# Assignment 1: CS 754, Advanced Image Processing

Due: 9th Feb before 11:55 pm

**Remember the honor code while submitting this (and every other) assignment. All members of the group should work on and *understand* all parts of the assignment. Exchange of answers between groups is not allowed. We will adopt a zero-tolerance policy against any violation, and we will expressly check for plagiarism.**

**Submission instructions:** You should ideally type out all the answers in Latex or else in MS Word with the equation editor. In either case, prepare a pdf file. Create a single zip or rar file containing the report, code and sample outputs and name it as follows: A1-IdNumberOfFirstStudent-IdNumberOfSecondStudent.zip. (If you are doing the assignment alone, the name of the zip file is A1-IdNumber.zip). Upload the file on moodle BEFORE 11:55 pm on 9th Feb. Beyond the cutoff time of 7 am on 10th Feb, no assignments will be accepted. Note that only one student per group should upload their work on moodle. Please preserve a copy of all your work until the end of the semester. *If you have difficulties, please do not hesitate to seek help from me.*

1. Let $\boldsymbol{\theta}^\star$ be the result of the following minimization problem (BP): $\min \|\boldsymbol{\theta}\|_1$ such that $\|\boldsymbol{y} - \boldsymbol{\Phi}\boldsymbol{\Psi}\boldsymbol{\theta}\|_2 \leq \varepsilon$, where $\boldsymbol{y}$ is an $m$-element measurement vector, $\boldsymbol{\Phi}$ is a $m \times n$ measurement matrix $(m < n)$, $\boldsymbol{\Psi}$ is a $n \times n$ orthonormal basis in which $n$-element signal $\boldsymbol{x}$ has a sparse representation of the form $\boldsymbol{x} = \boldsymbol{\Psi}\boldsymbol{\theta}$. Notice that $\boldsymbol{y} = \boldsymbol{\Phi}\boldsymbol{x} + \boldsymbol{\eta}$ and $\varepsilon$ is an upper bound on the magnitude of the noise vector $\boldsymbol{\eta}$.

   Theorem 3 we studied in class states the following: If $\boldsymbol{\Phi}$ obeys the restricted isometry property with isometry constant $\delta_{2s} < \sqrt{2} - 1$, then we have $\|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2 \leq C_1 s^{-1/2} \|\boldsymbol{\theta} - \boldsymbol{\theta_s}\|_1 + C_2\varepsilon$ where $C_1$ and $C_2$ are functions of only $\delta_{2s}$ and where $\forall i \in \mathcal{S}, \boldsymbol{\theta_{si}} = \theta_i; \forall i \notin \mathcal{S}, \boldsymbol{\theta_{si}} = 0$. Here $\mathcal{S}$ is a set containing the $s$ largest magnitude elements of $\boldsymbol{\theta}$.

   A curious student asks the following questions: '(1) It appears that the upper bound on $\|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2$ is reduced as $s$ increases, which goes against the very premise of compressed sensing. How do we address this apparent discrepancy? (2) It also appears that the error bound is independent of $m$. How do you address this? (3) Now consider that I gave you another theorem (called Theorem 3A), which is the same as Theorem 3 except that it requires that $\delta_{2s} < 0.3$. Out of Theorem 3 and Theorem 3A, which is the more useful theorem? Why? (4) It appears that if I set $\varepsilon = 0$ in BP, I can always reduce the upper bound on the error even if the noise vector $\boldsymbol{\eta}$ has non-zero magnitude. What am I missing?'
   Your job is to answer all four of the student's questions. [5+5+5+5=20 points]
   **Solution:** (1) The error bound actually does not decrease. First of all, $C_1, C_2$ are monotonically increasing functions of $\delta_{2s}$. As $s$ increases, so does $\delta_{2s}$ and hence $C_1, C_2$ both increase. Moreover as $s$ increases, the minimum number of measurements for the RIP to hold ($m \geq \mathcal{O}(s \log n)$), also increases. (2) The error bound is not independent of $m$. RIP for random matrices will typically require $m \geq \mathcal{O}(s \log n)$ measurements to guarantee reconstruction of $s$-sparse signals. (3) Theorem 3A is less powerful, as its guarantees hold for a smaller set of values of $\delta_{2s}$ than Theorem 3 which requires $\delta_{2s} < 0.414$. (4) The error bounds for Theorem 3 are *conditional* on $\|\boldsymbol{y} - \boldsymbol{Ax}\|_2$ actually following below a sensible value of $\varepsilon$, provided based on the noise level. Setting $\varepsilon$ to 0 would violate that condition and therefore the error bounds would no longer hold.
   **Marking scheme:** No marks for guesswork. A convincing reason needs to be provided in all cases.

2. In class, we studied a video compressive sensing architecture from the paper 'Video from a single exposure coded snapshot' published in ICCV 2011 (See http://www.cs.columbia.edu/CAVE/projects/single_shot_video/). Such a video camera acquires a 'coded snapshot' $E_u$ in a single exposure time interval $u$. This coded snapshot is the superposition of the form $E_u = \sum_{t=1}^{T} C_t \cdot F_t$ where $F_t$ is the image of the scene

at instant $t$ within the interval $u$ and $C_t$ is a randomly generated binary code at that time instant, which modulates $F_t$. Note that $E_u$, $F_t$ and $C_t$ are all 2D arrays. Also, the binary code generation as well as the final summation all occur within the hardware of the camera. Your task here is as follows:

(a) Read the 'cars' video in the homework folder in MATLAB using the 'mmread' function which has been provided in the homework folder and convert it to grayscale. Extract the first $T = 3$ frames of the video.

(b) Generate a $H \times W \times T$ random code pattern whose elements lie in $\{0, 1\}$. Compute a coded snapshot using the formula mentioned and add zero mean Gaussian random noise of standard deviation 2 to it. Display the coded snapshot in your report.

(c) Given the coded snapshot and assuming full knowledge of $C_t$ for all $t$ from 1 to $T$, your task is to estimate the original video sequence $F_t$. For this you should rewrite the aforementioned equation in the form $\boldsymbol{Ax} = \boldsymbol{b}$ where $\boldsymbol{x}$ is an unknown vector (vectorized form of the video sequence). Mention clearly what $\boldsymbol{A}$ and $\boldsymbol{b}$ are, in your report.

(d) You should perform the reconstruction using Orthogonal Matching Pursuit (OMP). For computational efficiency, we will do this reconstruction patchwise. Write an equation of the form $\boldsymbol{Ax} = \boldsymbol{b}$ where $\boldsymbol{x}$ represents the $i^{th}$ patch from the video and having size (say) $8 \times 8 \times T$ and mention in your report what $\boldsymbol{A}$ and $\boldsymbol{b}$ stand for. For perform the reconstruction, assume that each $8 \times 8$ slice in the patch is sparse or compressible in the 2D-DCT basis. Carefully work out the error term in the OMP algorithm, and explain this in your report!

(e) Repeat the reconstruction for all overlapping patches and average across the overlapping pixels to yield the final reconstruction. Display the reconstruction and mention the relative mean squared error between reconstructed and original data, in your report as well as in the code.

(f) Repeat this exercise for $T = 5, T = 7$ and mention the mention the relative mean squared error between reconstructed and original data again.

(g) **Note: To save time, extract a portion of about $120 \times 240$ around the lowermost car in the cars video and work entirely with it. In fact, you can show all your results just on this part. Some sample results are included in the homework folder.**

(h) Repeat the experiment with any consecutive 5 frames of the 'flame' video from the homework folder. [35 points = 18 points for successful OMP implementation + 7 points for carefully presenting error term bound + 10 points for displaying of all results]

**Solution and Marking scheme:** The values of a Gaussian random variable with mean 0 and standard deviation $\sigma$ lie in the range $[-3\sigma, +3\sigma]$ with 0.99 probability. Hence a reasonable bound for the error for OMP would be $\|\boldsymbol{b} - \boldsymbol{Ax}\|_2^2 \leq 3m\sigma^2$ for $m$ measurements. A more rigorous treatment would involve tail bounds of the chi-square distribution with $m$ degrees of freedom. We would have $P(\|\boldsymbol{b} - \boldsymbol{Ax}\|_2^2 \leq \sigma^2(m + 2m(\sqrt{t} + t))) \geq 1 - \exp(-tm)$. See https://math.stackexchange.com/questions/2864188/chi-squared-distribution-tail-bound and references therein.
The first solution with $3\sigma^2$ is less rigorous but acceptable and will be awarded full 7 marks. For a basic working OMP which is stub-tested on simple matrices, we will award 9 points out of 18. The remaining 9 are for correct implementation on *this* particular video CS problem. 10 points for display of appropriate results. For every result not displayed in the report, 2 points will be deducted. If the result is displayed in the code, but not the report, only 1 point will be deducted. For the flame video, the five consecutive frames should show *some* motion, otherwise two points are to be deducted.

3. We will prove why the value of the coherence between $m \times n$ measurement matrix $\boldsymbol{\Phi}$ (with all rows normalized to unit magnitude) and $n \times n$ orthonormal representation matrix $\boldsymbol{\Psi}$ must lie within the range $[1, \sqrt{n}]$ (both 1 and $\sqrt{n}$ inclusive). Recall that the coherence is given by the formula
$\mu(\boldsymbol{\Phi}, \boldsymbol{\Psi}) = \sqrt{n}\max_{i \in \{0,1,...,m-1\}, j \in \{0,1,...,n-1\}}|\boldsymbol{\Phi}^{i^t}\boldsymbol{\Psi}_j|$. Proving the upper bound should be very easy for you. To prove the lower bound, proceed as follows. Consider a unit vector $\boldsymbol{g} \in \mathbb{R}^n$. We know that it can be expressed as $\boldsymbol{g} = \sum_{k=1}^{n} \alpha_k \boldsymbol{\Psi_k}$ as $\boldsymbol{\Psi}$ is an orthonormal *basis*. Now prove that $\mu(\boldsymbol{g}, \boldsymbol{\Psi}) = \sqrt{n}\max_{i \in \{0,1,...,n-1\}} \frac{|\alpha_i|}{\sum_{j=1}^{n} \alpha_j^2}$.

Exploiting the fact that $g$ is a unit vector, prove that the minimal value of coherence is attained when $g = \sqrt{1/n}\sum_{k=1}^{n}\Psi_k$ and that hence the minimal value of coherence is 1. [10 points]

**Solution:** The upper bound is achieved when there exists a row of $\Phi$ which is equal to (or the element-wise negative of) some column of $\Psi$. This produces $\mu = \sqrt{n}$. As given in the question, we have $\mu(\Psi, g) = \max_{i \in \{0,1,...,n-1\}} \dfrac{|\alpha_i|}{\sqrt{\sum_{j=1}^{n}\alpha_j^2}}$. As $g$ is a unit vector and $\Psi$ is orthonormal, we must have $\|g\|_2^2 = \sum_{j=1}^{n}\alpha_j^2 = 1$.

Hence $\mu(\Psi, g) = \max_{i \in \{0,1,...,n-1\}}|\alpha_i|$, and we want the value of $\mu(\Psi, g)$ to be the least possible. Thus $\mu^{min}(\Psi, g) = \sqrt{n}\,\mathrm{argmin}_{\alpha_0,\alpha_1,...,\alpha_{n-1}}\max_{i \in \{0,1,...,n-1\}}|\alpha_i|$. Since $\sum_j \alpha_j^2 = 1$, the maximum value of $|\alpha_i|$ will be the least when $\forall i, \forall j \neq i, |\alpha_i| = |\alpha_j|$, which happens when $\forall i, \alpha_i = \dfrac{\pm 1}{\sqrt{n}}$. This yields $\mu^{min}(\Psi, g) = \sqrt{n}\left|\dfrac{\pm 1}{\sqrt{n}}\right| = 1$.

**Marking scheme:** 3 points for upper bound. 7 points for lower bound.

4. Compressive sensing reconstructions involve estimating a sparse signal $x \in \mathbb{R}^n, n \gg 2$ from a vector $y \in \mathbb{R}^m$ ($m \ll n$) of compressed measurements of the form $y = \Phi x$ where $\Phi \in \mathbb{R}^{m \times n}$ is the measurement matrix (assume there is no noise). Now answer the following questions, from first principles. **Do not merely quote theorems or algorithms.**

   (a) If it is known that $x$ has only 1 non-zero element and that the other elements are zero, can you uniquely estimate $x$ if $m = 1$? If yes, how? If not, why not? Now further suppose, you knew beforehand the index (but not the value) of the non-zero element of $x$? Does this help you any further? If yes, how? If not, why not?

   (b) If it is known that $x$ has only 1 non-zero element and that the other elements are zero, can you uniquely estimate $x$ if $m = 2$? If yes, how? If not, why not?

   (c) If it is known that $x$ has only 2 non-zero elements and that the other elements are zero, can you uniquely estimate $x$ if $m = 3$? If yes, describe an algorithm that is guaranteed to estimate it accurately. If not, explain why not, and explain whether there are any special instances of $\Phi$ for which unique estimation is possible?

   (d) Repeat part (c) with $m = 4$. [1+2+3+4=10 points]

**Solution:** (a) You cannot estimate $x$ if $m = 1$ because $\dfrac{y}{\Phi_j}$ could pass of as the non-zero value in $x$ for any index $j \in \{1, 2, ..., n\}$. If you knew the index with the non-zero element (say $k$), then you can uniquely determine $x$ as long $\Phi_k \neq 0$.
(b) In this case, we have $y_1 = \Phi_{1k}x_k$ and $y_2 = \Phi_{2k}x_k$. In this case, we can compute $\dfrac{y_1}{\Phi_{1j}}$ for all indices $j$ and choose the one for which $y_2 = \Phi_{2j}\dfrac{y_1}{\Phi_{1j}}$. The solution is unique if no two columns of $\Phi$ are parallel to each other.

(c) In this case, we have $\begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} \Phi_{1j} & \Phi_{1k} \\ \Phi_{2j} & \Phi_{2k} \\ \Phi_{3j} & \Phi_{3k} \end{pmatrix} \begin{pmatrix} x_j \\ x_k \end{pmatrix}$. A unique solution can be obtained only if you knew $j$ and $k$, the two non-zero indices. Otherwise a unique solution is not possible because to uniquely recover a 2-sparse vector, it is necessary and sufficient that any 4 columns of $\Phi$ be linearly independent (see proof for uniqueness of solution of problem P0 in the lecture slides). This is impossible because $\Phi$ has only 3 rows.
(d) In this case, a unique solution for the 2-sparse vector $x$ may be obtained if $\Phi$ satisfies the requirement that any four of its columns must be linearly independent. This is easy to satisfy with high probability if $\Phi$ is picked randomly from a Gaussian or a Bernoulli distribution. A brute-force method of estimating $x$ is to search over all column-subsets of $\Phi$ of size 4, and compute a putative solution for each by simple matrix inverse. The solution which is the sparsest is the one we want. Uniqueness of the solution is guaranteed by the assumption on $\Phi$ - see the lecture slides for proof of uniqueness of solution of problem P0. Another method: A brute-force method of estimating $x$ is to (1) search over all column-subsets of $\Phi$ of size 2 to yield indices $j_1, j_2$ as the possible support of the 2-sparse vector $x$, (2) compute a putative solution $\hat{x}_{j_1,j_2}$ for each

such pair $j_1, j_2$ by simple matrix inverse using any two measurements $y_1, y_2$, and (3) retain only that solution $\hat{\boldsymbol{x}}_{j_1,j_2}$ for which $y_3 = \boldsymbol{\Phi}^3\hat{\boldsymbol{x}}_{j_1,j_2}$ and $y_4 = \boldsymbol{\Phi}^4\hat{\boldsymbol{x}}_{j_1,j_2}$.

5. Read the paper 'Coded aperture compressive temporal imaging' which is a type of video compressive sensing camera developed at Duke University. The paper can be accessed here: `https://www.osapublishing.org/oe/fulltext.cfm?uri=oe-21-9-10526&id=253002`. A local copy is also placed in the homework folder. Answer the following questions:

   (a) Explain the similarities and differences between this architecture and the video compressed sensing architecture using coded snapshots (Hitomi et al, ICCV 2011) that we studied in class.

   (b) What cost function is the reconstruction algorithm in the paper based on? Explain the meaning and dimensions of every term in the cost function. (Note: you are not expected to describe the algorithms that minimize this cost function, but only mention the cost function itself) [8+7=15 points]

   **Solution and Marking Scheme:** For part (a): The overall architecture for video compressive sensing is very similar in both cases and the hardware acquisition is represented in the form $I(x, y) = \sum_{t=1}^{T} E(x, y, t)S(x, y, t)$ where $I(x, y)$ is the pixel at location $(x, y)$ in the coded snapshot, $T$ is the number of subframes of the underlying video given as $E(x, y, t)$ and $S(x, y, t)$ stands for the binary code value at pixel $(x, y)$ at time instant $t$. However in the Hitomi camera, the code $S(x, y, t)$ is implemented by means of a DMD array which changes its code $T$ times within a single exposure period. On the other hand, in the CACTI camera, the code $S(x, y, t)$ is implemented by means of a coded aperture mask which is moved $T$ times (at random) during a single exposure period. The movement of this coded aperture is due to a piezo-electric mechanism.
   **Marking scheme:** 4 points for basic similarity (equation of some form preferable but not necessary if the answer is clear). 4 points for the difference.
   Part (b): The cost function is given in equation 10 of the paper and has the form $J(\boldsymbol{f}) = \|\boldsymbol{g} - \boldsymbol{H}\boldsymbol{f}\|_2^2 + \lambda TV(\boldsymbol{f})$, where $\boldsymbol{g}$ is a vectorized form of the snapshot image $I$ of size $N \times N$, $\boldsymbol{f}$ is the underlying video of size $N \times N \times T$ and $\boldsymbol{H}$ is a forward model matrix of size $N^2 \times N^2 T$. $\boldsymbol{H}$ is a column-wise concatenation of diagonal matrices representing the code at each time $t$.
   Also $TV(\boldsymbol{f}) := \sum_{t=1}^{T}\sum_{i,j}\sqrt{(f(i+1, j, t) - f(i, j, t))^2 + (f(i, j+1, t) - f(i, j, t))^2}$ stands for the total variation of $\boldsymbol{f}$.
   **Marking scheme:** Clear definition of $\boldsymbol{g}, \boldsymbol{H}, \boldsymbol{f}$ in terms of the coded snapshot, aperture codes and the underlying video is essential. A clear definition of the TV norm must also be included in the answer (even if it is there in the paper being referred to).

6. Here is our mandatory Google search question. Note that this is the only question for which you can perform a google search to get the answer. A very interesting application of compressed sensing is in the area of optical or electron microscopy. Your task is to search for a research paper which applies compressed sensing in this application. Answer the following questions briefly:

   (a) Mention the title of the paper, where and when it was published and include a link to it.

   (b) Very briefly describe the hardware architecture used in the paper. You may refer to figures from the paper itself.

   (c) What reconstruction techhnique or cost function does the paper adopt for the sake of compressive reconstruction in this application? [3+4+4=10 points]

   **Solution:** I will give one example, referring to `https://www.pnas.org/content/109/26/E1679`, 'Compressive fluorescence microscopy for biological and hyperspectral imaging'. The hardware architecture for biological imaging in the paper is very similar to the Rice single pixel camera. The sensing matrix is random binary and implemented via a DMD system. The representation matrix proposed is either wavelets of Fourier. The reconstruction technique is same as problem P1 (or the LASSO problem in equation 6 – a relaxed version of P1).