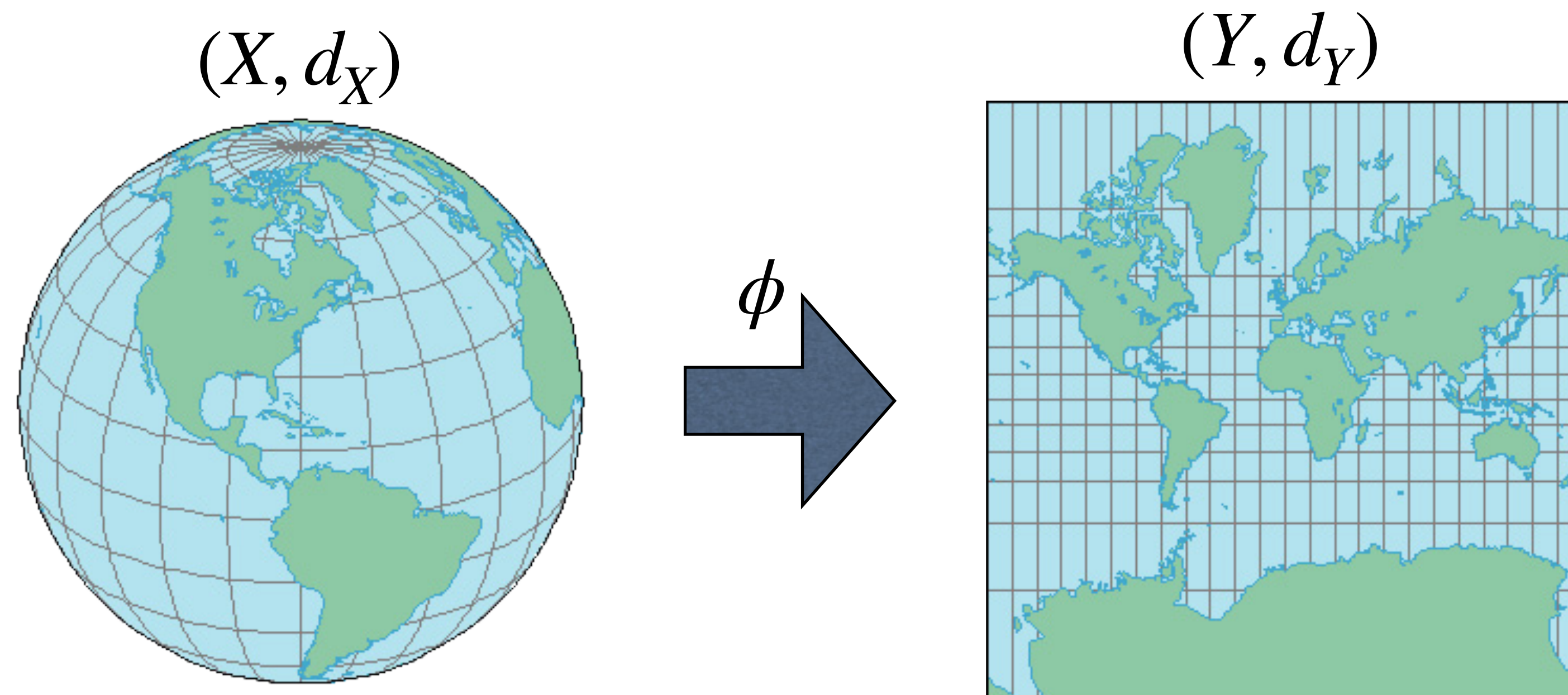# Advanced Algorithms

**Dimensionality Reduction**

刘明谋  **Nanjing University, Suzhou, 2025**

# Metric Embedding

- Two metric spaces: $(X, d_X)$ and $(Y, d_Y)$

$(X, d_X)$ $\qquad\qquad\qquad\qquad$ $(Y, d_Y)$



$\phi$

low-distortion: for small $\alpha \geq 1$

$$\forall x_1, x_2 \in X : \quad \frac{1}{\alpha} d_X(x_1, x_2) \leq d_Y(\phi(x_1), \phi(x_2)) \leq \alpha d_X(x_1, x_2)$$

# Dimension Reduction

Input: $n$ points $\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_n \in \mathbb{R}^{\color{red}d}$
Output: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \mathbb{R}^{\color{red}k}$ s.t. $\forall 1 \leq i, j \leq n$ :

$$(1 - \epsilon)\|\boldsymbol{x}_i - \boldsymbol{x}_j\| \leq \|\boldsymbol{y}_i - \boldsymbol{y}_j\| \leq (1 + \epsilon)\|\boldsymbol{x}_i - \boldsymbol{x}_j\|$$

- Usually we want $k \ll d$.

- How small can $k$ be?

- For what distance $\| \cdot \|$?
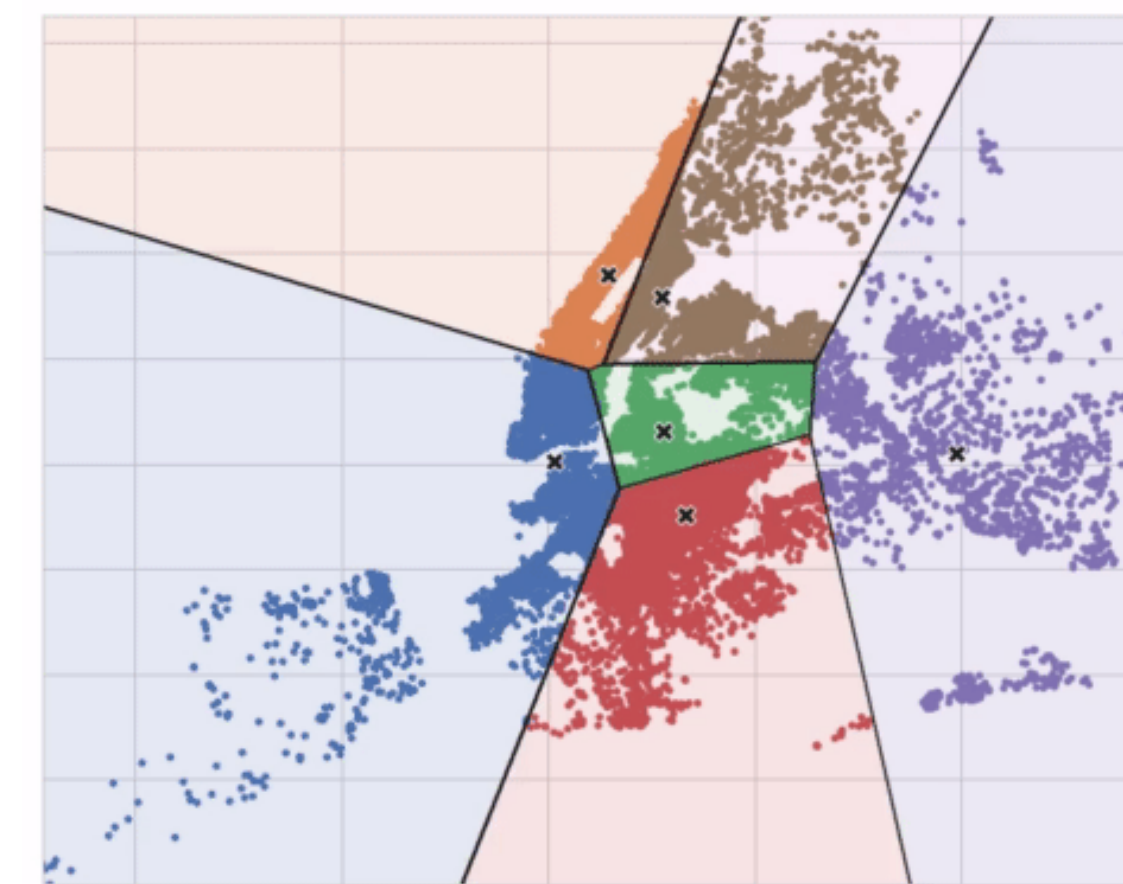
- The embedding should be efficiently constructible.

# $k$-means Clustering

**Input:** database $X = \{x_1, \ldots, x_n\} \subset \mathbb{R}^d$ and parameter $k \in \mathbb{N}^+$

**Output:** $y_1, \ldots, y_k \in \mathbb{R}^d$ minimizing

$$\sum_i \min_{j \in [k]} \|x_i - y_i\|_2^2$$

- Equivalently, $\min_{k\text{-parts } \mathscr{P}} \min_{y_1 \in \mathscr{P}_1, \ldots} \sum_{j \in [k]} \sum_{i \in \mathscr{P}_j} \|x_i - y_j\|_2^2$

- Fix $\mathscr{P}$, minimizer $y_j = \sum_{i \in \mathscr{P}_j} x_i \Big/ |\mathscr{P}_j|$ is the centroid

- Plug in, $\min_{k\text{-parts } \mathscr{P}} \sum_{j \in [k]} \sum_{i \neq i' \in \mathscr{P}_j} \|x_i - x_{i'}\|_2^2$



Voronoi partition

# Johnson-Lindenstrauss Theorem/Transformation (JLT)

# Johnson-Lindenstrauss Theorem

## (Johnson-Lindenstrauss 1984)

"In **Euclidian** space, it is always possible to embed
a set of $n$ points in *arbitrary* dimension to
$O(\log n)$ dimension with constant distortion."

**Theorem** (Johnson-Lindenstrauss 1984):

$\forall 0 < \epsilon < 1$, for any set $S$ of $n$ points from $\mathbb{R}^d$, there is a
$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ with $k = O(\epsilon^{-2} \log n)$, such that $\forall \boldsymbol{x}, \boldsymbol{y} \in S$ :

$$(1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|\phi(\boldsymbol{x}) - \phi(\boldsymbol{y})\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2$$

# Johnson-Lindenstrauss Theorem

**(Johnson-Lindenstrauss 1984)**

"In **Euclidian** space, it is always possible to embed
a set of $n$ points in *arbitrary* dimension to
$O(\log n)$ dimension with constant distortion."

**Theorem** (Johnson-Lindenstrauss 1984):

$\forall 0 < \epsilon < 1$, for any set $S$ of $n$ points from $\mathbb{R}^d$, there is a
$A \in \mathbb{R}^{k \times d}$ with $k = O(\epsilon^{-2} \log n)$, such that $\forall \boldsymbol{x}, \boldsymbol{y} \in S$ :

$$(1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \le \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \le (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2$$

- The **probabilistic method**: for random $A \in \mathbb{R}^{k \times d}$

$$\Pr\left[\forall \boldsymbol{x}, \boldsymbol{y} \in S : (1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \le \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \le (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2\right] = 1 - O\left(\frac{1}{n}\right)$$

w.h.p.

> **Theorem** (Johnson-Lindenstrauss 1984):
>
> $\forall 0 < \epsilon < 1$, for any set $S$ of $n$ points from $\mathbb{R}^d$, there is a $A \in \mathbb{R}^{k \times d}$ with $k = O(\epsilon^{-2} \log n)$, such that $\forall \boldsymbol{x}, \boldsymbol{y} \in S$ :
>
> $$(1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2$$

- The **probabilistic method**: for random $A \in \mathbb{R}^{k \times d}$

$$\Pr\left[\forall \boldsymbol{x}, \boldsymbol{y} \in S : (1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2\right] = 1 - O\left(\frac{1}{n}\right)$$

- **Efficient construction** of random $A \in \mathbb{R}^{k \times d}$:

  - projection onto uniform random $k$-dimensional subspace; (Johnson-Lindenstrauss; Dasgupta-Gupta)

  - independent Gaussian entries; (Indyk-Motwani)

  - i.i.d. -1/+1 entries; (Achlioptas)

# Dimension Reduction

**Input**: $n$ points $\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_n \in \mathbb{R}^{\textcolor{red}{d}}$

**Output**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \mathbb{R}^{\textcolor{red}{k}}$ s.t. $\forall 1 \leq i, j \leq n$ :

$$(1 - \epsilon)\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2 \leq \|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2$$

- for some suitable $k = O(\epsilon^{-2} \log n)$:

**J-L Transformation** (i.i.d. Gaussian entries)**:**

Entries of $A \in \mathbb{R}^{k \times d}$ are chosen i.i.d. from $\mathcal{N}(0, 1/k)$;
  (Gaussian distribution with mean 0 and variance $1/k$)

For $i = 1, 2, \ldots, n$: let $\boldsymbol{y}_i = A\boldsymbol{x}_i$;

- **Gaussian random variable $X \sim \mathcal{N}(\mu, \sigma^2)$:**

$$\Pr[X \leq t] = \int_{-\infty}^{t} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \, \mathrm{d}x \qquad \mathbb{E}[X] = \mu$$

$$\mathbf{Var}[X] = \sigma^2$$

# Norm Preservation

- $\forall 0 \le \epsilon \le 1, \forall$ set $S$ of $n$ points from $\mathbb{R}^d$
- Random matrix $A \in \mathbb{R}^{k \times d}$ with $k = (\epsilon^{-2} \log n)$:

**Johnson-Lindenstrauss Theorem**:

With high probability ($\ge 1 - O(1/n)$): $\forall \boldsymbol{x}, \boldsymbol{y} \in S$,

$$(1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \le \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \le (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2$$

$$1 - \epsilon \le \left\| A \frac{(\boldsymbol{x} - \boldsymbol{y})}{\|\boldsymbol{x} - \boldsymbol{y}\|_2} \right\|_2^2 \le 1 + \epsilon$$

**unit vector!**

union bound over all $O(n^2)$ pairs of $\boldsymbol{x}, \boldsymbol{y} \in S$

for any **unit vector** $\boldsymbol{u} \in \mathbb{R}^d$:

$$\Pr\left[ \left| \|A\boldsymbol{u}\|_2^2 - 1 \right| > \epsilon \right] < \frac{1}{n^3}$$

$A \in \mathbb{R}^{k \times d}$ : each entry of $A$ is chosen i.i.d. from $\mathcal{N}\left(0, \frac{1}{k}\right)$

for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :
$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] < \frac{1}{n^3}$$

$$\|A\boldsymbol{u}\|_2^2 = \sum_{i=1}^{k} (A\boldsymbol{u})_i^2 \quad \xrightarrow{\text{linearity of expectation}} \quad \mathbb{E}\left[\|A\boldsymbol{u}\|_2^2\right] = \sum_{i=1}^{k} \mathbb{E}\left[(A\boldsymbol{u})_i^2\right]$$

$$(A\boldsymbol{u})_i^2 = \left(\sum_{j=1}^{d} A_{ij} u_j\right)^2 \quad \text{each} \ A_{ij} \sim \mathcal{N}\left(0, \frac{1}{k}\right) \ \text{i.i.d.}$$

recall:

$$X \sim \mathcal{N}\left(\mu_1, \sigma_1^2\right), Y \sim \mathcal{N}\left(\mu_2, \sigma_2^2\right) \implies X + Y \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

$$\Longrightarrow \quad (A\boldsymbol{u})_i = \sum_{j=1}^{d} A_{ij} u_j \sim \mathcal{N}\left(0, \frac{\sum_{j=1}^{d} u_j^2}{k}\right) = \mathcal{N}\left(0, \frac{1}{k}\right)$$

independently!

$A \in \mathbb{R}^{k \times d}$ : each entry of $A$ is chosen i.i.d. from $\mathcal{N}\left(0, \frac{1}{k}\right)$

for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :

$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] < \frac{1}{n^3}$$

$$\|A\boldsymbol{u}\|_2^2 = \sum_{i=1}^{k}(A\boldsymbol{u})_i^2 \quad \xrightarrow{\text{linearity of expectation}} \quad \mathbb{E}\left[\|A\boldsymbol{u}\|_2^2\right] = \sum_{i=1}^{k}\mathbb{E}\left[(A\boldsymbol{u})_i^2\right]$$

$$(A\boldsymbol{u})_i \sim \mathcal{N}\left(0, \frac{1}{k}\right) \text{ i.i.d.}$$

$$\Longrightarrow \quad \mathbb{E}\left[(A\boldsymbol{u})_i^2\right] = \mathbf{Var}[(A\boldsymbol{u})_i] + \mathbb{E}\left[(A\boldsymbol{u})_i\right]^2 = \frac{1}{k}$$

$$\Longrightarrow \quad \mathbb{E}\left[\|A\boldsymbol{u}\|_2^2\right] = \sum_{i=1}^{k}\mathbb{E}\left[(A\boldsymbol{u})_i^2\right] = 1$$

$A \in \mathbb{R}^{k \times d}$ : each entry of $A$ is chosen i.i.d. from $\mathcal{N}\left(0, \frac{1}{k}\right)$

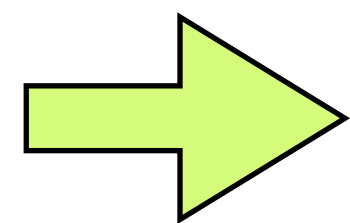for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :
$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] < \frac{1}{n^3}$$

$$\|A\boldsymbol{u}\|_2^2 = \sum_{i=1}^{k} (A\boldsymbol{u})_i^2$$

$(A\boldsymbol{u})_i \sim \mathcal{N}\left(0, \frac{1}{k}\right)$ i.i.d.

$$\mathbb{E}\left[\|A\boldsymbol{u}\|_2^2\right] = 1$$

$\Longrightarrow$ for i.i.d. $Y_1, Y_2, \ldots, Y_k \sim \mathcal{N}\left(0, \frac{1}{k}\right)$

$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] = \Pr\left[\left|\sum_{i=1}^{k} Y_i^2 - \mathbb{E}\left[\sum_{i=1}^{k} Y_i^2\right]\right| > \epsilon\right]$$

$A \in \mathbb{R}^{k \times d}$ : each entry of $A$ is chosen i.i.d. from $\mathcal{N}\left(0, \frac{1}{k}\right)$

for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :

$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] < \frac{1}{n^3}$$

for i.i.d. $Y_1, Y_2, \ldots, Y_k \sim \mathcal{N}\left(0, \frac{1}{k}\right)$    consider $X_i = \sqrt{k} \cdot Y_i$

$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] = \Pr\left[\left|\sum_{i=1}^{k} Y_i^2 - \mathbb{E}\left[\sum_{i=1}^{k} Y_i^2\right]\right| > \epsilon\right]$$

**Chernoff bound for $\chi^2$-distributions:**

For independent
$X_1, \ldots, X_k \in \mathcal{N}(0,1) \implies \begin{cases} \Pr\left[\sum_{i=1}^{k} X_i^2 > (1+\epsilon)k\right] < e^{-\epsilon^2 k/8} \\ \Pr\left[\sum_{i=1}^{k} X_i^2 < (1-\epsilon)k\right] < e^{-\epsilon^2 k/8} \end{cases}$

$A \in \mathbb{R}^{k \times d}$ : each entry of $A$ is chosen i.i.d. from $\mathcal{N}\left(0, \frac{1}{k}\right)$

for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :
$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] < \frac{1}{n^3}$$

for i.i.d. $X_1, X_2, \ldots, X_k \sim \mathcal{N}(0,1)$

$$\Pr\left[\left|\|A\boldsymbol{u}\|_2^2 - 1\right| > \epsilon\right] = \Pr\left[\sum_{i=1}^{k} X_i^2 > (1+\epsilon)k \text{ or } \sum_{i=1}^{k} X_i^2 < (1-\epsilon)k\right]$$

$< \frac{1}{n^3}$ for suitable $k = O(\epsilon^{-2} \log n)$

**Chernoff bound for $\chi^2$-distributions:**

For independent
$X_1, \ldots, X_k \in \mathcal{N}(0,1) \implies$
$$\begin{cases} \Pr\left[\sum_{i=1}^{k} X_i^2 > (1+\epsilon)k\right] < e^{-\epsilon^2 k/8} \\ \Pr\left[\sum_{i=1}^{k} X_i^2 < (1-\epsilon)k\right] < e^{-\epsilon^2 k/8} \end{cases}$$

For independent
$X_1, \ldots, X_k \in \mathcal{N}(0,1) \implies$
$$\begin{cases} \Pr\left[ \sum_{i=1}^{k} X_i^2 > (1+\epsilon)k \right] < e^{-\epsilon^2 k/8} \\ \Pr\left[ \sum_{i=1}^{k} X_i^2 < (1-\epsilon)k \right] < e^{-\epsilon^2 k/8} \end{cases}$$

for all $\lambda$>0: $\Pr\left[ \sum_{i=1}^{k} X_i^2 > \right. \left. \right]$

$$\leq e^{-(1+\epsilon)\lambda k} \cdot \mathbb{E}\left[ e^{\lambda \sum_{i=1}^{k} } \right]$$

$$X \sim \mathcal{N}(0,1)$$

$$\mathbb{E}\left[ e^{sX^2} \right] = \frac{1}{\sqrt{1-2s}}$$

$$\mathbb{E}\left[ e^{+sX^2} \right]$$
$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{sx^2} \cdot e^{-x^2/2} \, dx$$
$$= \frac{1}{\sqrt{2\pi}} \int e^{-(1-2s)x^2/2} \, dx$$

Set $y = \sqrt{1-2s}\, x$

$$= \frac{1}{\sqrt{1-2s}} \underbrace{\frac{1}{\sqrt{2\pi}} \int e^{-y^2/2} \, dy}_{= 1} = \frac{1}{\sqrt{1-2s}}$$

**Chernoff bound for $\chi^2$-distributions:**

For independent
$X_1, \ldots, X_k \in \mathcal{N}(0,1)$ $\implies$
$$\begin{cases} \Pr\left[ \sum_{i=1}^{k} X_i^2 > (1+\epsilon)k \right] < e^{-\epsilon^2 k/8} \\ \Pr\left[ \sum_{i=1}^{k} X_i^2 < (1-\epsilon)k \right] < e^{-\epsilon^2 k/8} \end{cases}$$

for all $\lambda > 0$: $\Pr\left[ \sum_{i=1}^{k} X_i^2 > (1+\epsilon)k \right] = \Pr\left[ e^{\lambda \sum_{i=1}^{k} X_i^2} > e^{(1+\epsilon)\lambda k} \right]$

$\leq e^{-(1+\epsilon)\lambda k} \cdot \mathbb{E}\left[ e^{\lambda \sum_{i=1}^{k} X_i^2} \right] = e^{-(1+\epsilon)\lambda k} \cdot \prod_{i=1}^{k} \mathbb{E}\left[ e^{\lambda X_i^2} \right]$

$X \sim \mathcal{N}(0,1)$ ⟹

$\mathbb{E}\left[ e^{sX^2} \right] = \dfrac{1}{\sqrt{1-2s}}$

$\leq e^{-\epsilon \lambda k} \left( \dfrac{e^{-\lambda}}{\sqrt{1-2\lambda}} \right)^k$

$\leq e^{-\epsilon \lambda k + 2\lambda^2 k}$    for $\lambda < 1/4$

$= e^{-\epsilon^2 k/8}$    choosing $\lambda = \epsilon/4$

# Johnson-Lindenstrauss Theorem

**(Johnson-Lindenstrauss 1984)**

"In **Euclidian** space, it is always possible to embed
a set of $n$ points in *arbitrary* dimension to
$O(\log n)$ dimension with constant distortion."

**Theorem** (Johnson-Lindenstrauss 1984):

$\forall 0 < \epsilon < 1$, for any set $S$ of $n$ points from $\mathbb{R}^d$, there is a
$A \in \mathbb{R}^{k \times d}$ with $k = O(\epsilon^{-2} \log n)$, such that $\forall \boldsymbol{x}, \boldsymbol{y} \in S :$

$$(1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2$$

- The **probabilistic method**: for random $A \in \mathbb{R}^{k \times d}$

$$\Pr\left[\forall \boldsymbol{x}, \boldsymbol{y} \in S : (1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2\right] = 1 - O\left(\frac{1}{n}\right)$$

w.h.p.

**Theorem** (Johnson-Lindenstrauss 1984):

$\forall 0 < \epsilon < 1$, for any set $S$ of $n$ points from $\mathbb{R}^d$, there is a $A \in \mathbb{R}^{k \times d}$ with $k = O(\epsilon^{-2} \log n)$, such that $\forall x, y \in S$:

$$(1 - \epsilon)\|x - y\|_2^2 \le \|Ax - Ay\|_2^2 \le (1 + \epsilon)\|x - y\|_2^2$$

- The probabilistic method: for random $A \in \mathbb{R}^{k \times d}$

$$\Pr\left[\forall x, y \in S : (1 - \epsilon)\|x - y\|_2^2 \le \|Ax - Ay\|_2^2 \le (1 + \epsilon)\|x - y\|_2^2\right] = 1 - O\left(\frac{1}{n}\right)$$

- **Efficient construction** of random $A \in \mathbb{R}^{k \times d}$:

  - projection onto uniform random $k$-dimensional subspace; (Johnson-Lindenstrauss; Dasgupta-Gupta)

  - independent Gaussian entries; (Indyk-Motwani)

  - i.i.d. -1/+1 entries; (Achlioptas)

# Dimension Reduction

**Input**: $n$ points $\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_n \in \mathbb{R}^{\textcolor{red}{d}}$

**Output**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \mathbb{R}^{\textcolor{red}{k}}$ s.t. $\forall 1 \leq i, j \leq n$ :

$$(1 - \epsilon)\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2 \leq \|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2$$
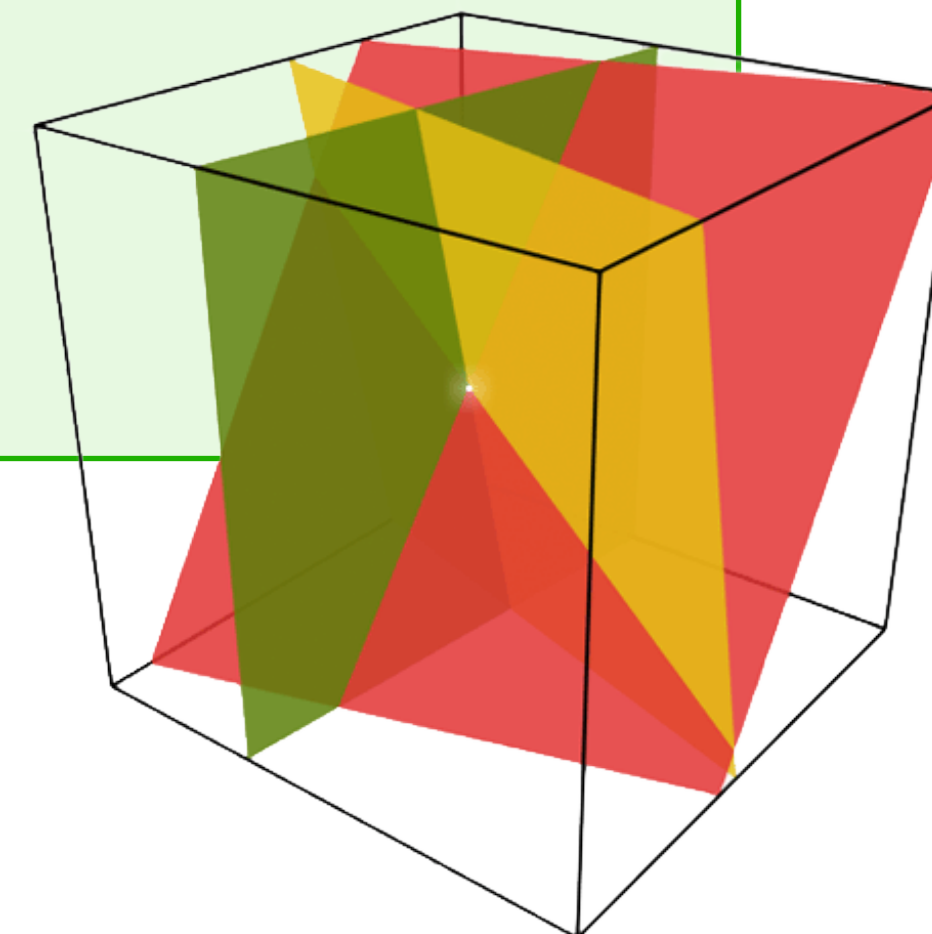
- for some suitable $k = O(\epsilon^{-2} \log n)$:

**J-L Transformation** (uniform $k$-dim subspace)**:**

The $k$ rows $A_1, \ldots, A_k$ of $A \in \mathbb{R}^{k \times d}$ are **orthogonal unit** vectors $\in \mathbb{R}^d$ chosen uniform at random;

For $i = 1, 2, \ldots, n$: let $\boldsymbol{y}_i = \sqrt{\dfrac{d}{k}} A \boldsymbol{x}_i$;
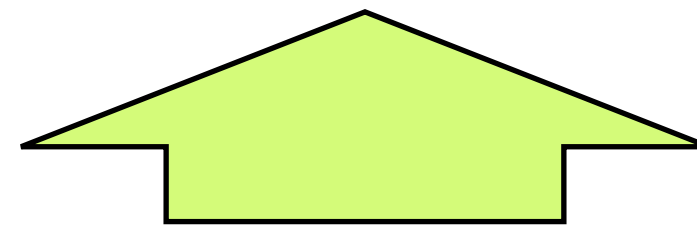


- $A \in \mathbb{R}^{k \times d}$ : projection onto a uniform $k$-dimensional subspace

$A \in \mathbb{R}^{k \times d}$ : projection onto uniform random $k$-dim subspace

for any unit vector $u \in \mathbb{R}^d$ :

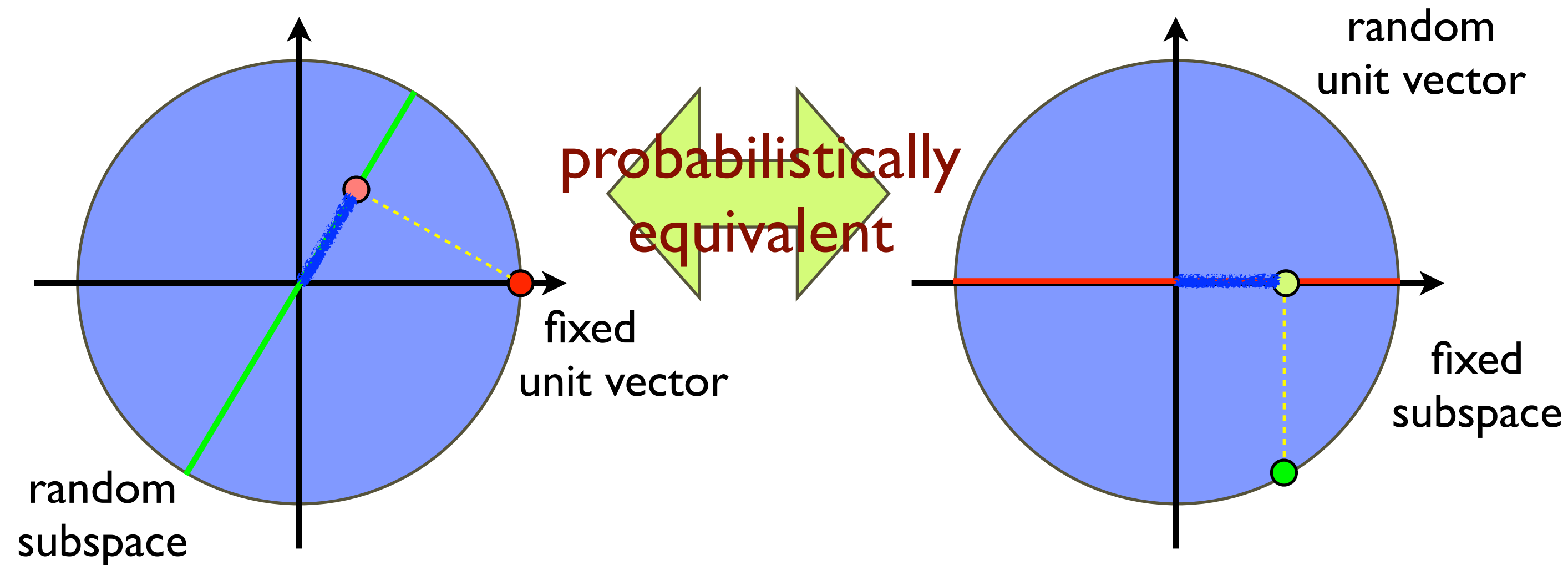$$\Pr\left[\left|\left\|\sqrt{\frac{d}{k}}Au\right\|_2^2 - 1\right| > \epsilon\right] < \frac{1}{n^3}$$

$$\Pr\left[\|Au\|_2^2 > (1+\epsilon)\frac{k}{d} \text{ or } \|Au\|_2^2 < (1-\epsilon)\frac{k}{d}\right] < \frac{1}{n^3}$$

$A \in \mathbb{R}^{k \times d}$ : projection onto uniform random $k$-dim subspace

for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :

$$\Pr\left[\ \|A\boldsymbol{u}\|_2^2 > (1+\epsilon)\frac{k}{d}\ \right] < \frac{1}{2n^3}$$

$$\Pr\left[\ \|A\boldsymbol{u}\|_2^2 < (1-\epsilon)\frac{k}{d}\ \right] < \frac{1}{2n^3}$$



probabilistically equivalent

random subspace

fixed unit vector

random unit vector

fixed subspace

*"inner-products are invariant under rotations"*

$A \in \mathbb{R}^{k \times d}$ : projection onto uniform random $k$-dim subspace

for any unit vector $\boldsymbol{u} \in \mathbb{R}^d$ :

$$\Pr\left[ \| A\boldsymbol{u} \|_2^2 > (1 + \epsilon)\frac{k}{d} \right] < \frac{1}{2n^3}$$
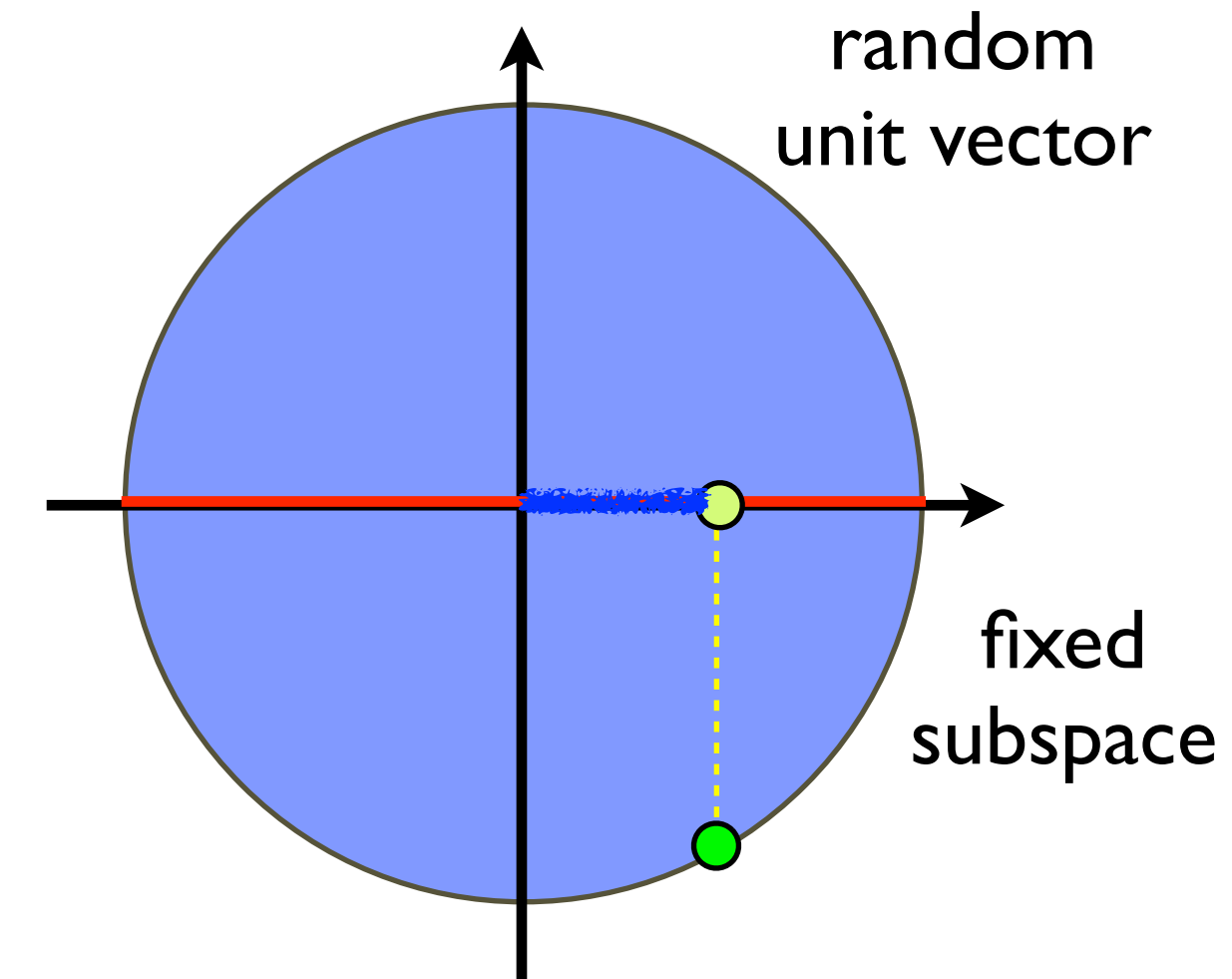
$$\Pr\left[ \| A\boldsymbol{u} \|_2^2 < (1 - \epsilon)\frac{k}{d} \right] < \frac{1}{2n^3}$$

uniform random unit vector $\in \mathbb{R}^d$:

$$Y = (Y_1, \ldots, Y_k, Y_{k+1}, \ldots, Y_d)$$

fixed $k$-dimensional subspace:

$$Z = (Y_1, \ldots, Y_k)$$



random unit vector

fixed subspace

$\|A\boldsymbol{u}\|$ is identically distributed as $\|Z\|$

uniform random unit vector $\in \mathbb{R}^d$:     $Y = (Y_1, \ldots, Y_d)$

sample $X = (X_1, \ldots, X_d) \in \mathbb{R}^d$ where $X_i \sim \mathcal{N}(0,1)$ *i.i.d.*;

$$\text{let } Y = \frac{X}{\|X\|};$$

density:
$$\Pr[X = x] = \prod_{i=1}^{d} \frac{1}{\sqrt{2\pi}} e^{-x_i^2/2} = (2\pi)^{-d/2} e^{-\|x\|_2^2/2}$$

Spherically symmetric!
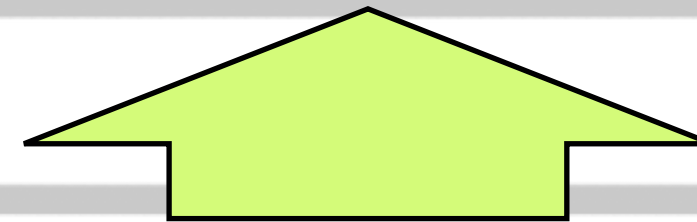
for some suitable $k = O(\varepsilon^{-2}\log n)$:

$$\Pr\left[\sum_{i=1}^{k} Y_i^2 > (1 + \epsilon)\frac{k}{d}\right] < \frac{1}{2n^3}$$

$$\Pr\left[\sum_{i=1}^{k} Y_i^2 < (1 - \epsilon)\frac{k}{d}\right] < \frac{1}{2n^3}$$

i.i.d. Gaussian random variables $X_1, X_2, \ldots, X_d \sim \mathcal{N}(0,1)$

for some suitable $k = \mathrm{O}(\varepsilon^{-2}\log n)$:

$$\Pr\left[\sum_{i=1}^{k} X_i^2 > (1+\epsilon)\frac{k}{d}\sum_{i=1}^{d} X_i^2\right] < \frac{1}{2n^3}$$

$$\Pr\left[\sum_{i=1}^{k} X_i^2 < (1-\epsilon)\frac{k}{d}\sum_{i=1}^{d} X_i^2\right] < \frac{1}{2n^3}$$

$$\Pr\left[(d-(1+\epsilon)k)\sum_{i=1}^{k} X_i^2 - (1+\epsilon)k\sum_{i=k+1}^{d} X_i^2 > 0\right] < \frac{1}{2n^3}$$

$$\Pr\left[(d-(1-\epsilon)k)\sum_{i=1}^{k} X_i^2 - (1-\epsilon)k\sum_{i=k+1}^{d} X_i^2 < 0\right] < \frac{1}{2n^3}$$

$$\mathbb{E}\left[e^{+sX^2}\right]$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{sx^2} \cdot e^{-x^2/2}\, dx$$

$$= \frac{1}{\sqrt{2\pi}} \int e^{-(1-2s)x^2/2}\, dx$$

set $y = \sqrt{1-2s}\, x$

$$= \frac{1}{\sqrt{1-2s}} \cdot \underbrace{\frac{1}{\sqrt{2\pi}} \int e^{-y^2/2}\, dy}_{=\,1} = \frac{1}{\sqrt{1-2s}}$$

$X_1, X_2, \ldots, X_d \sim \mathcal{N}(0,1)$

$n$):

$$\left. \epsilon)k \sum_{i=k+1}^{d} X_i^2 > 0 \right] < \frac{1}{2n^3}$$

$$\left. \epsilon)k \sum_{i=k+1}^{d} X_i^2 < 0 \right] < \frac{1}{2n^3}$$

$$\blacksquare = \Pr\left[ \exp\left\{ \lambda\left( ((1-\epsilon)k - d) \sum_{i=1}^{k} X_i^2 + (1-\epsilon)k \sum_{i=k+1}^{d} X_i^2 \right) \right\} > 1 \right] \quad \text{(arbitrary } \lambda > 0\text{)}$$

$$\leq \mathbb{E}\left[ \exp\left\{ \lambda\left( ((1-\epsilon)k - d) \sum_{i=1}^{k} X_i^2 + (1-\epsilon)k \sum_{i=k+1}^{d} X_i^2 \right) \right\} \right] \quad \text{(Markov's inequality)}$$

$$= \prod_{i=1}^{k} \mathbb{E}\left[ e^{\lambda((1-\epsilon)k - d)X_i^2} \right] \prod_{i=k+1}^{d} \mathbb{E}\left[ e^{\lambda(1-\epsilon)kX_i^2} \right]$$

i.i.d. Gaussian random variables $X_1, X_2, \ldots, X_d \sim \mathcal{N}(0,1)$

for some suitable $k = O(\varepsilon^{-2} \log n)$:

$$\Pr\left[(d - (1+\epsilon)k)\sum_{i=1}^{k} X_i^2 - (1+\epsilon)k\sum_{i=k+1}^{d} X_i^2 > 0\right] < \frac{1}{2n^3}$$

$$\Pr\left[(d - (1-\epsilon)k)\sum_{i=1}^{k} X_i^2 - (1-\epsilon)k\sum_{i=k+1}^{d} X_i^2 < 0\right] < \frac{1}{2n^3}$$

$$= \Pr\left[\exp\left\{\lambda\left(((1-\epsilon)k - d)\sum_{i=1}^{k} X_i^2 + (1-\epsilon)k\sum_{i=k+1}^{d} X_i^2\right)\right\} > 1\right] \quad \text{(arbitrary } \lambda > 0)$$

$$\leq (1 - 2\lambda((1-\epsilon)k - d))^{-\frac{k}{2}}(1 - 2\lambda(1-\epsilon)k)^{-\frac{d-k}{2}}$$

$$\leq (1-\epsilon)^{-\frac{k}{2}}\left(1 - \frac{\epsilon k}{d-k}\right)^{\frac{d-k}{2}} \leq \exp\left(-\frac{\epsilon^2 k}{4}\right)$$

set:
$$\lambda = \frac{\epsilon}{2(1-\epsilon)(d - (1-\epsilon)k)}$$

**Theorem** (Johnson-Lindenstrauss 1984):

$\forall 0 < \epsilon < 1$, for any set $S$ of $n$ points from $\mathbb{R}^d$, there is a $A \in \mathbb{R}^{k \times d}$ with $k = O(\epsilon^{-2} \log n)$, such that $\forall \boldsymbol{x}, \boldsymbol{y} \in S$:

$$(1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2$$

- The **probabilistic method**: for random $A \in \mathbb{R}^{k \times d}$

$$\Pr\left[ \forall \boldsymbol{x}, \boldsymbol{y} \in S : (1 - \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \leq \|A\boldsymbol{x} - A\boldsymbol{y}\|_2^2 \leq (1 + \epsilon)\|\boldsymbol{x} - \boldsymbol{y}\|_2^2 \right] = 1 - O\left(\frac{1}{n}\right)$$

- **Efficient construction** of random $A \in \mathbb{R}^{k \times d}$:

  - projection onto uniform random $k$-dimensional subspace; (Johnson-Lindenstrauss; Dasgupta-Gupta)

  - independent Gaussian entries; (Indyk-Motwani)

  - i.i.d. -1/+1 entries; (Achlioptas)

# Tug-Of-War Algorithm

**Count Sketch:** $z$

**Upon** each $x_i$: $z \leftarrow z + \sigma(x_i)$

**Query:** return $z^2$

**Chebyshev's Inequality**

For random variable $X$, for any $t > 0$,

$$\Pr\left[\,|X - \mathbb{E}[X]| \geq t\,\right] \leq \frac{\mathbf{Var}[X]}{t^2}$$

for any unit vector $\mathbf{u} \in \mathbb{R}^d$ :

$$\Pr\left[\,\left|\,\|A\mathbf{u}\|_2^2 - 1\,\right| \geq \epsilon\,\right] \leq \delta$$

- Unbiased: $\mathbb{E}[z^2] = \|x\|_2^2$
- Claim: $\mathbf{Var}(z^2) = O(\|x\|_2^2)$ for $A \in \{-1, +1\}^{1 \times d}$
- $\mathbf{Var}(z_+^2) = O(\|x\|_2^2/k)$ for $A \in \left\{-1/\sqrt{k}, +1/\sqrt{k}\right\}^{k \times d}$  *Any issue?*
- $\Pr\left[\,\left|\,\|A\mathbf{u}\|_2^2 - 1\,\right| \geq \epsilon\,\right] \leq \mathbf{Var}(z_+^2)/\epsilon^2 = O(1/k\epsilon^2) =: \delta$

**Chernoff bound for $\chi^2$-distributions:**

For independent
$X_1, \ldots, X_k \in \mathcal{N}(0,1) \implies$
$$\begin{cases} \Pr\left[ \sum_{i=1}^{k} X_i^2 > (1+\epsilon)k \right] < e^{-\epsilon^2 k/8} \\[2mm] \Pr\left[ \sum_{i=1}^{k} X_i^2 < (1-\epsilon)k \right] < e^{-\epsilon^2 k/8} \end{cases}$$

for all $\lambda > 0$: $\Pr\left[ \sum_{i=1}^{k} X_i^2 > (1+\epsilon)k \right] = \Pr\left[ e^{\lambda \sum_{i=1}^{k} X_i^2} > e^{(1+\epsilon)\lambda k} \right]$

$\leq e^{-(1+\epsilon)\lambda k} \cdot \mathbb{E}\left[ e^{\lambda \sum_{i=1}^{k} X_i^2} \right] = e^{-(1+\epsilon)\lambda k} \cdot \prod_{i=1}^{k} \mathbb{E}\left[ e^{\lambda X_i^2} \right]$

$$\boxed{\begin{array}{c} X \sim \mathcal{N}(0,1) \implies \\[3mm] \mathbb{E}\left[ e^{sX^2} \right] = \dfrac{1}{\sqrt{1-2s}} \end{array}}$$

$\leq e^{-\epsilon \lambda k} \left( \dfrac{e^{-\lambda}}{\sqrt{1-2\lambda}} \right)^k$

$\leq e^{-\epsilon \lambda k + 2\lambda^2 k}$    for $\lambda < 1/4$

$= e^{-\epsilon^2 k/8}$    choosing $\lambda = \epsilon/4$

# Count Sketch versus Gaussian

**Count Sketch+:** $CS[k]$ (initialized to all 0's)

Upon each $x_i$: $CS[j] \leftarrow CS[j] + \sigma_j(x_i)$, for all $j \leq k$

Query: return $z_+^2 = \sum CS[j]^2 / k$

$$M_{X^2}(t) \triangleq \mathbb{E}\left[\exp(tX^2)\right] = \sum_{k \geq 0} \frac{t^k}{k!} \mathbb{E}\left[X^{2k}\right]$$

- Count Sketch is not worse than Gaussian. Let $\mathbf{w} \triangleq (1,\ldots,1), \mathbf{Z} \sim \mathcal{N}^d(0,1), \mathbf{r} \sim \{-1, +1\}^d$

  - Claim: $\mathbb{E}\left[\langle \mathbf{x}, \mathbf{r} \rangle^{2k}\right] \leq \mathbb{E}\left[\langle \mathbf{w}, \mathbf{r} \rangle^{2k}\right] \leq \mathbb{E}\left[\langle \mathbf{w}, \mathbf{Z} \rangle^{2k}\right]$ for all $k \in \mathbb{N}, \mathbf{x} \in \mathbb{R}^d$ with $\|\mathbf{x}\|_2^2 = d$

- Eq(1): intuitively provable by symmetry.

- Eq(2): $\mathbb{E}\left[\langle \mathbf{w}, \mathbf{r} \rangle^{2k}\right] = \sum_{i_1,\ldots,i_{2k}} \mathbb{E}[r_{i_1} \cdot \ldots \cdot r_{i_{2k}}]$ is zero if any multiplicity is odd

  - Suffices to consider only $\sum_{i_1,\ldots,i_{2k}} \mathbb{E}[r_{i_1}^{2\ell_1}] \cdot \ldots \cdot \mathbb{E}[r_{i_{2k}}^{2\ell_{2k}}]$, so be $\mathbb{E}\left[\langle \mathbf{w}, \mathbf{Z} \rangle^{2k}\right]$.

# Count Sketch versus Gaussian

**Count Sketch+:** $\text{CS}[k]$ (initialized to all 0's)

**Upon** each $x_i$: $\text{CS}[j] \leftarrow \text{CS}[j] + \sigma_j(x_i)$, for all $j \leq k$

**Query**: return $z_+^2 = \sum \text{CS}[j]^2 / k$

$$M_{X^2}(t) \triangleq \mathbb{E}\left[\exp(tX^2)\right] = \sum_{k \geq 0} \frac{t^k}{k!} \mathbb{E}\left[X^{2k}\right]$$

- Count Sketch is not worse than Gaussian. Let $\mathbf{w} \triangleq (1,\ldots,1), \mathbf{Z} \sim \mathcal{N}^d(0,1), \mathbf{r} \sim \{-1,+1\}^d$

  - Claim: $\mathbb{E}\left[\langle \mathbf{x}, \mathbf{r}\rangle^{2k}\right] \leq \mathbb{E}\left[\langle \mathbf{w}, \mathbf{r}\rangle^{2k}\right] \leq \mathbb{E}\left[\langle \mathbf{w}, \mathbf{Z}\rangle^{2k}\right]$ for all $k \in \mathbb{N}, \mathbf{x} \in \mathbb{R}^d$ with $\|\mathbf{x}\|_2^2 = d$

- Eq(2): Suffices to consider only $\displaystyle\sum_{i_1,\ldots,i_{2k}} \mathbb{E}\left[r_{i_1}^{2\ell_1}\right] \cdot \ldots \cdot \mathbb{E}\left[r_{i_{2k}}^{2\ell_{2k}}\right]$, so be $\mathbb{E}\left[\langle \mathbf{w}, \mathbf{Z}\rangle^{2k}\right]$.

  - Claim: $\mathbb{E}\left[r_{i_1}^{2\ell_1}\right] \leq \mathbb{E}\left[Z_{i_1}^{2\ell_1}\right]$

  - $\mathbb{E}\left[r_{i_1}^{2\ell_1}\right] = 1, \mathbb{E}\left[Z_{i_1}^{2\ell_1}\right] = \frac{(2\ell_1)!}{(\ell!2^\ell)} \geq 1$

# Nearest Neighbor Search (NNS)

# Nearest Neighbor Search (NNS)

- Metric space $(X, \text{dist})$:

**Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in X$

**Query**: a point $\boldsymbol{x} \in X$

Find the datapoint $\boldsymbol{y}_i$ that is **closest** to $\boldsymbol{x}$.

## Applications in:

- database systems
- pattern matching
- machine learning
- image processing
- bioinformatics
- … …

# Nearest Neighbor Search (NNS)

**Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in [N]^d$

**Query**: a point $\boldsymbol{x} \in [N]^d$

Find the datapoint $\boldsymbol{y}_i$ that is **closest** to $\boldsymbol{x}$.

when dimension $d$ is small:



*k-d* tree



Voronoi diagram

# Nearest Neighbor Search (NNS)

- Hamming space $\{0,1\}^d$:

<div style="border:1px solid orange; background:#fdf3df; padding:1em;">

**Data**:  $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \{0,1\}^d$

**Query**:  a point $\boldsymbol{x} \in \{0,1\}^d$

Find the datapoint $\boldsymbol{y}_i$ that is **closest** to $\boldsymbol{x}$.

</div>

when **dimension** $d$ is high:

$$\text{say } d \gg \log n$$

**Curse of dimensionality**:

It is conjectured that to solve NNS in high dimension requires either super-polynomial($n$) space or super-polynomial($d$) time.

**Blessing**:   **randomization  +  approximation**

# Approximate Near Neighbor (ANN)

- Metric space $(X, \text{dist})$:

**Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in X$

**Query**: a point $\boldsymbol{x} \in X$

$c$-**ANN** (*Approximate **Nearest** Neighbor*):

Find a $\boldsymbol{y}_i$ such that $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot \min_{1 \leq j \leq n} \text{dist}(\boldsymbol{x}, \boldsymbol{y}_j)$

$(c, r)$-**ANN** (*Approximate **Near** Neighbor*):

return a $\boldsymbol{y}_i$ that $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot r$ if $\exists \boldsymbol{y}_j$ s.t. $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_j) \leq r$

"no" if $\forall \boldsymbol{y}_i$ , $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) > c \cdot r$

arbitrary if otherwise



$$r_0 = D_{\min} = \min_{1 \leq i < j \leq n} \text{dist}(\boldsymbol{y}_i, \boldsymbol{y}_j)$$

$$r_k = \sqrt{c} \cdot r_{k-1}$$

$$r_{\log_c(D_{\max}/D_{\min})} = D_{\max} = \max_{1 \leq i < j \leq n} \text{dist}(\boldsymbol{y}_i, \boldsymbol{y}_j)$$

- Metric space $(X, \text{dist})$:

**Data**: $n$ points $y_1, y_2, \ldots, y_n \in X$

**Query**: a point $x \in X$
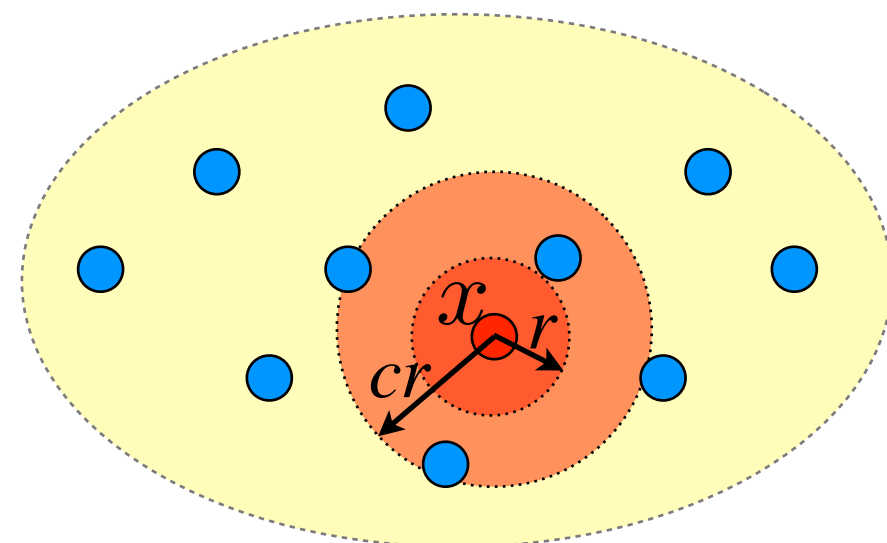
$c$-**ANN** (*Approximate **Nearest** Neighbor*):

Find a $y_i$ such that $\text{dist}(x, y_i) \leq c \cdot \min_{1 \leq j \leq n} \text{dist}(x, y_j)$

$(c, r)$-**ANN** (*Approximate **Near** Neighbor*):

return a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ if $\exists y_j$ s.t. $\text{dist}(x, y_j) \leq r$

"no" if $\forall y_i$, $\text{dist}(x, y_i) > c \cdot r$

arbitrary if otherwise

let $R = \dfrac{D_{\max}}{D_{\min}}$

$D_{\max} = \max_{1 \leq i < j \leq n} \text{dist}(y_i, y_j)$

$D_{\min} = \min_{1 \leq i < j \leq n} \text{dist}(y_i, y_j)$

$\forall r : (\sqrt{c}, r)$-**ANN** can be solved with space $s$ and query time $t$

$\Longrightarrow$

$c$-**ANN** can be solved within space $\text{O}(s \log_c R)$ and query time $\text{O}(t \log\log_c R)$

- Hamming space $\{0,1\}^d$:

**Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \{0,1\}^d$

**Query**: a point $\boldsymbol{x} \in \{0,1\}^d$

$(c, r)$-**ANN** (*Approximate **Near** Neighbor*):
  return a $\boldsymbol{y}_i$ that $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot r$ if $\exists \boldsymbol{y}_j$ s.t. $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_j) \leq r$

  answer "**no**" if $\forall \boldsymbol{y}_i$, $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_i) > c \cdot r$
  **arbitrary** if otherwise

- High dimension: $d \gg \log n$

**Dimension Reduction:**

$$z_i(j) = (A\boldsymbol{y}_i)_j = \left( \sum_{\ell=1}^{d} A_{j\ell} \boldsymbol{y}_i(\ell) \right) \mathrm{mod}\ 2$$

Let $k, p$ and $s$ to be fixed later;

sample $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in \mathrm{Bernoulli}(p)$;

for $i = 1, 2, \ldots, n$:  let $\boldsymbol{z}_i = A\boldsymbol{y}_i \in \{0,1\}^k$ on finite field $\mathrm{GF}(2)$;

store all $s$-balls $B_s(\boldsymbol{u}) = \{\boldsymbol{y}_i \mid \mathrm{dist}(\boldsymbol{u}, \boldsymbol{z}_i) \leq s\}$ for all $\boldsymbol{u} \in \{0,1\}^k$;

- **Hamming space** $\{0,1\}^d$:

  > **Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \{0,1\}^d$
  >
  > **Query**: a point $\boldsymbol{x} \in \{0,1\}^d$

- **High dimension**: $d \gg \log n$

  $$z_i(j) = (A\boldsymbol{y}_i)_j = \left( \sum_{\ell=1}^d A_{j\ell}\boldsymbol{y}_i(\ell) \right) \mathrm{mod}\ 2$$

  **Dimension Reduction:**

  Let $k$, $p$ and $s$ to be fixed later;

  sample $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in \mathrm{Bernoulli}(p)$;

  for $i = 1,2,\ldots,n$:   let $z_i = A\boldsymbol{y}_i \in \{0,1\}^k$ on finite field $\mathrm{GF}(2)$;

  store all $s$-balls $B_s(\boldsymbol{u}) = \{\boldsymbol{y}_i \mid \mathrm{dist}(\boldsymbol{u}, z_i) \leq s\}$ for all $\boldsymbol{u} \in \{0,1\}^k$;

  To answer **query** $\boldsymbol{x} \in \{0,1\}^d$:   retrieve $B_s(A\boldsymbol{x})$;

  $\qquad\qquad\qquad\qquad\qquad$ if $B_s(A\boldsymbol{x}) = \varnothing$ return "**no**"

  $\qquad\qquad\qquad\qquad\qquad$ else return any $\boldsymbol{y}_i \in B_s(A\boldsymbol{x})$

**space**: $\mathrm{O}(n2^k)$     **query time**: $\mathrm{O}(kd)$ computation $+ \mathrm{O}(1)$ memory access

- **Hamming space** $\{0,1\}^d$:

> **Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \{0,1\}^d$
>
> **Query**: a point $\boldsymbol{x} \in \{0,1\}^d$

- **High dimension**: $d \gg \log n$

**Dimension Reduction:**

$$\boldsymbol{z}_i(j) = (A\boldsymbol{y}_i)_j = \left( \sum_{\ell=1}^{d} A_{j\ell} \boldsymbol{y}_i(\ell) \right) \bmod 2$$
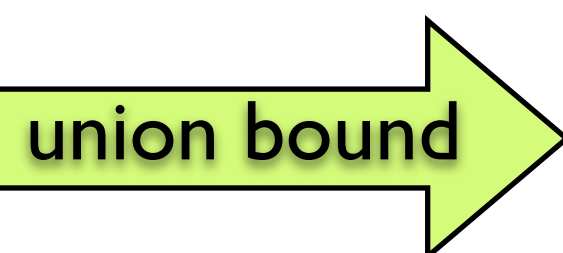
Let $k$, $p$ and $s$ to be fixed later;

sample $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in \mathrm{Bernoulli}(p)$;

for $i = 1,2,\ldots,n$:  let $\boldsymbol{z}_i = A\boldsymbol{y}_i \in \{0,1\}^k$ on finite field $\mathrm{GF}(2)$;

for suitable $k = \mathrm{O}(\log n)$, $p$ and $s$:

$\forall \boldsymbol{x}, \boldsymbol{y} \in \{0,1\}^d$:    $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}) \leq r$    $\Rightarrow$  $\Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s] < 1/n^2$

$\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}) > c \cdot r$  $\Rightarrow$  $\Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \leq s] < 1/n^2$

**union bound** ⟹ $(c, r)$-**ANN** is solved w.h.p.

random $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in$ Bernoulli($p$);

computation on GF(2):     $(A\boldsymbol{x})_i = \left( \sum_{j=1}^{d} A_{ij} x_i \right) \bmod 2$

for suitable $p$ and $s$:

$$\forall \boldsymbol{x}, \boldsymbol{y} \in \{0,1\}^d : \quad \mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}) \leq r \quad \Rightarrow \quad \Pr[\, \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s\,] < e^{-\Omega(k)}$$
$$\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\, \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \leq s\,] < e^{-\Omega(k)}$$

row vector $A_{i\cdot}$:     *i.i.d.* entries $\in$ Bernoulli($p$)

$$\Pr[(A\boldsymbol{x})_i \neq (A\boldsymbol{y})_i] = \Pr[\langle A_{i\cdot}, \boldsymbol{x} \rangle \neq \langle A_{i\cdot}, \boldsymbol{y} \rangle] = \frac{1}{2}\left(1 - (1 - 2p)^{\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y})}\right)$$

Why?

For uniform $\boldsymbol{u} \in \{0,1\}^d$:     $\Pr[\langle \boldsymbol{u}, \boldsymbol{x} \rangle \neq \langle \boldsymbol{u}, \boldsymbol{y} \rangle] = \frac{1}{2}$

generate $A_{i\cdot}$ as:

1. each $j \in [d]$ joins $D \subseteq [d]$ independently with probability $2p$;
2. for each $j \in D$: samples a uniform and independent $A_{ij} \in \{0,1\}$;
3. for each $j \notin D$: $A_{ij} = 0$;

$A_{i\cdot}$ restricted on $D$ is a uniform Boolean vector!

random $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in$ Bernoulli($p$);

computation on GF(2): $\quad (A\boldsymbol{x})_i = \left( \sum_{j=1}^{d} A_{ij} x_i \right) \bmod 2$

for suitable $p$ and $s$:

$\forall \boldsymbol{x},\boldsymbol{y} \in \{0,1\}^d: \quad \text{dist}(\boldsymbol{x},\boldsymbol{y}) \le r \quad \Rightarrow \quad \Pr[\, \text{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s\,] < e^{-\Omega(k)}$

$\text{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\, \text{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \le s\,] < e^{-\Omega(k)}$

row vector $A_{i\cdot}$: $\quad$ *i.i.d.* entries $\in$ Bernoulli($p$)

$$\Pr[(A\boldsymbol{x})_i \ne (A\boldsymbol{y})_i] = \Pr[\langle A_{i\cdot}, \boldsymbol{x}\rangle \ne \langle A_{i\cdot}, \boldsymbol{y}\rangle] = \frac{1}{2}\left(1 - (1-2p)^{\text{dist}(\boldsymbol{x},\boldsymbol{y})}\right)$$

choose $p$ to satisfy $(1-2p) = 2^{-1/r}$

$\text{dist}(\boldsymbol{x},\boldsymbol{y}) \le r \quad \Rightarrow \quad \Pr[\, (A\boldsymbol{x})_i \ne (A\boldsymbol{y})_i\,] \le 1/4$

$\text{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\, (A\boldsymbol{x})_i \ne (A\boldsymbol{y})_i\,] > 1/2 - 2^{-(c+1)}$

random $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in$ Bernoulli($p$);

computation on GF(2):  $(A\boldsymbol{x})_i = \left( \sum_{j=1}^{d} A_{ij} x_i \right) \mod 2$

for suitable $p$ and $s$:

$\forall \boldsymbol{x}, \boldsymbol{y} \in \{0,1\}^d :$   $\text{dist}(\boldsymbol{x}, \boldsymbol{y}) \leq r$   $\Rightarrow$   $\Pr[\text{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s] < e^{-\Omega(k)}$

$\text{dist}(\boldsymbol{x}, \boldsymbol{y}) > c \cdot r$   $\Rightarrow$   $\Pr[\text{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \leq s] < e^{-\Omega(k)}$

choose $p$ to satisfy $(1-2p) = 2^{-1/r}$

$\text{dist}(\boldsymbol{x}, \boldsymbol{y}) \leq r$   $\Rightarrow$   $\Pr[(A\boldsymbol{x})_i \neq (A\boldsymbol{y})_i] \leq 1/4$

$\text{dist}(\boldsymbol{x}, \boldsymbol{y}) > c \cdot r$   $\Rightarrow$   $\Pr[(A\boldsymbol{x})_i \neq (A\boldsymbol{y})_i] > 1/2 - 2^{-(c+1)}$

$\text{dist}(A\boldsymbol{x}, A\boldsymbol{y}) = X = \sum_{i=1}^{k} X_i$   where  $X_i = \begin{cases} 1 & \text{if } (A\boldsymbol{x})_i \neq (A\boldsymbol{y})_i \\ 0 & \text{otherwise} \end{cases}$

independent trials

random $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in$ Bernoulli($p$);

computation on GF(2): $\quad (A\boldsymbol{x})_i = \left( \sum_{j=1}^{d} A_{ij} x_i \right) \bmod 2$

for suitable $p$ and $s$:

$$\forall \boldsymbol{x}, \boldsymbol{y} \in \{0,1\}^d : \quad \mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) \leq r \quad \Rightarrow \quad \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s] < e^{-\Omega(k)}$$
$$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \leq s] < e^{-\Omega(k)}$$

choose $p$ to satisfy $(1-2p) = 2^{-1/r}$

$$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) \leq r \quad \Rightarrow \quad \Pr[\ X_i = 1\ ] \leq 1/4 \quad\quad \Rightarrow \quad \mathbf{E}[X] \leq k/4$$

$$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\ X_i = 1\ ] > 1/2 - 2^{-(c+1)} \quad \Rightarrow \quad \mathbf{E}[X] \leq (1/2 - 2^{-(c+1)})k$$

$$\mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) = X = \sum_{i=1}^{k} X_i \quad \textbf{where} \ \ X_i = \begin{cases} 1 & \text{if } (A\boldsymbol{x})_i \neq (A\boldsymbol{y})_i \\ 0 & \text{otherwise} \end{cases}$$

choose $s = (1/4 + 1/2 - 2^{-(c+1)})k/2 = (3/8 - 2^{-(c+2)})k$

$$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) \leq r \quad \Rightarrow \quad \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s] \leq \Pr[\ X > \mathbf{E}X + (1/8 - 2^{-(c+2)})k\ ]$$
$$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \leq s] \leq \Pr[\ X < \mathbf{E}X - (1/8 - 2^{-(c+2)})k\ ]$$

# Chernoff-Hoeffding Bound

**Chernoff Bound**:

For $X = \displaystyle\sum_{i=1}^{n} X_i$, where $X_1, \ldots, X_n \in \{0,1\}$ are *independent*

(or *negatively associated*),

for any $t > 0$:

$$\Pr\left[X \geq \mathbb{E}[X] + t\right] \leq \exp\left(-\frac{2t^2}{n}\right)$$

$$\Pr\left[X \leq \mathbb{E}[X] - t\right] \leq \exp\left(-\frac{2t^2}{n}\right)$$

random $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in$ Bernoulli($p$);

computation on GF(2): $\quad (A\boldsymbol{x})_i = \left( \sum_{j=1}^{d} A_{ij} x_i \right) \mod 2$

for suitable $p$ and $s$:

$$\forall \boldsymbol{x}, \boldsymbol{y} \in \{0,1\}^d: \quad \mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) \le r \quad \Rightarrow \quad \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s] < \mathrm{e}^{-\Omega(k)}$$
$$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \le s] < \mathrm{e}^{-\Omega(k)}$$

choose $p$ to satisfy $(1-2p) = 2^{-1/r}$

$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) \le r \quad \Rightarrow \quad \Pr[\ X_i = 1\ ] \le 1/4 \qquad \Rightarrow \quad \mathbf{E}[X] \le k/4$

$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \quad \Rightarrow \quad \Pr[\ X_i = 1\ ] > 1/2 - 2^{-(c+1)} \quad \Rightarrow \quad \mathbf{E}[X] \le (1/2 - 2^{-(c+1)})k$

$\mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) = X = \sum_{i=1}^{k} X_i \quad$ **where** $\ X_i = \begin{cases} 1 & \text{if } (A\boldsymbol{x})_i \ne (A\boldsymbol{y})_i \\ 0 & \text{otherwise} \end{cases}$

choose $s = (1/4 + 1/2 - 2^{-(c+1)})k/2 = (3/8 - 2^{-(c+2)})k$

$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) \le r \Rightarrow \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) > s] \le \Pr[\ X > \mathbf{E}X + (1/8 - 2^{-(c+2)})k\ ] \quad < \exp(-2(1/8 - 2^{-(c+2)})^2 k)$

$\mathrm{dist}(\boldsymbol{x},\boldsymbol{y}) > c \cdot r \Rightarrow \Pr[\ \mathrm{dist}(A\boldsymbol{x}, A\boldsymbol{y}) \le s] \le \Pr[\ X < \mathbf{E}X - (1/8 - 2^{-(c+2)})k\ ] \quad < \exp(-2(1/8 - 2^{-(c+2)})^2 k)$

- **Hamming space** $\{0,1\}^d$:

  > **Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in \{0,1\}^d$
  >
  > **Query**: a point $\boldsymbol{x} \in \{0,1\}^d$

- **High dimension**: $d \gg \log n$

> **Dimension Reduction:**
>
> Let $k = \frac{\ln n}{(1/8 - 2^{-(c+2)})^2}$, $p = \frac{1}{2} - 2^{-1-1/r}$, $s = \left(\frac{3}{8} - 2^{-(c+2)}\right)k$
>
> sample $k \times d$ Boolean matrix $A$ with *i.i.d.* entries $\in$ Bernoulli($p$);
>
> for $i = 1, 2, \ldots, n$: let $z_i = A\boldsymbol{y}_i \in \{0,1\}^k$ on finite field GF(2);
>
> store all $s$-balls $B_s(\boldsymbol{u}) = \{\boldsymbol{y}_i \mid \text{dist}(\boldsymbol{u}, z_i) \leq s\}$ for all $\boldsymbol{u} \in \{0,1\}^k$;
>
> To answer **query** $\boldsymbol{x} \in \{0,1\}^d$:  retrieve $B_s(A\boldsymbol{x})$;
>
>                             if $B_s(A\boldsymbol{x}) = \varnothing$ return "**no**"
>
>                             else return any $\boldsymbol{y}_i \in B_s(A\boldsymbol{x})$

**space**: $n^{O(1)}$      **query time**: $O(d \ln n)$      **solve** $(c, r)$-ANN w.h.p.

# Locality Sensitive Hashing (LSH)

# Locality-Sensitive Hashing (LSH)

(Indyk-Motwani 1998)

- Metric space $(X, \text{dist})$:

**Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in X$

**Query**: a point $\boldsymbol{x} \in X$

$(c, r)$-**ANN** (*Approximate **Near** Neighbor*):

return a $\boldsymbol{y}_i$ that $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot r$ if $\exists \boldsymbol{y}_j$ s.t. $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_j) \leq r$

"no" if $\forall \boldsymbol{y}_i$, $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) > c \cdot r$

arbitrary if otherwise

**Locality-sensitive hashing (LSH):**

A random $h : X \to U$ is an $(r, cr, p, q)$-**LSH** if $\forall x, y \in X$:

$$\text{dist}(x, y) \leq r \implies \Pr[h(x) = h(y)] \geq p$$

$$\text{dist}(x, y) > c \cdot r \implies \Pr[h(x) = h(y)] \leq q$$

# Locality-Sensitive Hashing (LSH)

- Metric space $(X, \text{dist})$:

**Locality-sensitive hashing (LSH)**:

A random $h : X \to U$ is an $(r, cr, p, q)$**-LSH** if $\forall x, y \in X$:

$$\text{dist}(x, y) \leq r \implies \Pr[h(x) = h(y)] \geq p$$

$$\text{dist}(x, y) > c \cdot r \implies \Pr[h(x) = h(y)] \leq q$$

**Proposition** (*tensorization*):

$\exists$ an $(r, cr, p, q)$-LSH $\qquad\qquad \exists$ an $\left(r, cr, p^k, q^k\right)$-LSH

$\qquad h : X \to U \qquad \implies \qquad\qquad h : X \to U^k$

draw independent $h_1, \ldots, h_k$ according to distribution of $h$

$$g(x) = (h_1(x), h_2(x), \ldots, h_k(x)) \in U^k$$

- Metric space $(X, \text{dist})$:

**Data**: $n$ points $y_1, y_2, \ldots, y_n \in X$

**Query**: a point $x \in X$

return a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ if $\exists y_j$ s.t. $\text{dist}(x, y_j) \leq r$

"no" if $\forall y_i$, $\text{dist}(x, y_i) > c \cdot r$

suppose we have $(r, cr, p^*, 1/n)$-LSH $g: X \rightarrow U$

$\forall x, y \in X$: $\quad \text{dist}(x, y) \leq r \quad \Rightarrow \Pr[\, g(x) = g(y) \,] \geq p^*$

$\text{dist}(x, y) > c \cdot r \Rightarrow \Pr[\, g(x) = g(y) \,] \leq 1/n$

**Data structure**: a dictionary for all key-value pairs $\langle g(y_i), y_i \rangle$;

Upon **query** $x \in X$: find all $y_i$ with $g(x) = g(y_i)$;

if encounter a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ return this $y_i$;

else return "no";

if real answer is "no": always correct

if real answer is not "no": correct with probability $\geq p^*$

(Use FKS) space: $O(n)$    time: $O(1) + O(1)$ in expectation

- Metric space $(X, \mathrm{dist})$:

**Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in X$

**Query**: a point $\boldsymbol{x} \in X$

    return a $\boldsymbol{y}_i$ that $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot r$ if $\exists \boldsymbol{y}_j$ s.t. $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_j) \leq r$

                  "no" if $\forall \boldsymbol{y}_i$ , $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_i) > c \cdot r$

suppose we have $(r, cr, p^*, 1/n)$-LSH $g: X \rightarrow U$

**Hash functions**: *i.i.d.* instances $g_1, \ldots, g_\ell$ of $g$, where $\ell = 1/p^*$;

**Data structure**: $\ell$ dictionaries, where the $j$th dictionary stores

          all key-value pairs $\langle g_j(\boldsymbol{y}_i), y_i \rangle$;

Upon **query $\boldsymbol{x} \in X$**:

    find $\leq 10\ell$ such $\boldsymbol{y}_i$ that $\exists j, g_j(\boldsymbol{x}) = g_j(\boldsymbol{y}_i)$;

    if encounter a $\boldsymbol{y}_i$ that $\mathrm{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot r$ then return this $\boldsymbol{y}_i$;

    else return "no";

metric space $(X, \text{dist})$        $(r, cr, p^*, 1/n)$-LSH $g: X \to U$

**Data**:  $n$ points $y_1, y_2, ..., y_n \in X$        **Query**:  $x \in X$

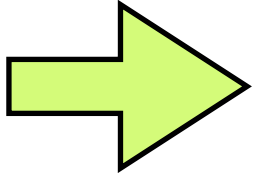Hash functions: *i.i.d.* instances $g_1, \ldots, g_\ell$ of $g$, where $\ell = 1/p^*$;

Data structure: $\ell$ dictionaries, where the $j$th dictionary stores
        all key-value pairs $\langle g_j(y_i), y_i \rangle$;

Upon **query $x \in X$**:
    find $\leq 10\ell$ such $y_i$ that $\exists j, g_j(x) = g_j(y_i)$;

    if encounter a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ then return this $y_i$;
    else return "no";

space:  $\text{O}(nl) = \text{O}(n/p^*)$    time:  $\text{O}(l) = \text{O}(1/p^*)$    (use FKS)

if real answer is "no" :  $\forall y_i$ , $\text{dist}(x, y_i) > c \cdot r$  $\Rightarrow$  always correct

if $\exists y_s$ s.t. $\text{dist}(x, y_s) \leq r$

        $\text{Pr}[\text{ answer "no" }] \leq ?$

metric space $(X, \text{dist})$ $\qquad$ $(r, cr, p^*, 1/n)$-LSH $g: X \to U$

**Data**: $n$ points $y_1, y_2, ..., y_n \in X$ $\qquad$ **Query**: $x \in X$

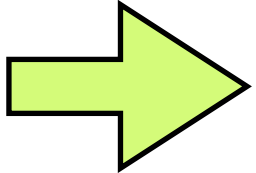Hash functions: *i.i.d.* instances $g_1, \ldots, g_\ell$ of $g$, where $\ell = 1/p^*$;

Data structure: $\ell$ dictionaries, where the $j$th dictionary stores
all key-value pairs $\langle g_j(y_i), y_i \rangle$;
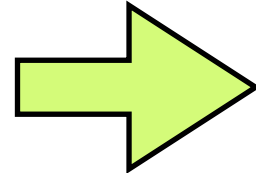
Upon **query** $x \in X$:
find $\leq 10\ell$ such $y_i$ that $\exists j, g_j(x) = g_j(y_i)$;

if encounter a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ then return this $y_i$;
else return "no";

if $\exists y_s$ s.t. $\text{dist}(x, y_s) \leq r$

$$\Pr[\text{ answer "no" }] \leq \boxed{\Pr[\forall j, g_j(x) \neq g_j(y_s)]} \leq (1 - p^*)^\ell \leq 1/e$$

$$+ \Pr[>10l \text{ bad } y_i \text{ that } \text{dist}(x, y_i) > c \cdot r \text{ but } \exists j \text{ s.t. } g_j(x) = g_j(y_i)]$$

Markov
inequality $\leq \mathbf{E}[\text{ \# of such bad } y_i] / 10l \quad \leq \frac{\ell \cdot n \cdot (1/n)}{10\ell} \leq 0.1$

linearity of expectation

metric space $(X, \text{dist})$ $\qquad$ $(r, cr, p^*, 1/n)$-LSH $g: X \to U$

**Data**: $n$ points $y_1, y_2, ..., y_n \in X$ $\qquad$ **Query**: $x \in X$

Hash functions: *i.i.d.* instances $g_1, \ldots, g_\ell$ of $g$, where $\ell = 1/p^*$;

Data structure: $\ell$ dictionaries, where the $j$th dictionary stores
all key-value pairs $\langle g_j(y_i), y_i \rangle$;

Upon **query** $x \in X$:

find $\leq 10\ell$ such $y_i$ that $\exists j, g_j(x) = g_j(y_i)$;

if encounter a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ then return this $y_i$;
else return "no";

space: $O(nl) = O(n/p^*)$ $\qquad$ time: $O(l) = O(1/p^*)$

if real answer is "no" : $\forall y_i, \text{dist}(x, y_i) > c \cdot r$ $\implies$ always correct

if $\exists y_s$ s.t. $\text{dist}(x, y_s) \leq r$

$\implies \Pr[\text{ answer "no" }] \leq 1/e + 0.1 < 0.5$
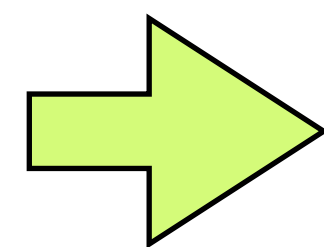
- $(c, r)$-ANN in metric space $(X, \text{dist})$:

> **Data**: $n$ points $\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_n \in X$
>
> **Query**: a point $\boldsymbol{x} \in X$
>
> return a $\boldsymbol{y}_i$ that $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) \leq c \cdot r$ if $\exists \boldsymbol{y}_j$ s.t. $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_j) \leq r$
>
> "no" if $\forall \boldsymbol{y}_i$, $\text{dist}(\boldsymbol{x}, \boldsymbol{y}_i) > c \cdot r$

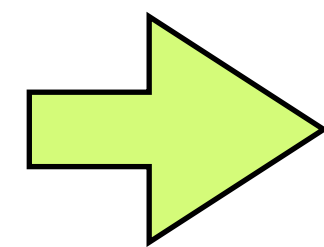suppose we have $(r, cr, p, q)$-LSH $h: X \rightarrow U$

⟹   we have $(r, cr, p^k, 1/n)$-LSH $g: X \rightarrow U^k$

for $k = \log_{(1/q)} n$     so   $p^k = p^{\log_{1/q} n} = n^{-\rho}$

where   $$\rho = \frac{\log p}{\log q}$$

⟹   solve $(c, r)$-ANN with space $O(n^{1+\rho})$

query time $\text{O}(n^\rho \cdot \log n)$ and one-sided error <0.5

- $(c, r)$-**ANN** in metric space $(X, \text{dist})$:

> **Data**: $n$ points $y_1, y_2, \ldots, y_n \in X$
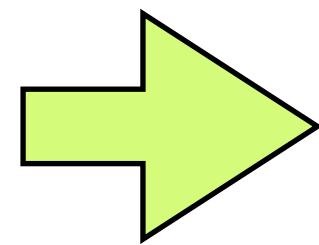>
> **Query**: a point $x \in X$
>
>   return a $y_i$ that $\text{dist}(x, y_i) \leq c \cdot r$ if $\exists y_j$ s.t. $\text{dist}(x, y_j) \leq r$
>
>   "no" if $\forall y_i$, $\text{dist}(x, y_i) > c \cdot r$

suppose we have $(r, cr, p, q)$-**LSH** $h: X \rightarrow U$

$$\rho = \frac{\log p}{\log q}$$

➡ solve $(c, r)$-**ANN** with space $O(n^{1+\rho})$

query time $O(n^{\rho} \cdot \log n)$ and one-sided error <0.5

- $(c, r)$-ANN in Hamming space $\{0,1\}^d$:

> **Data**: $n$ points $y_1, y_2, \ldots, y_n \in \{0,1\}^d$
>
> **Query**: a point $x \in \{0,1\}^d$
>   return a $y_i$ that $\mathrm{dist}(x, y_i) \leq c \cdot r$ if $\exists y_j$ s.t. $\mathrm{dist}(x, y_j) \leq r$
>                     "no" if $\forall y_i$, $\mathrm{dist}(x, y_i) > c \cdot r$

$$\forall\, x \in \{0, 1\}^d: \qquad h(x) = x_i \text{ for uniform random } i \in [d]$$

$$\mathrm{dist}(x, y) \leq r \quad \Rightarrow \Pr[\, h(x) = h(y)\, ] \geq 1\text{-}r/d$$

$$\mathrm{dist}(x, y) > c{\cdot}r \quad \Rightarrow \Pr[\, h(x) = h(y)\, ] \leq 1\text{-}cr/d$$

$h: \{0, 1\}^d \rightarrow \{0,1\}$ is an $(r, cr, 1\text{-}r/d, 1\text{-}cr/d)$-LSH

$$\rho = \frac{\log(1 - r/d)}{\log(1 - cr/d)} \leq \frac{1}{c}$$

➡ solve $(c, r)$-**ANN** in Hamming space with space $O(n^{1+1/c})$

query time $\mathrm{O}(n^{1/c}{\cdot}\log n)$ and one-sided error <0.5

## COMPUTER SCIENCE

# A neural algorithm for a fundamental computing problem

**Sanjoy Dasgupta,[1] Charles F. Stevens,[2,3] Saket Navlakha[4]***

Similarity search—for example, identifying similar images in a database or similar documents on the web—is a fundamental computing problem faced by large-scale information retrieval systems. We discovered that the fruit fly olfactory circuit solves this problem with a variant of a computer science algorithm (called locality-sensitive hashing). The fly circuit assigns similar neural activity patterns to similar odors, so that behaviors learned from one odor can be applied when a similar odor is experienced. The fly algorithm, however, uses three computational strategies that depart from traditional approaches. These strategies can be translated to improve the performance of computational similarity searches. This perspective helps illuminate the logic supporting an important sensory function and provides a conceptually new algorithm for solving a fundamental computational problem.

An essential task of many neural circuits is to generate neural activity patterns in response to input stimuli, so that different inputs can be specifically identified. We studied the circuit used to process odors in the fruit fly olfactory system and uncovered computational strategies for solving a fundamental machine learning problem: approximate similarity (or nearest-neighbors) search.

The fly olfactory circuit generates a "tag" for each odor, which is a set of neurons that fire when that odor is presented (1). This tag is critical for learning behavioral responses to different odors (2). For example, if a reward (e.g., sugar water) or a punishment (e.g., electric shock) is associated with an odor, that odor becomes attractive (a fly

pendence is removed (7, 8); that is, the distribution of firing rates across the 50 PN types is exponential, with close to the same mean for all odors and all odor concentrations (1). Thus, the first step in the circuit essentially "centers the mean"—a standard preprocessing step in many computational pipelines—using a technique called divisive normalization (8). This step is important so that the fly does not mix up odor intensity with odor type.

The second step, where the main algorithmic insight begins, involves a 40-fold expansion in the number of neurons: Fifty PNs project to 2000 Kenyon cells (KCs), connected by a sparse, binary random connection matrix (9). Each KC receives and sums the firing rates from about six randomly

out of the billions of images on the web—that look most similar to your elephant image. This is called the nearest-neighbors search problem, which is of fundamental importance in information retrieval, data compression, and machine learning (10). Each image is typically represented as a $d$-dimensional vector of feature values. (Each odor that a fly processes is a 50-dimensional feature vector of firing rates.) A distance metric is used to compute the similarity between two images (feature vectors), and the goal is to efficiently find the nearest neighbors of any query image. If the web contained only a few images, then brute force linear search could easily be used to find the exact nearest neighbors. If the web contained many images, but each image was represented by a low-dimensional vector (e.g., 10 or 20 features), then space-partitioning methods (12) would similarly suffice. However, for large databases with high-dimensional data, neither approach scales (11).

In many applications, returning an approximate set of nearest neighbors that are "close enough" to the query is adequate, so long as they can be found quickly. This has motivated an approach for finding approximate nearest neighbors by LSH (10). For the fly, as noted, the locality-sensitive property states that two odors that generate similar ORN responses will be represented by two tags that are themselves similar (Fig. 1B). Likewise, for image search, the tag of an elephant image will be more similar to the tag of another elephant image than to the tag of a skyscraper image.

Unlike a traditional (non-LSH) hash function, where the input points are scattered randomly and uniformly over the range, a LSH function provides a distance-preserving embedding of points from $d$-dimensional space into $m$-dimensional
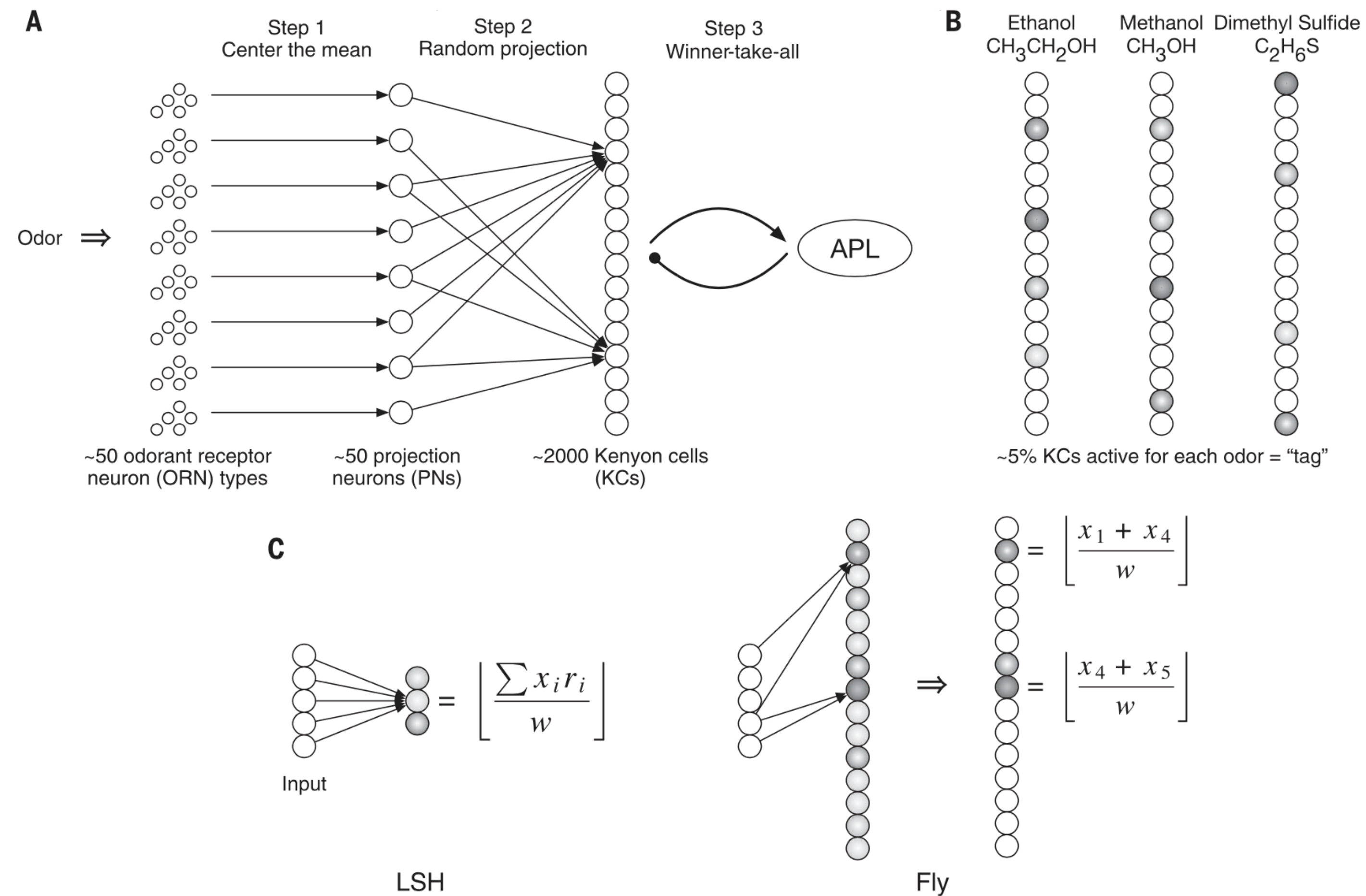
**Fig. 1. Mapping between the fly olfactory circuit and locality-sensitive hashing (LSH).** (**A**) Schematic of the fly olfactory circuit. In step 1, 50 ORNs in the fly's nose send axons to 50 PNs in the glomeruli; as a result of this projection, each odor is represented by an exponential distribution of firing rates, with the same mean for all odors and all odor concentrations. In step 2, the PNs expand the dimensionality, projecting to 2000 KCs connected by a sparse, binary random projection matrix. In step 3, the KCs receive feedback inhibition from the anterior paired lateral (APL) neuron, which leaves only the top 5% of KCs to remain firing spikes for the odor. This 5% corresponds to the tag (hash) for the odor. (**B**) Illustrative odor responses. Similar pairs of odors (e.g., methanol and ethanol) are assigned more similar tags than are dissimilar odors. Darker shading denotes higher activity. (**C**) Differences between conventional LSH and the fly algorithm. In the example, the computational complexity for LSH and the fly are the same. The input dimensionality $d = 5$. LSH computes $m = 3$ random projections, each of which requires 10 operations (five multiplications plus five additions). The fly computes $m = 15$ random projections, each of which requires two addition operations. Thus, both require 30 total operations. $\mathbf{x}$, input feature vector; $r$, Gaussian random variable; $w$, bin width constant for discretization.

# Fast Johnson-Lindenstrauss Transformation (FJLT)

# Subspace Embedding