

A tail inequality for quadratic forms of subgaussian random vectors

Daniel Hsu* Sham M. Kakade† Tong Zhang‡

Abstract

This article proves an exponential probability tail inequality for positive semidefinite quadratic forms in a subgaussian random vector. The bound is analogous to one that holds when the vector has independent Gaussian entries.

Keywords: Tail inequality; quadratic form; subgaussian random vectors; subgaussian chaos.

AMS MSC 2010: 60F10.

Submitted to ECP on June 11, 2012, final version accepted on October 29, 2012.

Supersedes arXiv:1110.2842.

1 Introduction

Suppose that $x = (x_1, \dots, x_n)$ is a random vector. Let $A \in \mathbb{R}^{n \times n}$ be a fixed matrix. A natural quantity that arises in many settings is the quadratic form $\|Ax\|^2 = x^\top (A^\top A)x$. Throughout $\|v\|$ denotes the Euclidean norm of a vector v , and $\|M\|$ denotes the spectral (operator) norm of a matrix M . We are interested in how close $\|Ax\|^2$ is to its expectation.

Consider the special case where x_1, \dots, x_n are independent standard Gaussian random variables. The following proposition provides an (upper) tail bound for $\|Ax\|^2$.

Proposition 1.1. *Let $A \in \mathbb{R}^{n \times n}$ be a matrix, and let $\Sigma := A^\top A$. Let $x = (x_1, \dots, x_n)$ be an isotropic multivariate Gaussian random vector with mean zero. For all $t > 0$,*

$$\Pr \left[\|Ax\|^2 > \text{tr}(\Sigma) + 2\sqrt{\text{tr}(\Sigma^2)t} + 2\|\Sigma\|t \right] \leq e^{-t}.$$

The proof, given in Appendix A.2, is straightforward given the rotational invariance of the multivariate Gaussian distribution, together with a tail bound for linear combinations of χ^2 random variables from [2]. We note that a slightly weaker form of Proposition 1.1 can be proved directly using Gaussian concentration [3].

In this note, we consider the case where $x = (x_1, \dots, x_n)$ is a *subgaussian* random vector. By this, we mean that there exists a $\sigma \geq 0$, such that for all $\alpha \in \mathbb{R}^n$,

$$\mathbb{E} [\exp(\alpha^\top x)] \leq \exp(\|\alpha\|^2 \sigma^2 / 2).$$

We provide a sharp upper tail bound for this case analogous to one that holds in the Gaussian case (indeed, the same as Proposition 1.1 when $\sigma = 1$).

*Microsoft Research New England, USA. E-mail: dahsu@microsoft.com

†Microsoft Research New England, USA. E-mail: skakade@microsoft.com

‡Department of Statistics, Rutgers University, USA. E-mail: tzhang@stat.rutgers.edu

Tail inequalities for sums of random vectors

One motivation for our main result comes from the following observations about sums of random vectors. Let a_1, \dots, a_n be vectors in a Euclidean space, and let $A = [a_1 | \dots | a_n]$ be the matrix with a_i as its i th column. Consider the squared norm of the random sum

$$\|Ax\|^2 = \left\| \sum_{i=1}^n a_i x_i \right\|^2 \tag{1.1}$$

where $x := (x_1, \dots, x_n)$ is a martingale difference sequence with $\mathbb{E}[x_i | x_1, \dots, x_{i-1}] = 0$ and $\mathbb{E}[x_i^2 | x_1, \dots, x_{i-1}] = \sigma^2$. Under mild boundedness assumptions on the x_i , the probability that the squared norm in (1.1) is much larger than its expectation

$$\mathbb{E}[\|Ax\|^2] = \sigma^2 \sum_{i=1}^n \|a_i\|^2 = \sigma^2 \operatorname{tr}(A^\top A)$$

falls off exponentially fast. This can be shown, for instance, using the following lemma by taking $u_i = a_i x_i$ (see Appendix A.1).

Proposition 1.2. *Let u_1, \dots, u_n be a martingale difference vector sequence, i.e.,*

$$\mathbb{E}[u_i | u_1, \dots, u_{i-1}] = 0, \quad \text{for all } i = 1, \dots, n,$$

such that

$$\sum_{i=1}^n \mathbb{E}[\|u_i\|^2 | u_1, \dots, u_{i-1}] \leq v \quad \text{and} \quad \|u_i\| \leq b$$

for all $i = 1, \dots, n$, almost surely. For all $t > 0$,

$$\Pr \left[\left\| \sum_{i=1}^n u_i \right\| > \sqrt{v} + \sqrt{8vt} + (4/3)bt \right] \leq e^{-t}.$$

After squaring the quantities in the stated probabilistic event, Proposition 1.2 gives the bound

$$\begin{aligned} \|Ax\|^2 \leq & \sigma^2 \cdot \operatorname{tr}(A^\top A) + \sigma^2 \cdot O \left(\operatorname{tr}(A^\top A)(\sqrt{t} + t) \right. \\ & \left. + \sqrt{\operatorname{tr}(A^\top A)} \max_i \|a_i\| (t + t^{3/2}) + \max_i \|a_i\|^2 t^2 \right) \end{aligned}$$

with probability at least $1 - e^{-t}$ when the x_i are almost surely bounded by 1 (or any constant).

Unfortunately, this bound obtained from Proposition 1.2 can be suboptimal when the x_i are subgaussian. For instance, if the x_i are Rademacher random variables, so $\Pr[x_i = +1] = \Pr[x_i = -1] = 1/2$, then it is known that

$$\|Ax\|^2 \leq \operatorname{tr}(A^\top A) + O \left(\sqrt{\operatorname{tr}((A^\top A)^2)} t + \|A\|^2 t \right) \tag{1.2}$$

with probability at least $1 - e^{-t}$. A similar result holds for any subgaussian distribution on the x_i [1]. This is an improvement over the previous bound because the deviation terms (i.e., those involving t) can be significantly smaller, especially for large t .

In this work, we give a simple proof of (1.2) with explicit constants that match the analogous bound when the x_i are independent standard Gaussian random variables.

2 Positive semidefinite quadratic forms

Our main theorem, given below, is a generalization of (1.2).

Theorem 2.1. *Let $A \in \mathbb{R}^{n \times n}$ be a matrix, and let $\Sigma := A^\top A$. Suppose that $x = (x_1, \dots, x_n)$ is a random vector such that, for some $\mu \in \mathbb{R}^n$ and $\sigma \geq 0$,*

$$\mathbb{E} [\exp(\alpha^\top (x - \mu))] \leq \exp(\|\alpha\|^2 \sigma^2 / 2) \quad (2.1)$$

for all $\alpha \in \mathbb{R}^n$. For all $t > 0$,

$$\Pr \left[\|Ax\|^2 > \sigma^2 \cdot \left(\text{tr}(\Sigma) + 2\sqrt{\text{tr}(\Sigma^2)t} + 2\|\Sigma\|t \right) + \text{tr}(\Sigma\mu\mu^\top) \cdot \left(1 + 2\left(\frac{\|\Sigma\|^2}{\text{tr}(\Sigma^2)} t \right)^{1/2} \right) \right] \leq e^{-t}.$$

Remark 2.2. *If $\mu = 0$, then the assumption (2.1) implies $\mathbb{E}[x] = 0$ and $\text{cov}(x) \preceq \sigma^2 I$. In this case,*

$$\mathbb{E}[\|Ax\|^2] = \text{tr}(\Sigma \text{cov}(x)) \leq \sigma^2 \text{tr}(\Sigma), \quad \text{var}(\|Ax\|^2) = O(\sigma^4 \text{tr}(\Sigma^2)),$$

so probability inequality may be interpreted as a Bernstein inequality. If $\mu = 0$ and $\sigma = 1$, then the probability inequality reads

$$\Pr \left[\|Ax\|^2 > \text{tr}(\Sigma) + 2\sqrt{\text{tr}(\Sigma^2)t} + 2\|\Sigma\|t \right] \leq e^{-t},$$

which is the same as Proposition 1.1.

Remark 2.3. *Our proof (via (2.2), (2.4), and (2.5)) actually establishes the following upper bounds on the moment generating function of $\|Ax\|^2$ for $0 \leq \eta < 1/(2\sigma^2\|\Sigma\|)$:*

$$\begin{aligned} \mathbb{E} [\exp(\eta\|Ax\|^2)] &\leq \mathbb{E} \left[\exp \left(\sigma^2 \|A^\top z\|^2 \eta + \mu^\top A^\top z \sqrt{2\eta} \right) \right] \\ &\leq \exp \left(\sigma^2 \text{tr}(\Sigma)\eta + \frac{\sigma^4 \text{tr}(\Sigma^2)\eta^2 + \|A\mu\|^2 \eta}{1 - 2\sigma^2\|\Sigma\|\eta} \right) \end{aligned}$$

where z is a vector of n independent standard Gaussian random variables.

Proof of Theorem 2.1. Let z be a vector of n independent standard Gaussian random variables (sampled independently of x). For any $\alpha \in \mathbb{R}^n$,

$$\mathbb{E} [\exp(z^\top \alpha)] = \exp(\|\alpha\|^2 / 2). \quad (2.2)$$

Thus, for any $\lambda \in \mathbb{R}$ and $\varepsilon \geq 0$, we have the following decoupling (which holds, in fact, for any random vector x):

$$\begin{aligned} \mathbb{E} [\exp(\lambda z^\top Ax)] &\geq \mathbb{E} \left[\exp(\lambda z^\top Ax) \mid \|Ax\|^2 > \varepsilon \right] \cdot \Pr [\|Ax\|^2 > \varepsilon] \\ &\geq \exp \left(\frac{\lambda^2 \varepsilon}{2} \right) \cdot \Pr [\|Ax\|^2 > \varepsilon]. \end{aligned} \quad (2.3)$$

Moreover, using (2.1),

$$\begin{aligned} \mathbb{E} [\exp(\lambda z^\top Ax)] &= \mathbb{E} \left[\mathbb{E} \left[\exp(\lambda z^\top A(x - \mu)) \mid z \right] \exp(\lambda z^\top A\mu) \right] \\ &\leq \mathbb{E} \left[\exp \left(\frac{\lambda^2 \sigma^2}{2} \|A^\top z\|^2 + \lambda \mu^\top A^\top z \right) \right]. \end{aligned} \quad (2.4)$$

A tail inequality for quadratic forms of subgaussian random vectors

Let USV^\top be a singular value decomposition of A ; where U and V are, respectively, matrices of orthonormal left and right singular vectors; and $S = \text{diag}(\sqrt{\rho_1}, \dots, \sqrt{\rho_m})$ is the diagonal matrix of corresponding singular values. Note that

$$\|\rho\|_1 = \sum_{i=1}^n \rho_i = \text{tr}(\Sigma), \quad \|\rho\|_2^2 = \sum_{i=1}^n \rho_i^2 = \text{tr}(\Sigma^2), \quad \text{and} \quad \|\rho\|_\infty = \max_i \rho_i = \|\Sigma\|.$$

By rotational invariance, $y := U^\top z$ is an isotropic multivariate Gaussian random vector with mean zero. Therefore $\|A^\top z\|^2 = z^\top US^2U^\top z = \rho_1 y_1^2 + \dots + \rho_n y_n^2$ and $\mu^\top A^\top z = \nu^\top y = \nu_1 y_1 + \dots + \nu_n y_n$, where $\nu := SV^\top \mu$ (note that $\|\nu\|^2 = \|SV^\top \mu\|^2 = \|A\mu\|^2$). Let $\gamma := \lambda^2 \sigma^2 / 2$. By Lemma 2.4,

$$\mathbb{E} \left[\exp \left(\gamma \sum_{i=1}^n \rho_i y_i^2 + \frac{\sqrt{2\gamma}}{\sigma} \sum_{i=1}^n \nu_i y_i \right) \right] \leq \exp \left(\|\rho\|_1 \gamma + \frac{\|\rho\|_2^2 \gamma^2 + \|\nu\|^2 \gamma / \sigma^2}{1 - 2\|\rho\|_\infty \gamma} \right) \quad (2.5)$$

for $0 \leq \gamma < 1/(2\|\rho\|_\infty)$. Combining (2.3), (2.4), and (2.5) gives

$$\Pr [\|Ax\|^2 > \varepsilon] \leq \exp \left(-\varepsilon \gamma / \sigma^2 + \|\rho\|_1 \gamma + \frac{\|\rho\|_2^2 \gamma^2 + \|\nu\|^2 \gamma / \sigma^2}{1 - 2\|\rho\|_\infty \gamma} \right)$$

for $0 \leq \gamma < 1/(2\|\rho\|_\infty)$ and $\varepsilon \geq 0$. Choosing

$$\varepsilon := \sigma^2 (\|\rho\|_1 + \tau) + \|\nu\|^2 \sqrt{1 + \frac{2\|\rho\|_\infty \tau}{\|\rho\|_2^2}} \quad \text{and} \quad \gamma := \frac{1}{2\|\rho\|_\infty} \left(1 - \sqrt{\frac{\|\rho\|_2^2}{\|\rho\|_2^2 + 2\|\rho\|_\infty \tau}} \right),$$

we have

$$\begin{aligned} & \Pr \left[\|Ax\|^2 > \sigma^2 (\|\rho\|_1 + \tau) + \|\nu\|^2 \sqrt{1 + \frac{2\|\rho\|_\infty \tau}{\|\rho\|_2^2}} \right] \\ & \leq \exp \left(-\frac{\|\rho\|_2^2}{2\|\rho\|_\infty^2} \left(1 + \frac{\|\rho\|_\infty \tau}{\|\rho\|_2^2} - \sqrt{1 + \frac{2\|\rho\|_\infty \tau}{\|\rho\|_2^2}} \right) \right) = \exp \left(-\frac{\|\rho\|_2^2}{2\|\rho\|_\infty^2} h_1 \left(\frac{\|\rho\|_\infty \tau}{\|\rho\|_2^2} \right) \right) \end{aligned}$$

where $h_1(a) := 1 + a - \sqrt{1 + 2a}$, which has the inverse function $h_1^{-1}(b) = \sqrt{2b} + b$. The result follows by setting $\tau := 2\sqrt{\|\rho\|_2^2 t + 2\|\rho\|_\infty t} = 2\sqrt{\text{tr}(\Sigma^2)t + 2\|\Sigma\|t}$. \square

The following lemma is a standard estimate of the logarithmic moment generating function of a quadratic form in standard Gaussian random variables, proved much along the lines of the estimate from [2].

Lemma 2.4. *Let z be a vector of n independent standard Gaussian random variables. Fix any non-negative vector $\alpha \in \mathbb{R}_+^n$ and any vector $\beta \in \mathbb{R}^n$. If $0 \leq \lambda < 1/(2\|\alpha\|_\infty)$, then*

$$\log \mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^n \alpha_i z_i^2 + \sum_{i=1}^n \beta_i z_i \right) \right] \leq \|\alpha\|_1 \lambda + \frac{\|\alpha\|_2^2 \lambda^2 + \|\beta\|_2^2 / 2}{1 - 2\|\alpha\|_\infty \lambda}.$$

Proof. Fix $\lambda \in \mathbb{R}$ such that $0 \leq \lambda < 1/(2\|\alpha\|_\infty)$, and let $\eta_i := 1/\sqrt{1 - 2\alpha_i \lambda} > 0$ for $i = 1, \dots, n$. We have

$$\begin{aligned} \mathbb{E} [\exp(\lambda \alpha_i z_i^2 + \beta_i z_i)] &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp(-z_i^2/2) \exp(\lambda \alpha_i z_i^2 + \beta_i z_i) dz_i \\ &= \eta_i \exp\left(\frac{\beta_i^2 \eta_i^2}{2}\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi \eta_i^2}} \exp\left(-\frac{1}{2\eta_i^2} (z_i - \beta_i \eta_i^2)^2\right) dz_i \end{aligned}$$

so

$$\log \mathbb{E} \left[\exp \left(\lambda \sum_{i=1}^n \alpha_i z_i^2 + \sum_{i=1}^n \beta_i z_i \right) \right] = \frac{1}{2} \sum_{i=1}^n \beta_i^2 \eta_i^2 + \frac{1}{2} \sum_{i=1}^n \log \eta_i^2.$$

The right-hand side can be bounded using the inequalities

$$\frac{1}{2} \sum_{i=1}^n \log \eta_i^2 = -\frac{1}{2} \sum_{i=1}^n \log(1 - 2\alpha_i \lambda) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^{\infty} \frac{(2\alpha_i \lambda)^j}{j} \leq \|\alpha\|_1 \lambda + \frac{\|\alpha\|_2^2 \lambda^2}{1 - 2\|\alpha\|_{\infty} \lambda}$$

and

$$\frac{1}{2} \sum_{i=1}^n \beta_i^2 \eta_i^2 \leq \frac{\|\beta\|_2^2 / 2}{1 - 2\|\alpha\|_{\infty} \lambda}. \quad \square$$

Example: fixed-design regression with subgaussian noise

We give a simple application of Theorem 2.1 to fixed-design linear regression with the ordinary least squares estimator.

Let x_1, \dots, x_n be fixed design vectors in \mathbb{R}^d . Let the responses y_1, \dots, y_n be random variables for which there exists $\sigma > 0$ such that

$$\mathbb{E} \left[\exp \left(\sum_{i=1}^n \alpha_i (y_i - \mathbb{E}[y_i]) \right) \right] \leq \exp \left(\sigma^2 \sum_{i=1}^n \alpha_i^2 \right)$$

for any $\alpha_1, \dots, \alpha_n \in \mathbb{R}$. This condition is satisfied, for instance, if

$$y_i = \mathbb{E}[y_i] + \varepsilon_i$$

for independent subgaussian zero-mean noise variables $\varepsilon_1, \dots, \varepsilon_n$. Let $\Sigma := \sum_{i=1}^n x_i x_i^\top / n$, which we assume is invertible without loss of generality. Let

$$\beta := \Sigma^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i \mathbb{E}[y_i] \right)$$

be the coefficient vector of minimum expected squared error (i.e., $\mathbb{E}[n^{-1} \sum_{i=1}^n (x_i^\top \beta - y_i)^2] = \min!$). The ordinary least squares estimator is given by

$$\hat{\beta} := \Sigma^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right).$$

The excess loss $R(\hat{\beta})$ of $\hat{\beta}$ is the difference between the expected squared error of $\hat{\beta}$ and that of β :

$$R(\hat{\beta}) := \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (x_i^\top \hat{\beta} - y_i)^2 \right] - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (x_i^\top \beta - y_i)^2 \right].$$

It is easy to see that

$$R(\hat{\beta}) = \|\Sigma^{1/2}(\hat{\beta} - \beta)\|^2 = \left\| \sum_{i=1}^n (\Sigma^{-1/2} x_i) (y_i - \mathbb{E}[y_i]) \right\|^2.$$

By Theorem 2.1,

$$\Pr \left[R(\hat{\beta}) > \frac{\sigma^2 (d + 2\sqrt{dt} + 2t)}{n} \right] \leq e^{-t}.$$

Note that in the case that $\mathbb{E}[(y_i - \mathbb{E}[y_i])^2] = \sigma^2$ for each i , then

$$\mathbb{E}[R(\hat{\beta})] = \frac{\sigma^2 d}{n};$$

so the tail inequality above is essentially tight when the y_i are independent Gaussian random variables.

A Standard tail inequalities

A.1 Martingale tail inequalities

The following is a standard form of Bernstein's inequality stated for martingale difference sequences.

Lemma A.1 (Bernstein's inequality for martingales). *Let d_1, \dots, d_n be a martingale difference sequence with respect to random variables x_1, \dots, x_n (i.e., $\mathbb{E}[d_i | x_1, \dots, x_{i-1}] = 0$ for all $i = 1, \dots, n$) such that $|d_i| \leq b$ and $\sum_{i=1}^n \mathbb{E}[d_i^2 | x_1, \dots, x_{i-1}] \leq v$. For all $t > 0$,*

$$\Pr \left[\sum_{i=1}^n d_i > \sqrt{2vt} + (2/3)bt \right] \leq e^{-t}.$$

Proposition 1.2 is an immediate consequence of the following folklore results, together with Jensen's inequality. Lemma A.2 is a straightforward application of Bernstein's inequality to a Doob martingale, and Lemma A.3 is proved by a simple induction argument.

Lemma A.2. *Let u_1, \dots, u_n be random vectors such that $\sum_{i=1}^n \mathbb{E}[\|u_i\|^2 | u_1, \dots, u_{i-1}] \leq v$ and $\|u_i\| \leq b$ for all $i = 1, \dots, n$, almost surely. For all $t > 0$,*

$$\Pr \left[\left\| \sum_{i=1}^n u_i \right\| - \mathbb{E} \left[\left\| \sum_{i=1}^n u_i \right\| \right] > \sqrt{8vt} + (4/3)bt \right] \leq e^{-t}.$$

Lemma A.3. *If u_1, \dots, u_n is a martingale difference vector sequence (c.f. Proposition 1.2), then $\mathbb{E}[\| \sum_{i=1}^n u_i \|^2] = \sum_{i=1}^n \mathbb{E}[\|u_i\|^2]$.*

A.2 Gaussian quadratic forms and χ^2 tail inequalities

It is well-known that if $z \sim \mathcal{N}(0, 1)$ is a standard Gaussian random variable, then z^2 follows a χ^2 distribution with one degree of freedom. The following inequality from [2] gives a bound on linear combinations of χ^2 random variables.

Lemma A.4 (χ^2 tail inequality; [2]). *Let q_1, \dots, q_n be independent χ^2 random variables, each with one degree of freedom. For any vector $\gamma = (\gamma_1, \dots, \gamma_n) \in \mathbb{R}_+^n$ with non-negative entries, and any $t > 0$,*

$$\Pr \left[\sum_{i=1}^n \gamma_i q_i > \|\gamma\|_1 + 2\sqrt{\|\gamma\|_2^2 t} + 2\|\gamma\|_\infty t \right] \leq e^{-t}.$$

Proof of Proposition 1.1. Let $V\Lambda V^\top$ be an eigen-decomposition of $A^\top A$, where V is a matrix of orthonormal eigenvectors, and $\Lambda := \text{diag}(\rho_1, \dots, \rho_n)$ is the diagonal matrix of corresponding eigenvalues ρ_1, \dots, ρ_n . By the rotational invariance of the distribution, $z := V^\top x$ is an isotropic multivariate Gaussian random vector with mean zero. Thus, $\|Ax\|^2 = z^\top \Lambda z = \rho_1 z_1^2 + \dots + \rho_n z_n^2$, and the z_i^2 are independent χ^2 random variables, each with one degree of freedom. The claim now follows from a tail bound for χ^2 random variables (Lemma A.4). \square

References

- [1] D. L. Hanson and F. T. Wright, *A bound on tail probabilities for quadratic forms in independent random variables*, The Annals of Math. Stat. **42** (1971), no. 3, 1079–1083. MR-0279864
- [2] B. Laurent and P. Massart, *Adaptive estimation of a quadratic functional by model selection*, The Annals of Statistics **28** (2000), no. 5, 1302–1338. MR-1805785
- [3] G. Pisier, *The volume of convex bodies and banach space geometry*, Cambridge University Press, 1989. MR-1036275

Acknowledgments. We thank the anonymous reviewers for their helpful comments.