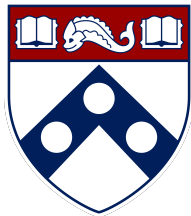


Phase Transition and Convex Geometry



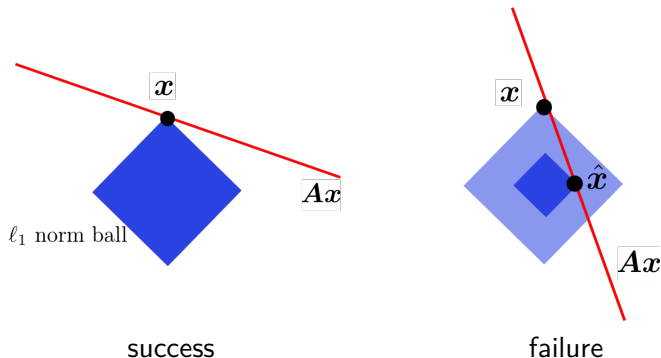
Yuxin Chen

Wharton Statistics & Data Science, Spring 2022

ℓ_1 minimization for sparse signal recovery

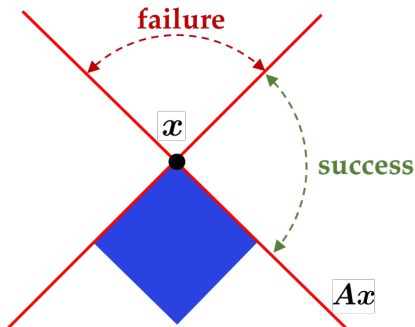
$$\begin{aligned} & \text{minimize}_{\mathbf{x} \in \mathbb{R}^p} && \|\mathbf{x}\|_1 \\ & \text{s.t.} && \mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^n \end{aligned}$$

What is the probability that ℓ_1 minimization succeeds in recovering \mathbf{x} ?



Probability of success of ℓ_1 minimization

Suppose \mathbf{A} is randomly generated



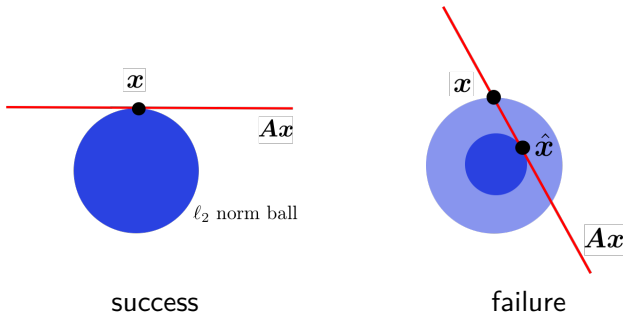
If $\mathbf{x} = [0, 1]^\top$ and $\mathbf{A} = [a_1, a_2] \in \mathbb{R}^{1 \times 2}$ (i.i.d. Gaussian), then

$$\mathbb{P}\{\ell_1\text{-min succeeds}\} = \frac{1}{2} \quad (\text{proved by figure})$$

How about ℓ_2 minimization?

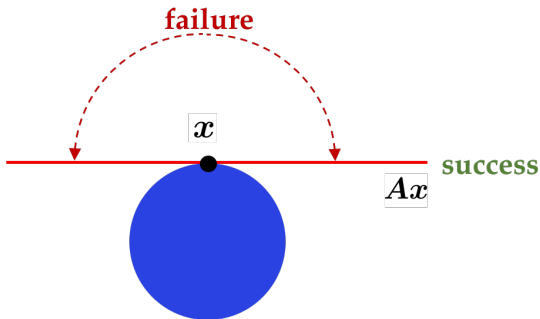
Suppose A is randomly generated (e.g. i.i.d. Gaussian), and consider

$$\begin{aligned} \text{minimize}_{\mathbf{x} \in \mathbb{R}^p} \quad & \|\mathbf{x}\|_2 \\ \text{s.t.} \quad & \mathbf{y} = A\mathbf{x} \in \mathbb{R}^n \end{aligned}$$



Probability of success of ℓ_2 minimization

Suppose \mathbf{A} is randomly generated



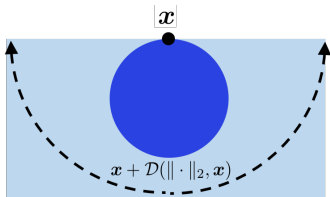
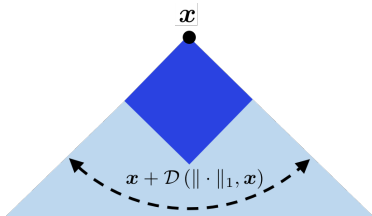
If $\mathbf{x} = [0, 1]^\top$ and $\mathbf{A} = [a_1, a_2] \in \mathbb{R}^{1 \times 2}$ (i.i.d. Gaussian), then

$$\mathbb{P}\{\ell_2\text{-min succeeds}\} = 0 \quad (\text{proved by figure})$$

Key metric: volume of descent cone

Suppose A is randomly generated, and consider

$$\begin{aligned} \text{minimize}_{\mathbf{x} \in \mathbb{R}^p} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{y} = A\mathbf{x} \in \mathbb{R}^n \end{aligned} \tag{10.1}$$



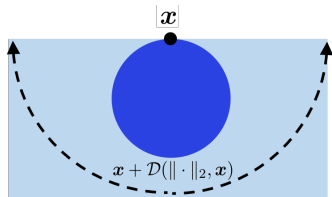
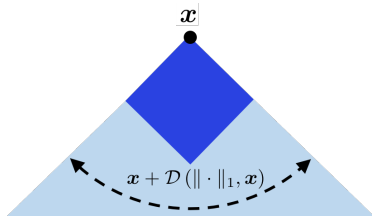
The success probability of (10.1) depends on the volume of the **descent cone**

$$\mathcal{D}(f, \mathbf{x}) := \{\mathbf{h} : \exists \epsilon > 0 \text{ s.t. } f(\mathbf{x} + \epsilon\mathbf{h}) \leq f(\mathbf{x})\}$$

Key metric: volume of descent cone

Suppose A is randomly generated, and consider

$$\begin{aligned} \text{minimize}_{\mathbf{x} \in \mathbb{R}^p} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{y} = A\mathbf{x} \in \mathbb{R}^n \end{aligned} \tag{10.1}$$



We need to compute the probability of 2 convex cones sharing a ray:

$$\mathbb{P}\left\{ (10.4) \text{ succeeds} \right\} = \mathbb{P}\left\{ \mathcal{D}(f, \mathbf{x}) \cap \{\mathbf{h} : A\mathbf{h} = \mathbf{0}\} = \{\mathbf{0}\} \right\}$$

Kinematic formula

Lemma 10.1 (Theorem 6.5.6, Schneider & Wolfgang '08)

Let $\mathcal{C}, \mathcal{K} \subseteq \mathbb{R}^d$ be convex cones, and Q a random orthogonal basis:

$$\mathbb{P}\{\mathcal{C} \cap Q\mathcal{K} \neq \{0\}\} = \sum_{i=0}^d \left(1 + (-1)^{i+1}\right) \sum_{j=i}^d \nu_i(\mathcal{C}) \nu_{d+i-j}(\mathcal{K}),$$

where $\nu_k \geq 0$ is called the k th intrinsic volume.

- Exact but not workable formula!
- Calls for a simpler expression

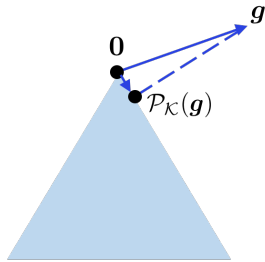
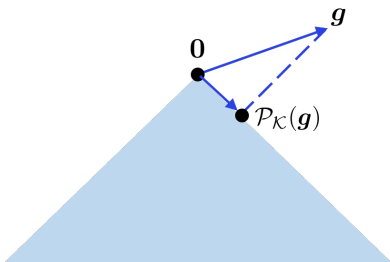
Statistical dimension and Gaussian width

Definition 10.2 (Statistical dimension)

For any **convex** cone \mathcal{K} , its statistical dimension is defined as

$$\text{stat-dim}(\mathcal{K}) := \mathbb{E}[\|\mathcal{P}_{\mathcal{K}}(\mathbf{g})\|_2^2]$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$; $\mathcal{P}_{\mathcal{K}}(\mathbf{g}) := \arg \min_{\mathbf{z} \in \mathcal{K}} \|\mathbf{g} - \mathbf{z}\|_2$: Euclidean projection



Statistical dimension and Gaussian width

Definition 10.2 (Statistical dimension)

For any **convex** cone \mathcal{K} , its statistical dimension is defined as

$$\text{stat-dim}(\mathcal{K}) := \mathbb{E}[\|\mathcal{P}_{\mathcal{K}}(\mathbf{g})\|_2^2]$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$; $\mathcal{P}_{\mathcal{K}}(\mathbf{g}) := \arg \min_{\mathbf{z} \in \mathcal{K}} \|\mathbf{g} - \mathbf{z}\|_2$: Euclidean projection

- If \mathcal{K} is k -dimensional subspace, then

$$\text{stat-dim}(\mathcal{K}) = k \quad (\text{so it is indeed a measure of "dimension"})$$

Statistical dimension and Gaussian width

Definition 10.2 (Statistical dimension)

For any **convex** cone \mathcal{K} , its statistical dimension is defined as

$$\text{stat-dim}(\mathcal{K}) := \mathbb{E}[\|\mathcal{P}_{\mathcal{K}}(\mathbf{g})\|_2^2]$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$; $\mathcal{P}_{\mathcal{K}}(\mathbf{g}) := \arg \min_{\mathbf{z} \in \mathcal{K}} \|\mathbf{g} - \mathbf{z}\|_2$: Euclidean projection

- A related definition: Gaussian width

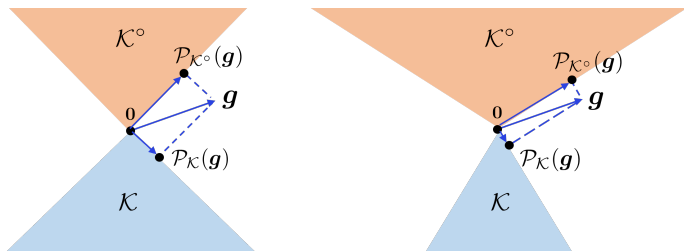
$$w(\mathcal{K}) := \mathbb{E} \left[\sup_{\mathbf{z} \in \mathcal{K}, \|\mathbf{z}\|_2=1} \langle \mathbf{z}, \mathbf{g} \rangle \right]$$

- (Homework) $w^2(\mathcal{K}) \leq \text{stat-dim}(\mathcal{K}) \leq w^2(\mathcal{K}) + 1$

Polar cone

The polar cone of a convex cone \mathcal{K} is defined as

$$\mathcal{K}^\circ = \{\mathbf{y} \mid \mathbf{y}^\top \mathbf{x} \leq 0, \quad \forall \mathbf{x} \in \mathcal{K}\}$$



- Moreau's decomposition:

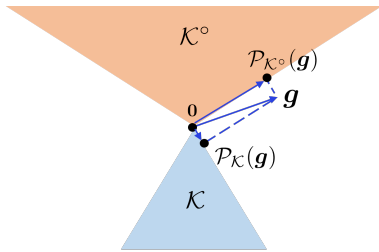
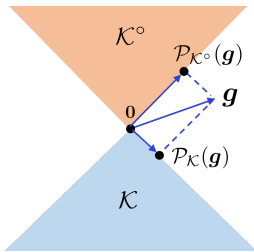
$$\mathbf{g} = \mathcal{P}_{\mathcal{K}}(\mathbf{g}) + \mathcal{P}_{\mathcal{K}^\circ}(\mathbf{g})$$

where $\langle \mathcal{P}_{\mathcal{K}}(\mathbf{g}), \mathcal{P}_{\mathcal{K}^\circ}(\mathbf{g}) \rangle = 0$

Polar cone

The polar cone of a convex cone \mathcal{K} is defined as

$$\mathcal{K}^\circ = \{\mathbf{y} \mid \mathbf{y}^\top \mathbf{x} \leq 0, \quad \forall \mathbf{x} \in \mathcal{K}\}$$



- \mathcal{K} is called *self-dual* if

$$\mathcal{K} = \underbrace{-\mathcal{K}^\circ}_{\text{dual cone}}$$

- If \mathcal{K} is self-dual in \mathbb{R}^d , then

$$\text{stat-dim}(\mathcal{K}) = d/2 \quad (\text{by Moreau's decomposition})$$

Examples

- Nonnegative orthant \mathbb{R}_+^d (self-dual)

$$\text{stat-dim}(\mathbb{R}_+^d) = \frac{1}{2}d$$

- Second-order cone $\mathbb{L}^d := \left\{ \begin{bmatrix} \mathbf{x} \\ \tau \end{bmatrix} \in \mathbb{R}^{d+1} : \|\mathbf{x}\|_2 \leq \tau \right\}$
(self-dual)

$$\text{stat-dim}(\mathbb{L}^d) = \frac{1}{2}(d+1)$$

- Set of symmetric positive semidefinite matrices $\mathbb{S}_+^{d \times d}$ (self-dual)

$$\text{stat-dim}(\mathbb{S}_+^{d \times d}) = \frac{1}{2} \cdot \frac{d(d+1)}{2}$$

Second-order cone is self-dual

The dual cone of \mathbb{L}^d is

$$\begin{aligned} \mathcal{C} &= \left\{ \begin{bmatrix} \mathbf{y} \\ \alpha \end{bmatrix} \mid 0 \leq \mathbf{y}^\top \mathbf{x} + \alpha\tau, \forall \begin{bmatrix} \mathbf{x} \\ \tau \end{bmatrix} \in \mathbb{L}^d \right\} \\ &= \left\{ \begin{bmatrix} \mathbf{y} \\ \alpha \end{bmatrix} \mid 0 \leq \inf_{\tau \geq 0} \inf_{\mathbf{x}: \|\mathbf{x}\|_2 \leq \tau} (\mathbf{y}^\top \mathbf{x} + \alpha\tau) \right\} \\ &= \left\{ \begin{bmatrix} \mathbf{y} \\ \alpha \end{bmatrix} \mid 0 \leq \inf_{\tau \geq 0} \inf_{\mathbf{x}: \|\mathbf{x}\|_2 \leq \tau} (-\|\mathbf{y}\|_2 \|\mathbf{x}\|_2 + \alpha\tau) \right\} \\ &= \left\{ \begin{bmatrix} \mathbf{y} \\ \alpha \end{bmatrix} \mid 0 \leq \inf_{\tau \geq 0} (\alpha - \|\mathbf{y}\|_2) \tau \right\} \end{aligned}$$

Since $\tau \geq 0$, one has $\begin{bmatrix} \mathbf{y} \\ \alpha \end{bmatrix} \in \mathcal{C}$ iff $\|\mathbf{y}\|_2 \leq \alpha$ and hence $\mathcal{C} = \mathbb{L}^d$

Polarity relation

The statistical dimension can be expressed in terms of the polar cone:

$$\text{stat-dim}(\mathcal{K}) = \mathbb{E} \left[\text{dist}^2(\mathbf{g}, \mathcal{K}^\circ) \right] := \mathbb{E} \left[\inf_{\mathbf{z} \in \mathcal{K}^\circ} \|\mathbf{g} - \mathbf{z}\|_2^2 \right]$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$

- A direct consequence:

$$\text{dist}(\mathbf{g}, \mathcal{K}) = \|\mathbf{g} - \mathcal{P}_{\mathcal{K}}(\mathbf{g})\|_2 \stackrel{\text{Moreau decomposition}}{=} \|\mathcal{P}_{\mathcal{K}^\circ}(\mathbf{g})\|_2$$

Approximate kinematic formula

$\mathcal{C}, \mathcal{K} \in \mathbb{R}^d$: convex cones; $Q \in \mathbb{R}^{d \times d}$: random orthogonal basis

Theorem 10.3 (Amelunxen, Lotz, McCoy & Tropp '13)

$$\text{stat-dim}(\mathcal{C}) + \text{stat-dim}(\mathcal{K}) \leq d - \Theta(\sqrt{d \log d})$$

$$\implies Q\mathcal{C} \cap \mathcal{K} = \{\mathbf{0}\} \quad \text{with high prob.}$$

$$\text{stat-dim}(\mathcal{C}) + \text{stat-dim}(\mathcal{K}) \geq d + \Theta(\sqrt{d \log d})$$

$$\implies Q\mathcal{C} \cap \mathcal{K} \neq \{\mathbf{0}\} \quad \text{with high prob.}$$

- 2 randomly rotated cones share a ray \iff their aggregate statistical dimension exceeds ambient dimension
- Sharp concentration: the fluctuation does not exceed $O(\sqrt{d \log d})$

Example: logistic regression

Suppose we obtain n independent *binary* samples:

$$y_i = \begin{cases} 1, & \text{with prob. } \frac{1}{1+\exp(-\mathbf{a}_i^\top \mathbf{x})} \\ -1, & \text{with prob. } \frac{1}{1+\exp(\mathbf{a}_i^\top \mathbf{x})} \end{cases} \quad 1 \leq i \leq n$$

where $\{\mathbf{a}_i \in \mathbb{R}^p\}$: known design vectors; $\mathbf{x} \in \mathbb{R}^p$: unknown signal

- the likelihood for each y_i :

$$\begin{aligned} \mathcal{L}(\mathbf{x}; y_i) &= \frac{1}{1 + \exp(-\mathbf{a}_i^\top \mathbf{x})} \mathbb{1}\{y_i = 1\} + \frac{1}{1 + \exp(\mathbf{a}_i^\top \mathbf{x})} \mathbb{1}\{y_i = -1\} \\ &= \frac{1}{1 + \exp(-y_i \mathbf{a}_i^\top \mathbf{x})} \end{aligned}$$

Example: logistic regression

Maximum likelihood estimation (logistic regression)

$$\text{minimize}_{\mathbf{x}} - \sum_{i=1}^n \log \mathcal{L}(\mathbf{x}; y_i) = \sum_{i=1}^n \log \left\{ 1 + \exp(-y_i \mathbf{a}_i^\top \mathbf{x}) \right\}$$

- Consider a simple case

true signal $\mathbf{x} = \mathbf{0}$ (global null); $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$

- We'd naturally hope the MLE $\hat{\mathbf{x}}$ to be small (since $\mathbf{x} = \mathbf{0}$)
- **Question:** is $\|\hat{\mathbf{x}}\|_2$ always small under the global null (i.e. $\mathbf{x} = \mathbf{0}$)?

Example: logistic regression

Fact 10.4 (Cover '65; Sur, Chen, Candes '17)

Suppose $x = \mathbf{0}$. If $p > n/2 - \Theta(\sqrt{n \log n})$, then $\|\hat{x}\|_2 = \infty$ w.h.p.

- $n = 2p$ is indeed a sharp boundary in the sense that $\|\hat{x}\|_2 \lesssim 1$ if $n/p > 2$ (Sur, Chen, Candes '17)

Proof of Fact 10.4

Note that $\log \{1 + \exp(-y_i \mathbf{a}_i^\top \mathbf{x})\} \geq 0$. Thus, if $\exists \hat{\mathbf{x}}$ s.t.

$$y_i \mathbf{a}_i^\top \hat{\mathbf{x}} = +\infty, \quad 1 \leq i \leq n, \quad (10.2)$$

then $\log \{1 + \exp(-y_i \mathbf{a}_i^\top \mathbf{x})\} = 0$ for all i , and hence $\hat{\mathbf{x}}$ must be the MLE. In this case, $\|\hat{\mathbf{x}}\|_2 = \infty$

Proof of Fact 10.4 (cont.)

It remains to check when we can find $\hat{\mathbf{x}}$ obeying (10.2), or equivalently, when we have

$$\underbrace{\left\{ \mathbf{u} \mid u_i = y_i \mathbf{a}_i^\top \mathbf{x}, \mathbf{x} \in \mathbb{R}^p \right\}}_{\text{QC: } p\text{-dimensional}} \cap \underbrace{\mathbb{R}_+^n}_{\mathcal{K}} \neq \{0\} \quad (10.3)$$

Note that \mathbf{y} is independent of \mathbf{A} when $\mathbf{x} = \mathbf{0}$. By Theorem 10.3, if

$$p + \underbrace{\text{stat-dim}(\mathbb{R}_+^n)}_{n/2} > n + \Theta(\sqrt{n \log n}) \quad (\text{or } p > n/2 + \Theta(\sqrt{n \log n})),$$

then (10.3) holds. This establishes the claim that $\hat{\mathbf{x}}$ is unbounded.

Phase transition for inverse problems

Suppose $\mathbf{A} \in \mathbb{R}^{n \times p}$ is i.i.d. Gaussian, and consider

$$\begin{aligned} & \text{minimize}_{\mathbf{x} \in \mathbb{R}^p} && f(\mathbf{x}) \\ & \text{s.t.} && \mathbf{y} = \mathbf{A}\mathbf{x} \end{aligned} \tag{10.4}$$

Key: convex geometry

(10.4) succeeds



$$\{\mathbf{h} : \mathbf{A}\mathbf{h} = \mathbf{0}\} \cap \mathcal{D}(f, \mathbf{x}) = \{\mathbf{0}\}$$



$$\underbrace{\text{stat-dim}(\{\mathbf{h} : \mathbf{A}\mathbf{h} = \mathbf{0}\})}_{= p-n} + \text{stat-dim}(\mathcal{D}(f, \mathbf{x})) < p \quad (\text{by Theorem 10.3})$$

Phase transition for inverse problems

Suppose $A \in \mathbb{R}^{n \times p}$ is i.i.d. Gaussian, and consider

$$\begin{aligned} \text{minimize}_{\mathbf{x} \in \mathbb{R}^p} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{y} = A\mathbf{x} \end{aligned} \tag{10.4}$$

Theorem 10.5 (Amelunxen, Lotz, McCoy & Tropp '13)

$$\begin{aligned} n > \text{stat-dim}(\mathcal{D}(f, \mathbf{x})) + \Theta(\sqrt{p \log p}) \\ \implies \quad & (10.4) \text{ succeeds with high prob.} \end{aligned}$$

$$\begin{aligned} n < \text{stat-dim}(\mathcal{D}(f, \mathbf{x})) - \Theta(\sqrt{p \log p}) \\ \implies \quad & (10.4) \text{ fails with high prob.} \end{aligned}$$

Computing statistical dimension of descent cones?

Example: the decent cone w.r.t. ℓ_∞ norm

$$\text{stat-dim}(\mathcal{D}(\|\cdot\|_\infty, \mathbf{x})) = d - s/2$$

where $s = \#\{i : |x_i| = \|\mathbf{x}\|_\infty\}$

Proof: WLOG, suppose $\mathbf{x} = [1, \dots, 1, x_{s+1}, \dots, x_d]^\top$ with $1 > x_{s+1} \geq \dots \geq x_d \geq 0$. Then

$$\mathcal{D}(\|\cdot\|_\infty, \mathbf{x}) = (\mathbb{R}_-^s) \times \mathbb{R}^{d-s}$$

$$\implies \text{stat-dim}(\mathcal{D}(\|\cdot\|_\infty, \mathbf{x})) = \frac{1}{2}s + d - s = d - \frac{1}{2}s$$

Computing statistical dimension of descent cone?

In general, there is a duality between descent cone and subdifferentials
set of subgradients

$$(\mathcal{D}(f, \mathbf{x}))^\circ = \text{cone}(\partial f(\mathbf{x})) := \bigcup_{\tau \geq 0} \tau \partial f(\mathbf{x})$$

$$\Rightarrow \text{stat-dim}(\mathcal{D}(f, \mathbf{x})) = \mathbb{E} \left[\inf_{\tau \geq 0} \min_{\mathbf{u} \in \partial f(\mathbf{x})} \|\mathbf{g} - \tau \mathbf{u}\|_2^2 \right]$$

Computing statistical dimension of descent cone?

In general, there is a duality between descent cone and subdifferentials
set of subgradients

$$(\mathcal{D}(f, \mathbf{x}))^\circ = \text{cone}(\partial f(\mathbf{x})) := \bigcup_{\tau \geq 0} \tau \partial f(\mathbf{x})$$

Lemma 10.6 (informal, Amelunxen, Lotz, McCoy & Tropp '13)

$$\text{stat-dim}(\mathcal{D}(f, \mathbf{x})) \approx \inf_{\tau \geq 0} \mathbb{E} \left[\min_{\mathbf{u} \in \partial f(\mathbf{x})} \|\mathbf{g} - \tau \mathbf{u}\|_2^2 \right]$$

Example: ℓ_1 minimization

WLOG, suppose $x_1, \dots, x_k > 0$, $x_{k+1} = \dots = x_p = 0$.

$$\begin{aligned}\mathbb{E} \left[\min_{\mathbf{u} \in \partial \|\mathbf{x}\|_1} \|\mathbf{g} - \tau \mathbf{u}\|_2^2 \right] &= \mathbb{E} \left[\sum_{i=1}^k (g_i - \tau)^2 + \sum_{i=k+1}^p \min_{|u_i| \leq 1} (g_i - \tau u_i)^2 \right] \\ &= k (1 + \tau^2) + (p - k) \mathbb{E} \left[\min_{|u_i| \leq 1} (g_i - \tau u_i)^2 \right] \\ &= k (1 + \tau^2) + (p - k) \mathbb{E} \left[(|g_i| - \tau)_+^2 \right]\end{aligned}$$

By Lemma 10.6,

$$\begin{aligned}\text{stat-dim}(\mathcal{D}(\|\cdot\|_1, \mathbf{x})) &\approx \inf_{\tau \geq 0} \mathbb{E} \left[\min_{\mathbf{u} \in \partial \|\mathbf{x}\|_1} \|\mathbf{g} - \tau \mathbf{u}\|_2^2 \right] \\ &= \inf_{\tau \geq 0} \left\{ k (1 + \tau^2) + (p - k) \sqrt{\frac{2}{\pi}} \int_{\tau}^{\infty} (z - \tau)^2 e^{-z^2} dz \right\}\end{aligned}$$

Numerical phase transition

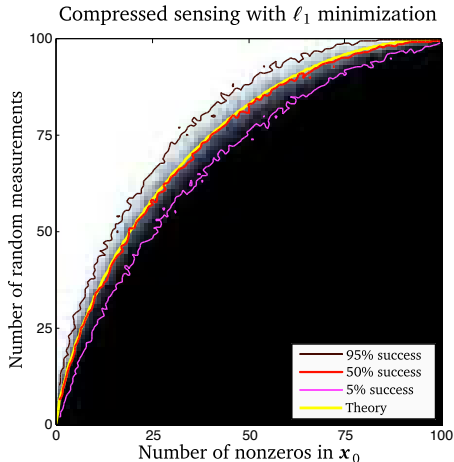


Figure credit: Amelunxen, Lotz, McCoy, & Tropp '13

Reference

- *“Living on the edge: Phase transitions in convex programs with random data,”* D. Amelunxen, M. Lotz, M. McCoy, J. Tropp, *Information and Inference*, 2014.
- *“Neighborliness of randomly projected simplices in high dimensions,”* D. Donoho and J. Tanner, *PNAS*, 2005.
- *“Stochastic and integral geometry,”* R. Schneider and W. Wolfgang, *Springer Science & Business Media*, 2008.
- *“Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition,”* T. Cover, *IEEE trans. on electronic computers*, 1965.
- *“The likelihood ratio test in high-dimensional logistic regression is asymptotically a rescaled chi-square,”* P. Sur, Y. Chen, and E. Candes, 2017.