \quad (5)$$

$$+ E[(\hat{\beta} - \beta_0)'X_{n+1}(\tilde{X}'_{n+1} - X'_{n+1})\hat{\beta}]. \quad (6)$$

For incentive compatibility to hold in large samples, we need the sum of the right-hand side terms to be greater than or equal to zero. The first term on the right side (4) is always non-negative. Hence, if we prove that (5) and (6) converge to zero, we establish asymptotic incentive compatibility. However, the size of terms in (5) and (6) will depend on the mean of higher-order estimation error of Lasso.

To bound these error, we prove new oracle inequalities, which are different from those that are given in the literature for $\|\hat{\beta} - \beta_0\|_1$. These inequalities will serve an important role in proving our main result in the next section (Theorem 3). Besides, they are of independent interest as they extend previous results on sub-Gaussian data to *heteroskedastic* (conditionally) data sets that are commonly used in econometrics. Our proof technique will also look at a less conservative bound compared with Jankova and van de Geer (2018). Hence, our new inequalities contribute to the literature on high-dimensional econometrics where they can be used for proving generalized semiparametric efficiency of Lasso-type-estimators (as, e.g., in Jankova and van de Geer (2018)).

Our first result in this section is a k -th moment bound for the l_1 norm of the Lasso bias. A key concept used in this result is the *exception probability* for the event $\mathcal{F} := \{\mathcal{A}_1 \cap \mathcal{A}_2\}$, where \mathcal{A}_1 and \mathcal{A}_2 are defined in (A.6) and (A.9), which represent the empirical process-noise, and the eigenvalue condition, respectively. The exception probability is the complement of the event \mathcal{F} , and is denoted by $P(\mathcal{F}^c)$. An explicit upper bound for the exception probability is calculated in Lemma A.4.

Theorem 1. *Under Assumptions 1-2, if n is sufficiently large and $\lambda_n \geq \frac{P(\mathcal{F}^c)^{1/4k}}{s_0^{1/2}}$, then*

$$[E\|\hat{\beta} - \beta_0\|_1^k]^{1/k} = O(s_0\lambda_n).$$

This result is valid uniformly over $\mathcal{B}_{l_0}(s_0) = \{\|\beta_0\|_{l_0} \leq s_0\}$.

If we set $k = 1$ we can learn whether the Lasso estimator is unbiased. By the above Theorem, Assumption

2 and (A.15) imply $s_0\lambda_n \rightarrow 0$. Hence, in large samples, we have unbiasedness in the large λ_n case. Next, we provide the k -th moment bound for l_1 norm for the Lasso estimator.

Theorem 2. *Under Assumptions 1-2, if n is sufficiently large n and $\lambda_n \geq P(\mathcal{F}^c)^{1/2k}/s_0^{1/2}$, then*

$$[E\|\hat{\beta}\|_1^k]^{1/k} = O(s_0^{1/2}).$$

This result is valid uniformly over $\mathcal{B}_{l_0}(s_0) = \{\|\beta_0\|_{l_0} \leq s_0\}$.

This is a new result and a simple extension of Theorem 1 above. The rate in Theorem diverges to infinity if $s_0 \rightarrow \infty$ as $n \rightarrow \infty$.

4 Incentive Compatibility of Lasso

Our main result, which is new in the literature on penalized regressions, establishes that the Lasso estimator is incentive compatible for a sufficiently large sample size. In other words, we show that when $n \rightarrow \infty$

$$E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 \geq E[X'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2.$$

for all \tilde{X}'_{n+1} and for every β_0 , where the expectation is taken with respect to the reporting user's attributes X'_{n+1} (this is our ex-ante notion of incentive-compatibility that we explained in Section 2.1) and with respect to the statistician's realized sample (since the reporting user does not observe this sample).

The next theorem is our main result, which provides sufficient conditions for incentive compatibility. Its proof makes use of the following notation.

$$M_3 := \max_{1 \leq j \leq p} |X_{n+1,j}|,$$

$$M_4 := \max_{1 \leq j \leq p} |\tilde{X}_{n+1,j} - X_{n+1,j}|.$$

Note that M_4 is nothing more than the absolute magnitude of the misreport on a given variable j by the $n+1$ user.

Theorem 3. *Under Assumptions 1 and 2, the Lasso estimator is incentive compatible in large samples ($n \rightarrow \infty$) if the following conditions hold:*

$$\lambda_n \geq P(\mathcal{F}^c)^{1/16}/s_0^{1/2} \tag{7}$$

and

$$s_0^{3/2} \sqrt{\frac{\ln p}{n}} [EM_3^4]^{1/4} [EM_4^4]^{1/4} \rightarrow 0. \tag{8}$$

Furthermore, incentive compatibility is valid uniformly over $\mathcal{B}_{l_0}(s_0) = \{\|\beta_0\|_{l_0} \leq s_0\}$.

Remarks.

1. Theorem 3 establishes that a *sufficient* condition for incentive compatibility is that the tuning parameter λ_n needs to be large “enough”. A simple way to choose λ_n to satisfy (7) is to use the upper bound of the exception probability

$$\lambda_n := \text{upperbound}(P(\mathcal{F}^c)^{1/16}),$$

in Lemma A.4. The simulations in the next section address the issue of whether such a bound is feasible.

2. The typical concern with Lasso is the consistency of the estimator ($\|\hat{\beta} - \beta_0\|_1 = o_p(1)$), which can be achieved by making sure that λ_n goes to zero at a relatively fast rate (as Lemma A.1 in Appendix A shows, this rate is $s_0\lambda_n \rightarrow 0$). However, if λ_n gets too small, the Lasso estimator may admit many nonzero variables incorrectly (i.e., it creates an *overfit*). Consequently, when the number of regressors p is very large, the expectation of the sum of l_1 errors ($E\|\hat{\beta} - \beta_0\|_1$) can grow arbitrarily large, and incentive compatibility may be violated. Put differently, *consistency does not imply incentive compatibility in large samples*.

To see this formally, let $\{A\}$ denote the event A . Then for all $\epsilon > 0$,

$$\begin{aligned} E\|\hat{\beta} - \beta_0\|_1 &= P(\|\hat{\beta} - \beta_0\|_1 \leq \epsilon)\{\|\hat{\beta} - \beta_0\|_1 \leq \epsilon\} + P(\|\hat{\beta} - \beta_0\|_1 > \epsilon)\{\|\hat{\beta} - \beta_0\|_1 > \epsilon\} \\ &\leq P(\|\hat{\beta} - \beta_0\|_1 \leq \epsilon)\epsilon + P(\|\hat{\beta} - \beta_0\|_1 > \epsilon)\{\|\hat{\beta} - \beta_0\|_1 > \epsilon\}, \end{aligned} \tag{9}$$

By consistency, the first term in (9) will go to zero when ϵ approaches zero. However, the second term may be large and can dominate all the expectation in large dimensions (this is also discussed on p.2339 of Jankova and van de Geer (2018)) In other words, just using the l_1 estimator bound on its own does *not* imply a bound for the *expectation* of l_1 error. It requires a non-trivial proof. Consistency does not imply unbiasedness, and hence, it does not imply incentive compatibility.

Why is overfitting a significant issue for incentive compatibility? The intuition is as follows. Suppose the tuning parameter is sufficiently small so that given the user’s prior on the true coefficients, she expects that many irrelevant variables will be included in the estimator. To correct this bias, she can report that these variables are equal to zero.

3. The second sufficient condition (8) allows the distance between the user’s report \tilde{X}_{n+1} and the truth X_{n+1} to be of any magnitude since $EM_4 \equiv E\|\tilde{X}_{n+1} - X_{n+1}\|_\infty$ can be arbitrarily large. Since the above conditions are sufficient but not necessary, it remains an open question whether incentive compatibility can be achieved with a tuning parameter that is lower than the threshold in (7) without restricting the magnitude of the deviation between the user’s reported and true attributes.

4. Note that (8) requires stricter sparsity than Assumption 2. If $EM_3^4 = O(1)$ and $EM_4^4 = O(1)$, then condition (8) amounts to $s_0^{3/2} \sqrt{\frac{\ln p}{n}} \rightarrow 0$, which is a sparsity requirement still stronger than Assumption 2(ii). In addition, if we let $EM_4^4 = O(\ln n)$ and $EM_3^4 = O(\ln n)$, then $s_0^{3/2} \sqrt{\frac{\ln p}{n}} (\ln n)^{1/2} \leq s_0^{3/2} \sqrt{\frac{2 \ln p}{n}} \rightarrow 0$ is needed to get incentive compatibility with $n \leq p$.

5. A natural question that arises is whether condition (7) is compatible with the l_1 norm consistency of Lasso? In other words, consistency requires a small λ_n , but incentive compatibility requires a large λ_n , so

are they compatible with each other? When we select a large λ_n to satisfy incentive compatibility, we should not sacrifice consistency, i.e. we need $s_0\lambda_n \rightarrow 0$. To verify whether this is possible, we can take the lower bound on the tuning parameter in (7) and see whether we can achieve consistency. Note that

$$s_0\lambda_n = s_0 \frac{P(\mathcal{F}^c)^{1/16}}{s_0^{1/2}} = s_0^{1/2} P(\mathcal{F}^c)^{1/16}, \quad (10)$$

From (A.22) in the Appendix, an upper bound on this exception probability is:

$$P(\mathcal{F}^c) \leq \frac{2}{p^{C_1}} + \frac{K[EM_1^2 + EM_2^2]}{nl np}, \quad (11)$$

where C_1 and K are positive constants. With $l = 1, 2$, it therefore follows from (10) and (11) that we need

$$s_0^8/p^{C_1} \rightarrow 0, \quad s_0^8 \max_l EM_l^2/nl np \rightarrow 0,$$

to have consistency. These two conditions are not unreasonable in the sense that they are consistent with (n, p) increasing to infinity. Also they are compatible with moments satisfying condition (8) in Theorem 3.

6. Finally, note that $\lambda_n = O(\sqrt{\frac{\ln p}{n}})$ represents an upper bound in terms of rates for λ_n , whereas (7) represents a lower bound. We can then take for a positive constant $C > 0$

$$C \frac{\sqrt{\ln p}}{\sqrt{n}} \geq \lambda_n \geq \frac{P(\mathcal{F}^c)^{1/16}}{s_0^{1/2}}.$$

The question is, are there suitable combinations of n and p that satisfy these inequalities? By using algebra and the upper bound for exception probability (A.22), we obtain the requirement that,

$$C s_0^{1/2} \geq \left[\frac{2n}{p^{C_1}} + \frac{K[EM_1^2 + EM_2^2]}{nl np} \right]^{1/16} \frac{\sqrt{n}}{\sqrt{\ln p}},$$

which is plausible for $p > n$ and large n since the left hand side may diverge and the right side may go to zero. This may be the case for example when p is exponential in n .

5 Simulations

This section has two objectives. First, it illustrates how in practice the tuning parameter can be chosen to ensure incentive compatibility of the Lasso estimator. Second, it demonstrates that by appropriately choosing the tuning parameter (in line with the conditions in Theorem 3), incentive compatibility is satisfied regardless of the magnitude of the ‘‘lie’’ (i.e., the distance between the true and reported attributes)

We provide a simple simulation setup. We model

$$y_i = X_i' \beta_0 + u_i,$$

where $\beta_0 = (1, 0'_{p-s_0}, 1'_{s_0-1})'$, 0_{p-s_0} is a $p - s_0$ column vector of all zero elements, and 1_{s_0-1} is a $s_0 - 1$ dimensional column vector of all ones. Let s_0 represent the sparsity of the above model and set $s_0 = 20$.

We present three different simulations. In Design 1, we choose X_i to be a $p \times 1$ vector of a t distribution with five degrees of freedom. The new $n + 1$ user has the same distribution for her attributes but is independent of the first n users. The errors u_i are also chosen from a t_5 distribution but independently of the regressors. Tables 1-3 displays these results. For the second simulation (Design 2), we only change the distribution of the attributes for the $n + 1$ user to a t distribution with three degrees of freedom. In Design 2, we keep the same distribution for the errors and the same attribute distribution for the first n users from Design 1. The results are displayed in Tables 4-6. In the third simulation (Design 3), we change only the following in Design 2. We introduce a multivariate normal distribution for the attributes of users $i = 1, \dots, n$, such that the covariance between the j and m -th random variables are governed by

$$\Sigma_{j,m} = 0.5^{|j-m|},$$

for $j = 1, \dots, p$ and $m = 1, \dots, p$. Thus, the correlation between the adjacent random variables is 0.5, and this declines when the random variables are further apart. This Toeplitz type structure is commonly used in the high dimensional literature (see Caner and Kock (2018)). In Design 3, we keep the distribution of the $n + 1$ user and the errors from Design 2. The results are presented in Tables 6-9.

We aim to demonstrate that with a “large” tuning parameter as in Theorem 3, incentive compatibility can be achieved when the sample size n is large enough. As mentioned in the previous section, one possible choice of a tuning parameter that satisfies Theorem 3 is the upper bound on the exception probability,

$$\lambda_n \geq \text{upperbound}(P(\mathcal{F}^c)^{1/16}).$$

The issue is to make the exception probability, $P(\mathcal{F}^c)$ operational and usable. Note that an upper bound on this probability is (with positive constants $C_1 > 0, C_2 > 0, K > 0$)

$$P(\mathcal{F}^c) \leq \frac{2}{p^{C_1}} + \frac{K[EM_1^2 + EM_2^2]}{nl np} \leq \frac{2}{p^{C_1}} + \frac{C_2}{(lnp)^2}, \quad (12)$$

by observing that for $l = 1, 2$

$$\begin{aligned} \frac{K \max_l EM_l^2}{nl np} &= \left[\frac{K^{1/2} \sqrt{\max_l EM_l^2}}{\sqrt{n} \sqrt{lnp}} \right]^2 \\ &= \left[\frac{K^{1/2} \sqrt{\max_l EM_l^2} \sqrt{lnp}}{\sqrt{n}} \right]^2 \left(\frac{1}{lnp} \right)^2 \\ &\leq \frac{C_2}{(lnp)^2}, \end{aligned}$$

where we use Assumption 2(i). Hence, we can write the upper bound of the exception probability by using $p \geq 1$

$$\frac{2}{p^{C_1}} + \frac{C_2}{(lnp)^2} \leq 2 + \frac{C_2}{(lnp)^2}.$$

The tuning parameter is as follows

$$\lambda_n := \left[2 + \frac{C_2}{(\ln p)^2}\right]^{1/16}, \quad (13)$$

where C_2 can start from a small positive value and stop at a large positive value, and we select the optimal C_2 and λ_n according to the Generalized Information Criterion as in Caner and Kock (2018), which gives consistent model selection with weighted Lasso choices in the least squares framework. Therefore, our choice of λ_n is *above* a lower bound, which prevents overfitting (this is the novel insight of Theorem 3). On the other hand, to prevent a very large λ_n and ensure consistency of Lasso, the lower bound inversely depends on p .

Define

$$\lambda_n^* := \operatorname{argmin}_{\lambda_n \in \Lambda} [\ln(\hat{\sigma}^2(\lambda_n)) + \frac{\hat{s}(\lambda_n)}{n} \ln(n) \ln(\ln(p))],$$

where $\hat{s}(\lambda_n)$ is the number of nonzero elements in the Lasso estimator, given a choice of λ_n in a grid Λ , and $\hat{\sigma}^2(\lambda_n)$ is the mean squared residuals from the Lasso regression, given a choice of λ_n in a grid Λ . We form Λ as follows: We take C_2 in a grid of values $[2 + \frac{C_2}{(\ln p)^2}]$ as in (13). Let $C_2 := [0.1, 0.5, 1, 2, 10, 20, 50, 100]$, so Λ is the grid of values of λ_n depending on C_2 . The number of iterations is 1,000.

The ‘‘Report’’ column in Tables 1-3 display $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ as the mean squared error from a false report by the user. ‘‘Truth’’ refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$. The difference between $\tilde{X}_{n+1} - X_{n+1}$ is kept at three levels: 5, 2 and 0.2 (for all p variables), which represent large, medium, and small deviations from the truth. We have $p = 100, 250, 500$, and for each p level we analyze $n = 100, 200, 400$.

The numbers in each cell of the tables correspond to the disutility of the user (i.e., the mean square difference between the statistician’s estimate and the optimal action). Hence, smaller numbers correspond to higher payoffs. Let us compare the three tables when $p = 500$ and $n = 400$. In Table 1, which corresponds to a large magnitude of a lie, the user’s disutility from reporting the truth is 24.36, while the disutility from lying is 334.01. Hence, the $n + 1$ user prefers to be truthful. In Table 2, for a medium magnitude of lies, truth-telling induces a disutility of 25.03, while lying induces a higher disutility of 71.63. Finally, in Table 3, if the lie is ‘‘close’’ to the truth, the disutility from truth-telling is 24.36, while the disutility from lying is 24.73. Similar comparisons hold in the tables’ remaining cells, suggesting that that Lasso’s incentive compatibility is achieved. In Tables 4-9, the same message from Tables 1-3 carries over: Lasso is incentive-compatible with our tuning parameter choice. Tables 6 and 9 show that with a minor lie, Lasso is still incentive-compatible and the difference between the truth and lie in MSE sense is larger compared with Table 3 of Design 1.

6 Conclusion

The growing reliance on machine learning in automating decisions previously made by people raises the question of how people would interact with these automated systems. In particular, would people have an

Table 1: Design 1-Incentive Compatibility Scenarios: Difference 5

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	25.93	385.25	23.61	338.91	25.89	321.55
$p = 250$	27.75	353.90	25.25	353.08	26.41	337.91
$p = 500$	26.71	305.53	25.55	333.97	24.36	334.01

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 2: Design 1-Incentive Compatibility Scenarios: Difference 2

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	26.12	80.11	25.88	76.10	23.56	71.87
$p = 250$	25.88	75.64	24.25	77.89	25.20	72.40
$p = 500$	28.02	72.21	25.20	76.73	25.03	71.63

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 3: Design 1-Incentive Compatibility Scenarios: Difference 0.2

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	25.60	25.90	25.73	25.93	24.61	25.03
$p = 250$	26.06	26.87	23.90	24.27	25.43	25.84
$p = 500$	28.34	28.98	24.94	25.62	24.36	24.73

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 4: Design 2-Incentive Compatibility Scenarios: Difference 5

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	45.62	392.71	39.25	368.91	40.22	345.82
$p = 250$	47.11	374.02	45.55	355.38	46.59	342.67
$p = 500$	45.08	326.11	43.77	370.69	47.70	350.90

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 5: Design 2-Incentive Compatibility Scenarios: Difference 2

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	48.63	97.92	59.58	94.09	58.65	85.68
$p = 250$	46.19	89.30	44.12	98.91	43.25	95.22
$p = 500$	55.57	105.84	41.61	90.43	40.95	92.95

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 6: Design 2-Incentive Compatibility Scenarios: Difference 0.2

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	50.48	51.03	44.02	45.72	44.42	45.07
$p = 250$	45.98	47.33	48.88	49.47	42.56	43.68
$p = 500$	48.84	48.92	44.08	44.79	41.45	42.18

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 7: Design 3-Incentive Compatibility Scenarios: Difference 5

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	26.22	2766.29	18.84	3004.57	16.77	3136.86
$p = 250$	24.77	2725.52	21.41	3000.16	20.35	3130.54
$p = 500$	34.32	2722.20	19.68	2981.42	17.98	3139.44

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 8: Design 3-Incentive Compatibility Scenarios: Difference 2

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	27.43	464.41	19.82	493.65	17.16	515.18
$p = 250$	23.49	454.13	19.76	488.40	16.14	507.25
$p = 500$	25.25	448.99	38.48	503.34	14.56	509.49

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

Table 9: Design 3-Incentive Compatibility Scenarios: Difference 0.2

	$n = 100$		$n = 200$		$n = 400$	
Dimension	Truth	Report	Truth	Report	Truth	Report
$p = 100$	24.16	28.51	16.77	20.07	14.35	19.63
$p = 250$	25.80	29.93	19.41	25.04	16.14	22.63
$p = 500$	27.11	30.47	19.97	25.11	18.24	22.52

Note: "Truth" refers to $E[X'_{n+1}(\hat{\beta} - \beta_0)]^2$ and "Report" refers to $E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2$ in Incentive Compatibility Definition. Smaller values of these average squared errors are desirable.

incentive to act strategically in order to manipulate such automated systems? This strategic interaction will become particularly important when these automated systems start playing a more prominent role in medical decision-making or even in driving.

This paper takes only a small preliminary step towards addressing this question by studying whether a user would want to lie to an automated system that uses Lasso to predict that user's ideal outcome based on her reported attributes. Our main contribution is showing that truthful reporting can be ensured by appropriately adjusting the tuning parameter to be larger than what is required for consistency. Our result is also significant from a pure econometrics point of view: Just concentrating on oracle inequalities and post-selection inference can lead to a small tuning parameter, which in turn, can lead to model overfitting, which then introduces an incentive to misreport. If users have an incentive to provide false input to algorithms used for estimation and prediction, then it is no longer clear that one can rely on the output of these algorithms.

In the next part, Appendix A considers the proofs when $p > n$, and Appendix B considers the case $p \leq n$, and relaxing Assumption 2(iii).

A Appendix A

A.1 Notation

In this section, we show some results that will help us in proofs. Define random vector of variables $F_i := (F_{i1}, \dots, F_{ij}, \dots, F_{ip})'$. Also define $\sigma_F^2 := n(\max_{1 \leq j \leq p} \text{var} F_{ij})$, and $M_F := \max_{1 \leq i \leq n} \max_{1 \leq j \leq p} |F_{ij} - EF_{ij}|$. Note that $\hat{\mu}_j := n^{-1} \sum_{i=1}^n F_{ij}$, and $\mu_j := EF_{ij}$.

A.2 Maximal Inequalities

We use two assumptions that will provide us maximal inequalities.

Assumption A.1. Assume F_i are iid random vectors across $i = 1, 2, \dots, n$ with $\max_{1 \leq j \leq p} \text{var} F_{ij}$ bounded away from infinity uniformly in n .

Assumption A.2. Assume

$$\frac{\sqrt{EM_F^2} \sqrt{\ln p}}{\sqrt{n}} \rightarrow 0.$$

We use the following maximal inequality. With Assumption A.1, Lemma E.2(ii) of Chernozhukov et al. (2017) is: (see (A.2) of Caner and Kock (2019))

$$P \left[\max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| \geq 2E \max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| + \frac{t}{n} \right] \leq \exp(-t^2/3\sigma_F^2) + K \frac{EM_F^2}{t^2}, \quad (\text{A.1})$$

for a constant $K > 0$. With Assumptions A.1-A.2 here, Caner and Kock (2019) or Lemma E.1 of Chernozhukov et al. (2017) provides

$$\begin{aligned} E \max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| &\leq K \left[\frac{\sqrt{\ln p}}{\sqrt{n}} + \frac{\sqrt{EM_F^2 \ln p}}{n} \right] \\ &= O\left(\frac{\sqrt{\ln p}}{\sqrt{n}}\right). \end{aligned} \quad (\text{A.2})$$

Define the sequence $\kappa_n = \ln p$. Set $t = t_n = (n\kappa_n)^{1/2}$ to have (A.1) as

$$\begin{aligned} P \left[\max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| \geq 2E \max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| + \frac{\sqrt{\kappa_n}}{\sqrt{n}} \right] &\leq \exp(-C_1 \kappa_n) + K \frac{EM_F^2}{n\kappa_n} \\ &= \frac{1}{p^{C_1}} + \frac{KEM_F^2}{n \ln p} \end{aligned} \quad (\text{A.3})$$

where $C_1 > 0$, is a positive constant.

Now combine (A.2) with (A.3) to have

$$\begin{aligned}
P(\max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| \geq 2K[\frac{\sqrt{lnp}}{\sqrt{n}} + \frac{(EM_F^2)^{1/2}lnp}{n}] + \frac{\sqrt{lnp}}{\sqrt{n}}) \\
\leq \frac{1}{p^{C_1}} + \frac{KEM_F^2}{n(lnp)} = o(1),
\end{aligned} \tag{A.4}$$

by Assumptions A1-A.2. This shows also that, since EM_F^2 is nondecreasing in n

$$\max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| = O_p(\sqrt{lnp}/\sqrt{n}). \tag{A.5}$$

A.2.1 Events

Before the assumptions, we need to define events that will be helpful. The first event is:

$$\mathcal{A}_1 = \left\{ 2 \left\| \frac{u'X}{n} \right\|_{\infty} \leq \lambda_n \right\}, \tag{A.6}$$

which controls the noise. This is the maximal correlation between regressors and errors. We want this to be bounded with probability approaching one, and this upper bound, λ_n , itself is converging to zero in our proofs. We show that in Lemma A.2. So in large samples, this proof technique amounts to verification of exogeneity of regressors. This is standard in high dimensional econometrics, for a recent analysis see Lemma A.4 of Caner and Kock (2018).

We start with defining first population counterparts of restricted eigenvalue conditions and then show the empirical version also. These are standard in high dimensional econometrics and statistics and can be seen from Assumption 1 of Caner and Kock (2018).

We define the population adaptive restricted eigenvalue of Σ

$$\phi_{\Sigma}^2(s) = \min \left\{ \frac{\delta' \Sigma \delta}{\|\delta_S\|_2^2} : \delta \in R^p - \{0\}, \|\delta_{S^c}\|_1 \leq 3\sqrt{s}\|\delta_S\|_2, |S| \leq s \right\}. \tag{A.7}$$

Note that if $\Sigma = EX_iX_i'$ has full rank, the population adaptive restricted eigenvalue being positive is satisfied by Assumption 1. Also instead of minimizing all over R^p , we minimize vectors that satisfy $\|\delta_S^c\|_1 \leq 3\|\delta_S\|_1$. Even in the cases that Σ does not have full rank, it is possible that minimal adaptive restricted eigenvalue condition is satisfied due to optimization over a restricted set. The parameter δ will be related to structural parameter β in the proofs.

First define the empirical adaptive restricted eigenvalue condition, which is empirical counterpart of the population version in Assumption 1:

$$\hat{\phi}_{\Sigma}^2(s) = \min \left\{ \frac{\delta' \hat{\Sigma} \delta}{\|\delta_S\|_2^2} : \delta \in R^p - \{0\}, \|\delta_{S^c}\|_1 \leq 3\sqrt{s}\|\delta_S\|_2, |S| \leq s \right\}. \tag{A.8}$$

We are interested in behavior of the minimal empirical adaptive restricted eigenvalue condition evaluated for set S_0 at cardinality s_0 . The second event is:

$$\mathcal{A}_2 = \left\{ \hat{\phi}_{\Sigma}^2(s_0) \geq \phi_{\Sigma}^2(s_0)/2 \right\}. \tag{A.9}$$

Empirical adaptive restricted eigenvalue condition is needed since in case of $p > n$, $X'X$ is singular and the minimal eigenvalue of $X'X$ is zero. Empirical adaptive eigenvalue is over a restricted set which we prove to be positive, with probability approaching one, in Lemma A.3. This is also standard in high dimensional econometrics, see Lemma A.6 of Caner and Kock (2018). Set $\mathcal{F} = \mathcal{A}_1 \cap \mathcal{A}_2$, and the complement event as \mathcal{F}^c .

A.2.2 Proofs of Lemmata

The following four Lemmata are the intermediate results that are used for Theorems.

Lemma A.1. *Under the joint event $\mathcal{F} := \{\mathcal{A}_1 \cap \mathcal{A}_2\}$ we have*

$$\|\hat{\beta} - \beta_0\|_1 \leq \frac{24\lambda_n s_0}{\phi_{\Sigma}^2(s_0)}.$$

This is also valid uniformly over $\mathcal{B}_{l_0}(s_0) = \{\|\beta_0\|_{l_0} \leq s_0\}$.

Proof of Lemma A.1. Using $\hat{\beta}$ definition

$$\|Y - X\hat{\beta}\|_n^2 + 2\lambda_n \sum_{j=1}^p |\hat{\beta}_j| \leq \|Y - X\beta_0\|_n^2 + 2\lambda_n \sum_{j=1}^p |\beta_{0,j}|.$$

Use the model $Y = X\beta_0 + u$ on the first left side term as well as the first right side term to simplify the inequality above combining with Holder's Inequality

$$\begin{aligned} \|X(\hat{\beta} - \beta_0)\|_n^2 + 2\lambda_n \sum_{j=1}^p |\hat{\beta}_j| &\leq 2 \left| \frac{u'X}{n} (\hat{\beta} - \beta_0) \right| + 2\lambda_n \sum_{j=1}^p |\beta_{0,j}| \\ &\leq 2 \left\| \frac{u'X}{n} \right\|_{\infty} \|\hat{\beta} - \beta_0\|_1 + 2\lambda_n \sum_{j=1}^p |\beta_{0,j}| \end{aligned}$$

On the right side assuming we are on the event \mathcal{A}_1

$$2 \left\| \frac{u'X}{n} \right\|_{\infty} \|\hat{\beta} - \beta_0\|_1 \leq \lambda_n \|\hat{\beta} - \beta_0\|_1.$$

So we have

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + 2\lambda_n \sum_{j=1}^p |\hat{\beta}_j| \leq \lambda_n \|\hat{\beta} - \beta_0\|_1 + 2\lambda_n \sum_{j=1}^p |\beta_{0,j}|.$$

Use $\|\hat{\beta}\|_1 = \|\hat{\beta}_{S_0}\|_1 + \|\hat{\beta}_{S_0^c}\|_1$ on the second term for the left side of the inequality immediately above

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + 2\lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq \lambda_n \|\hat{\beta} - \beta_0\|_1 + 2\lambda_n \sum_{j=1}^p |\beta_{0,j}| - 2\lambda_n \sum_{j \in S_0} |\hat{\beta}_j|.$$

By assumption of sparsity $\sum_{j \in S_0^c} |\beta_{0,j}| = 0$, and using the reverse triangle inequality we have

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + 2\lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq \lambda_n \|\hat{\beta} - \beta_0\|_1 + 2\lambda_n \sum_{j \in S_0} |\hat{\beta}_j - \beta_{0,j}|.$$

Next by $\|\hat{\beta} - \beta_0\|_1 = \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_1 + \|\hat{\beta}_{S_0^c}\|_1$ for the first term on the right side of the inequality immediately above

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + \lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq 3\lambda_n \sum_{j \in S_0} |\hat{\beta}_j - \beta_{0,j}|.$$

Use $\|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_1 \leq \sqrt{s_0} \|\hat{\beta} - \beta_{0,S_0}\|_2$ above on the right side to have

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + \lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq 3\lambda_n \sqrt{s_0} \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_2. \quad (\text{A.10})$$

Ignoring the first term on the left of (A.10), (A.10) shows that we satisfy the restricted set condition in empirical adaptive restricted eigenvalue condition, so we have

$$\|\hat{\beta}_{S_0^c}\|_1 \leq 3\sqrt{s_0} \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_2.$$

Using $\delta = \hat{\beta} - \beta_0$ in the empirical adaptive restricted eigenvalue condition (A.8) in (A.10)

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + \lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq 3\lambda_n \sqrt{s_0} \frac{\|X'(\hat{\beta} - \beta_0)\|_n}{\hat{\phi}_\Sigma(s_0)}.$$

Then use $3uv \leq u^2/2 + 9v^2/2$ with $u = \lambda_n \sqrt{s_0}/\hat{\phi}_\Sigma(s_0)$, $v = \|X(\hat{\beta} - \beta_0)\|_n$ to get

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + \lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq \frac{\|X(\hat{\beta} - \beta_0)\|_n^2}{2} + \frac{9}{2} \frac{\lambda_n^2 s_0}{\hat{\phi}_\Sigma^2(s_0)}.$$

Simplify above

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + 2\lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq \frac{9\lambda_n^2 s_0}{\hat{\phi}_\Sigma^2(s_0)}.$$

Use the event \mathcal{A}_2 we get the following

$$\|X(\hat{\beta} - \beta_0)\|_n^2 + 2\lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| \leq \frac{18\lambda_n^2 s_0}{\hat{\phi}_\Sigma^2(s_0)}.$$

This implies the oracle inequality

$$\|X(\hat{\beta} - \beta_0)\|_n^2 \leq \frac{18\lambda_n^2 s_0}{\hat{\phi}_\Sigma^2(s_0)}. \quad (\text{A.11})$$

To get to the l_1 bound ignore the first term in (A.10) and add both sides $\lambda_n \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_1$ to have

$$\lambda_n \sum_{j \in S_0^c} |\hat{\beta}_j| + \lambda_n \sum_{j \in S_0} |\hat{\beta}_j - \beta_{0,j}| = \lambda_n \|\hat{\beta} - \beta_0\|_1 \leq \lambda_n \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_1 + 3\lambda_n \sqrt{s_0} \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_2,$$

by seeing also $\sum_{j \in S_0^c} |\beta_{0,j}| = 0$. Now use the norm inequality $\|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_1 \leq \sqrt{s_0} \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_2$ to have

$$\lambda_n \|\hat{\beta} - \beta_0\|_1 \leq 4\lambda_n \sqrt{s_0} \|\hat{\beta}_{S_0} - \beta_{0,S_0}\|_2.$$

Use the empirical adaptive restricted eigenvalue condition with $\delta = \hat{\beta} - \beta_0$

$$\|\hat{\beta} - \beta_0\|_1 \leq 4\sqrt{s_0} \frac{\|X(\hat{\beta} - \beta_0)\|_n}{\hat{\phi}_\Sigma(s_0)}.$$

Use (A.11) and the event \mathcal{A}_2 to have

$$\begin{aligned} \|\hat{\beta} - \beta_0\|_1 &\leq 4\sqrt{s_0} \left[\frac{3\sqrt{2}\lambda_n\sqrt{s_0}}{\phi_\Sigma(s_0)} \right] \left[\frac{1}{\hat{\phi}_\Sigma(s_0)} \right] \\ &\leq \frac{24\lambda_n s_0}{\phi_\Sigma^2(s_0)}. \end{aligned} \quad (\text{A.12})$$

Note that uniformity over $\mathcal{B}_{l_0}(s_0)$ follows since the upper bound in (A.12) depends on β_0 only through s_0 .

Q.E.D

Lemma A.2. (i). Under Assumption 1, and since $\kappa_n = \ln p$

$$P(\mathcal{A}_1) \geq 1 - \exp(-C_1\kappa_n) - \frac{KEM_1^2}{(n\kappa_n)} = 1 - \frac{1}{p^{C_1}} - \frac{KEM_1^2}{n\ln p}$$

(ii). Under added Assumption 2 to Assumption 1, $P(\mathcal{A}_1) \rightarrow 1$.

(iii). Under added Assumption 2 to Assumption 1, $\lambda_n = O(\sqrt{\ln p/n})$.

Proof of Lemma A.2. (i). Establish the probability bound on \mathcal{A}_1 via Assumption 1, using (A.3)(A.4) with $F_i = X_i u_i$ there and $\kappa_n = \ln p$, we have

$$P(\mathcal{A}_1) \geq 1 - \exp(-C_1\kappa_n) - K \frac{EM_1^2}{(n\kappa_n)} = 1 - \frac{1}{p^{C_1}} - \frac{KEM_1^2}{n\ln p}, \quad (\text{A.13})$$

with

$$\lambda_n = K \left[\sqrt{\frac{\ln p}{n}} + \frac{\sqrt{EM_1^2 \ln p}}{n} \right] + \sqrt{\frac{\ln p}{n}}. \quad (\text{A.14})$$

(ii). By Assumption 2, we have the proof.

(iii). By Assumption 2, we have

$$\lambda_n = O(\sqrt{\ln p/n}). \quad (\text{A.15})$$

Q.E.D.

Lemma A.3. Under Assumptions 1, 2, $\kappa_n = \ln p$

$$P(\mathcal{A}_2) \geq 1 - \exp(-C_1\kappa_n) - \frac{KEM_2^2}{(n\kappa_n)} = 1 - \frac{1}{p^{C_1}} - \frac{KEM_2^2}{n\ln p} = 1 - o(1).$$

Proof of Lemma A.3. Start with

$$\begin{aligned} \left| \delta' \frac{X'X}{n} \delta \right| &= \left| \delta' \left(\frac{X'X}{n} - \Sigma + \Sigma \right) \delta \right| \\ &\geq |\delta' \Sigma \delta| - |\delta' (\hat{\Sigma} - \Sigma) \delta|. \end{aligned} \quad (\text{A.16})$$

The second term on the right side of (A.16) can be bounded by repeated application of Holders inequality

$$|\delta' (\hat{\Sigma} - \Sigma) \delta| \leq \|\delta\|_1^2 \|\hat{\Sigma} - \Sigma\|_\infty.$$

So (A.16) becomes

$$|\delta' \hat{\Sigma} \delta| \geq |\delta' \Sigma \delta| - \|\delta\|_1^2 \|\hat{\Sigma} - \Sigma\|_\infty. \quad (\text{A.17})$$

Now we digress a bit to simplify (A.17). Note that we have the restriction set definition

$$\|\delta_{S_0^c}\|_1 \leq 3\sqrt{s_0}\|\delta_{S_0}\|_2,$$

where we add $\|\delta_{S_0}\|_1$ to both sides

$$\begin{aligned} \|\delta\|_1 &\leq 3\sqrt{s_0}\|\delta_{S_0}\|_2 + \|\delta_{S_0}\|_1 \\ &\leq 3\sqrt{s_0}\|\delta_{S_0}\|_2 + \sqrt{s_0}\|\delta_{S_0}\|_2 \\ &= 4\sqrt{s_0}\|\delta_{S_0}\|_2, \end{aligned}$$

where we used the norm inequality $\|\delta_{S_0}\|_1 \leq \sqrt{s_0}\|\delta_{S_0}\|_2$ in the second inequality above. So we get

$$\frac{\|\delta\|_1^2}{\|\delta_{S_0}\|_2^2} \leq 16s_0.$$

Now divide (A.17) by $\|\delta_{S_0}\|_2^2 > 0$ to have

$$\frac{|\delta' \hat{\Sigma} \delta|}{\|\delta_{S_0}\|_2^2} \geq \frac{|\delta' \Sigma \delta|}{\|\delta_{S_0}\|_2^2} - 16s_0 \|\hat{\Sigma} - \Sigma\|_\infty.$$

Minimize over δ on the both sides

$$\hat{\phi}_\Sigma^2(s_0) \geq \phi_\Sigma^2(s_0) - 16s_0 \|\hat{\Sigma} - \Sigma\|_\infty. \quad (\text{A.18})$$

So if we can prove that with probability approaching one, $16s_0 \|\hat{\Sigma} - \Sigma\|_\infty \leq \phi_\Sigma^2(s_0)/2$, that will imply of $\hat{\phi}_\Sigma^2(s_0) \geq \phi_\Sigma^2(s_0)/2$ with probability approaching one. Define $\epsilon_n = 16s_0 t_1$, where

$$t_1 = K \left[\sqrt{\frac{\ln p^2}{n}} + \frac{\sqrt{EM_2^2 \ln p^2}}{n} \right] + \sqrt{\frac{\ln p}{n}}. \quad (\text{A.19})$$

By (A.3)(A.4), via Assumption 1

$$\begin{aligned} P[16s_0 \|\hat{\Sigma} - \Sigma\|_\infty > \epsilon_n] &= P[\|\hat{\Sigma} - \Sigma\|_\infty > t_1] \\ &\leq \exp(-C_1 \ln p) + \frac{KEM_2^2}{(n \ln p)} \\ &\rightarrow 0, \end{aligned} \quad (\text{A.20})$$

where we use Assumption 2 for the probability tail converging to zero. Also see that by Assumption 2, $\epsilon_n \rightarrow 0$ since $s_0 \sqrt{\ln p/n} \rightarrow 0$. So we get, with probability approaching one, $16s_0 \|\hat{\Sigma} - \Sigma\|_\infty \leq \epsilon_n \leq \phi_\Sigma^2(s_0)/2$, since left side of that inequality converges to zero in probability, and the right side is constant. Then by (A.18)(A.20)

$$\begin{aligned} P[\hat{\phi}_\Sigma^2(s_0) \geq \phi_\Sigma^2(s_0)/2] &\geq 1 - \exp(-C_1 \kappa_n) - \frac{KEM_2^2}{(n \kappa_n)} \\ &= 1 - \frac{1}{p^{C_1}} - \frac{KEM_2^2}{n \ln p} \\ &= 1 - o(1). \end{aligned} \quad (\text{A.21})$$

Q.E.D.

We need the following Lemma for the exception set $\mathcal{F}^c := \{A_1 \cap A_2\}^c$ upper bound probability.

Lemma A.4. Under Assumptions 1, 2, with $\kappa_n = \ln p$

$$\begin{aligned} P(\mathcal{F}^c) &\leq 2\exp(-C_1\kappa_n) + \frac{K[EM_1^2 + EM_2^2]}{(n\kappa_n)} \\ &= \frac{2}{p^{C_1}} + \frac{K(EM_1^2 + EM_2^2)}{n\ln p} = o(1). \end{aligned}$$

Proof of Lemma A.4.

Now we provide an upper bound for the probability $P(\mathcal{F}^c)$ in our case under Assumptions 1, 2, by using Lemmata A.2-A.3

$$\begin{aligned} P(\mathcal{F}^c) &= P(\mathcal{A}_1 \cap \mathcal{A}_2)^c = P(\mathcal{A}_1^c \cup \mathcal{A}_2^c) \leq P(\mathcal{A}_1^c) + P(\mathcal{A}_2^c) \\ &\leq 2\exp(-C_1\kappa_n) + \frac{K[EM_1^2 + EM_2^2]}{(n\kappa_n)} \\ &= \frac{2}{p^{C_1}} + \frac{K[EM_1^2 + EM_2^2]}{n\ln p} \\ &\rightarrow 0. \end{aligned} \tag{A.22}$$

Q.E.D.

A.2.3 New Oracle Inequality Proofs

We start with proof of Theorems 1-2, where they are used as inputs to proof of Theorem 3. Theorems 1-2 consider the new oracle inequalities.

Proof of Theorem 1. We proceed in several steps.

Denote the joint event $\mathcal{F} = \{\mathcal{A}_1 \cap \mathcal{A}_2\}$. \mathcal{F}^c is \mathcal{F} 's complement. See that

$$E\|\hat{\beta} - \beta_0\|_1^k = E\|\hat{\beta} - \beta_0\|_1^k 1_{\{\mathcal{F}\}} + E\|\hat{\beta} - \beta_0\|_1^k 1_{\{\mathcal{F}^c\}}. \tag{A.23}$$

We want to form rates for the right side terms in (A.23).

Step 1. Note that by Lemma A.1, the first term on the right side of (A.23) is:

$$E\|\hat{\beta} - \beta_0\|_1^k 1_{\{\mathcal{F}\}} = O(s_0^k \lambda_n^k). \tag{A.24}$$

Now we want to evaluate the second term on the right side of (A.23). But before that we need the following intermediate step.

Step 2. Use Nemirowski's moment inequality, Lemma 14.24 in Buhlmann and van de Geer (2011), with for all $k \geq 1$, for the first inequality, and for the second inequality by Loeve's c_r inequality, and for the equality we use u_i being iid, also the definition of $\sigma^2 := Eu_i^2$,

$$\begin{aligned} E \left| \frac{\sum_{i=1}^n u_i^2 - \sigma^2}{n} \right|^k &\leq [8\ln(2)]^{k/2} E \left[\frac{\sum_{i=1}^n (u_i^4)}{n^2} \right]^{k/2} \\ &\leq \frac{Cn^{(k/2)-1}}{n^k} \sum_{i=1}^n Eu_i^{2k} \\ &= C[Eu_i^{2k}]n^{-k/2} = O(n^{-k/2}) = o(1), \end{aligned}$$

by Assumption 1. Before the next result we provide the inequality,

$$|x + y|^k \leq 2^{k-1}(|x|^k + |y|^k), \quad (\text{A.25})$$

for $k \geq 1$, and x, y being generic scalars, and σ^2 being bounded above by Assumption 1 and using (A.25)

$$\begin{aligned} E \left| \frac{1}{n} \sum_{i=1}^n u_i^2 \right|^k &= E \left| \frac{1}{n} \sum_{i=1}^n (u_i^2 - \sigma^2) + \sigma^2 \right|^k \\ &\leq 2^{k-1} \left[E \left| \frac{1}{n} \sum_{i=1}^n (u_i^2 - \sigma^2) \right|^k + (\sigma^2)^k \right] \\ &= O(n^{-k/2}) + O(1) = O(1). \end{aligned} \quad (\text{A.26})$$

Step 3. Now we have to form another l_1 expectation bound for lasso that will be key to the second right side term analysis in (A.23). This step 3 modifies the proof of Theorem 1, supplement, p.4 of Jankova and van de Geer (2018). We extend their proof to non-sub-Gaussian case and show that their bound is very conservative, and we provide a new less conservative bound. Start with the definition of lasso.

$$\|Y - X\hat{\beta}\|_n^2 + 2\lambda_n \|\hat{\beta}\|_1 \leq \|Y - X\beta_0\|_n^2 + 2\lambda_n \|\beta_0\|_1.$$

Ignore the first term and use the model $u = Y - X\beta_0$ to have

$$\|\hat{\beta}\|_1 \leq \frac{\|u\|_n^2}{2\lambda_n} + \|\beta_0\|_1.$$

Then use triangle inequality and then the inequality above

$$\|\hat{\beta} - \beta_0\|_1 \leq \|\hat{\beta}\|_1 + \|\beta_0\|_1 \leq \frac{\|u\|_n^2}{2\lambda_n} + 2\|\beta_0\|_1. \quad (\text{A.27})$$

Next taking the k th moment of the sampling error in l_1 norm, and using (A.25) by taking expectations there for the second inequality below

$$E \|\hat{\beta} - \beta_0\|_1^k \leq E \left[\frac{\|u\|_n^2}{2\lambda_n} + 2\|\beta_0\|_1 \right]^k \leq 2^{k-1} \{ E \left[\frac{\|u\|_n^2}{2\lambda_n} \right]^k + 2\|\beta_0\|_1^k \} \quad (\text{A.28})$$

We use the assumption $\|\beta_0\|_2 = O(1)$ to have

$$\|\beta_0\|_1^k \leq (\sqrt{s_0} \|\beta_0\|_2)^k = O(s_0^{k/2}). \quad (\text{A.29})$$

Then use the last equation with (A.26) in (A.28) to have

$$E \left[\frac{\|u\|_n^2}{2\lambda_n} \right]^k + 2\|\beta_0\|_1^k = O(\lambda_n^{-k}) + O(s_0^{k/2}) = O(\max(s_0^{k/2}, \lambda_n^{-k})). \quad (\text{A.30})$$

Note that proof of Jankova and van de Geer (2018) use $s_0^{k/2} \lambda_n^{-k}$ but this is very conservative upper bound since both two terms in multiplication is diverging with n . But a better bound is $\max(s_0^{k/2}, \lambda_n^{-k})$.

We get the rough bound for expectation using (A.30) in (A.28)

$$E\|\hat{\beta} - \beta_0\|_1^k = O(\max(s_0^{k/2}, \lambda_n^{-k})). \quad (\text{A.31})$$

Note that rates in (A.24)(A.31) are different and the last rate in this step is a rough bound which will be helpful in the next step. The rate in (A.31) is diverging to infinity.

Step 4. Rewrite the expectation using event $\mathcal{F}, \mathcal{F}^c$.

$$\begin{aligned} E\|\hat{\beta} - \beta_0\|_1^k &= E\|\hat{\beta} - \beta_0\|_1^k 1_{\{\mathcal{F}\}} + E\|\hat{\beta} - \beta_0\|_1^k 1_{\{\mathcal{F}^c\}} \\ &\leq O(s_0^k \lambda_n^k) + \sqrt{E\|\hat{\beta} - \beta_0\|_1^{2k}} \sqrt{E 1_{\{\mathcal{F}^c\}}} \\ &= O(s_0^k \lambda_n^k) + O(\max(s_0^{k/2}, \lambda_n^{-k})) \sqrt{P(\mathcal{F}^c)} \end{aligned} \quad (\text{A.32})$$

where we use (A.24) and Cauchy-Schwartz inequality for the first inequality, and the second equality is by (A.31).

First possibility of a rate is (jointly holding):

$$s_0^k \lambda_n^k \geq s_0^{k/2} P(\mathcal{F}^c)^{1/2}. \quad (\text{A.33})$$

$$s_0^k \lambda_n^k \geq \lambda_n^{-k} P(\mathcal{F}^c)^{1/2}. \quad (\text{A.34})$$

By (A.32)(A.33)(A.34)

$$E\|\hat{\beta} - \beta_0\|_1^k = O(s_0^k \lambda_n^k).$$

We can simplify further (A.33)(A.34), respectively they are

$$\lambda_n \geq P(\mathcal{F}^c)^{1/2k} / s_0^{1/2}, \quad (\text{A.35})$$

and

$$\lambda_n \geq P(\mathcal{F}^c)^{1/4k} / s_0^{1/2}. \quad (\text{A.36})$$

Since $P(\mathcal{F}^c)^{1/4k} \geq P(\mathcal{F}^c)^{1/2k}$, $k \geq 1$ (A.34) implies (A.33) or (A.36) implies (A.35). So if $\lambda_n \geq P(\mathcal{F}^c)^{1/4k} / s_0^{1/2}$ then

$$E\|\hat{\beta} - \beta_0\|_1^k = O(s_0^k \lambda_n^k). \quad (\text{A.37})$$

Of course there is another possibility-subcase that provides the rate in (A.37). That is when

$$s_0^k \lambda_n^k \geq \lambda_n^{-k} P(\mathcal{F}^c)^{1/2}, \quad (\text{A.38})$$

jointly holding with

$$\lambda_n^{-k} P(\mathcal{F}^c)^{1/2} \geq s_0^{k/2} P(\mathcal{F}^c)^{1/2}. \quad (\text{A.39})$$

This results in the same sufficient condition (A.36) via only (A.38) since by Assumption 2, $s_0\lambda_n \rightarrow 0$ which results in by rewriting (A.39): $1 \geq (s_0\lambda_n)^k/s_0^{k/2}$. So (A.39) is always satisfied with sufficiently large n . Note also that joint inequalities $s_0^k\lambda_n^k \geq s_0^{k/2}P(\mathcal{F}^c)^{1/2}$ jointly holding with

$$s_0^{k/2}P(\mathcal{F}^c)^{1/2} \geq \lambda_n^{-k}P(\mathcal{F}^c)^{1/2} \quad (\text{A.40})$$

is not possible since (A.40) implies $s_0^{k/2}\lambda_n^k \geq 1$ which is equivalent to $(s_0\lambda_n)^k/(s_0^{k/2}) \geq 1$. This last inequality cannot hold given $s_0\lambda_n \rightarrow 0$ Assumption 2 in large n .

To combine all the results for the k th moment of the estimation error, for values of $\lambda_n \geq P(\mathcal{F}^c)^{1/4k}/s_0^{1/2}$,

$$E\|\hat{\beta} - \beta_0\|_1^k = O(s_0^k\lambda_n^k).$$

The uniformity over $\mathcal{B}_{l_0}(s_0)$ follows since the rates in (A.24)(A.31)-(A.34) depends on β_0 only by s_0 .

Q.E.D.

Remark. Proof of Theorem 1 in Jankova and van de Geer (2018), in their appendix, p.5, shows that they use assumption:

$$\lambda_n \geq \frac{P(\mathcal{F}^c)^{1/4k}}{s_0^{1/4}}, \quad (\text{A.41})$$

which is equivalent to the following condition as shown in p.3 of proof of Theorem 1 in Jankova and van de Geer (2018)

$$\tau^2 > 2kln[(\sqrt{s_0}\lambda_n^2)^{-1}]/lnp,$$

given that $\lambda_n \geq C\tau\sqrt{lnp/n}$ and $C > 0, \tau > 1$ with

$$P(\mathcal{F}^c) \leq \frac{2}{(2p)\tau^{2/2}} \quad (\text{A.42})$$

by Lemma 7 in appendix of Jankova and van de Geer (2018). Our result and theirs are not comparable in terms of λ_n since they assume sub-Gaussian data, and ours is more general.

Proof of Theorem 2.

We start with

$$E\|\hat{\beta}\|_1^k = E\|\hat{\beta}\|_1^k 1_{\{\mathcal{F}\}} + E\|\hat{\beta}\|_1^k 1_{\{\mathcal{F}^c\}} \leq E\|\hat{\beta}\|_1^k 1_{\{\mathcal{F}\}} + \sqrt{E\|\hat{\beta}\|_1^{2k}} \sqrt{P(\mathcal{F}^c)}, \quad (\text{A.43})$$

by using Cauchy-Schwartz inequality. Then use triangle inequality on set \mathcal{F} and by Lemma A.1, and norm inequality to have

$$\begin{aligned} \|\hat{\beta}\|_1 &\leq \|\hat{\beta} - \beta_0\|_1 + \|\beta_0\|_1 \\ &\leq \frac{24\lambda_n s_0}{\phi_\Sigma^2(s_0)} + \sqrt{s_0}\|\beta_0\|_2 \\ &= O_p(\sqrt{s_0}), \end{aligned}$$

by Assumptions 1, 2. This last rate shows that

$$E\|\hat{\beta}\|_1^k 1_{\{\mathcal{F}\}} = O(s_0^{k/2}). \quad (\text{A.44})$$

To handle the second right side term in (A.43) we start with the second inequality in (A.27) and ignore $\|\beta_0\|_1$ in the middle to have

$$\|\hat{\beta}\|_1 \leq \frac{\|u\|_n^2}{2\lambda_n} + \|\beta_0\|_1.$$

then follow (A.30) to get

$$\begin{aligned} \sqrt{E\|\hat{\beta}\|_1^{2k} P(\mathcal{F}^c)^{1/2}} &= O(\max(s_0^{k/2}, \lambda_n^{-k})) P(\mathcal{F}^c)^{1/2} \\ &= O(\lambda_n^{-k} P(\mathcal{F}^c)^{1/2}), \end{aligned} \quad (\text{A.45})$$

and to get the second equality by Assumption 2(ii) $(s_0\lambda_n)^k/s_0^{k/2} \leq 1$ since the ratio on the left converges to zero, so this means $s_0^{k/2} \leq \lambda_n^{-k}$ with sufficiently large n .

Now use (A.44) with (A.45) in (A.43)

$$E\|\hat{\beta}\|_1^k = O(s_0^{k/2}) + O(\lambda_n^{-k} P(\mathcal{F}^c)^{1/2}). \quad (\text{A.46})$$

If $\lambda_n \geq P(\mathcal{F}^c)^{1/2k}/s_0^{1/2}$ it is clear that

$$s_0^{k/2} \geq \lambda_n^{-k} P(\mathcal{F}^c)^{1/2}, \quad (\text{A.47})$$

So by (A.47) in (A.46) we have the desired result. **Q.E.D.**

Q.E.D.

A.2.4 Main Theorem Proof: Incentive Compatibility

Proof of Theorem 3.

By Theorem 1 and 2 we can choose the larger of λ_n in those theorems, with $s_0 \geq 1$, and since it is nondecreasing with n ,

$$\lambda_n \geq \frac{P(\mathcal{F}^c)^{1/4k}}{s_0^{1/2}} \geq \frac{P(\mathcal{F}^c)^{1/2k}}{s_0^{1/2}} \quad (\text{A.48})$$

Add and subtract $X'_{n+1}\hat{\beta}$ inside the right hand side of the incentive compatibility definition:

$$\begin{aligned} E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 &= E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\hat{\beta} + X'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 \\ &= E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\hat{\beta}]^2 + E[X'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 \\ &\quad + E[\hat{\beta}'(\tilde{X}_{n+1} - X_{n+1})X'_{n+1}(\hat{\beta} - \beta_0)] \\ &\quad + E[(\hat{\beta} - \beta_0)'X_{n+1}(\tilde{X}'_{n+1} - X'_{n+1})\hat{\beta}]. \end{aligned} \quad (\text{A.49})$$

Using the definition of incentive compatibility, with defining $D_{n+1} := \tilde{X}_{n+1} - X_{n+1}$, we have

$$E[\tilde{X}'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 - E[X'_{n+1}\hat{\beta} - X'_{n+1}\beta_0]^2 = E[\hat{\beta}'D_{n+1}D'_{n+1}\hat{\beta}] \quad (\text{A.50})$$

$$+ E[\hat{\beta}'D_{n+1}X'_{n+1}(\hat{\beta} - \beta_0)] \quad (\text{A.51})$$

$$+ E[(\hat{\beta} - \beta_0)'X_{n+1}D'_{n+1}\hat{\beta}]. \quad (\text{A.52})$$

Now analyze (A.51), the analysis of (A.52) is the same and thus omitted. See that

$$\begin{aligned}
\hat{\beta}' D_{n+1} X'_{n+1} (\hat{\beta} - \beta_0) &\leq |\hat{\beta}' D_{n+1} X'_{n+1} (\hat{\beta} - \beta_0)| \\
&\leq |\hat{\beta}' D_{n+1}| |X'_{n+1} (\hat{\beta} - \beta_0)| \\
&\leq \|\hat{\beta}\|_1 \|D_{n+1}\|_\infty \|X_{n+1}\|_\infty \|\hat{\beta} - \beta_0\|_1,
\end{aligned} \tag{A.53}$$

where we use Holder's inequality. Then

$$E[\hat{\beta}' D_{n+1} X'_{n+1} (\hat{\beta} - \beta_0)] \leq E \left[\|\hat{\beta}\|_1 \|D_{n+1}\|_\infty \|X_{n+1}\|_\infty \|\hat{\beta} - \beta_0\|_1 \right] \tag{A.54}$$

$$\leq [E\|\hat{\beta}_1^4\|^{1/4} [E\|D_{n+1}\|_\infty^4]^{1/4} [E\|X_{n+1}\|_\infty^4]^{1/4} [E\|\hat{\beta} - \beta_0\|_1^4]^{1/4}] \tag{A.55}$$

$$= [E\|\hat{\beta}_1^4\|^{1/4} [EM_4^4]^{1/4} [EM_3^4]^{1/4} [E\|\hat{\beta} - \beta_0\|_1^4]^{1/4}] \tag{A.56}$$

where we apply (A.53) for the first inequality and Holder's Inequality in the second inequality above, and the last equality comes from M_3, M_4 definitions. Then we apply Theorems 1-2 with $k = 4$. We assume $\lambda_n \geq P(\mathcal{F}^c)^{1/16}/s_0^{1/2}$ and if

$$s_0^{3/2} \sqrt{\frac{\ln p}{n}} [EM_3^4]^{1/4} [EM_4^4]^{1/4} \rightarrow 0, \tag{A.57}$$

we see that (A.56) goes to zero, by Theorems 1-2, and $\lambda_n = O(\sqrt{\frac{\ln p}{n}})$.

So looking at incentive compatibility definition and (A.50)-(A.52)

$$E[\tilde{X}'_{n+1} \hat{\beta} - X'_{n+1} \beta_0]^2 - E[X'_{n+1} \hat{\beta} - X'_{n+1} \beta_0]^2 = E[\hat{\beta}' D_{n+1} D'_{n+1} \hat{\beta}] + o(1), \tag{A.58}$$

where the first right side term in (A.58) is nonnegative and the other terms are negligible in large samples by (A.57).

The uniformity over $\mathcal{B}_{l_0}(s_0)$ goes through since Theorems 1, 2 depend on β_0 only through s_0 , and they are the main ingredient in the proof.

Q.E.D.

B Appendix B

Here we consider results when $p \leq n$, and relaxing Assumption 2(iii).

B.1 When $p \leq n$

There are minor modifications in the proofs compared to $p > n$. We consider them here. One major change is since $p \leq n$, we set $\kappa_n = lnn$. Change Assumption 2(ii) so that $s_0 \sqrt{\ln n/n} \rightarrow 0$.

We provide the maximal inequality here. Now take the case of $p \leq n$, and combine (A.2) with (A.3) to have with $\kappa_n = lnn$ in that case

$$\begin{aligned}
P\left(\max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| \geq 2K \left[\frac{\sqrt{\ln p}}{\sqrt{n}} + \frac{(EM_F^2)^{1/2} \ln p}{n} \right] + \frac{\sqrt{\ln n}}{\sqrt{n}} \right) \\
\leq \frac{1}{n^{C_1}} + \frac{EM_F^2}{n(lnn)} = o(1),
\end{aligned} \tag{B.1}$$

by Assumptions A1-A.2. To see this point

$$\frac{EM_F^2}{nl\kappa_n} = \left[\left(\frac{(EM_F^2)^{1/2} \sqrt{\ln p}}{\sqrt{n}} \right) \frac{1}{\sqrt{\ln n} \sqrt{\ln p}} \right]^2 = o(1). \quad (\text{B.2})$$

This shows also that

$$\max_{1 \leq j \leq p} |\hat{\mu}_j - \mu_j| = O_p(\sqrt{\ln n}/\sqrt{n}). \quad (\text{B.3})$$

Lemma A.1 will be the same. Lemma A.2(i) lower bound probability has $\kappa_n = l\kappa_n$ now. Lemma A.2(ii) is the same. Lemma A.2(iii) will change to $\lambda_n = O(\sqrt{\ln n}/\sqrt{n})$. Lemma A.3 use $\kappa_n = l\kappa_n$, so (A.19) becomes

$$t_1 = K \left[\frac{\sqrt{\ln p^2}}{\sqrt{n}} + \frac{\sqrt{EM_2^2 \ln p^2}}{n} \right] + \frac{\sqrt{\ln n}}{n}.$$

Lemma A.4 is the same with $\kappa_n = l\kappa_n$.

Given these results, the proof of Theorem 1 is the same with $\lambda_n = O(\sqrt{\frac{\ln n}{n}})$. Theorem 2 does not change. Theorem 3 condition will be changing to

$$s_0^{3/2} \sqrt{\frac{\ln n}{n}} [EM_3^4]^{1/4} [EM_4^4]^{1/4} \rightarrow 0,$$

B.2 Relaxing Assumption 2(iii)

In this subsection we relax Assumption 2(iii) from $\|\beta_0\|_2 = O(1)$ to $\|\beta_0\|_2 = O(\sqrt{s_0})$ and we explain the logic and meaning of this new assumption.

Assumption 2(iv).

$$\|\beta_0\|_2 = O(\sqrt{s_0}).$$

Assumption 2(iii) which is suggested by Jankova and van de Geer (2018) and simplifies their paper in semiparametric efficient estimators. Our Assumption 2(iv) here generalizes that assumption and in the case of s_0 being constant becomes Assumption 2(iii). The implication of Assumption 2(iv) is that all nonzero coefficients can be constant and none of them has to be local to zero.

$$\|\beta_0\|_2 = \sqrt{\sum_{j=1}^p \beta_{0,j}^2} = \sqrt{\sum_{j \in S_0} \beta_{0,j}^2} = O(\sqrt{s_0}).$$

In terms of Section 2 discussion after Assumption 2, this implies $S_0 = F_1$, and F_2 is an empty set. So Assumption 2(iv) can simultaneously allow s_0 increasing with n , and all large nonzero coefficients in S_0 . Previously in Assumption 2(iii), there can be only a fixed number of large coefficients, and increasing $(s_0 - f_1)$ number of local to zero (small) coefficients.

We proceed in a way that we only change the proofs in Appendix A, when necessary. All lemmata in Appendix A goes through, there is no usage of Assumption 2(iii) there. The first change comes in step 3

of Theorem 1 proof. First (A.29) changes to $\|\beta_0\|_1^k = O(s_0^k)$ under Assumption 2(iv) instead of Assumption 2(iii). Then (A.30) becomes

$$E \left[\frac{\|u\|_n^2}{2\lambda_n} \right]^k + 2\|\beta_0\|_1^k = O(\max(s_0^k, \lambda_n^{-k})). \quad (\text{B.4})$$

Then (A.32) changes to following

$$E\|\hat{\beta} - \beta_0\|_1^k = O(s_0^k \lambda_n^k) + O(\max(s_0^k, \lambda_n^{-k}) \sqrt{P(\mathcal{F}^c)}). \quad (\text{B.5})$$

Instead of (A.33)(A.34) we have the following conditions, to establish the rate for the oracle inequality (i.e. mean l_1 norm bound to k th order)

$$s_0^k \lambda_n^k \geq s_0^k P(\mathcal{F}^c)^{1/2}. \quad (\text{B.6})$$

$$s_0^k \lambda_n^k \geq \lambda_n^{-k} P(\mathcal{F}^c)^{1/2}. \quad (\text{B.7})$$

Using (B.5)-(B.7)

$$E\|\hat{\beta} - \beta_0\|_1^k = O(s_0^k \lambda_n^k). \quad (\text{B.8})$$

The conditions (B.6)(B.7) can be written as

$$\lambda_n \geq \max(P(\mathcal{F}^c)^{1/2k}, P(\mathcal{F}^c)^{1/4k} / s_0^{1/2}), \quad (\text{B.9})$$

where the tuning parameter choice under Assumption 2(iv) which is (B.9) is larger than or equal to choice by Assumption 2(iii), which is the second component in the max on the right side of (B.9). The discussion after this in step 4 is the same, given Assumption 2(i)-(ii). So we have the following result:

Corollary B.1. *Under Assumptions 1, 2(i)(ii)(iv), with*

$$\lambda_n \geq \max(P(\mathcal{F}^c)^{1/2k}, P(\mathcal{F}^c)^{1/4k} / s_0^{1/2}).$$

we have

$$[E\|\hat{\beta} - \beta_0\|_1^k]^{1/k} = O(s_0 \lambda_n).$$

The result is also uniform over l_0 ball \mathcal{B}_{l_0}

Now we modify the proof of Theorem 2. In that respect, by Assumption 2(iv) the rate after (A.43) becomes

$$\|\hat{\beta}\|_1 = O_p(s_0). \quad (\text{B.10})$$

Then (A.46) changes to

$$E\|\hat{\beta}\|_1^k = O(s_0^k) + O(\lambda_n^{-k} P(\mathcal{F}^c)^{1/2}). \quad (\text{B.11})$$

We can show that

$$s_0^k \geq \lambda_n^{-k} P(\mathcal{F}^c)^{1/2}, \quad (\text{B.12})$$

if we have

$$\lambda_n \geq P(\mathcal{F}^c)^{1/2k} / s_0. \quad (\text{B.13})$$

Then given (B.13), using (B.12) in (B.11) we have

$$E\|\hat{\beta}\|_1^k = O(s_0^k).$$

So we established the following Corollary to Theorem 2. The result is different from Theorem 2 and the k th moment of l_1 error grows faster here in Corollary B.2 if s_0 increases with n . So relaxed assumption comes with a cost that will affect main incentive compatibility condition.

Corollary B.2. *Under Assumptions 1, 2(i)(ii)(iv), with*

$$\lambda_n \geq P(\mathcal{F}^c)^{1/2k}/s_0.$$

we have

$$[E\|\hat{\beta}\|_1^k]^{1/k} = O(s_0).$$

The result is also uniform over l_0 ball \mathcal{B}_{l_0}

Now we follow the proof of Theorem 3 and substitute Assumption 2(iv) instead of Assumption 2(iii). Note that our λ_n choice must choose the maximum of the ones in Corollary B.1 and B.2. Clearly Corollary B.1 tuning parameter is larger than the one in Corollary B.2. The only place we have to change there is (A.57). Given $\lambda_n \geq \max(P(\mathcal{F}^c)^{1/8}, \frac{P(\mathcal{F}^c)^{1/16}}{s_0^{1/2}})$ we need

$$s_0^2 \sqrt{\frac{\ln p}{n}} [EM_3^4]^{1/4} [EM_4^4]^{1/4} \rightarrow 0,$$

to have Incentive Compatibility in large samples. So we have the following counterpart to Theorem 3.

Corollary B.3. *Under Assumptions 1, 2(i)(ii)(iv) and*

$$\lambda_n \geq \max(P(\mathcal{F}^c)^{1/8}, \frac{P(\mathcal{F}^c)^{1/16}}{s_0^{1/2}}),$$

and

$$s_0^2 \sqrt{\frac{\ln p}{n}} [EM_3^4]^{1/4} [EM_4^4]^{1/4} \rightarrow 0,$$

lasso is Incentive Compatible. The result is also uniform over l_0 ball \mathcal{B}_{l_0} .

Clearly, there are two differences between Theorem 3 and Corollary B.3 here. First, we need a tuning parameter in Corollary B.3 which may be larger than or equal to the one in Theorem 3. Then, incentive compatibility of lasso is more difficult to achieve, due to sparsity, s_0 , having exponent of 2 here instead of 3/2 in Theorem 3.

References

Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80, 2369–2429.

- Belloni, A., V. Chernozhukov, and C. Hansen (2014). Inference on treatment effects after selection among high dimensional controls. *Review of Economic Studies* 81, 608–650.
- Buhlmann, P. and S. van de Geer (2011). Statistics for high-dimensional data. *Springer Verlag*.
- Cai, Y., C. Daskalakis, and C. Papadimitrou (2015). Optimum statistical estimation with strategic data sources. *Proceedings of the 28 th Conference on Learning Theory* 40, 1–40.
- Caner, M. and A. B. Kock (2018). Asymptotically honest confidence regions for high dimensional parameters by the desparsified conservative lasso. *Journal of Econometrics* 203, 143–168.
- Caner, M. and A. B. Kock (2019). High dimensional linear gmm. *arXiv:1811.08779*.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2017). Central limit theorems and bootstrap in high dimensions. *Annals of Probability* 45, 2309–2452.
- Chernozhukov, V., M. Goldman, V. Semenova, and M. Taddy (2018). Orthogonal machine learning for demand estimation: High dimensional causal inference in dynamic panels. *arXiv:1712.09988*.
- Chiang, H. (2020). Many average partial effects: with an application to text regression. *Working Paper*.
- Chiang, H. and Y. Sasaki (2019). Causal inference by quantile regression kink designs. *Journal of Econometrics* 210, 405–433.
- Cummings, R., S. Ioannidis, and K. Ligett (2015). Truthful linear regression. *Conference on Learning Theory* 40, 448–483.
- Dekel, O., F. Fischer, and A. Procaccia (2010). Incentive compatible regression learning. *Journal of Computer System and Sciences* 76, 759–77.
- Eliasz, K. and R. Spiegler (2019). The model selection curse. *American Economic Review-Insights* 1, 127–140.
- Eliasz, K. and R. Spiegler (2020). On incentive compatible estimators. *Working Paper-Tel Aviv University*.
- Gao, C., A. Van der Vaart, and H. Zhou (2015). A general framework for bayes structured linear models. *arXiv:1506.02174*.
- Hardt, M., N. Megiddo, C. Papadimitrou, and M. Wooters (2016). Strategic classification. *Proceedings of the ACM Conference on Innovations in Theoretical Computer Science*, 111–122.
- Jankova, J. and S. van de Geer (2018). Semi-parametric efficiency bounds for high-dimensional models. *Annals of Statistics* 46, 2336–2359.
- Kock, A. (2016). Oracle inequalities, variable selection and uniform inference in high-dimensional correlated random effects panel data models. *Journal of Econometrics* 195, 71–85.

- Kock, A. and H. Tang (2019). Inference in high-dimensional dynamic panel data models. *Econometric Theory* 35, 295–359.
- Meir, R., A. Procaccia, and J. Rosenschein (2012). Algorithms for strategyproof classification. *Artificial Intelligence* 186, 123–156.
- Perte, J. and J. Perote-Pena (2004). Strategy-proof estimators for simple regression. *Mathemtical Social Sciences* 47, 153–176.
- Shaywitz, D. (2020). "the alignment problem" review: When machines miss the point. *The Wall Street Journal*, A25,25 October.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of The Royal Statistical Society Series B* 58, 267–288.
- van de Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Annals of Statistics*.