

Lecture 12. Dimension Reduction

Bao Wang

Department of Mathematics

Scientific Computing and Imaging Institute

University of Utah

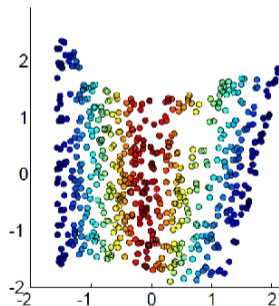
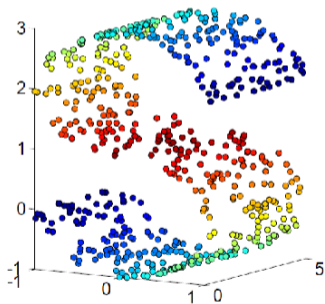
Math 5750/6880, Fall 2021

Outline

- Principal component analysis (PCA)
- Random projection
- Compressed sensing

Dimension reduction

Dimension reduction: mapping data in a high dimensional space into a new space whose dimensionality is much smaller.



Linear dimension reduction

Linear dimension reduction: if the original data is in \mathbb{R}^d and we want to embed it into $\mathbb{R}^n (n < d)$ then we would like to find a matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$ that induces the mapping $\mathbf{x} \rightarrow \mathbf{W}\mathbf{x}$. A natural criterion for choosing \mathbf{W} is in a way that will enable a reasonable recovery of \mathbf{x} from $\mathbf{W}\mathbf{x}$ is possible.

Principal Component Analysis (PCA)

Principal component analysis (PCA) – Motivation

- Let $\mathbf{x}_1, \dots, \mathbf{x}_m$ be m vectors in \mathbb{R}^d . We would like to reduce the dimensionality of these vectors using a linear transformation.
- A matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$, where $n < d$, induces a mapping $\mathbf{x} \rightarrow \mathbf{W}\mathbf{x}$, where $\mathbf{W}\mathbf{x} \in \mathbb{R}^n$ is the lower dimensionality representation of \mathbf{x} . Then, a second matrix $\mathbf{U} \in \mathbb{R}^{d \times n}$ can be used to (approximately) recover each original vector \mathbf{x} from its compressed version.
- That is, for a compressed vector $\mathbf{y} = \mathbf{W}\mathbf{x}$, where \mathbf{y} is in the low dimensional space \mathbb{R}^n , we can construct $\tilde{\mathbf{x}} = \mathbf{U}\mathbf{y}$, so that $\tilde{\mathbf{x}}$ is the recovered version of \mathbf{x} and resides in the original high dimensional space \mathbb{R}^d .
- We want $\|\mathbf{x} - \tilde{\mathbf{x}}\|$ to be small.

Principal component analysis (PCA) – Formulation

PCA: find the compression matrix \mathbf{W} and the recovering matrix \mathbf{U} so that the total squared distance between the original and recovered vectors is minimal; namely, we aim at solving the problem

$$\arg \min_{\mathbf{W} \in \mathbb{R}^{n \times d}, \mathbf{U} \in \mathbb{R}^{d \times n}} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{UWx}_i\|^2. \quad (1)$$

To solve the problem (1), we first show that the optimal solution takes a specific form.

Lemma 1. Let (\mathbf{U}, \mathbf{W}) be a solution to (1). Then the columns of \mathbf{U} are orthonormal (namely, $\mathbf{U}^\top \mathbf{U}$ is the identity matrix of \mathbb{R}^n) and $\mathbf{W} = \mathbf{U}^\top$.

Proof. Fix any \mathbf{U}, \mathbf{W} and consider the mapping $\mathbf{x} \rightarrow \mathbf{UWx}$. The range of this mapping, $R = \{\mathbf{UWx} : \mathbf{x} \in \mathbb{R}^d\}$, is an n dimensional linear subspace of \mathbb{R}^d . Let $\mathbf{V} \in \mathbb{R}^{d \times n}$ be a matrix whose columns form an orthonormal basis of this subspace, namely, the range of \mathbf{V} is R and $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$. Therefore, each vector in R can be written as \mathbf{Vy} where $\mathbf{y} \in \mathbb{R}^n$.

Also, for every $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} \in \mathbb{R}^n$, we have

$$\|\mathbf{x} - \mathbf{Vy}\|^2 = \|\mathbf{x}\|^2 + \mathbf{y}^\top \mathbf{V}^\top \mathbf{Vy} - 2\mathbf{y}^\top \mathbf{V}^\top \mathbf{x} = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\mathbf{y}^\top (\mathbf{V}^\top \mathbf{x}),$$

where we used the fact that $\mathbf{V}^\top \mathbf{V}$ is the identity matrix of \mathbb{R}^n .

Minimizing the preceding expression with respect to \mathbf{y} by comparing the gradient with respect to \mathbf{y} to zero gives that $\mathbf{y} = \mathbf{V}^\top \mathbf{x}$. Therefore, for each \mathbf{x} we have that

$$\mathbf{V}\mathbf{V}^\top \mathbf{x} = \arg \min_{\tilde{\mathbf{x}} \in R} \|\mathbf{x} - \tilde{\mathbf{x}}\|^2.$$

In particular this holds for $\mathbf{x}_1, \dots, \mathbf{x}_m$ and therefore we can replace \mathbf{U}, \mathbf{W} by $\mathbf{V}, \mathbf{V}^\top$ and by that do not increase the objective

$$\sum_{i=1}^m \|\mathbf{x}_i - \mathbf{U}\mathbf{W}\mathbf{x}_i\|_2^2 \geq \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{V}\mathbf{V}^\top \mathbf{x}_i\|_2^2.$$

Since this holds for every \mathbf{U}, \mathbf{W} the proof of the lemma follows.

According to the preceding lemma, we can rewrite the optimization problem given in (1) as follows:

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{d \times n}, \mathbf{U}^\top \mathbf{U} = \mathbf{I}} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_i\|_2^2. \quad (2)$$

We further simplify the optimization problem by using the following elementary algebraic manipulations. For every $\mathbf{x} \in \mathbb{R}^d$ and a matrix $\mathbf{U} \in \mathbb{R}^{d \times n}$ such that $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ we have

$$\begin{aligned}\|\mathbf{x} - \mathbf{U}\mathbf{U}^\top \mathbf{x}\|^2 &= \|\mathbf{x}\|^2 - 2\mathbf{x}^\top \mathbf{U}\mathbf{U}^\top \mathbf{x} + \mathbf{x}^\top \mathbf{U}\mathbf{U}^\top \mathbf{U}\mathbf{U}^\top \mathbf{x} \\ &= \|\mathbf{x}\|^2 - \mathbf{x}^\top \mathbf{U}\mathbf{U}^\top \mathbf{x} = \|\mathbf{x}\|^2 - \text{trace}(\mathbf{U}^\top \mathbf{x}\mathbf{x}^\top \mathbf{U}),\end{aligned}\tag{3}$$

where the trace of a matrix is the sum of its diagonal entries. Since the trace is a linear operator, this allows us to rewrite (2) as follows:

$$\arg \max_{\mathbf{U} \in \mathbb{R}^{d \times n}: \mathbf{U}^\top \mathbf{U} = \mathbf{I}} \text{trace}\left(\mathbf{U}^\top \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top \mathbf{U}\right).\tag{4}$$

$$\text{trace}(\mathbf{A} + \mathbf{B}) = \text{trace}(\mathbf{A}) + \text{trace}(\mathbf{B})$$

- Let $\mathbf{A} = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top$.
- The matrix \mathbf{A} is symmetric (all eigenvalues are real and \mathbf{A} is diagonalizable) and therefore it can be written using its spectral decomposition as $\mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{V}^\top$, where \mathbf{D} is diagonal and $\mathbf{V}^\top \mathbf{V} = \mathbf{V} \mathbf{V}^\top = \mathbf{I}$. Here, the elements on the diagonal of \mathbf{D} are the eigenvalues of \mathbf{A} and the columns of \mathbf{V} are the corresponding eigenvectors.
- We assume without of generality that $D_{1,1} \geq D_{2,2} \geq \dots \geq D_{d,d}$. Since \mathbf{A} is positive semidefinite it also holds that $D_{d,d} \geq 0$. We claim that the solution to (4) is the matrix \mathbf{U} whose columns are the n eigenvectors of \mathbf{A} corresponding to the largest n eigenvalues.

Theorem. Let $\mathbf{x}_1, \dots, \mathbf{x}_m$ be arbitrary vectors in \mathbb{R}^d , let $\mathbf{A} = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top$, and let $\mathbf{u}_1, \dots, \mathbf{u}_n$ be eigenvectors of the matrix \mathbf{A} corresponding to the largest n eigenvalues of \mathbf{A} . Then, the solution to the PCA optimization problem given in (1) is the set \mathbf{U} to be the matrix whose columns are $\mathbf{u}_1, \dots, \mathbf{u}_n$ and to set $\mathbf{W} = \mathbf{U}^\top$.

Proof. Let \mathbf{VDV}^\top be the spectral decomposition of \mathbf{A} . Fix some matrix $\mathbf{U} \in \mathbb{R}^{d \times n}$ with orthonormal columns and let $\mathbf{B} = \mathbf{V}^\top \mathbf{U}$. Then, $\mathbf{VB} = \mathbf{VV}^\top \mathbf{U} = \mathbf{U}$. It follows that

$$\mathbf{U}^\top \mathbf{AU} = \mathbf{B}^\top \mathbf{V}^\top \mathbf{VDV}^\top \mathbf{VB} = \mathbf{B}^\top \mathbf{DB},$$

and therefore

$$\text{trace}(\mathbf{U}^\top \mathbf{AU}) = \text{trace}(\mathbf{B}^\top \mathbf{DB}) = \sum_{j=1}^d D_{j,j} \sum_{i=1}^n B_{j,i}^2.$$

Note that $\mathbf{B}^\top \mathbf{B} = \mathbf{U}^\top \mathbf{VV}^\top \mathbf{U} = \mathbf{U}^\top \mathbf{U} = \mathbf{I}$. Therefore, the columns of \mathbf{B} are also orthonormal, which implies that $\sum_{j=1}^d \sum_{i=1}^n B_{j,i}^2 = n$. In addition, let $\tilde{\mathbf{B}} \in \mathbb{R}^{d \times d}$ be a matrix such that its first n columns are the columns of \mathbf{B} and in addition $\tilde{\mathbf{B}}^\top \tilde{\mathbf{B}} = \mathbf{I}$. Then, for every j we have $\sum_{i=1}^d \tilde{B}_{j,i}^2 = 1$, implying that $\sum_{i=1}^n B_{j,i}^2 \leq 1$.

It follows that

$$\text{trace}(\mathbf{U}^\top \mathbf{A} \mathbf{U}) \leq \max_{\beta \in [0,1]^d: \|\beta\|_1 \leq n} \sum_{j=1}^d D_{j,j} \beta_j.$$

It is not hard to verify that the right-hand side equals to $\sum_{j=1}^n D_{j,j}$ (note that $D_{j,j}$ has been sorted and by setting $\beta_1 = \dots = \beta_n = 1$ and $\beta_{n+1} = \dots = \beta_d = 0$).

We have therefore shown that for every matrix $\mathbf{U} \in \mathbb{R}^{d \times n}$ with orthonormal columns it hold that $\text{trace}(\mathbf{U}^\top \mathbf{A} \mathbf{U}) \leq \sum_{j=1}^n D_{j,j}$. On the other hand, if we set \mathbf{U} to be the matrix whose columns are the n leading eigenvectors of \mathbf{A} we obtain that

$$\text{trace}(\mathbf{U}^\top \mathbf{A} \mathbf{U}) = \sum_{j=1}^n D_{j,j}, \text{ using spectral decomposition of } \mathbf{A},$$

and this concludes our proof.

Remark. The proof of the above theorem tells us that the value of the objective of (4) is $\sum_{i=1}^n D_{i,i}$; combine this with (3) and note that $\sum_{i=1}^m \|\mathbf{x}_i\|^2 = \text{trace}(\mathbf{A}) = \sum_{i=1}^d D_{i,i}$ we obtain that the optimal objective value of (1) is $\sum_{i=n+1}^d D_{i,i}$.

Remark. It is a common practice to "center" the examples before applying PCA. That is, we first calculate $\mu = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$ and then apply PCA on the vectors $(\mathbf{x}_1 - \mu), \dots, (\mathbf{x}_m - \mu)$. This is also related to the interpretation of PCA as variance maximization.

A More Efficient Solution for the Case $d \gg m$

In some situations the original dimensionality of the data is much larger than the number of examples m . The computational complexity of calculating the PCA solution as described previously is $O(d^3)$ (for calculating eigenvalues of \mathbf{A}) plus $O(md^2)$ (for constructing the matrix \mathbf{A}).

How to calculate PCA solution efficiently when $d \gg m$?

A More Efficient Solution for the Case $d \gg m$

- Recall that the matrix \mathbf{A} is defined to be $\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top$.
- We can rewrite $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$ where $\mathbf{X} \in \mathbb{R}^{m \times d}$ is a matrix whose i -th row is \mathbf{x}_i^\top .
- Consider the matrix $\mathbf{B} = \mathbf{X}\mathbf{X}^\top$. That is, $\mathbf{B} \in \mathbb{R}^{m \times m}$ is the matrix whose i, j element equals $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$.
- Suppose that \mathbf{u} is an eigenvector of \mathbf{B} : That is, $\mathbf{B}\mathbf{u} = \lambda\mathbf{u}$ for some $\lambda \in \mathbb{R}$.
- Multiplying the equality by \mathbf{X}^\top and using the definition of \mathbf{B} we obtain $\mathbf{X}^\top \mathbf{X}\mathbf{X}^\top \mathbf{u} = \lambda \mathbf{X}^\top \mathbf{u}$.
- But using the definition of \mathbf{A} , we get that $\mathbf{A}(\mathbf{X}^\top \mathbf{u}) = \lambda(\mathbf{X}^\top \mathbf{u})$. Thus, $\frac{\mathbf{X}^\top \mathbf{u}}{\|\mathbf{X}^\top \mathbf{u}\|}$ is an eigenvector of \mathbf{A} with eigenvalue of λ .

A More Efficient Solution for the Case $d \gg m$

We can therefore calculate the PCA solution by calculating the eigenvalues of \mathbf{B} instead of \mathbf{A} . The complexity is $O(m^3)$ (for calculating eigenvalues of \mathbf{B}) and m^2d (for constructing the matrix \mathbf{B}).

PCA

Input A matrix of m examples $\mathbf{X} \in \mathbb{R}^{m \times d}$, number of components n .

if ($m > d$)

$\mathbf{A} = \mathbf{X}^\top \mathbf{X}$, let $\mathbf{u}_1, \dots, \mathbf{u}_n$ be the eigenvectors of \mathbf{A} with largest eigenvalues

else

$$\mathbf{B} = \mathbf{X}\mathbf{X}^\top$$

Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be the eigenvectors of \mathbf{B} with largest eigenvalues for $i = 1, \dots, n$ set $\mathbf{u}_i = \frac{1}{\|\mathbf{X}^\top \mathbf{v}_i\|} \mathbf{X}^\top \mathbf{v}_i$

output $\mathbf{u}_1, \dots, \mathbf{u}_n$

To illustrate how PCA works, let us generate vectors in \mathbb{R}^2 that approximately reside on a line, namely, on a one dimensional subspace of \mathbb{R}^2 . For example, suppose that each example is of the form $(x, x + y)$ where x is chosen uniformly at random from $[-1, 1]$ and y is sampled from a Gaussian distribution with mean 0 and standard deviation of 0.1. Suppose we apply PCA on this data. Then, the eigenvector corresponding to the largest eigenvalue will be chosen to be the vector $(1/\sqrt{2}, 1/\sqrt{2})$. When projecting a point $(x, x + y)$ on this principal component we will obtain the scalar $(2x + y)/\sqrt{2}$. The reconstruction of the original vector will be $(1/\sqrt{2}, 1/\sqrt{2}) * (2x + y)/\sqrt{2} = ((x + y/2), (x + y/2))$. In Figure 1, we depict the original versus reconstructed data.

Note that $\mathbf{A} = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^\top \sim \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ with $\lambda_1 = 0$ ($\mathbf{v}_1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$) and $\lambda_2 = 2$ ($\mathbf{v}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$).

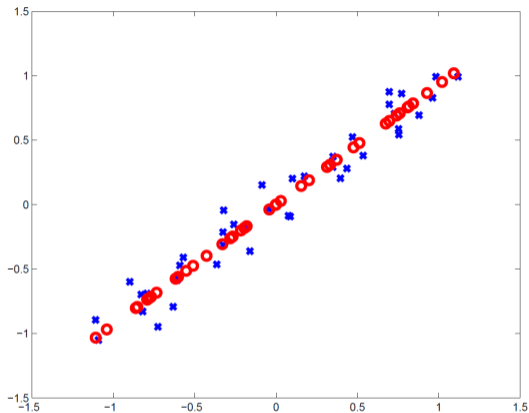


Figure: A set of vectors in \mathbb{R}^2 (blue) and their reconstruction after dimensionality reduction to \mathbb{R}^1 using PCA (red).

Random Projections

We show that reducing the dimension by using a random linear transformation leads to a simple compression scheme with a surprisingly low distortion. The transformation $\mathbf{x} \rightarrow \mathbf{W}\mathbf{x}$, when \mathbf{W} is a random matrix (each $W_{i,j}$ is an independent normal random variable), is often referred to as a random projection.

Let $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$. A matrix \mathbf{W} does not distort too much the distance between \mathbf{x}_1 and \mathbf{x}_2 if the ratio

$$\frac{\|\mathbf{W}\mathbf{x}_1 - \mathbf{W}\mathbf{x}_2\|}{\|\mathbf{x}_1 - \mathbf{x}_2\|}$$

is close to 1. In other words, the distances between \mathbf{x}_1 and \mathbf{x}_2 before and after the transformation are almost the same. To show that $\|\mathbf{W}\mathbf{x}_1 - \mathbf{W}\mathbf{x}_2\|$ is not too far away from $\|\mathbf{x}_1 - \mathbf{x}_2\|$ it suffices to show that \mathbf{W} does not distort the norm of the difference vector $\mathbf{x} = \mathbf{x}_1 - \mathbf{x}_2$. Therefore, from now on we focus on the ratio $\frac{\|\mathbf{W}\mathbf{x}\|}{\|\mathbf{x}\|}$. We start with analyzing the distortion caused by applying a random projection to a single vector.

Random projections

Lemma 2. Fix some $\mathbf{x} \in \mathbb{R}^d$. Let $\mathbf{W} \in \mathbb{R}^{n \times d}$ be a random matrix such that each $W_{i,j}$ is an independent normal random variable. Then, for every $\epsilon \in (0, 3)$ we have

$$\mathbb{P} \left[\left| \frac{\|(1/\sqrt{n})\mathbf{W}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} - 1 \right| > \epsilon \right] \leq 2e^{-\epsilon^2 n/6}.$$

Lemma: Concentration of χ^2 Variables

Let X_1, \dots, X_k be k independent normally distributed random variables. That is, for all i , $X_i \sim N(0, 1)$. The distribution of the random variable X_i^2 is called χ^2 and the distribution of the random variable $Z = X_1^2 + \dots + X_k^2$ is called χ_k^2 . Clearly, $\mathbb{E}[X_i^2] = 1$ and $\mathbb{E}[Z] = k$. The following lemma states that χ_k^2 is concentrated around its mean.

Lemma. Let $Z \sim \chi_k^2$. Then, for all $\epsilon > 0$ we have

$$\mathbb{P}[Z \leq (1 - \epsilon)k] \leq e^{-\epsilon^2 k/6},$$

and for all $\epsilon \in (0, 3)$ we have

$$\mathbb{P}[Z \geq (1 + \epsilon)k] \leq e^{-\epsilon^2 k/6}.$$

Finally, for all $\epsilon \in (0, 3)$,

$$\mathbb{P}[(1 - \epsilon)k \leq Z \leq (1 + \epsilon)k] \geq 1 - 2e^{-\epsilon^2 k/6}.$$

Proof of Lemma 2. Without loss of generality, we can assume that $\|\mathbf{x}\|^2 = 1$. Therefore, an equivalent inequality is

$$\mathbb{P}\left[(1 - \epsilon)n \leq \|\mathbf{W}\mathbf{x}\|^2 \leq (1 + \epsilon)n\right] \geq 1 - 2e^{-\epsilon^2 n/6}.$$

Let \mathbf{w}_i be the i -th row of \mathbf{W} . The random variable $\langle \mathbf{w}_i, \mathbf{x} \rangle$ is a weighted sum of d independent normal random variables and therefore it is normally distributed with zero mean and variance $\sum_j x_j^2 = \|\mathbf{x}\|^2 = 1$ (the variance of $(\mathbf{w}_i)_i * x_i$ is x_i^2 , also $\text{var}(r_1 + r_2) = \text{var}(r_1) + \text{var}(r_2)$, i.e., $\langle \mathbf{w}_i, \mathbf{x} \rangle$ is the standard Gaussian). Therefore, the random variable $\|\mathbf{W}\mathbf{x}\|^2 = \sum_{i=1}^n (\langle \mathbf{w}_i, \mathbf{x} \rangle)^2$ has a χ_n^2 distribution. The claim now follows directly from a measure concentration property of χ^2 random variables.

The Johnson-Lindenstrauss lemma follows from this using a simple union bound argument.

Johnson-Lindenstrauss Lemma

Lemma 3. [Johnson-Lindenstrauss Lemma] Let Q be a finite set of vectors in \mathbb{R}^d . Let $\delta \in (0, 1)$ and n be an integer such that

$$\epsilon = \sqrt{\frac{6 \log(2|Q|/\delta)}{n}} \leq 3.$$

Then, with probability at least $1 - \delta$ over a choice of a random matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$ such that each element of \mathbf{W} is distributed normally with zero mean and variance of $1/n$ we have

$$\sup_{\mathbf{x} \in Q} \left| \frac{\|\mathbf{W}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} - 1 \right| < \epsilon.$$

Proof. Combining Lemma 2 and the union bound ($P(A \cup B) \leq P(A) + P(B)$) we have that for every $\epsilon \in (0, 3)$:

$$\mathbb{P} \left[\sup_{\mathbf{x} \in Q} \left| \frac{\|\mathbf{W}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} - 1 \right| > \epsilon \right] \leq 2|Q|e^{-\epsilon^2 n/6}.$$

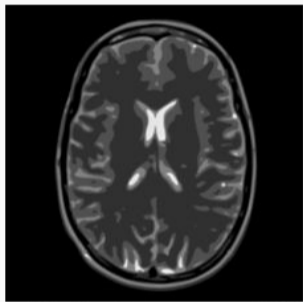
Let δ denote the right-hand-side of the inequality; thus we obtain that

$$\epsilon = \sqrt{\frac{6 \log(2|Q|/\delta)}{n}}.$$

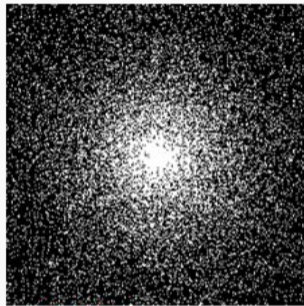
Remark. Interestingly, the bound given in the JL lemma does not depend on the original dimension of \mathbf{x} . In fact, the bound holds even if \mathbf{x} is in an infinite dimensional Hilbert space.

Compressed Sensing

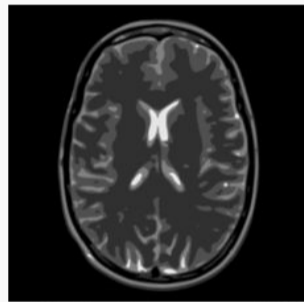
Compressed sensing is a dimensionality reduction technique which utilizes a prior assumption that the original vector is sparse in some basis.



Ground truth



25% subsampling in k -space (MRI)



Reconstruction

To motivate compressed sensing, consider a vector $\mathbf{x} \in \mathbb{R}^d$ that has at most s nonzero elements. That is,

$$\|\mathbf{x}\|_0 := |\{i : x_i \neq 0\}| \leq s.$$

Clearly, we can compress \mathbf{x} by representing it using s (index, value) pairs. This compression is lossless – we can reconstruct \mathbf{x} exactly from the s (index, value) pairs.

Now, let us take one step forward and assume that $\mathbf{x} = \mathbf{U}\boldsymbol{\alpha}$ where $\boldsymbol{\alpha}$ is a sparse vector, $\|\boldsymbol{\alpha}\|_0 \leq s$, and \mathbf{U} is a fixed orthonormal matrix. That is, \mathbf{x} has a sparse representation in another basis. It turns out that many natural vectors are (at least approximately) sparse in some representation. – For instance, natural image in wavelet representation.

Can we still compress \mathbf{x} into roughly s numbers?

Compressed sensing

One simple way to do this is to multiply \mathbf{x} by \mathbf{U}^\top , which yields the sparse vector α , and then represent α by its s (index,value) pairs. However, this requires us to first to “sense” \mathbf{x} , to store it, and then to multiply it by \mathbf{U}^\top .

This raises a very natural question: *Why go to so much effort to acquire all the data when most of what we get will be thrown away? Cannot we just directly measure the part that will not end up being thrown away?*

Compressed sensing

- Compressed sensing is a technique that simultaneously acquires and compresses the data. The key result is that a random linear transformation can compress \mathbf{x} without losing information.

- The number of measurements needed is order of $s \log(d)$ rather than order d . That is, we roughly acquire only the important information about the signal.

Compressed sensing

- As we will see later, the price we pay is a slower reconstruction phase. In some situations, it makes sense to save time in compression even at the price of a slower reconstruction. For example, a security camera should sense and compress a large amount of images while most of the time we do not need to decode the compressed data at all.
- Furthermore, in many practical applications, compression by a linear transformation is advantageous because it can be performed efficiently in hardware. An important application of compressed sensing is medical imaging, in which requiring fewer measurements translates to less radiation for the patient.

Premise of compressed sensing

- > It is possible to reconstruct any sparse signal fully if it was compressed by $\mathbf{x} \rightarrow \mathbf{W}\mathbf{x}$, where \mathbf{W} is a matrix which satisfies a condition called the **Restricted Isoperimetric Property (RIP)**. A matrix that satisfies this property is guaranteed to have a low distortion of the norm of any sparse representable vector.
- > The reconstruction can be calculated in polynomial time by solving a linear program.
- > A random $n \times d$ ($d > n$) matrix is likely to satisfy the RIP condition provided that n is greater than an order of $s \log(d)$.

Definition. [RIP] A matrix $\mathbf{W} \in \mathbb{R}^{n \times d}$ is (ϵ, s) -RIP if for all $\mathbf{x} \neq 0$ s.t. $\|\mathbf{x}\|_0 \leq s$ we have

$$\left| \frac{\|\mathbf{W}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} - 1 \right| \leq \epsilon.$$

Compressed sensing

Theorem 1. Let $\epsilon < 1$ and let \mathbf{W} be a $(\epsilon, 2s)$ -RIP matrix. Let \mathbf{x} be a vector s.t. $\|\mathbf{x}\|_0 \leq s$, let $\mathbf{y} = \mathbf{W}\mathbf{x}$ be the compression of \mathbf{x} , and let

$$\tilde{\mathbf{x}} \in \arg \min_{\mathbf{v}: \mathbf{W}\mathbf{v}=\mathbf{y}} \|\mathbf{v}\|_0,$$

be a reconstructed vector. Then, $\tilde{\mathbf{x}} = \mathbf{x}$.

Remark. We can exactly recovery \mathbf{x} from the compression \mathbf{y} by solving an optimization problem.

Proof. We assume, by contradiction, that $\tilde{\mathbf{x}} \neq \mathbf{x}$. Since \mathbf{x} satisfies the constraints in the optimization problem that defines $\tilde{\mathbf{x}}$ we clearly have that $\|\tilde{\mathbf{x}}\|_0 \leq \|\mathbf{x}\|_0 \leq s$. Therefore, $\|\mathbf{x} - \tilde{\mathbf{x}}\|_0 \leq 2s$ and we can apply the RIP inequality on the vector $\mathbf{x} - \tilde{\mathbf{x}}$. But, since $\mathbf{W}(\mathbf{x} - \tilde{\mathbf{x}}) = 0$ we get that $|0 - 1| \leq \epsilon$, which leads to a contradiction.

Remark. The theorem above establishes that RIP matrices yield a lossless compression scheme for sparse vectors. It also provides a (nonefficient) reconstruction scheme.

Compressed sensing

The reconstruction scheme given in Theorem 1 seems to be nonefficient because we need to minimize a combinatorial objective (the sparsity of \mathbf{v}). Quite surprisingly, it turns out that we can replace the combinatorial objective, $\|\mathbf{v}\|_0$, with a convex objective, $\|\mathbf{v}\|_1$, which leads to a linear programming problem that can be solved efficiently.

Theorem 2. Assume that the conditions of Theorem 1 holds and that $\epsilon < \frac{1}{1+\sqrt{2}}$. Then

$$\mathbf{x} = \arg \min_{\mathbf{v}: \mathbf{W}\mathbf{v}=\mathbf{y}} \|\mathbf{v}\|_0 = \arg \min_{\mathbf{v}: \mathbf{W}\mathbf{v}=\mathbf{y}} \|\mathbf{v}\|_1.$$

When a matrix is RIP

- A random matrix with $n \geq \Omega(s \log(d))$ are likely to be RIP.
- In fact, multiplying a random matrix by an orthonormal matrix also provides an RIP matrix.
- This is important for compressing signals of the form $\mathbf{x} = \mathbf{U}\alpha$ where \mathbf{x} is not sparse but α is sparse. In that case, if \mathbf{W} is a random matrix and we compress using $\mathbf{y} = \mathbf{W}\mathbf{x}$ then this is the same as compressing α by $\mathbf{y} = (\mathbf{W}\mathbf{U})\alpha$ and since $\mathbf{W}\mathbf{U}$ is also RIP we can reconstruct α (and thus also \mathbf{x}) from \mathbf{y} .

More precisely, we have:

Theorem 3. Let \mathbf{U} be an arbitrary fixed $d \times d$ orthonormal matrix, let ϵ, δ be scalars on $(0, 1)$, let s be an integer in $[d]$, and let n be an integer that satisfies

$$n \geq 100 \frac{s \log(40d/(\delta\epsilon))}{\epsilon^2}.$$

Let $\mathbf{W} \in \mathbb{R}^{n \times d}$ be a matrix s.t. each element of \mathbf{W} is distributed normally with zero mean and variance of $1/n$. Then, with probability at least $1 - \delta$ over the choice of \mathbf{W} , the matrix \mathbf{WU} is (ϵ, s) -RIP.