# Randomized Dimensionality Reduction in Machine Learning

**Sebastian Mayer**

Institute for Numerical Simulation, University of Bonn

universität**bonn**

Institute for Numerical Simulation
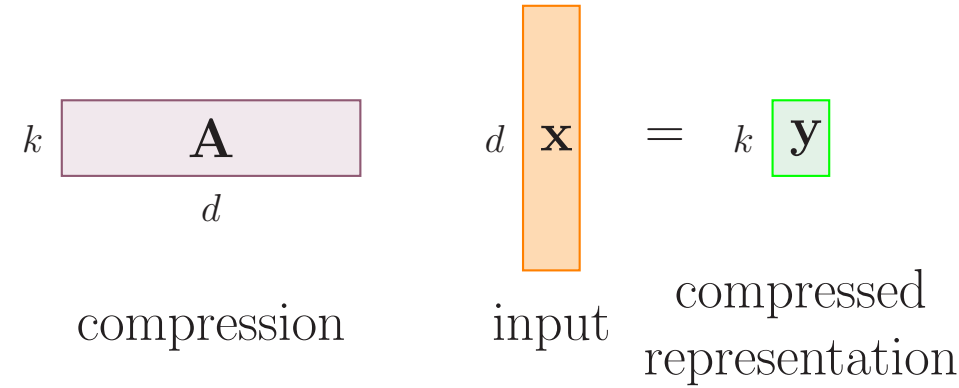Rheinische Friedrich-Wilhelms-Universität Bonn

## Introduction

Typical situation in machine learning:

- Feature space $\mathcal{X} \subseteq \mathbb{R}^d$ with $d$ **very large** (e.g. $d = 40k$):
  - text represented as *bag-of-word* (BOW),
  - image represented as *histogram of gradients* (HOG),
  - property list in *computational biology*.
- **Subproblem:** computing distances or dot products
- Working on $\mathcal{X}$ difficult $\Longrightarrow$ **data compression** $y = Ax$.

$$k \quad \boxed{\mathbf{A}} \quad d \quad \boxed{\mathbf{x}} \ = \ k \ \boxed{\mathbf{y}}$$
$$d$$

compression    input    compressed representation

Here:

**Compression = randomized dimensionality reduction**

Motivations for (randomized) dimensionality reduction:

- computation speed-up
- memory or storage constraints
- trick in algorithm design
- better compression guarantees

**Can we compress while preserving important properties?**

## Johnson-Lindenstrauss embeddings (JLE): Basic version [JL '84; Dasgupta, Gupta '99]

Observation:

- Let $\mathcal{A}$ be a $(k \times d)$-**random matrix** with *i.i.d.* entries $\mathcal{A}_{ij} \sim \mathcal{N}(0, 1/k)$.
- For fixed $x \in \mathbb{R}^d$, we have $\mathbb{E}_{\mathcal{A}} \|\mathcal{A}x\|_2^2 = \|x\|_2^2$.

Given:

- **finite** $S \subset \mathbb{R}^d$,
- **Accuracy parameter** $\varepsilon \in (0, 1)$,
- Failure probability $\delta \in (0, 1)$,

**Theorem:** If
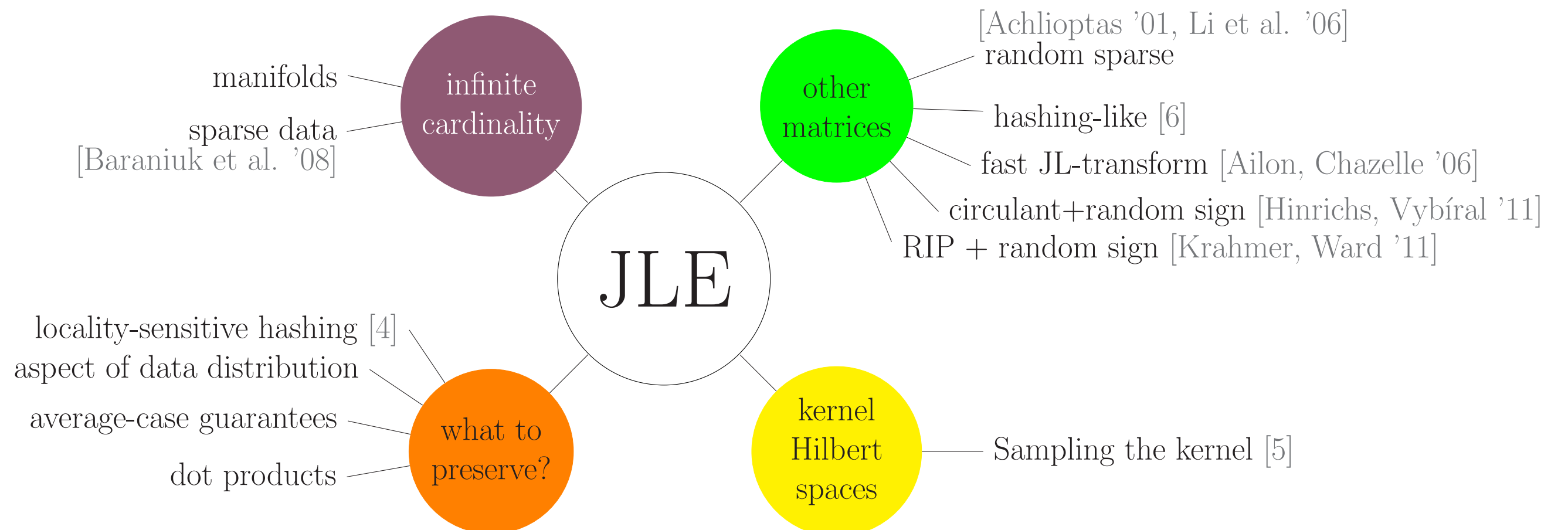$$k = \mathcal{O}\left(\varepsilon^{-2} \ln(|S|/\delta)\right),$$
then a realization $A = \text{draw}(\mathcal{A})$ is with probability $1 - \delta$ an $\varepsilon$-**isometry** on $S$:
$$(1 - \varepsilon)\|x - y\|_2^2 \leq \|A(x - y)\|_2^2 \leq (1 + \varepsilon)\|x - y\|_2^2,$$
for all $x, y \in S$.

**Preservation of *relative* distances!**

## JLE: Extensions and Generalizations

infinite cardinality
- manifolds
- sparse data [Baraniuk et al. '08]

other matrices
- [Achlioptas '01, Li et al. '06] random sparse
- hashing-like [6]
- fast JL-transform [Ailon, Chazelle '06]
- circulant+random sign [Hinrichs, Vybíral '11]
- RIP + random sign [Krahmer, Ward '11]

JLE

what to preserve?
- locality-sensitive hashing [4]
- aspect of data distribution
- average-case guarantees
- dot products

kernel Hilbert spaces
- Sampling the kernel [5]

## Example 1: Nearest neighbour problems [e.g. Indyk '01]

Given: *finite* set of $n$ points $S \subset \mathbb{R}^d$.

$r$-similarity graph: Compute $G = (S, E)$ such that for $p, q \in S$,
$$(p, q) \in E \iff d(p, q) = \|p - q\|_2 \leq r.$$

Computation time: $\mathcal{O}(dn^2)$

**Can we speed up computations?**

**Relax problem**:

$(r, R, \delta)$-**similarity graph**: Compute $G' = (S, E')$ such that with probability $1 - \delta$,
$$(p, q) \in E' \implies d(p, q) \leq r$$
$$(p, q) \notin E' \implies d(p, q) \geq R.$$

Computation time: *JLE with sparse matrix:* $C = \left(\frac{R^2 + r^2}{R^2 - r^2}\right)^2$
$$\underbrace{\mathcal{O}\left(C \log n \, n^2\right)}_{\text{dist. comp.}} + \underbrace{\mathcal{O}\left(C \sqrt{d} \, n \, \log n\right)}_{\text{preprocessing}}$$

## Example 2: Learning mixtures of Gaussians [Dasgupta '99]

Given: $n$ independent, $\mathbb{R}^d$-valued Gaussians $\mathcal{N}(\mu_1, \Sigma), \ldots, \mathcal{N}(\mu_n, \Sigma)$.

Goal: Estimate $\mu_1, \ldots, \mu_n$ from observations $X_1, \ldots, X_m$ originating from *any* Gaussian with equal probability.

- $\gamma$-**separated Gaussians:**
$$\|\mu_i - \mu_j\|_2 \geq \gamma \sqrt{d \, \lambda_{\max}(\Sigma)}.$$
- **Eccentricity:**
$$\text{ecc}(\Sigma) = \sqrt{\lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma)}.$$

*If the Gaussians are sufficiently separated and have low eccentricity, then there is a relatively simple algorithm which needs $\mathcal{O}(\rho^{-d})$ observations to estimate with accuracy $\rho \sqrt{d \, \lambda_{max}(\Sigma)}$.*

**Can we overcome the curse of dimensionality and even handle higher eccentricity?**
$$A = \text{draw}(\mathcal{A})$$
$$\Longrightarrow$$
**auxiliary Gaussians**
$$\tilde{\mu}_i = A\mu_i, \ \tilde{\Sigma} = A\Sigma A^T$$

Observation: If $X \sim \mathcal{N}(\mu, \Sigma)$, then
$$\mathbb{E}_{\mathcal{A}} \|\mathbb{E}_X \mathcal{A}X\|_2^2 = \|\mu\|_2^2$$
$$\mathbb{E}_{\mathcal{A}} \text{Cov}(\mathcal{A}X) = \frac{\text{trace}(\Sigma)}{k} \, \text{Id}_k.$$

**Theorem:** If
$$k = \mathcal{O}\left(\varepsilon^{-2} \log(2n/\delta)\right)$$
and
$$\text{ecc}(\Sigma) = \mathcal{O}\left(\frac{\tau\sqrt{d}}{\log(k/(\delta\tau))}\right),$$
then with probability $1 - \delta$

- **Preservation of separability:**
$$\|\tilde{\mu}_i - \tilde{\mu}_j\|_2^2 \geq (1 - \varepsilon)(1 - \tau)\gamma^2 k \lambda_{\max}(\tilde{\Sigma}),$$
- **Reduction of eccentricity:**
$$(1 - \tau)\lambda \leq |v^T \tilde{\Sigma} v| \leq (1 + \tau)\lambda$$
for all $v \in \mathbb{S}_2^{k-1}$ and $\lambda = \text{trace}(\Sigma)/k$.

## Example 3: Learning a linear classifier

Related work: [1, 2, 3]

Given: Classification problem $(\mathcal{D}, \ell)$, where
- (unknown) distribution $\mathcal{D}$ over $\mathcal{X} \subseteq \{x \in \mathbb{R}^d : \|x\|_2 \leq R\}$,
- labelling function $\ell : \mathcal{X} \to \{-1, 1\}$.

Properties:
- $(\mathcal{D}, \ell)$ **linearly separable**: there is $w \in \mathbb{R}^d$
$$\mathcal{R}_{0,\mathcal{D}}(w) := \mathbb{P}_{X \sim \mathcal{D}}(\ell(X)\langle w, X \rangle) \geq 0) = 1.$$
- $(\mathcal{D}, \ell)$ **linearly separable at margin** $\gamma$:
$$(\exists w \in \mathbb{S}_2^{d-1}) \quad \mathbb{P}_{X \sim \mathcal{D}}(\ell(X)\langle w, X \rangle) \geq \gamma) = 1.$$

**Hinge risk:**
- $\mathcal{R}_{\gamma,\mathcal{D}}(w) = \mathbb{E}_{X \sim \mathcal{D}}\left(1 - \frac{\ell(X)\langle w, X \rangle}{\gamma}\right)_+$
- convex, $1/\gamma$-Lipschitz
- *Surrogate* for misclassification rate: $\mathcal{R}_{0,\mathcal{D}}(w) \leq \mathcal{R}_{\gamma,\mathcal{D}}(w).$

Learning task: Draw *i.i.d.* samples $X_1, \ldots, X_m \sim \mathcal{D}$, learn $\hat{w} \in \mathbb{R}^d$ with minimal **misclassification rate**:
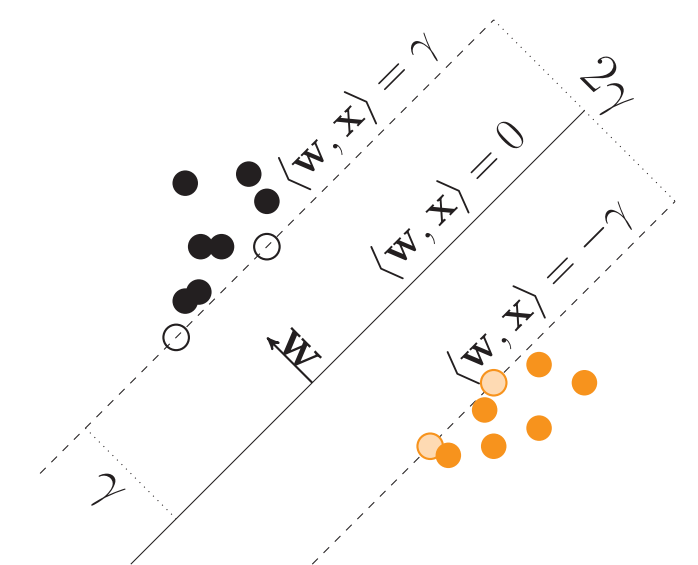$$\hat{w} = \text{argmin}_{w \in \mathbb{R}^d} \mathcal{R}_{0,\mathcal{D}}(w).$$
- Required amount of samples $m$ increases with decreasing $\gamma$.

**Can we learn and predict on compressed data?**
$$A = \text{draw}(\mathcal{A})$$
$$\Longrightarrow$$
**compressed classification problem** $(A\mathcal{D}, \ell)$

## Learning a linear classifier: Some results [M. '13]

**Theorem (Preservation of linear separability):**
Fix $w \in \mathbb{R}^d$ with $\|w\|_2 \leq R$. If
$$k = \mathcal{O}\left(\left(\frac{1 + R^2}{\varepsilon\gamma}\right)^2 \log(1/\delta)\right),$$
then a realization $A = \text{draw}(\mathcal{A})$ fulfils
$$|\mathcal{R}_{\gamma,A\mathcal{D}}(Aw) - \mathcal{R}_{\gamma,\mathcal{D}}(w)| \leq \varepsilon$$
with probability $1 - \delta$.

**Theorem (Preservation of sample complexity):**
Let $\mathcal{X} \cong \mathbb{R}_s^d$ and $w \in \mathbb{S}_2^{d-1}$. If
$$k = \mathcal{O}\left(\left(\frac{1 + R^2}{\tau\gamma}\right)^2 s \log(d/s) + \log(1/\delta))\right),$$
then a realization $A = \text{draw}(\mathcal{A})$ fulfils
$$\mathbb{P}_{X \sim \mathcal{D}}(\ell(X)\langle Aw, AX \rangle) \geq (1 - \tau)\gamma) = 1$$
with probability $1 - \delta$.

**Numerical experiment: Support Vector Machines (SVM)**

- Model: $\mathcal{D}$ uniform distribution over $\mathcal{C}_1 \cup \mathcal{C}_2$, where
$$\mathcal{C}_1 := (R + \gamma)e_1 + [-R, R]^{50}, \quad \mathcal{C}_2 := -(R + \gamma)e_1 + [-R, R]^{50}.$$
- Training: fixed number of training samples; train both in ambient and projected space.
- Evaluation: Accuracy (=1-misclass. rate) in projected space relative to accuracy in ambient space.

- Interpretation:

**Good:**
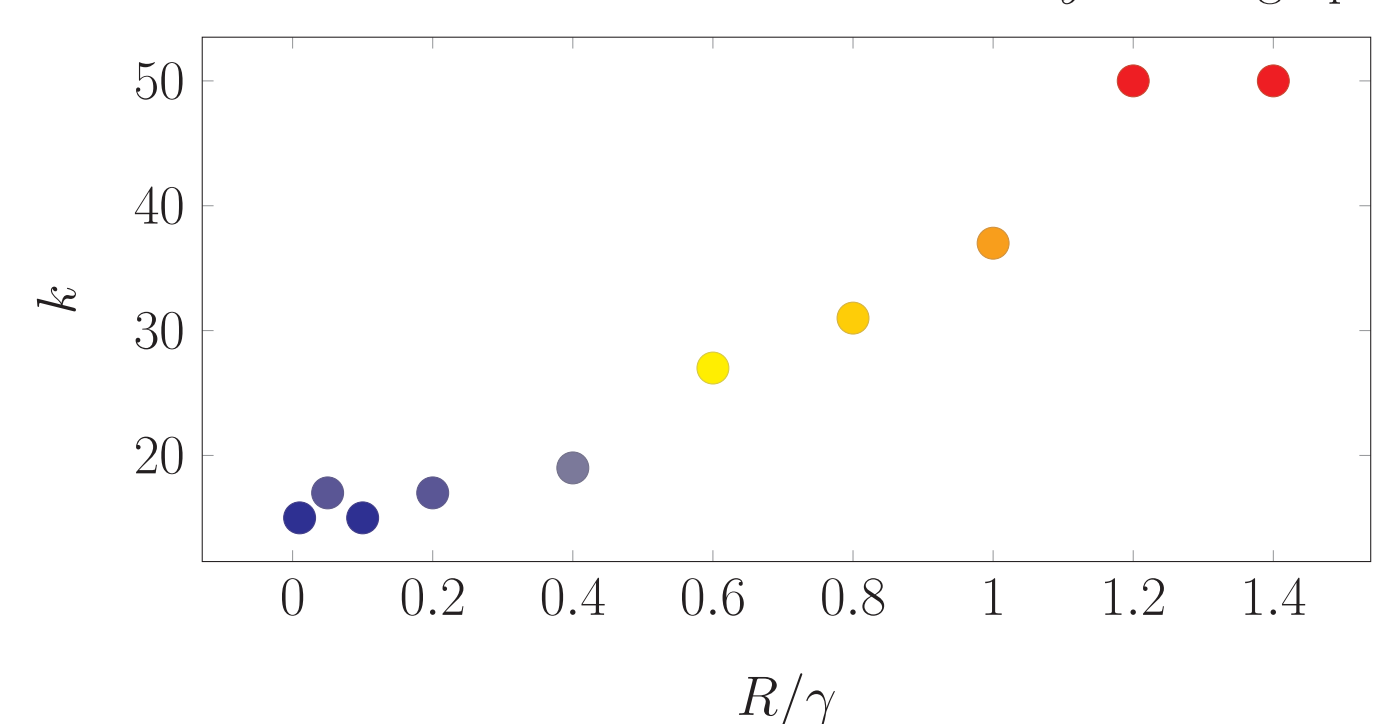- easy learning problems,
- labelling cheap.

**Bad:**
- harder learning problems,
- labelling expensive.

Min. dimension to attain 0.95 *relative accuracy* with high probability

## References

[1] Arriaga, Vempala (2006): *An algorithmic theory of learning: Robust concepts and random projection.*

[2] Balcan, Blum, Vempala (2006): *Kernels as features: On kernels, margins, and low-dimensional mappings.*

[3] Calderbank, Jafarpour (2012): *Finding needles in compressed haystacks.*

[4] Gionis, Indyk, Motwani (1999): *Similarity search via hashing in high dimensions.*

[5] Rahimi, Recht (2007): *Random features for large-scale kernel machines.*

[6] Weinberger, Dasgupta, Langford, Smola, Attenberg (2009): *Feature hashing for large scale multitask learning.*