#### **SIGTACS** Seminar Series

#### Metric Embeddings and Applications in Computer Science

Presented by : Purushottam Kar

January 10, 2009



# Outline

### 1 Introduction

- **2** Embeddings into Normed Spaces
- **3** Dimensionality Reduction
- 4 The JL Lemma





#### Definition (Metric)

A Metric is a structure  $(X, \rho)$  where  $\rho$  is a distance measure  $\rho : X \times X \to \mathbb{R}$  which is non-negative, symmetric and satisfies the triangle inequality.

#### **Definition (Embedding Distortion)**

An embedding  $f : X \to Y$  from a metric space  $(X, \rho)$  to another metric space  $(Y, \sigma)$  is said to have a distortion D if  $D = \sup_{x,y \in X} \frac{\sigma(f(x), f(y))}{\rho(x, y)} \cdot \sup_{x,y \in X} \frac{\rho(x, y)}{\sigma(f(x), f(y))}.$ 

Such embeddings are also called *bi-Lipschitz* embeddings.

### **Embeddings**

• Various criterion used to evaluate embeddings



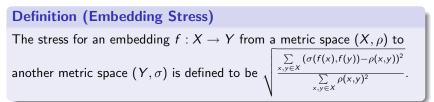
### **Embeddings**

- Various criterion used to evaluate embeddings
- Distortion, Stress, Residual Variance ...

# **Definition (Embedding Stress)** The stress for an embedding $f : X \to Y$ from a metric space $(X, \rho)$ to another metric space $(Y, \sigma)$ is defined to be $\sqrt{\frac{\sum\limits_{x,y \in X} (\sigma(f(x), f(y)) - \rho(x, y))^2}{\sum\limits_{x,y \in X} \rho(x, y)^2}}$ .

### **Embeddings**

- Various criterion used to evaluate embeddings
- Distortion, Stress, Residual Variance ...



• Lead to very interesting algorithmic questions



• Started out as a branch of functional analysis



- Started out as a branch of functional analysis
- Algorithmic applications



- Started out as a branch of functional analysis
- Algorithmic applications
  - Metric Embeddings for datasets operating with a non-metric



- Started out as a branch of functional analysis
- Algorithmic applications
  - Metric Embeddings for datasets operating with a non-metric
  - Dimensionality reduction to reduce storage space costs, processing time

- Started out as a branch of functional analysis
- Algorithmic applications
  - Metric Embeddings for datasets operating with a non-metric
  - Dimensionality reduction to reduce storage space costs, processing time
  - Facilitate pruning procedures in database searches

- Started out as a branch of functional analysis
- Algorithmic applications
  - Metric Embeddings for datasets operating with a non-metric
  - Dimensionality reduction to reduce storage space costs, processing time
  - Facilitate pruning procedures in database searches
  - Preserve residual variance (PCA), inter-point similarity (Random Projections), Stress (MDS)

- Started out as a branch of functional analysis
- Algorithmic applications
  - Metric Embeddings for datasets operating with a non-metric
  - Dimensionality reduction to reduce storage space costs, processing time
  - Facilitate pruning procedures in database searches
  - Preserve residual variance (PCA), inter-point similarity (Random Projections), Stress (MDS)
- Streaming Algorithms

**Theorem (Frétchet's Embedding)** 

Every n-point metric can be isometrically embedded into  $I_\infty$ 

• Fréchet's Embedding technique - non-expansive

#### Theorem (Frétchet's Embedding)

Every n-point metric can be isometrically embedded into  $I_\infty$ 

- Fréchet's Embedding technique non-expansive
- Choose coordinates as projections onto some fixed sets

#### Theorem (Frétchet's Embedding)

Every n-point metric can be isometrically embedded into  $I_\infty$ 

- Fréchet's Embedding technique non-expansive
- Choose coordinates as projections onto some fixed sets
- Triangle inequality ensures contractive embeddings

#### Theorem (Frétchet's Embedding)

Every n-point metric can be isometrically embedded into  $I_\infty$ 

- Fréchet's Embedding technique non-expansive
- Choose coordinates as projections onto some fixed sets
- Triangle inequality ensures contractive embeddings
- Choice of "landmark" sets gives other algorithms

#### Theorem (Frétchet's Embedding)

Every n-point metric can be isometrically embedded into  $I_\infty$ 

- Fréchet's Embedding technique non-expansive
- Choose coordinates as projections onto some fixed sets
- Triangle inequality ensures contractive embeddings
- Choice of "landmark" sets gives other algorithms
- Embedding dimension can be reduced to O(qn<sup><sup>1</sup>/<sub>q</sub></sup> ln n) by tolerating a distortion of 2q − 1.

### **Embedding into** *l*<sub>2</sub>

#### Theorem (Bourgain's Embedding)

Every n-point metric can be  $O(\log n)$ -embedded into  $I_2$ 

• Uses a random selection of the landmark sets

## **Embedding into** *l*<sub>2</sub>

#### Theorem (Bourgain's Embedding)

Every n-point metric can be  $O(\log n)$ -embedded into  $l_2$ 

- Uses a random selection of the landmark sets
- Tight The graph metric of a constant degree expander has  $\Omega(\log n)$  distortion into any Euclidean space

SIGTACS Seminar Series

### **Embedding into** *l*<sub>2</sub>

#### Theorem (Bourgain's Embedding)

Every n-point metric can be  $O(\log n)$ -embedded into  $l_2$ 

- Uses a random selection of the landmark sets
- Tight The graph metric of a constant degree expander has  $\Omega(\log n)$  distortion into any Euclidean space
- Any embedding of the Hamming cube into  $l_2$  incurs  $\Omega\left(\sqrt{\log n}\right)$  distortion

# **Dimensionality Reduction in** *l*<sub>1</sub>

• Impossible - A *D*-embedding of *n* points may require  $n^{\Omega(1/D^2)}$  dimensions

# **Dimensionality Reduction in** $I_1$

- Impossible A *D*-embedding of *n* points may require  $n^{\Omega(1/D^2)}$  dimensions
- No "flattening" results known for other *l<sub>p</sub>* metrics either ...

# **Dimensionality Reduction in** $I_1$

- Impossible A *D*-embedding of *n* points may require  $n^{\Omega(1/D^2)}$  dimensions
- No "flattening" results known for other *l<sub>p</sub>* metrics either ...
- Except for p = 2

#### Theorem (The JL-Lemma)

Given  $\epsilon > 0$  and integer n, let  $k \ge k_0 = \mathcal{O}(\epsilon^{-2} \log n)$ . For every set P of n points in  $\mathbb{R}^d$  there exists  $f : \mathbb{R}^d \longrightarrow \mathbb{R}^k$  such that for all  $u, v \in P$ 

$$(1-\epsilon) \|u-v\|^2 \le \|f(u)-f(v)\|^2 \le (1+\epsilon) \|u-v\|^2$$

• Implementation as a randomized algorithm

#### Theorem (The JL-Lemma)

Given  $\epsilon > 0$  and integer n, let  $k \ge k_0 = \mathcal{O}(\epsilon^{-2} \log n)$ . For every set P of n points in  $\mathbb{R}^d$  there exists  $f : \mathbb{R}^d \longrightarrow \mathbb{R}^k$  such that for all  $u, v \in P$ 

$$(1-\epsilon)||u-v||^2 \le ||f(u)-f(v)||^2 \le (1+\epsilon)||u-v||^2$$

- Implementation as a randomized algorithm
- Equivalent interpretations random projection vs. random rotation

#### Theorem (The JL-Lemma)

Given  $\epsilon > 0$  and integer n, let  $k \ge k_0 = \mathcal{O}(\epsilon^{-2} \log n)$ . For every set P of n points in  $\mathbb{R}^d$  there exists  $f : \mathbb{R}^d \longrightarrow \mathbb{R}^k$  such that for all  $u, v \in P$ 

$$(1-\epsilon)||u-v||^2 \le ||f(u)-f(v)||^2 \le (1+\epsilon)||u-v||^2$$

- Implementation as a randomized algorithm
- Equivalent interpretations random projection vs. random rotation
- Various Proofs known [IM98], [DG99], [AV99], [A01]

#### Theorem (The JL-Lemma)

Given  $\epsilon > 0$  and integer n, let  $k \ge k_0 = \mathcal{O}(\epsilon^{-2} \log n)$ . For every set P of n points in  $\mathbb{R}^d$  there exists  $f : \mathbb{R}^d \longrightarrow \mathbb{R}^k$  such that for all  $u, v \in P$ 

$$(1-\epsilon)||u-v||^2 \le ||f(u)-f(v)||^2 \le (1+\epsilon)||u-v||^2$$

- Implementation as a randomized algorithm
- Equivalent interpretations random projection vs. random rotation
- Various Proofs known [IM98], [DG99], [AV99], [A01]
- Common Technique

Point Drafting  $\longrightarrow$  Set Drafting  $\stackrel{\text{Union Bound}}{\longrightarrow}$  Set Embedding

• Instead of choosing from an uncountably infinite domain, can we choose vectors from a finite set of vectors ?

- Instead of choosing from an uncountably infinite domain, can we choose vectors from a finite set of vectors ?
- Achlioptas: In fact 'choosing' from the *d*-dimensional Hamming Cube {1, -1}<sup>d</sup> works.



- Instead of choosing from an uncountably infinite domain, can we choose vectors from a finite set of vectors ?
- Achlioptas: In fact 'choosing' from the *d*-dimensional Hamming Cube  $\{1, -1\}^d$  works.
- Consider a random vector  $R = (X_1, X_2, ..., X_d)$ , where each  $X_i$  is chosen from one of the two distributions:

$$D_1 = rac{1}{\sqrt{d}} \left\{ egin{array}{cc} -1 & ext{ with probability } 1/2 \ 1 & ext{ with probability } 1/2 \end{array} 
ight.$$

$$D_2 = \frac{1}{\sqrt{d}} \begin{cases} -\sqrt{3} \\ 0 \\ \sqrt{3} \end{cases}$$

with probability 1/6with probability 2/3with probability 1/6

Discussion

### **Enter Achlioptas**

• Pick k such random vectors  $R_1, R_2, \ldots R_k$ .

- Pick k such random vectors  $R_1, R_2, \ldots R_k$ .
- For a given unit vector α = (α<sub>1</sub>, α<sub>2</sub>,..., α<sub>d</sub>), the low (k-)dimensional vector corresponding to α is

$$f(\alpha) = \sqrt{\frac{d}{k}} (\langle \alpha, R_1 \rangle, \langle \alpha, R_2 \rangle, \dots, \langle \alpha, R_k \rangle)$$

- Pick k such random vectors  $R_1, R_2, \ldots R_k$ .
- For a given unit vector α = (α<sub>1</sub>, α<sub>2</sub>,..., α<sub>d</sub>), the low (k-)dimensional vector corresponding to α is

$$f(\alpha) = \sqrt{\frac{d}{k}} (\langle \alpha, R_1 \rangle, \langle \alpha, R_2 \rangle, \dots, \langle \alpha, R_k \rangle)$$

• Advantage: Simple and can be implemented as SQL queries.

### Main Theorem

• Let 
$$S = \langle \alpha, R_1 \rangle^2 + \langle \alpha, R_2 \rangle^2 + \cdots \langle \alpha, R_k \rangle^2$$

#### Theorem (Main Theorem)

For every d-dimensional unit vector  $\alpha$ , integer  $k \ge 1$  and  $\epsilon > 0$ 

$$\Pr\left[S \ge (1 \pm \epsilon)\frac{k}{d} \cdot 1\right] \le e^{\frac{-k}{2}\left(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3}\right)}$$

### Main Theorem

• Let 
$$S = \langle \alpha, R_1 \rangle^2 + \langle \alpha, R_2 \rangle^2 + \cdots \langle \alpha, R_k \rangle^2$$

#### Theorem (Main Theorem)

For every d-dimensional unit vector  $\alpha$ , integer  $k \ge 1$  and  $\epsilon > 0$ 

$$\Pr\left[S \ge (1 \pm \epsilon)\frac{k}{d} \cdot 1\right] \le e^{\frac{-k}{2}\left(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3}\right)}$$

• Hence, if  $k \ge \frac{4+2\beta}{\epsilon^2/2-\epsilon^3/3} \log n$ , this probability becomes smaller than  $\frac{2}{n^{2+\beta}}$  which is inverse polynomial w.r.t n.

# **Expected Value of** $||f(\alpha)||^2$

 $\bullet\,$  On expectation the length of a unit vector  $\alpha$  is preserved.

$$E\left[\|f(\alpha)\|^2\right] = E\left[\sum_{i=1}^k \frac{d}{k} \left(\sum_{j=1}^d X_j \alpha_j\right)^2\right]$$
$$= \frac{d}{k} \sum_{i=1}^k \left(\sum_{j=1}^d E[X_j^2] \alpha_j^2 + \sum_{j
$$= \frac{d}{k} \sum_{i=1}^k \frac{1}{d} = 1 = \|\alpha\|^2$$$$

## Deviation from Expectation: Proof of Main Theorem

• By Markov inequality,

$$\Pr\left[S > (1+\epsilon)\frac{k}{d}\right] < E\left[e^{hS}\right]e^{-(1+\epsilon)\frac{hk}{d}}$$
$$\Pr\left[S < (1-\epsilon)\frac{k}{d}\right] < E\left[e^{-hS}\right]e^{(1-\epsilon)\frac{hk}{d}}$$

## Deviation from Expectation: Proof of Main Theorem

• By Markov inequality,

$$\Pr\left[S > (1+\epsilon)\frac{k}{d}\right] < E\left[e^{hS}\right]e^{-(1+\epsilon)\frac{hk}{d}}$$
$$\Pr\left[S < (1-\epsilon)\frac{k}{d}\right] < E\left[e^{-hS}\right]e^{(1-\epsilon)\frac{hk}{d}}$$

• Since the vectors  $R'_i s$  are all chosen independently we can rewrite the above as

$$\Pr\left[S > (1+\epsilon)\frac{k}{d}\right] < \left(E\left[e^{hQ_1^2}\right]\right)^k e^{-(1+\epsilon)\frac{hk}{d}}$$
  
$$\Pr\left[S < (1-\epsilon)\frac{k}{d}\right] < \left(E\left[e^{-hQ_1^2}\right]\right)^k e^{(1-\epsilon)\frac{hk}{d}}$$

where  $Q_1 = \langle lpha, R_1 
angle$ 

## **Proof of Main Theorem**

• By Taylor's Expansion,

$$\Pr\left[S < (1-\epsilon)\frac{k}{d}\right] < \left(E\left[1-hQ_1^2 + \frac{hQ_1^4}{2}\right]\right)^k e^{-(1+\epsilon)\frac{hk}{d}}$$
$$= \left(1-\frac{h}{d} + \frac{h^2 E[Q_1^4]}{2}\right)^k e^{(1-\epsilon)\frac{hk}{d}}$$

#### Lemma

For  $h \in [0, d/2)$  and all  $d \ge 1$ ,

$$E\left[e^{hQ_1^2}\right] \leq \frac{1}{\sqrt{1-2h/d}}$$
(1)  
$$E\left[Q_1^4\right] \leq \frac{3}{d^2}$$
(2)

# Proof of Main Theorem using Inequalities (1) and (2)

• If we take  $h = \frac{d\epsilon}{2(1+\epsilon)}$ , for the upper bound we have the following:

$$\Pr\left[S > (1+\epsilon)\frac{k}{d}\right] < \left(\frac{1}{\sqrt{1-2h/d}}\right)^k e^{-(1+\epsilon)\frac{hk}{d}}$$
$$= ((1+\epsilon)e^{-\epsilon})^{k/2} < e^{\frac{-k}{2}(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3})}.$$

## Proof of Main Theorem using Inequalities (1) and (2)

• If we take  $h = \frac{d\epsilon}{2(1+\epsilon)}$ , for the upper bound we have the following:

$$\Pr\left[S > (1+\epsilon)\frac{k}{d}\right] < \left(\frac{1}{\sqrt{1-2h/d}}\right)^k e^{-(1+\epsilon)\frac{hk}{d}}$$
$$= ((1+\epsilon)e^{-\epsilon})^{k/2} < e^{\frac{-k}{2}(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3})}.$$

• For the same value of *h*, for the lower bound we get:

$$\Pr\left[S < (1-\epsilon)\frac{k}{d}\right] < \left(1-h/d + \frac{3h^2}{2d^2}\right)^k e^{(1-\epsilon)\frac{hk}{d}} < e^{\frac{-k}{2}(\frac{\epsilon^2}{2} - \frac{\epsilon^3}{3})}.$$

Introduction

Discussion

## **Proof of Inequality (2)**

• For inequality (2)  

$$E[Q_{1}^{4}] = \left(\sum_{i=1}^{d} X_{i} \alpha_{i}\right)^{4} = \sum_{i} E[X_{i}^{4}] \alpha_{i}^{4} + \left(\binom{4}{1,3} \sum_{i < j} E[X_{i}^{3}] E[X_{j}] \alpha_{i}^{3} \alpha_{i} + \binom{4}{2,2} \sum_{i < j} E[X_{i}^{2}] E[X_{j}^{2}] \alpha_{i}^{2} \alpha_{j}^{2} + \left(\binom{4}{2,1,1} \sum_{i < j < k} E[X_{i}^{2}] E[X_{j}] E[X_{j}] E[X_{k}] \alpha_{i}^{2} \alpha_{j} \alpha_{k} + \left(\binom{4}{1,1,1,1} \sum_{i < j < k < l} E[X_{i}] E[X_{i}] E[X_{j}] E[X_{k}] E[X_{i}] \alpha_{i} \alpha_{j} \alpha_{k} \alpha_{l} \right) \\ = \frac{1}{d^{2}} (\alpha^{4} + 6 \sum_{i < j} \alpha_{i}^{2} \alpha_{j}^{2}) \le \frac{3}{d^{2}}.$$

# **Proof of Inequality (1)**

 The idea is to first make the random variable Q<sub>1</sub> independent of α and then compare the even moments of Q<sub>1</sub> with a properly scaled normal distribution.

#### Lemma (Worst Vector Lemma)

For all unit vectors 
$$\alpha$$
,  $E[Q_1^{2k}(\alpha)] \leq E[Q_1^{2k}(w)]$ , where  $w = \frac{1}{\sqrt{d}}(1, 1, \dots, 1)$  for  $k = 1, 2, \dots$ 

#### Lemma (Normal Bound Lemma)

If  $T \sim N(0, 1/d)$ , then  $E[Q_1^{2k}(w)] \leq E[T^{2k}]$ , where  $w = \frac{1}{\sqrt{d}}(1, 1, ..., 1)$ for k = 1, 2, ...

## **Proof of Inequality (1)**

$$E\left[e^{hT^{2}}\right] = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{\lambda^{2}/2} e^{h\lambda^{2}/d} d\lambda$$
  

$$= \frac{1}{\sqrt{1-2h/d}}$$
  

$$= E\left[\sum_{k=0}^{\infty} \frac{h^{k}T^{2k}}{k!}\right] \qquad \text{(using MCT)}$$
  

$$= \sum_{k=0}^{\infty} \frac{h^{k}E\left[T^{2k}\right]}{k!}$$
  

$$\geq \sum_{k=0}^{\infty} \frac{h^{k}E\left[Q_{1}^{2k}(w)\right]}{k!} = E\left[e^{hQ_{1}(w)^{2}}\right] \geq E\left[e^{hQ_{1}(\alpha)^{2}}\right]$$

#### **Proving the Worst Vector Lemma**

• Let  $r_1$  and  $r_2$  be i.i.d. r.v. distributed as  $\{-1, +1\}$  with equal probability. Furthermore let a, b, T be any reals and  $c = \sqrt{(a^2 + b^2)/2}$  and k > 0 be any integer, then

$$E\left[(T+ar_1+br_2)^{2k}\right] \leq E\left[(T+cr_1+cr_2)^{2k}\right]$$

#### **Proving the Worst Vector Lemma**

• Let 
$$r_1$$
 and  $r_2$  be i.i.d. r.v. distributed as  $\{-1, +1\}$  with equal probability. Furthermore let  $a, b, T$  be any reals and  $c = \sqrt{(a^2 + b^2)/2}$  and  $k > 0$  be any integer, then

$$E\left[(T + ar_1 + br_2)^{2k}\right] \le E\left[(T + cr_1 + cr_2)^{2k}\right]$$

• Let  $R_1 = \frac{1}{\sqrt{d}}(r_1, r_2, \dots, r_d)$ . Thus we have

$$E\left[Q_{1}(\alpha)^{2k}\right] = \frac{1}{d^{k}} \sum_{R} E\left[(R + \alpha_{1}r_{1} + \alpha_{2}r_{2})^{2k}\right] \Pr\left[\sum_{i=3}^{d} \alpha_{i}r_{i} = \frac{R}{\sqrt{d}}\right]$$

$$\leq \frac{1}{d^{k}} \sum_{R} E\left[(R + cr_{1} + cr_{2})^{2k}\right] \Pr\left[\sum_{i=3}^{d} \alpha_{i}r_{i} = \frac{R}{\sqrt{d}}\right]$$

$$= E\left[Q_{1}(\theta)^{2k}\right]$$
where  $c = \sqrt{(\alpha_{1}^{2} + \alpha_{2}^{2})/2}$ 

#### **Proving the Worst Vector Lemma**

• Let 
$$r_1$$
 and  $r_2$  be i.i.d. r.v. distributed as  $\{-1, +1\}$  with equal probability. Furthermore let  $a, b, T$  be any reals and  $c = \sqrt{(a^2 + b^2)/2}$  and  $k > 0$  be any integer, then

$$E\left[(T + ar_1 + br_2)^{2k}\right] \le E\left[(T + cr_1 + cr_2)^{2k}\right]$$

• Let  $R_1 = \frac{1}{\sqrt{d}}(r_1, r_2, \dots, r_d)$ . Thus we have

$$E\left[Q_{1}(\alpha)^{2k}\right] = \frac{1}{d^{k}} \sum_{R} E\left[\left(R + \alpha_{1}r_{1} + \alpha_{2}r_{2}\right)^{2k}\right] \Pr\left[\sum_{i=3}^{d} \alpha_{i}r_{i} = \frac{R}{\sqrt{d}}\right]$$

$$\leq \frac{1}{d^{k}} \sum_{R} E\left[\left(R + cr_{1} + cr_{2}\right)^{2k}\right] \Pr\left[\sum_{i=3}^{d} \alpha_{i}r_{i} = \frac{R}{\sqrt{d}}\right]$$

$$= E\left[Q_{1}(\theta)^{2k}\right]$$

where  $c = \sqrt{(\alpha_1^2 + \alpha_2^2)/2}$ 

•  $\theta$  is a more "uniform" unit vector than  $\alpha$ .

#### **Proving the Normal Bound Lemma**

• Let  $\{T_i\}_{i=1}^d$  be i.i.d. normal r.v.. By stability of normal distribution

$$T = \frac{1}{d}\sum_{i=1}^{d} T_i \sim N(0, 1/d)$$

۲

#### Discussion

### **Proving the Normal Bound Lemma**

• Let  $\{T_i\}_{i=1}^d$  be i.i.d. normal r.v.. By stability of normal distribution

$$T=rac{1}{d}\sum\limits_{i=1}^d T_i\sim N(0,1/d)$$
 We also have  $Q_1(w)=rac{1}{d}\sum\limits_{i=1}^d r_1$ 

$$E[Q_1^{2k}(w)] = \frac{1}{d^{2k}} \sum_{i_1=1}^d \dots \sum_{i_{2k}=1}^d E[r_{i_1} \dots r_{i_{2k}}]$$
$$E[T^{2k}] = \frac{1}{d^{2k}} \sum_{i_1=1}^d \dots \sum_{i_{2k}=1}^d E[T_{i_1} \dots T_{i_{2k}}]$$

## **Proving the Normal Bound Lemma**

• Let  $\{T_i\}_{i=1}^d$  be i.i.d. normal r.v.. By stability of normal distribution

$$T = rac{1}{d}\sum_{i=1}^{d}T_i \sim N(0, 1/d)$$

• We also have 
$$Q_1(w) = rac{1}{d}\sum_{i=1}^d r_1$$

$$E[Q_1^{2k}(w)] = \frac{1}{d^{2k}} \sum_{i_1=1}^d \dots \sum_{i_{2k}=1}^d E[r_{i_1} \dots r_{i_{2k}}]$$
$$E[T^{2k}] = \frac{1}{d^{2k}} \sum_{i_1=1}^d \dots \sum_{i_{2k}=1}^d E[T_{i_1} \dots T_{i_{2k}}]$$

• For each index assignment we have

$$E[r_{i_1}\ldots r_{i_{2k}}] \leq E[T_{i_1}\ldots T_{i_{2k}}]$$











### **Open questions**

- Plenty !
- No-flattening results for other  $l_p$  metrics, non metrics



## **Open questions**

- Plenty !
- No-flattening results for other  $l_p$  metrics, non metrics
- Embeddability of non-metrics into metric spaces useful in databases, learning

SIGTACS Seminar Series

22 / 23



## **Open questions**

- Plenty !
- No-flattening results for other  $l_p$  metrics, non metrics
- Embeddability of non-metrics into metric spaces useful in databases, learning
- Information Theoretic Metrics KL, Bhattacharyya, Mahalanobis widely used



## THANK YOU

