

# CIS 700: “algorithms for Big Data”

## Lecture 5: Dimension Reduction

Slides at <http://grigory.us/big-data-class.html>

**Grigory Yaroslavtsev**

<http://grigory.us>



# $L_p$ -norm Estimation

- Stream:  $m$  updates  $(x_i, \Delta_i) \in [n] \times \mathbb{R}$  that define vector  $f$  where  $f_j = \sum_{i:x_i=j} \Delta_i$ .
- **Example:** For  $n = 4$

$$\langle (1,3), (3, 0.5), (1,2), (2, -2), (2,1), (1, -1), (4,1) \rangle$$
$$f = (4, -1, 0.5, 1)$$

- $L_p$ -norm:  $\|f\|_p = (\sum_i |f_i|^p)^{\frac{1}{p}}$

# $L_p$ -norm Estimation

- $L_p$ -norm:  $\|f\|_p = (\sum_i |f|^p)^{\frac{1}{p}}$
- Two lectures ago:
  - $\|f\|_0 = F_0$ -moment
  - $\|f\|_2^2 = F_2$ -moment (via AMS sketching)
- Space:  $O\left(\frac{\log n}{\epsilon^2} \log \frac{1}{\delta}\right)$
- Technique: linear sketches
  - $\|f\|_0$ :  $\sum_{i \in S} f_i$  for random sets  $S$
  - $\|f\|_2^2$ :  $\sum_i \sigma_i f_i$  for random signs  $\sigma_i$

# AMS as dimensionality reduction

- Maintain a “linear sketch” vector

$$\mathbf{Z} = (Z_1, \dots, Z_k) = Rf$$

$$Z_i = \sum_{j \in [n]} \sigma_{ij} f_j, \text{ where } \sigma_{ij} \in_R \{-1, 1\}$$

- Estimator  $\mathbf{Y}$  for  $\|f\|_2^2$ :

$$\frac{1}{k} \sum_{i=1}^k Z_i^2 = \frac{\|Rf\|_2^2}{k}$$

- “Dimensionality reduction”:  $x \rightarrow Rx$ , “heavy” tail

$$\Pr \left[ \left| \mathbf{Y} - \|f\|_2^2 \right| \geq c \left( \frac{2}{k} \right)^{\frac{1}{2}} \|f\|_2^2 \right] \leq \frac{1}{c^2}$$

# Normal Distribution

- Normal distribution  $N(0,1)$ 
  - Range:  $(-\infty, +\infty)$
  - Density:  $\mu(x) = (2\pi)^{-\frac{1}{2}} e^{-\frac{x^2}{2}}$
  - Mean = 0, Variance = 1
- Basic facts:
  - If  $X$  and  $Y$  are independent r.v. with normal distribution then  $X + Y$  has normal distribution
  - $Var[cX] = c^2 Var[X]$
  - If  $X, Y$  are independent, then  $Var[X + Y] = Var[X] + Var[Y]$

# Johnson-Lindenstrauss Transform

- Instead of  $\pm 1$  let  $\sigma_i$  be i.i.d. random variables from normal distribution  $N(0,1)$

$$Z = \sum_i \sigma_i f_i$$

- We still have  $\mathbb{E}[Z^2] = \sum_i f_i^2 = \|f\|_2^2$  because:
  - $\mathbb{E}[\sigma_i]\mathbb{E}[\sigma_j] = 0$ ;  $\mathbb{E}[\sigma_i^2] = \text{“variance of } \sigma_i \text{”} = 1$
- Define  $\mathbf{Z} = (Z_1, \dots, Z_k)$  and define:

$$\|\mathbf{Z}\|_2^2 = \sum_j Z_j^2 \quad \left( \mathbb{E}[\|\mathbf{Z}\|_2^2] = k\|f\|_2^2 \right)$$

- **JL Lemma:** There exists  $C > 0$  s.t. for small enough  $\epsilon > 0$ :

$$\Pr \left[ \left| \|\mathbf{Z}\|_2^2 - k\|f\|_2^2 \right| > \epsilon k\|f\|_2^2 \right] \leq \exp(-C\epsilon^2 k)$$

# Proof of JL Lemma

- **JL Lemma:**  $\exists C > 0$  s.t. for small enough  $\epsilon > 0$ :  
$$\Pr \left[ \left| \|\mathbf{Z}\|_2^2 - k \|f\|_2^2 \right| > \epsilon k \|f\|_2^2 \right] \leq \exp(-C \epsilon^2 k)$$
- Assume  $\|f\|_2^2 = 1$ .
- We have  $\mathbf{Z}_i = \sum_j \sigma_{ij} f_j$  and  $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_k)$   
$$\mathbb{E} \left[ \|\mathbf{Z}\|_2^2 \right] = k \|f\|_2^2 = k$$
- **Alternative form of JL Lemma:**  
$$\Pr \left[ \|\mathbf{Z}\|_2^2 > k(1 + \epsilon)^2 \right] \leq \exp(-\epsilon^2 k + O(k \epsilon^3))$$

# Proof of JL Lemma

- Alternative form of JL Lemma:

$$\Pr \left[ \|\mathbf{Z}\|_2^2 > k(1 + \epsilon)^2 \right] \leq \exp(-\epsilon^2 k + O(k \epsilon^3))$$

- Let  $Y = \|\mathbf{Z}\|_2^2$  and  $\alpha = k(1 + \epsilon)^2$
- For every  $s > 0$  we have:

$$\Pr[Y > \alpha] = \Pr[e^{sY} > e^{s\alpha}]$$

- By Markov and independence of  $\mathbf{Z}'_i$ s:

$$\Pr[e^{sY} > e^{s\alpha}] \leq \frac{\mathbb{E}[e^{sY}]}{e^{s\alpha}} = e^{-s\alpha} \mathbb{E} \left[ e^{s \sum_i Z_i^2} \right] = e^{-s\alpha} \prod_{i=1}^k \mathbb{E} \left[ e^{s Z_i^2} \right]$$

- We have  $Z_i \in N(0,1)$ , hence:

$$\mathbb{E} \left[ e^{s Z_i^2} \right] = (2\pi)^{-\frac{1}{2}} \int_{-\infty}^{\infty} e^{s t^2} e^{-\frac{t^2}{2}} dt = \frac{1}{\sqrt{1 - 2s}}$$



# Proof of JL Lemma

- Alternative form of JL Lemma:

$$\Pr \left[ \|\mathbf{Z}\|_2^2 > k(1 + \epsilon)^2 \right] \leq \exp(-\epsilon^2 k + O(k \epsilon^3))$$

- For every  $\mathbf{s} > 0$  we have:

$$\Pr[\mathbf{Y} > \alpha] \leq e^{-\mathbf{s}\alpha} \prod_{i=1}^k \mathbb{E} \left[ e^{\mathbf{s}Z_i^2} \right] = e^{-\mathbf{s}\alpha} (1 - 2\mathbf{s})^{-\frac{k}{2}}$$

- Let  $\mathbf{s} = \frac{1}{2} \left( 1 - \frac{k}{\alpha} \right)$  and recall that  $\alpha = k(1 + \epsilon)^2$
- A calculation finishes the proof:

$$\Pr[\mathbf{Y} > \alpha] \leq \exp(-\epsilon^2 k + O(k \epsilon^3))$$

# Johnson-Lindenstrauss Transform

- Single vector:  $k = O\left(\frac{\log \frac{1}{\delta}}{\epsilon^2}\right)$ 
  - Tight:  $k = \Omega\left(\frac{\log \frac{1}{\delta}}{\epsilon^2}\right)$  [Woodruff'10]
- $n$  vectors simultaneously:  $k = O\left(\frac{\log n \log \frac{1}{\delta}}{\epsilon^2}\right)$ 
  - Tight:  $k = \Omega\left(\frac{\log n \log \frac{1}{\delta}}{\epsilon^2}\right)$  [Molinaro, Woodruff, Y. '13]
- Distances between  $n$  vectors =  $O(n^2)$  vectors:

$$k = O\left(\frac{\log n \log \frac{1}{\delta}}{\epsilon^2}\right)$$

# Random Variables and Norms

- For a random variable  $\mathbf{X}$  and  $p \geq 1$  let:

$$\|\mathbf{X}\|_p = \mathbb{E}[X^p]^{1/p}$$

Facts:

- For any  $c$ :  $\|c\mathbf{X}\|_p = c\|\mathbf{X}\|_p$
- $\|\cdot\|_p$  is a norm (Minkowski's inequality)
- $\|\cdot\|_p \leq \|\cdot\|_q$  for  $p \leq q$  (Monotonicity of norms)
- Jensen's inequality (used a lot for  $F = |x|^p$ ):  
If  $F$  is convex then  $F(\mathbb{E}[\mathbf{X}]) \leq \mathbb{E}[F(\mathbf{X})]$

# Khintchine Inequality

- [Khintchine] For  $p \geq 1, x \in \mathbb{R}^n$  and  $\sigma_i$  i.i.d. Rademachers:

$$\left\| \sum_i \sigma_i x_i \right\|_p \leq \sqrt{p} \|x\|_2$$

- For  $r_i$  (either  $\sigma_i$  or  $g_i \sim N(0,1)$ ) expand  $\mathbb{E}[(\sum_i r_i x_i)^p]$
- All odd powers of  $r_i$  are zero
- All even moments for  $\sigma_i$  are 1, and for  $g_i$  are  $\geq 1$
- $\left\| \sum_i \sigma_i x_i \right\|_p \leq \left\| \sum_i g_i x_i \right\|_p$
- $\sum_i g_i x_i \sim N(0, \|x\|_2^2) \Rightarrow \left\| \sum_i g_i x_i \right\|_p \leq \sqrt{p} \|x\|_2$

# Symmetrization

- [Symmetrization]: If  $Z_1, \dots, Z_n$  are independent and  $\sigma_i$  are i.i.d. Rademachers:

$$\left\| \sum_i Z_i - \mathbb{E} \sum_i Z_i \right\|_p \leq 2 \left\| \sum_i \sigma_i Z_i \right\|_p$$

- Let  $Y_1 \dots Y_n$  be independent with the same distribution as  $Z_i$
- $\left\| \sum_i Z_i - \mathbb{E} \sum_i Z_i \right\|_p = \left\| \sum_i Z_i - \mathbb{E}_Y \sum_i Y_i \right\|_p$   
 $\leq \left\| \sum_i (Z_i - Y_i) \right\|_p$  (Jensen)  
 $= \left\| \sum_i \sigma_i (Z_i - Y_i) \right\|_p$  ( $Z_i - Y_i$  are independent and symmetric)  
 $\leq 2 \left\| \sum_i \sigma_i Z_i \right\|_p$  (triangle inequality)

# Decoupling

- Let  $x_1, \dots, x_n$  be independent with mean 0 and  $x'_1, \dots, x'_n$  identically distributed as  $x_i$  and independent of them. For any  $a_{ij}$  and  $p \geq 1$ :

$$\left\| \sum_{i \neq j} a_{ij} x_i x_j \right\|_p \leq 4 \left\| \sum_{i,j} a_{ij} x_i x'_j \right\|_p$$

- Let  $\eta_1, \dots, \eta_n$  be i.i.d. Bernoullis (0/1 w.p. 1/2):

$$\begin{aligned} \left\| \sum_{i \neq j} a_{ij} x_i x_j \right\|_p &= 4 \left\| \mathbb{E}_\eta \sum_{i \neq j} a_{ij} x_i x_j \eta_i (1 - \eta_j) \right\|_p \\ &\leq 4 \left\| \sum_{i \neq j} a_{ij} x_i x_j \eta_i (1 - \eta_j) \right\|_p \text{ (Jensen)} \end{aligned}$$

- There exists  $\eta' \in \{0,1\}^n$  such that:

$$\left\| \sum_{i \neq j} a_{ij} x_i x_j \eta_i (1 - \eta_j) \right\|_p \leq \left\| \sum_{i \in S} \sum_{j \in \bar{S}} a_{ij} x_i x_j \right\|_p$$

where  $S = \{i : \eta' = 1\}$ .

# Decoupling (continued)

Let  $x_S$  be an  $S$ -dimensional vector of  $x_i$  for  $i \in S$ .

$$\begin{aligned}
 & \bullet \left\| \sum_{i \in S} \sum_{j \in \bar{S}} a_{ij} x_i x_j \right\|_p = \left\| \sum_{i \in S} \sum_{j \in \bar{S}} a_{ij} x_i x'_j \right\|_p \\
 &= \left\| \mathbb{E}_{x_{\bar{S}}} \mathbb{E}_{x'_S} \sum_{i,j} a_{ij} x_i x'_j \right\|_p \quad (\mathbb{E}[x_i] = \mathbb{E}[x'_i] = 0) \\
 &\leq \left\| \sum_{i,j} a_{ij} x_i x'_j \right\|_p \text{ (Jensen)}
 \end{aligned}$$

• Overall:

$$\left\| \sum_{i \neq j} a_{ij} x_i x_j \right\|_p \leq 4 \left\| \sum_{i,j} a_{ij} x_i x'_j \right\|_p$$

# Hanson-Wright Inequality

- For  $\sigma_1, \dots, \sigma_n$  independent Rademachers and  $A \in \mathbb{R}^{n \times n}$  real and symmetric for all  $p \geq 1$ :  
$$\left| \left| \sigma^T A \sigma - \mathbb{E}[\sigma^T A \sigma] \right| \right|_p \leq \sqrt{p} \|A\|_F + p \|A\|$$
- Recall:

$$- \|A\|_F = \sqrt{\sum_{ij} a_{ij}^2} = \sqrt{\text{Tr}(A^T A)}$$

$$- \|A\| = \sup_{\{v \neq 0\}} \frac{\|Av\|_2}{\|v\|_2}$$



# Hanson-Wright Inequality

- For  $\sigma_1, \dots, \sigma_n$  independent Rademachers and  $A \in \mathbb{R}^{n \times n}$  real and symmetric for all  $p \geq 1$ :

$$\left| \left| \sigma^T A \sigma - \mathbb{E}[\sigma^T A \sigma] \right| \right|_p \leq \sqrt{p} \|A\|_F + p \|A\|$$

$$\left| \left| \sigma^T A \sigma - \mathbb{E}[\sigma^T A \sigma] \right| \right|_p \leq \left| \left| \sigma^T A \sigma' \right| \right|_p \text{ (decoupling)}$$

$$\leq \sqrt{p} \left| \left| \left| A \sigma \right|_2 \right| \right|_p \text{ (Khintchine)}$$

$$= \sqrt{p} \left| \left| \left| A \sigma \right|_2^2 \right| \right|_{p/2}^{\frac{1}{2}}$$

$$\leq \sqrt{p} \left| \left| \left| A \sigma \right|_2^2 \right| \right|_p^{\frac{1}{2}} \text{ (monotonicity of norms)}$$

# Hanson-Wright (continued)

$$\begin{aligned}
 \sqrt{p} \left\| \left\| |A\sigma|_2 \right\| \right\|_p &\leq \dots \leq \sqrt{p} \left\| \left\| |A\sigma|_2^2 \right\| \right\|_p^{\frac{1}{2}} \\
 &\leq \sqrt{p} \left( \mathbb{E} \left[ \left\| |A\sigma|_2^2 \right\| \right] + \left\| \left\| |A\sigma|_2^2 \right\| - \mathbb{E} \left[ \left\| |A\sigma|_2^2 \right\| \right] \right\|_p \right)^{\frac{1}{2}} \text{ (triangle ineq.)} \\
 &= \sqrt{p} \left( \|A\|_F^2 + \left\| \left\| |A\sigma|_2^2 \right\| - \mathbb{E} \left[ \left\| |A\sigma|_2^2 \right\| \right] \right\|_p \right)^{\frac{1}{2}} \\
 &\leq \sqrt{p} \|A\|_F + \sqrt{p} \left\| \left\| |A\sigma|_2^2 \right\| - \mathbb{E} \left[ \left\| |A\sigma|_2^2 \right\| \right] \right\|_p^{\frac{1}{2}} \\
 &\preccurlyeq \sqrt{p} \|A\|_F + \sqrt{p} \left| \sigma^T A^T A \sigma' \right|_p^{\frac{1}{2}} \text{ (decoupling)} \\
 &\preccurlyeq \sqrt{p} \|A\|_F + p^{\frac{3}{4}} \left\| \left\| A^T A \sigma \right\|_2 \right\|_p^{1/2} \text{ (Khintchine)} \\
 &\preccurlyeq \sqrt{p} \|A\|_F + p^{\frac{3}{4}} \|A\|^{\frac{1}{2}} \left\| \left\| Ax \right\|_2 \right\|_p^{\frac{1}{2}}
 \end{aligned}$$

# Hanson-Wright (continued)

$$\sqrt{p} \left| \left| \left| A\sigma \right|_2 \right| \right|_p \leq \sqrt{p} \|A\|_F + p^{\frac{3}{4}} \|A\|^{\frac{1}{2}} \left| \left| \left| A\sigma \right|_2 \right| \right|_p^{\frac{1}{2}}$$

Let  $E = \left| \left| \left| Ax \right|_2 \right| \right|_p^{\frac{1}{2}}$  then  $E^2 - Cp^{\frac{1}{4}} \|A\|^{\frac{1}{2}} E - C \|A\|_F \leq 0$

- $E \leq$  larger root of the quadratic equation above
- $E^2 \leq \sqrt{p} \|A\|_F + p \|A\|$
- (Hanson-Wright) For  $\sigma_1, \dots, \sigma_n$  independent Rademachers and  $A \in \mathbb{R}^{n \times n}$  real and symmetric for all  $p \geq 1$ :

$$\left| \left| \sigma^T A \sigma - \mathbb{E}[\sigma^T A \sigma] \right| \right|_p \leq \sqrt{p} \|A\|_F + p \|A\|$$