

COMPSCI 514: ALGORITHMS FOR DATA SCIENCE

Cameron Musco

University of Massachusetts Amherst. Spring 2020.

Lecture 13

- **Midterm is on Thursday.**
- No calculators, cheatsheets, or other aids.
- Very important to do some practice problems and to try them first with no resources, to simulate the exam.
- Make sure you can recognize when to apply the fundamentals: union bound, linearity of expectation and variance, Markov's inequality, Chebyshev's inequality, indicator random variables.
- Understand the goal of each algorithm/data structure. I.e., what problem it solves with what guarantees. No need to memorize proofs.

Last Few Classes:

The Johnson-Lindenstrauss Lemma

- Reduce n data points in any dimension d to $O\left(\frac{\log n/\delta}{\epsilon^2}\right)$ dimensions and preserve (with probability $\geq 1 - \delta$) all pairwise distances up to $1 \pm \epsilon$.
- Compression is linear via multiplication with a random, data oblivious, matrix (linear compression)

High-Dimensional Geometry

- Why high-dimensional space is so different than low-dimensional space.
- How the JL Lemma can still work.

Next Few Classes: Low-rank approximation, the SVD, and principal component analysis (PCA).

- Reduce d -dimensional data points to a smaller dimension m .
- Like JL, **compression is linear** – by applying a matrix.
- Chose this matrix carefully, taking into account **structure of the dataset**.
- Can give better compression than random projection.

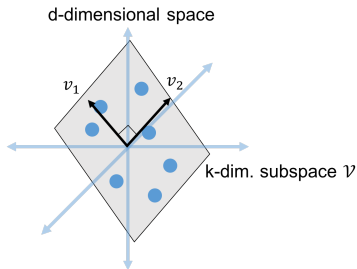
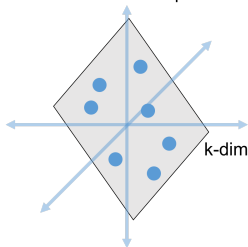
Will be using a fair amount of linear algebra: orthogonal basis, column/row span, eigenvectors, etc,

RANDOMIZED ALGORITHMS UNIT TAKEAWAYS

- Randomization is an important tool in working with large datasets.
- Lets us solve 'easy' problems that get really difficult on massive datasets. Fast/space efficient look up (hash tables and bloom filters), distinct items counting, frequent items counting, near neighbor search, etc.
- The analysis of randomized algorithms leads to complex output distributions, which we can't compute exactly.
- We use concentration inequalities to bound these distributions and behaviors like accuracy, space usage, and runtime.
- Concentration inequalities and probability tools used in randomized algorithms are also fundamental in statistics, machine learning theory, probabilistic modeling of complex systems, etc.

EMBEDDING WITH ASSUMPTIONS

Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie in any k -dimensional subspace \mathcal{V} of \mathbb{R}^d d-dimensional space



Claim: Let $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns. For all \vec{x}_i, \vec{x}_j :

$$\|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

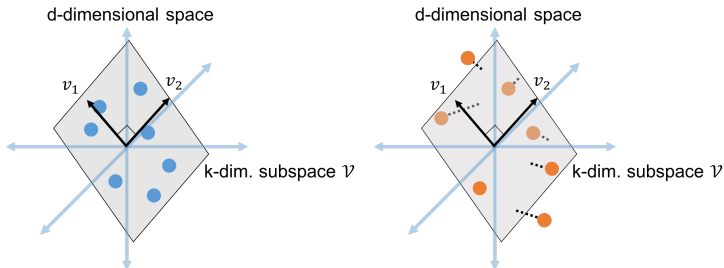
- $\mathbf{V}^T \in \mathbb{R}^{k \times d}$ is a linear embedding of $\vec{x}_1, \dots, \vec{x}_n$ into k dimensions with **no distortion**.
- An actual projection, analogous to a JL random projection $\mathbf{\Pi}$.

Claim: Let $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns. For all $\vec{x}_i, \vec{x}_j \in \mathcal{V}$:

$$\|\mathbf{V}^T \vec{x}_i - \mathbf{V}^T \vec{x}_j\|_2 = \|\vec{x}_i - \vec{x}_j\|_2.$$

EMBEDDING WITH ASSUMPTIONS

Main Focus of Today: Assume that data points $\vec{x}_1, \dots, \vec{x}_n$ lie close to any k -dimensional subspace \mathcal{V} of \mathbb{R}^d .



Letting $\vec{v}_1, \dots, \vec{v}_k$ be an orthonormal basis for \mathcal{V} and $\mathbf{V} \in \mathbb{R}^{d \times k}$ be the matrix with these vectors as its columns, $\mathbf{V}^T \vec{x}_i \in \mathbb{R}^k$ is still a good embedding for $x_i \in \mathbb{R}^d$. The key idea behind low-rank approximation and principal component analysis (PCA).

- How do we find \mathcal{V} and \mathbf{V} ?
- How good is the embedding?