

Compressed Least-Squares Regression

Authored by:

Rémi Munos
Odalric Maillard

Abstract

We consider the problem of learning, from K input data, a regression function in a function space of high dimension N using projections onto a random subspace of lower dimension M . From any linear approximation algorithm using empirical risk minimization (possibly penalized), we provide bounds on the excess risk of the estimate computed in the projected subspace (compressed domain) in terms of the excess risk of the estimate built in the high-dimensional space (initial domain). We apply the analysis to the ordinary Least-Squares regression and show that by choosing $M=O(\sqrt{K})$, the estimation error (for the quadratic loss) of the “Compressed Least Squares Regression is $O(1/\sqrt{K})$ up to logarithmic factors. We also discuss the numerical complexity of several algorithms (both in initial and compressed domains) as a function of N , K , and M .

1 Paper Body

We consider the problem of learning, from K data, a regression function in a linear space of high dimension N using projections onto a random subspace of lower dimension M . From any algorithm minimizing the (possibly penalized) empirical risk, we provide bounds on the excess risk of the estimate computed in the projected subspace (compressed domain) in terms of the excess risk of the estimate built in the high-dimensional space (initial domain). We show that solving the problem in the compressed domain instead of the initial domain reduces the estimation error at the price of an increased (but controlled) approximation error. We apply the analysis to Least-Squares (LS) regression and discuss the excess risk and numerical complexity of the resulting “Compressed Least Squares Regression” (CLSR) in terms of N , K , and M . When we choose $M = O(\sqrt{K})$, we show that CLSR has an estimation error of order $O(\log K / \sqrt{K})$.

1

Problem setting

We consider a regression problem where we observe data $DK = (\{x_k, y_k\}_{k=1}^K)$ (where $x_k \in X$ and $y_k \in R$) are assumed to be independently and

identically distributed (i.i.d.) from some distribution P , where $x_k \in \mathcal{X}$ and $y_k = f(x_k) + \epsilon_k(x_k)$, where f is the (unknown) target function, and ϵ_k a centered independent noise of variance $\sigma^2(x_k)$. For a given class of functions F , and $f \in F$, we define the empirical (quadratic) error def

$$L_K(f) = \frac{1}{K} \sum_{k=1}^K [y_k - f(x_k)]^2,$$

and the generalization (quadratic) error def

$$L(f) = E_{(X,Y) \sim P} [(Y - f(X))^2].$$

Our goal is to return a regression function $\hat{f} \in F$ with lowest possible generalization error $L(\hat{f})$. Notations: In the sequel we will make use of the following notations about norms: for $h : \mathcal{X} \rightarrow \mathbb{R}$, we write $\|h\|_P$ for the L_2 norm of h with respect to (w.r.t.) the measure P , $\|h\|_{PK}$ for the L_2 norm

$\|h\|_{PK} = \sqrt{\frac{1}{K} \sum_{i=1}^K h(x_i)^2}$ of h w.r.t. the empirical measure P_K , and for $u \in \mathbb{R}^n$, $\|u\|$ denotes by default $\|u\|_2 = \sqrt{\sum_{i=1}^n u_i^2}$. The measurable function minimizing the generalization error is f^* , but it may be the case that $f^* \notin F$. For any regression function \hat{f} , we define the excess risk $L(\hat{f}) - L(f^*) = \| \hat{f} - f^* \|_P^2$, which decomposes as the sum of the estimation error $L(\hat{f}) - \inf_{f \in F} L(f)$ and the approximation error $\inf_{f \in F} L(f) - L(f^*) = \inf_{f \in F} \|f - f^*\|_P^2$ which measures the distance between f^* and the function space F . 1

In this paper we consider a class of linear functions FN defined as the span of a set of N functions def $FN = \text{span} \{ \phi_n \}_{n=1}^N$ called features. Thus: $FN = \{ f = \sum_{n=1}^N \alpha_n \phi_n, \alpha_n \in \mathbb{R} \}$. When the number of data K is larger than the number of features N , the ordinary Least-Squares Regression (LSR) provides the LS solution \hat{f} which is the minimizer of the empirical risk $L_K(\hat{f})$ in FN . Note that here $L_K(\hat{f})$ rewrites $\frac{1}{K} \|Y - \Phi \hat{\alpha}\|^2$ where Φ is the $K \times N$ matrix with elements $(\phi_n(x_k))_{n=1, k=1}^N$ and Y the K -vector with components $(y_k)_{k=1}^K$. Usual results provide bound on the estimation error as a function of the capacity of the function space and the number of data. In the case of linear approximation, the capacity measures (such as covering numbers [23] or the pseudo-dimension [16]) depend on the number of features (for example the pseudo-dimension is at most $N + 1$). For example, let \hat{f} be a LS estimate (minimizer of L_K in FN), then (a more precise statement will be stated later in Subsection 3) the expected estimation error is bounded as:

$$N \log K E L(\hat{f}) \leq \inf_{f \in FN} L(f) + c \frac{1}{K}, \quad (1)$$

where c is a universal constant, $\sigma^2 = \sup_{x \in \mathcal{X}} \sigma^2(x)$, and the expectation is taken with respect to P . Now, the excess risk is the sum of this estimation error and the approximation error $\inf_{f \in FN} \|f - f^*\|_P^2$ of the class FN . Since the later usually decreases when the number of features N increases [13] (e.g. when FN is dense in $L_2(P)$), we see the usual tradeoff between small estimation error (low N) and small approximation error (large N). In this paper we are interested in the setting when N is large so that the approximation error is small. Whenever N is larger than K we face the overfitting problem since there are more parameters than actual data (more variables than constraints), which is illustrated in the bound (1) which provides no information

about the generalization ability of any LS estimate. In addition, there are many minimizers (in fact a vector space of same dimension as the null space of T) of the empirical risk. To overcome the problem, several approaches have been proposed in the literature: ℓ_2 LS solution with minimal norm: The solution is the minimizer of the empirical error with minimal (ℓ_1 or ℓ_2)-norm: $b = \arg \min_{\|b\|_2} \|Y - Tb\|_2$, (or a robust solution $\arg \min_{\|b\|_2} \|Y - Tb\|_2$). The choice of ℓ_2 -norm yields the ordinary LS solution. The choice of ℓ_1 -norm has been used for generating sparse solutions (e.g. the Basis Pursuit [10]), and assuming that the target function admits a sparse decomposition, the field of Compressed Sensing [9, 21] provides sufficient conditions for recovering the exact solution. However, such conditions (e.g. that T possesses a Restricted Isometric Property (RIP)) does not hold in general in this regression setting. On another aspect, solving these problems (both for ℓ_1 or ℓ_2 -norm) when N is large is numerically expensive. ℓ_1 Regularization. The solution is the minimizer of the empirical error plus a penalty term, for example $\hat{b} = \arg \min_{\|b\|_1} L_K(f) + \lambda \|b\|_1$, for $p = 1$ or 2 .

where λ is a parameter and usual choices for the norm are ℓ_2 (ridge-regression [20]) and ℓ_1 (LASSO [19]). A close alternative is the Dantzig selector [8, 5] which solves: $b = \arg \min_{\|b\|_1} \|Y - Tb\|_2$. The numerical complexity and generalization bounds of those methods depend on the sparsity of the target function decomposition in F_N . Now if we possess a sequence of function classes $(F_N)_N$ with increasing capacity, we may perform structural risk minimization [22] by solving in each model the empirical risk penalized by a term that depends on the size of the model: $\hat{b}_N = \arg \min_{f \in F_N} L_K(f) + \text{pen}(N, K)$, where the penalty term measures the capacity of the function space. In this paper we follow another approach where instead of searching in the large space F_N (where $N \leq K$) for a solution that minimizes the empirical error plus a penalty term, we simply search for the empirical error minimizer in a (randomly generated) lower dimensional subspace $G_M \subset F_N$ (where $M \leq K$). Our contribution: We consider a set of M random linear combinations of the initial N features and perform our favorite LS regression algorithm (possibly regularized) using those M compressed

features. This is equivalent to projecting the K points $\{x_k\}_{k=1}^K \subset \mathbb{R}^N$, $k = 1..K$ from the initial domain (of size N) onto a random subspace of dimension M , and then performing the regression in the compressed domain (i.e. span of the compressed features). This is made possible because random projections approximately preserve inner products between vectors (by a variant of the Johnson-Lindenstrauss Lemma stated in Proposition 1. Our main result is a bound on the excess risk of a linear estimator built in the compressed domain in terms of the excess risk of the linear estimator built in the initial domain (Section 2). We further detail the case of ordinary Least-Squares Regression (Section 3) and discuss, in terms of M , N , K , the different tradeoffs concerning the excess risk (reduced estimation error in the compressed domain versus increased approximation error introduced by the random projection) and the numerical complexity (reduced complexity of solving the LSR in the compressed domain versus the additional load of performing the projection). As a consequence,

we show that by choosing $M = O(K)$ projections we define a Compressed Least-Squares Regression which uses $O(NK^{3/2})$ elementary operations to compute a regression function with estimation error (relatively to the initial function space FN) of order $\log K/K$ up to a multiplicative factor which depends on the best approximation of f in FN . This is competitive with the best methods, up to our knowledge. Related works: Using dimension reduction and random projections in various learning areas has received considerable interest over the past few years. In [7], the authors use a SVM algorithm in a compressed space for the purpose of classification and show that their resulting algorithm has good generalization properties. In [25], the authors consider a notion of compressed linear regression. For data $Y = X\theta + \epsilon$, where ϵ is the target and ϵ a standard noise, they use compression of the set of data, thus considering $AY = AX\theta + A\epsilon$, where A has a Restricted Isometric Property. They provide an analysis of the LASSO estimator built from these compressed data, and discuss a property called sparsistency, i.e. the number of random projections needed to recover θ (with high probability) when it is sparse. These works differ from our approach in the fact that we do not consider a compressed (input and/or output) data space but a compressed feature space instead. In [11], the authors discuss how compressed measurements may be useful to solve many detection, classification and estimation problems without having to reconstruct the signal ever. Interestingly, they make no assumption about the signal being sparse, like in our work. In [6, 17], the authors show how to map a kernel $k(x, y) = \langle \phi(x), \phi(y) \rangle$ into a low-dimensional space, while still approximately preserving the inner products. Thus they build a low-dimensional feature space specific for (translation invariant) kernels.

2

Linear regression in the compressed domain

We remind that the initial set of features is $\{\phi_n : X \rightarrow \mathbb{R}, 1 \leq n \leq N\}$ and the initial domain $FN = \{f = \sum_{n=1}^N \alpha_n \phi_n, \alpha_n \in \mathbb{R}\}$ is the span of those features. We write $\phi(x)$ the N -vector of components $(\phi_n(x))_{n=1}^N$. Let us now define the random projection. Let A be a $M \times N$ matrix of i.i.d. elements drawn for some distribution \mathcal{D} . Examples of distributions are: \mathcal{D} Gaussian random variables $N(0, 1/M)$,

\mathcal{D} Bernoulli distributions, i.e. which takes values $\pm 1/M$ with equal probability $1/2$, \mathcal{D} Distribution taking values $\pm \sqrt{3}/M$ with probability $1/6$ and 0 with probability $2/3$. The following result (proof in the supplementary material) states the property that inner-product are approximately preserved through random projections (this is a simple consequence of the JohnsonLindenstrauss Lemma): Proposition 1 Let $(u_k)_{k=1}^K$ and v be vectors of \mathbb{R}^N . Let A be a $M \times N$ matrix of i.i.d. elements drawn from one of the previously defined distributions. For any $\epsilon > 0$, $\delta > 0$, for $M \geq \frac{1}{\delta^2} \frac{1}{\epsilon^2} \log \frac{4K}{\delta}$, we have, with probability at least $1 - \delta$, for all $k \leq K$,

?

6

— $Au_k - Av - \epsilon \|u_k - v\|$ — $\leq \epsilon \|u_k - v\|$ — $\leq \epsilon \|u_k - v\|$. 3
def

We now introduce the set of M compressed features $(\phi_m)_{1 \leq m \leq M}$ such that $\phi_m(x) = \frac{1}{M} \sum_{i=1}^M \langle x, \phi_i \rangle \phi_i$. We also write $\phi(x)$ the M -vector of components $(\phi_m(x))_{m=1}^M$. Thus $\phi(x) = \frac{1}{M} A_m \phi_m(x)$. PM def $\phi(x) = A \phi(x)$. We define the compressed domain $GM = \{g \mid g = \sum_{m=1}^M \phi_m \phi_m^T, \phi \in \mathbb{R}^M\}$ the span of the compressed features (vector space of dimension at most M). Note that each $\phi_m \in \mathbb{F}^N$, thus GM is a subspace of \mathbb{F}^N . 2.1

Approximation error

We now compare the approximation error assessed in the compressed domain GM versus in the initial space \mathbb{F}^N . This applies to the linear algorithms mentioned in the introduction such as ordinary LS regression (analyzed in details in Section 3), but also its penalized versions, e.g. LASSO and ridge regression. Define $\lambda = \arg \min_{\lambda} \frac{1}{2} \sum_{i=1}^N L(f_i - \lambda) = L(f - \lambda)$ the parameter of the best regression function in \mathbb{F}^N . Theorem 1 For any $\lambda \geq 0$, any $M \geq 15 \log(8K/\lambda)$, let A be a random $M \times N$ matrix defined like in Proposition 1, and GM be the compressed domain resulting from this choice of A . Then with probability at least $1 - \lambda$, r

$\frac{8 \log(8K/\lambda)}{M} + 2 \log \frac{4}{\lambda} \leq \frac{2}{M} \inf_{g \in GM} \sum_{i=1}^M (f_i - g_i)^2 \leq \frac{2}{M} \inf_{f \in \mathbb{F}^N} \sum_{i=1}^M (f_i - \lambda)^2 + 2 \sup_{x \in X} \sum_{i=1}^M (\phi_i(x) - \lambda)^2 + \inf_{f \in \mathbb{F}^N} \sum_{i=1}^M (f_i - \lambda)^2$. $g \in GM$ $f \in \mathbb{F}^N$ $M \geq 2K$ $x \in X$ (2) This theorem shows the tradeoff in terms of estimation and approximation errors for an estimator g_b obtained in the compressed domain compared to an estimator f_b obtained in the initial domain: λ Bounds on the estimation error of g_b in GM are usually smaller than that of f_b in \mathbb{F}^N when $M \leq N$ (since the capacity of \mathbb{F}^N is larger than that of GM). λ Theorem 1 says that the approximation error assessed in GM increases by at most $O(\log(K/\lambda))$ compared to that in \mathbb{F}^N . M def

def

Proof: Let us write $f + = f - \lambda = \arg \min_{f \in \mathbb{F}^N} \sum_{i=1}^M (f_i - \lambda)^2$ and $g + = g - \lambda$. The approximation error assessed in the compressed domain GM is bounded as $\inf_{g \in GM} \sum_{i=1}^M (g_i - \lambda)^2$

$g \in GM$

λ

$$\sum_{i=1}^M (g_i + \lambda - \lambda)^2 = \sum_{i=1}^M (g_i + \lambda - \lambda)^2 + \sum_{i=1}^M (f_i + \lambda - \lambda)^2 = \sum_{i=1}^M (g_i + \lambda - \lambda)^2 + \sum_{i=1}^M (f_i + \lambda - \lambda)^2,$$

(3)

since $f +$ is the orthogonal projection of $f -$ on \mathbb{F}^N and $g +$ belongs to \mathbb{F}^N . We now bound $\sum_{i=1}^M (g_i + \lambda - \lambda)^2$ def

def

$f + = \sum_{i=1}^M (f_i - \lambda) \phi_i$ using concentration inequalities. Define $Z(x) = A \phi(x) = \sum_{i=1}^M \phi_i \phi_i^T \phi(x)$. Define $\lambda_2 = \log(8K/\lambda) \frac{8}{M} \log(8K/\lambda)$. For $M \geq 15 \log(8K/\lambda)$ we have $\lambda \leq 3/4$ thus $M \geq 2 \lambda_2 / 3 \geq 6$. Proposition 1 applies and says that on an event E of probability at least $1 - \lambda/2$, we have for all $k \leq K$, def

$$\sum_{i=1}^M (Z(x_k) - \lambda)^2 \leq \sum_{i=1}^M (Z(x_k) - \lambda)^2 + \sum_{i=1}^M (Z(x_k) - \lambda)^2 \leq \sup_{x \in X} \sum_{i=1}^M (Z(x) - \lambda)^2 = C$$

(4)

$x \in X$

On the event E , we have with probability at least $1 - \lambda/2$, $+$

$$\sum_{i=1}^M (g_i - \lambda)^2$$

$$f + = \sum_{i=1}^M (f_i - \lambda)^2$$

(in dimension 1) with coefficients $b_i = \int_{\mathbb{R}} f(x) \psi_i(x) dx$. Thus here $\beta = 1$. Other classes (such as Sobolev spaces) lead to larger values of β related to the order of differentiability. $\beta \leq \frac{1}{2}$ By choosing $c_i = i^{-\beta/2}$, we have $\sum_{i=1}^N |b_i|^2 \leq C \sum_{i=1}^N i^{-\beta}$. Thus if $\beta > 1$, then this term is bounded by a constant that does not depend on N . If $\beta = 1$ then it is bounded by $O(\log N)$, and if $0 < \beta < 1$, then it is bounded by $O(N^{1-\beta})$. However any orthonormal basis, even rescaled, would not necessarily yield a small $\sum_{i=1}^N |b_i|^2$ term (this is all the more true when the dimension of X is large). The desired property that the coefficients (b_i) of the decomposition of f rapidly decrease to 0 indicates that hierarchical bases, such as wavelets, that would decompose the function at different scales, may be interesting. Wavelets: Consider an infinite family of wavelets in $[0, 1]$: $(\psi_{h,l})$ (indexed by $n \geq 1$ or equivalently by the scale $h \geq 0$ and translation $0 \leq l \leq 2^h - 1$) where $\psi_{h,l}(x) = 2^{h/2} \psi(2^h x - l)$ and ψ is the mother wavelet. Then consider $N = 2^h$ features $(\psi_{h,l})_{l=0}^{2^h-1}$ defined as the rescaled def

wavelets $\psi_{h,l} = 2^{h/2} \psi(2^h x - l)$, where $c_l \geq 0$ are some coefficients. Assume the mother wavelet ψ is C^p (for $p \geq 1$), has at least p vanishing moments, and that for all $h \geq 0$, $\sup_{l=0}^{2^h-1} \int_{\mathbb{R}} \psi_{h,l}(x) dx = 0$. Then the following proposition (proof in the supplementary material) provides a bound on $\sup_{x \in \mathbb{R}} \sum_{l=0}^{2^h-1} \psi_{h,l}(x)^2$ (thus on $E(\sum_{l=0}^{2^h-1} \psi_{h,l}(x)^2)$) by a constant independent of N : Proposition 2 Assume that f is (L, β) -Lipschitz (i.e. for all $v \in X$ there exists a polynomial p_v of degree β such that for all $u \in X$, $|f(u) - p_v(u)| \leq L \|u - v\|^\beta$) with $1/2 \leq \beta \leq p$. Then setting $R_1 = 2^{h/2} \int_{\mathbb{R}} \psi(x)^2 dx$, we have $\sup_{x \in \mathbb{R}} \sum_{l=0}^{2^h-1} \psi_{h,l}(x)^2 \leq R_1 2^{h/2} \int_{\mathbb{R}} \psi(x)^2 dx$, which is independent of N . Notice that the Haar wavelets has $p = 1$ vanishing moment but is not C^1 , thus the Proposition does not apply directly. However direct computations show that if f is L -Lipschitz (i.e. $\beta = 1$) then $\int_{\mathbb{R}} \psi_{h,l}(x)^2 dx \leq L^2 2^{h/2} \int_{\mathbb{R}} \psi(x)^2 dx$, and thus $\sup_{x \in \mathbb{R}} \sum_{l=0}^{2^h-1} \psi_{h,l}(x)^2 \leq 4(1/2) L^2 \int_{\mathbb{R}} \psi(x)^2 dx$ with $ch = 2^{h/2}$. 7

4.2

Comparison with other methods In the case when the factor $\sum_{l=0}^{2^h-1} \psi_{h,l}(x)^2$ does not depend on N (such as in the previous example), the bound (8) on the excess error (assessed in \mathbb{P} risk of CLSR states that the estimation \mathbb{P} terms of FN) of CLSR is $O(\log K / K)$. It is clear that whenever $N \leq K$ (which is the case of interest here), this is better than the ordinary LSR in the initial domain, whose estimation error is $O(N \log K / K)$. It is difficult to compare this result with LASSO (or the Dantzig selector that has similar properties [5]) for which an important aspect is to design sparse regression functions or to recover a solution assumed to be sparse. From [12, 15, 24] one deduces that under some assumptions, the estimation error of LASSO is of order $S \log KN$ where S is the sparsity (number of non-zero coefficients) of the \mathbb{P} best regressor f^* in FN . If $S \leq K$ then LASSO is more interesting than CLSR in terms of excess risk. Otherwise CLSR may be an interesting alternative although this method does not make any assumption about the sparsity of f^* and its goal is not to recover a possible sparse f^* but only to make good predictions. However, in some sense our method finds a sparse solution in the fact that the regression function g_L lies in a space GM of small dimension $M \ll N$ and can thus

be expressed using only M coefficients. Now in terms of numerical complexity, CLSR requires $O(N K^{3/2})$ operations to build the matrix and compute the regression function, whereas according to [18], the (heuristic) complexity of the LASSO algorithm is $O(N K^2)$ in the best cases (assuming that the number of steps required for convergence is $O(K)$, which is not proved theoretically). Thus CLSR seems to be a good and simple competitor to LASSO.

5

Conclusion

We considered the case when the number of features N is larger than the number of data K . The result stated in Theorem 1 enables to analyze the excess risk of any linear regression algorithm (LS or its penalized versions) performed in the compressed domain GM versus in the initial space FN . In the compressed domain the estimation error is reduced but an additional (controlled) approximation error (when compared to the best regressor in FN) comes into the picture. In the case of LS 2 has a mild dependency on N , then by choosing a regression, when the term $\|E(X) - E\|_2^2$ is small, CLSR has an estimation error (assessed in terms of FN) bounded by $O(\log K / K)$ and has numerical complexity $O(N K^{3/2})$. In short, CLSR provides an alternative to usual penalization techniques where one first selects a random subspace of lower dimension and then performs an empirical risk minimizer in this subspace. Further work needs to be done to provide additional settings (when the space X is of dimension ≥ 1) for which the term $\|E(X) - E\|_2^2$ is small. Acknowledgements: The authors wish to thank Laurent Jacques for numerous comments and Alessandro Lazaric and Mohammad Ghavamzadeh for exciting discussions. This work has been supported by French National Research Agency (ANR) through COSINUS program (project EXPLO-RA, ANR-08-COSI-004).

2 References

- [1] Dimitris Achlioptas. Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *Journal of Computer and System Sciences*, 66(4):671–687, June 2003.
- [2] Nir Ailon and Bernard Chazelle. Approximate nearest neighbors and the fast JohnsonLindenstrauss transform. In *STOC '06: Proceedings of the thirty-eighth annual ACM symposium on Theory of computing*, pages 557–563, New York, NY, USA, 2006. ACM.
- [3] Jean-Yves Audibert and Olivier Catoni. Risk bounds in linear regression through pac-bayesian truncation. Technical Report HAL : hal-00360268, 2009.
- [4] David Bau III and Lloyd N. Trefethen. *Numerical linear algebra*. Philadelphia: Society for Industrial and Applied Mathematics, 1997. 8
- [5] Peter J. Bickel, Ya'acov Ritov, and Alexandre B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. To appear in *Annals of Statistics*, 2008.
- [6] Avrim Blum. Random projection, margins, kernels, and feature-selection. *Subspace, Latent Structure and Feature Selection*, pages 52–68, 2006.
- [7] Robert Calderbank, Sina Jafarpour, and Robert Schapire. Compressed learn-

ing: Universal sparse dimensionality reduction and learning in the measurement domain. Technical Report, 2009. [8] Emmanuel Candes and Terence Tao. The Dantzig selector: Statistical estimation when p is much larger than n . *Annals of Statistics*, 35:2313, 2007. [9] Emmanuel J. Candes and Justin K. Romberg. Signal recovery from random projections. volume 5674, pages 76?86. SPIE, 2005. [10] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20:33?61, 1998. [11] Mark A. Davenport, Michael B. Wakin, and Richard G. Baraniuk. Detection and estimation with compressive measurements. Technical Report TREE 0610, Department of Electrical and Computer Engineering, Rice University, 2006. [12] E. Greenshtein and Y. Ritov. Persistency in high dimensional linear predictor-selection and the virtue of over-parametrization. *Bernoulli*, 10:971?988, 2004. [13] L. Gy?orfi, M. Kohler, A. Krzy?zak, and H. Walk. A distribution-free theory of nonparametric regression. Springer-Verlag, 2002. [14] Sham M. Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In Daphne Koller, Dale Schuurmans, Yoshua Bengio, and Leon Bottou, editors, *Neural Information Processing Systems*, pages 793? 800. MIT Press, 2008. [15] Yuval Nardi and Alessandro Rinaldo. On the asymptotic properties of the group Lasso estimator for linear models. *Electron. J. Statist.*, 2:605?633, 2008. [16] D. Pollard. *Convergence of Stochastic Processes*. Springer Verlag, New York, 1984. [17] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Neural Information Processing Systems*, 2007. [18] Saharon Rosset and Ji Zhu. Piecewise linear regularized solution paths. *Annals of Statistics*, 35:1012, 2007. [19] Robert Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58:267?288, 1994. [20] A. N. Tikhonov. Solution of incorrectly formulated problems and the regularization method. *Soviet Math Dokl* 4, pages 1035?1038, 1963. [21] Yaakov Tsaig and David L. Donoho. Compressed sensing. *IEEE Trans. Inform. Theory*, 52:1289?1306, 2006. [22] Vladimir N. Vapnik. *The nature of statistical learning theory*. Springer-Verlag New York, Inc., New York, NY, USA, 1995. [23] Tong Zhang. Covering number bounds of certain regularized linear function classes. *Journal of Machine Learning Research*, 2:527?550, 2002. [24] Tong Zhang. Some sharp performance bounds for least squares regression with L1 regularization. To appear in *Annals of Statistics*, 2009. [25] Shuheng Zhou, John D. Lafferty, and Larry A. Wasserman. Compressed regression. In John C. Platt, Daphne Koller, Yoram Singer, and Sam T. Roweis, editors, *Neural Information Processing Systems*. MIT Press, 2007.