# Robust Regression

Method of least absolute deviation ($l_1$ -regression)

- Find x* that minimizes $|Ax\text{-}b|_1 = \Sigma\ |b_i - <A_{i*}, x>|$

- Cost is less sensitive to outliers than least squares

- Can solve via linear programming

# Solving $l_1$-regression via Linear Programming

- Minimize $(1,\ldots,1) \cdot (\alpha^+ + \alpha^-)$
- Subject to:

$$A\,x + \alpha^+ - \alpha^- = b$$
$$\alpha^+, \alpha^- \geq 0$$

- Generic linear programming gives poly(nd) time

- Want much faster time using sketching!

# Well-Conditioned Bases

- For an n x d matrix A, can choose an n x d matrix U with orthonormal columns for which A = UW, and $|Ux|_2 = |x|_2$ for all x

- Can we find a U for which A = UW and $|Ux|_1 \approx |x|_1$ for all x?

- Let A = QW where Q has full column rank, and define $|z|_{Q,1} = |Qz|_1$
  - $|z|_{Q,1}$ is a norm

- Let C = $\{z \in R^d : |z|_{Q,1} \leq 1\}$ be the unit ball of $|.|_{Q,1}$

- C is a convex set which is symmetric about the origin
  - Lowner-John Theorem: can find an ellipsoid E such that: $E \subseteq C \subseteq \sqrt{d}E$, where E = $\{z \in R^d : z^T Fz \leq 1\}$
  - $\left(z^T Fz\right)^{.5} \leq |z|_{Q,1} \leq \sqrt{d}\left(z^T Fz\right)^{.5}$
  - $F = GG^T$ since F defines an ellipsoid

- Define $U = QG^{-1}$

# Well-Conditioned Bases

- Recall $U = QG^{-1}$ where

$$\left(z^T F z\right)^{.5} \leq |z|_{Q,1} \leq \sqrt{d}\left(z^T F z\right)^{.5} \text{ and } F = GG^T$$

- $|Ux|_1 = |QG^{-1}x|_1 = |Qz|_1 = |z|_{Q,1}$ where $z = G^{-1}x$

- $z^T F z \quad = \left(x^T (G^{-1})^T G^T G \, (G^{-1})x\right) = x^T x = |x|_2^2$

- So $|x|_2 \leq |Ux|_1 \leq \sqrt{d}|x|_2$

- So $\dfrac{|x|_1}{\sqrt{d}} \leq |x|_2 \leq |Ux|_1 \leq \sqrt{d}|x|_2 \leq \sqrt{d}|x|_1$

# Net for $\ell_1$ − Ball

- Consider the unit $\ell_1$-ball B = $\{x \in R^d : |x|_1 = 1\}$
- Subset N is a $\gamma$-net if for all $x \in B$, there is a $y \in N$, such that $|x - y|_1 \leq \gamma$
- Greedy construction of N
  - While there is a point $x \in B$ of distance larger than $\gamma$ from every point in N, include x in N
- The $\ell_1$-ball of radius $\gamma/2$ around every point in N is contained in the $\ell_1$-ball of radius 1+ $\gamma/2$ around $0^d$
- Further, all such ball are disjoint
- Ratio of volume of d-dimensional similar polytopes of radius 1+ $\gamma/2$ to radius $\gamma/2$ is $(1 + \gamma/2)^d/(\gamma/2)^d$, so $|N| \leq (1 + \gamma/2)^d/(\gamma/2)^d$

# Net for $\ell_1 -$ Subspace

- Let A = UW for a well-conditioned basis U
  - $|x|_1 \leq |Ux|_1 \leq d|x|_1$ for all x

- Let N be a $(\gamma/d) -$net for the unit $\ell_1$-ball B

- Let M = {Ux | x in N}, so $|M| \leq (1 + \gamma/(2d))^d/(\gamma/(2d))^d$

- Claim: For every x in B, there is a y in M for which $|Ux - y|_1 \leq \gamma$

- Proof: Let x' in B be such that $|x - x'|_1 \leq \gamma/d$
  Then $|Ux - Ux'|_1 \leq d|x - x'|_1 \leq \gamma$, using the
  well-conditioned basis property. Set y = Ux'

- $|M| \leq \left(\dfrac{d}{\gamma}\right)^{O(d)}$

# Rough Algorithm Overview

$$\min_{x \text{ in } R^d} |Ax-b|_1 = \min_{x \text{ in } R^d} |Ux - b'|_1$$

Sample $\text{poly}(d/\varepsilon)$ rows of $U \circ b'$ proportional to their $l_1$-norm.

**STOP** Compute poly(d)- approximation

Compute well-conditioned basis **STOP**

Find x' such that
$|Ax'-b|_1 \leq \text{poly}(d) \min_{x \text{ in } R^d}$
Let $b' = b-Ax'$ be the residual

Find a basis $A=UW$ so that for all x in $R^d$,
$|x|_1/\text{poly}(d) \leq |Ux|_1 \leq \text{poly}(d) |x|_1$

Takes nnz(A)

Now generic linear programming is efficient

Solve $l_1$-regression on the sample, obtaining vector x, and output x

Will focus on showing how to quickly compute

1. A poly(d)-approximation

2. A well-conditioned basis

# Sketching Theorem

## Theorem

- There is a probability space over (d log d) $\times$ n matrices R such that for any n$\times$d matrix A, with probability at least 99/100 we have for all x:

$$|Ax|_1 \leq |RAx|_1 \leq d \log d \cdot |Ax|_1$$

## Embedding

- is linear
- is independent of A
- preserves lengths of an infinite number of vectors

# Application of Sketching Theorem

Computing a d(log d)-approximation

- Compute RA and Rb

- Solve $x' = \text{argmin}_x |RAx - Rb|_1$

- Main theorem applied to $A \circ b$ implies $x'$ is a d log d – approximation

- RA, Rb have d log d rows, so can solve $l_1$-regression efficiently

# Application of Sketching Theorem

## Computing a well-conditioned basis

1. Compute RA

2. Compute W so that RAW is orthonormal (in the $l_2$-sense)

3. Output U = AW

## U = AW is well-conditioned because

$$|AWx|_1 \leq |RAWx|_1 \leq (d \log d)^{1/2} |RAWx|_2 = (d \log d)^{1/2} |x|_2 \leq (d \log d)^{1/2} |x|_1$$

and

$$|AWx|_1 \geq |RAWx|_1/(d \log d) \gtrsim |RAWx|_2/(d \log d) = |x|_2/(d \log d) \geq |x|_1 /(d^{3/2} \log d)$$

# Sketching Theorem

Theorem:

- There is a probability space over (d log d) $\times$ n matrices R such that for any n$\times$d matrix A, with probability at least 99/100 we have for all x:

$$|Ax|_1 \leq |RAx|_1 \leq d \log d \cdot |Ax|_1$$

A dense R that works:
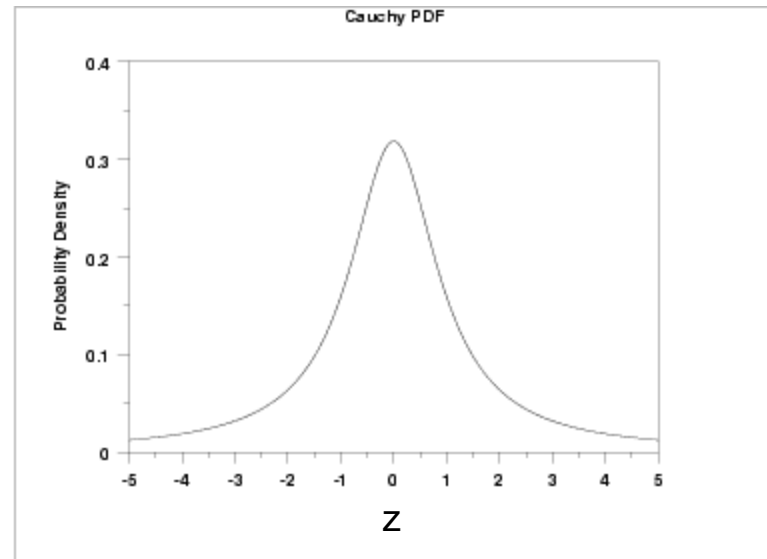
The entries of R are i.i.d. Cauchy random variables, scaled by 1/(d log d)
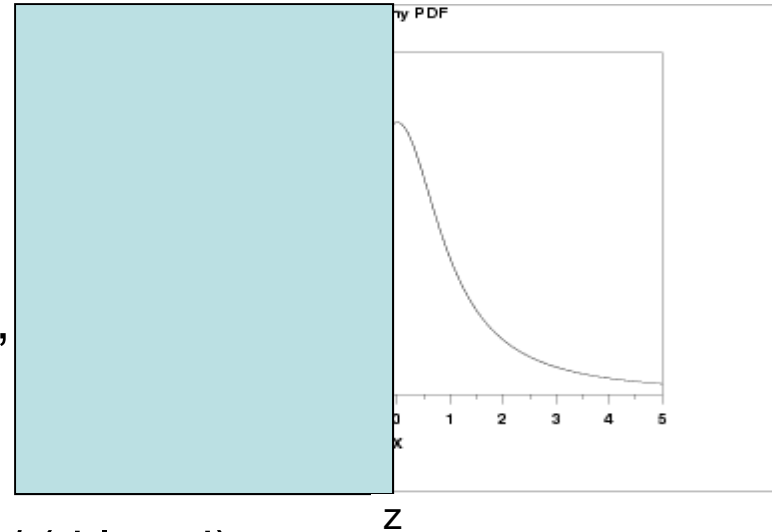
# Cauchy Random Variables

- pdf(z) = $1/(\pi(1+z^2))$ for z in (-4 , 4 )

- Undefined expectation and infinite variance



Cauchy PDF

- 1-stable:
    - If $z_1$, $z_2$, …, $z_n$ are i.i.d. Cauchy, then for a 5 $R^n$,

    $$a_1 \cdot z_1 + a_2 \cdot z_2 + … + a_n \cdot z_n \; \hat{\sim} \; |a|_1 \cdot z, \text{ where z is Cauchy}$$

- Can generate as the ratio of two standard normal random variables

# Proof of Sketching Theorem

- By 1-stability,
  - For all rows r of R,
    - $\langle r, Ax\rangle = |Ax|_1 \cdot Z / (d \log d)$,
      where Z is a Cauchy

$$z$$

- $RAx = (|Ax|_1 \cdot Z_1, \ldots, |Ax|_1 \cdot Z_{d \log d}) / (d \log d)$,
  where $Z_1, \ldots, Z_{d \log d}$ are i.i.d. Cauchy

- $|RAx|_1 = |Ax|_1 \sum_j |Z_j| / (d \log d)$
  - The $|Z_j|$ are half-Cauchy

- $\sum_j |Z_j| = \Omega(d \log d)$ with probability 1-exp(-d log d) by Chernoff

- But the $|Z_j|$ are heavy-tailed…

# Proof of Sketching Theorem

- $\sum_j |Z_j|$ is heavy-tailed, so $|RAx|_1 = |Ax|_1 \sum_j |Z_j| / (d \log d)$ may be large

- Each $|Z_j|$ has c.d.f. asymptotic to $1-\Theta(1/z)$ for z in [0, 4 )

- There *exists* a well-conditioned basis of A
  - Suppose w.l.o.g. the basis vectors are $A_{*1}, \ldots, A_{*d}$

- $|RA_{*i}|_1 \cong |A_{*i}|_1 f \sum_j |Z_{i,j}| / (d \log d)$

- Let $E_{i,j}$ be the event that $|Z_{i,j}| \leq d^3$
  - Define $Z'_{i,j} = |Z_{i,j}|$ if $|Z_{i,j}| \leq d^3$, and $Z'_{i,j} = d^3$ otherwise
  - $E[Z_{i,j} \mid E_{i,j}] = E[Z'_{i,j} \mid E_{i,j}] = O(\log d)$

- Let E be the event that for all i,j, $E_{i,j}$ occurs
  - $Pr[E] \geq 1 - \frac{\log d}{d}$
- What is $E[Z'_{i,j} \mid E]$?

# Proof of Sketching Theorem

- What is $E[Z'_{i,j} \mid E]$?

- $E[Z'_{i,j}|E_{i,j}] = E[Z'_{i,j}|E_{i,j}, E] \Pr[E \mid E_{i,j}] + E[Z'_{i,j}|E_{i,j}, \neg E] \Pr[\neg E \mid E_{i,j}]$

$$\geq E[Z'_{i,j}|E_{i,j}, E] \Pr[E \mid E_{i,j}]$$

$$= E[Z'_{i,j}|E] \cdot \left( \frac{\Pr[E_{i,j}|E] \Pr[E]}{\Pr[E_{i,j}]} \right)$$

$$\geq E[Z'_{i,j}|E] \cdot \left( 1 - \frac{\log d}{d} \right)$$

- So, $E[Z'_{i,j}|E] = O(\log d)$
- $|RA_{*i}|_1 \stackrel{\frown}{=} |A_{*i}|_1 \cdot \sum_{i,j} |Z_{i,j}| / (d \log d)$
- With constant probability, $\sum_i |RA_{*i}|_1 = O(\log d) \sum_i |A_{*i}|_1$

# Proof of Sketching Theorem

- With constant probability, $\sum_i |RA_{*i}|_1 = O(\log d) \sum_i |A_{*i}|_1$

- Recall $A_{*1}, \ldots, A_{*d}$ is a well-conditioned basis, and we showed the existence of such a basis earlier

- We will use the <span style="color:red">Auerbach basis</span> which always exists:
    - For all x, $|x|_4 \leq |Ax|_1$
    - $\sum_i |A_{*i}|_1 = d$

- $\sum_i |RA_{*i}|_1 = O(d \log d)$

- For all x, $|RAx|_1 \leq \sum_i |RA_{*i} x_i| \leq |x|_4 \sum_i |RA_{*i}|_1$
$$= |x|_4 \, O(d \log d)$$
$$= O(d \log d) \, |Ax|_1$$

# Where are we?

- Suffices to show for all x with $|x|_1 = 1$, that $|Ax|_1 \leq |RAx|_1 \leq d \log d \cdot |Ax|_1$
- We know

  - (1) there is a $\gamma$-net M, with $|M| \leq \left(\dfrac{d}{\gamma}\right)^{O(d)}$, of the set {Ax such that $|x|_1 = 1$}

  - (2) for any fixed x, $|RAx|_1 \geq |Ax|_1$ with probability $1 - \exp(-d \log d)$
  - (3) for all x, $|RAx|_1 = O(d \log d)|Ax|_1$

- Set $\gamma = 1/(d^3 \log d)$ so $|M| \leq d^{O(d)}$
  - By a union bound, for all y in M, $|Ry|_1 \geq |y|_1$

- Let x with $|x|_1 = 1$ be arbitrary. Let y in M satisfy $|Ax - y|_1 \leq \gamma = 1/(d^3 \log d)$

- $|RAx|_1 \geq |Ry|_1 - |R(Ax - y)|_1$
  $\geq |y|_1 - O(d \log d)|Ax - y|_1$
  $\geq |y|_1 - O(d \log d)\gamma$
  $\geq |y|_1 - O\left(\dfrac{1}{d^2}\right)$
  $\geq |y|_1/2$   (why?)

120

# Sketching to solve $l_1$-regression [CW, MM]

- Most expensive operation is computing R*A where R is the matrix of i.i.d. Cauchy random variables

- All other operations are in the "smaller space"

- Can speed this up by choosing R as follows:

$$
\begin{bmatrix}
0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\
1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & -1 & 1 & 0 & -1 & 0 \\
0 & -1 & 0 & 0 & 0 & 0 & 0 & 1
\end{bmatrix}
\cdot
\begin{bmatrix}
C_1 & & & & \\
& C_2 & & & \\
& & C_3 & & \\
& & & \ldots & \\
& & & & C_n
\end{bmatrix}
$$

# Further sketching improvements [WZ]

- Can show you need a fewer number of sampled rows in later steps if instead choose R as follows

- Instead of diagonal of Cauchy random variables, choose diagonal of reciprocals of exponential random variables

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1/E_1 & & & & \\ & 1/E_2 & & & \\ & & 1/E_3 & & \\ & & & \dots & \\ & & & & 1/E_n \end{bmatrix}$$

# Turnstile Streaming Model

- Underlying n-dimensional vector x initialized to $0^n$

- Long stream of updates $x_i \leftarrow x_i + \Delta_i$ for $\Delta_i$ in $\{-1,1\}$

- At end of the stream, x is promised to be in the set $\{-M, -M+1, \ldots, M-1, M\}^n$ for some bound $M \leq \text{poly}(n)$

- Output an approximation to f(x) whp

- Goal: use as little space (in bits) as possible
  - Massive data: stock transactions, weather data, genomes

# Example Problem: Norms

- Suppose you want $|x|_p^p = \Sigma_{i=1}^n |x_i|^p$

- Want Z for which $(1-\varepsilon) |x|_p^p \leq Z \leq (1+\varepsilon) |x|_p^p$ with probability $> 9/10$

# Example Problem: Euclidean Norm

- Want Z for which $(1-\varepsilon) \, |x|_2^2 \leq Z \leq (1+\varepsilon) \, |x|_2^2$

- Sample a random CountSketch matrix S with $1/\epsilon^2$ rows

- Can store S efficiently using limited independence

-  If $x_i \leftarrow x_i + \Delta_i$ in the stream, then $Sx \leftarrow Sx + \Delta_i S_{*,i}$

- At end of stream, output $|Sx|_2^2$

- With probability at least 9/10, $|Sx|_2^2 = (1 \pm \epsilon)|x|_2^2$

- Space complexity is $1/\epsilon^2$ words, each word is O(log n) bits

# Example Problem: 1-Norm

- Want Z for which $(1-\varepsilon)\,|x|_1 \le Z \le (1+\varepsilon)\,|x|_1$

- Sample a random Cauchy matrix S?

- Can store S with $\frac{1}{\epsilon}$ words of space [Kane, Nelson, W]

- If $x_i \leftarrow x_i + \Delta_i$ in the stream, then $Sx \leftarrow Sx + \Delta_i S_{*,i}$

- Space complexity is $1/\epsilon^2$ words, each word is $O(\log n)$ bits ?

- At end of stream, output $|Sx|_1$ ?

- *Cauchy random variables have no concentration…*

# 1-Norm Estimator

- Probability density function f(x) of |C| for a Cauchy random variable C is $f(x) = \dfrac{2}{\pi(1+x^2)}$

- Cumulative distribution function F(z):

$$F(z) = \int_0^z f(x)dx = \frac{2}{\pi}\arctan(z)$$

- Since $\tan(\pi/4) = 1$, F(1) = ½, so median(|C|) = 1

- If you take $r = \dfrac{\log\left(\frac{1}{\delta}\right)}{\epsilon^2}$ independent samples $X_1, \dots, X_r$ from F, and $X = \text{median}_i X_i$ , then F(X) in [1/2-$\epsilon$, 1/2+ $\epsilon$] with large probability

- $F^{-1}(X) = \tan\left(\dfrac{X\pi}{2}\right) \in [1 - 4\epsilon, 1 + 4\epsilon]$