# PHYS/ASTR 8150
## Fourier Transforms
## Linear Inverse Problems
## Applications to Imaging : Denoising and Deblurring

Fabien Baron

Georgia State University

*fbaron@gsu.edu*

Fall 2025

# DFT: Discrete Fourier Transform

- The Discrete Fourier Transform (DFT) transforms a sequence of $N$ complex numbers $\{\mathbf{x_n}\} := x_0, x_1, \ldots, x_{N-1}$ into another sequence of complex numbers, $\{\mathbf{X_k}\} := X_0, X_1, \ldots, X_{N-1}$, which is defined by
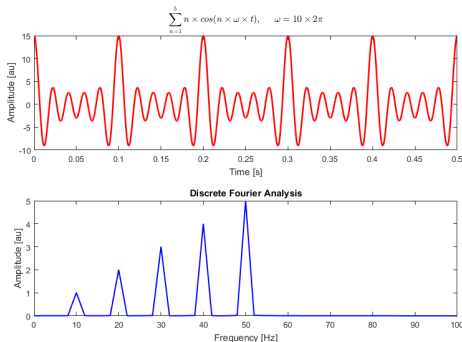
$$X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-\frac{i2\pi}{N}kn} = \sum_{n=0}^{N-1} x_n \cdot \left[ \cos\left(\frac{2\pi}{N}kn\right) - i \cdot \sin\left(\frac{2\pi}{N}kn\right) \right].$$

- The inverse transform is given by:

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k \cdot e^{i2\pi kn/N}$$

- The DFT is a linear transformation and thus can also be written in terms of a DFT matrix $\mathbf{F}$
- The transform is sometimes denoted by the symbol $\mathcal{F}$ or $\mathbf{F}$, as in $\mathbf{X} = \mathcal{F}\{\mathbf{x}\}$ or $\mathcal{F}(\mathbf{x})$ or $\mathbf{Fx}$.
- When scaled appropriately it becomes a unitary matrix $\mathcal{F}\mathcal{F}^{-1} = I$ and the $X_k$ thus be viewed as coefficients of $x$ in an orthonormal basis.

$$\sum_{n=1}^{5} n \times \cos(n \times \omega \times t), \quad \omega = 10 \times 2\pi$$

Discrete Fourier Analysis

- Evaluating this definition directly requires $O\left(N^2\right)$ operations: there are $N$ outputs $X_k$, and each output requires a sum of $N$ terms. An FFT is any method to compute the same results in $O(N \log N)$ operations.
- The Fast Fourier Transform is computed on a uniformly sampled array and returns a uniformly sampled **spatial frequency decomposition**.

# FFT for 2D imaging

- Check the fft_imaging.jl script
- What do frequencies represent? Low vs high frequencies?
- What does fftshift do?

# Phase of an image

- We often call the phase of an image the phase of its Fourier transform.
- The Fourier phase carries most of the structure location information, while the amplitude carries most of the structure power.

# Power Spectral Densities

- The PSD of an image describes how its energy is distributed across spatial frequencies. The Power Spectral Density of image *a* is:

$$PSD = |F(a)|^2$$

The square modulus in the PSDs offers better noise debiasing than modulus-only when under Gaussian noise conditions (e.g. trying to assess the PSD of an object by repeat measurements).

- Natural objects/scenes have a PSD that falls off approximately as $1/f^\alpha$ with spatial frequency, $\alpha \simeq 2$. Images with randomized Fourier phase but $1/f^2$ amplitude spectra look "cloudy" or "naturalistic".

# Understanding what Convolutions are

- Convolutions are used to model imaging instruments with linear responses.
- The response to an impulse function, the Point Spread Function (PSF), tells us what a point entering the instrument will looks like on the detector.
- One can then calculate the image resulting from looking at an object (= merely a collection of points) using a convolution.
- An image = the convolution of an object by a point-spread function
- In image space, each pixel of the original object becomes the PSF weighted by that pixel flux
- In Fourier space, a convolution acts as a filter. A convolution in image space becomes a Hadamard product (also known as "element by element multiplication").

# Convolution math

- Continuous functions:

$$(a * b)(t) := \int_{-\infty}^{\infty} a(\tau) b(t - \tau) d\tau$$

- Discrete functions:

$$(a * b)[n] = \sum_{m=-\infty}^{\infty} a[m] b[n - m]$$

- In practice in CS, while the convolution can be computed directly for a small number of pixels N, they involve $N^2$ operations. FFT methods are used for better $(N \log N)$ scaling.

$$c = a * b \rightarrow F(c) = F(a) \cdot F(b) \rightarrow c = \Re \left\{ F^{-1} \left( F(a) \cdot F(b) \right) \right\}$$

# Correlation

- Cross-correlation can be used to match templates *a* and *b* by identifying correlation peaks.
- Discrete version:

$$(a \star b)[n] = \sum_{m=-\infty}^{\infty} \overline{a[m-n]} b[m]$$

where $\overline{\text{something}}$ denotes the complex conjugate.

- $c = a \star b \rightarrow F(c) = F(a) \cdot \overline{F(b)} \rightarrow c = \Re \left\{ F^{-1} \left( F(a) \cdot \overline{F(b)} \right) \right\}$

- An autocorrelation is the correlation of an function/image/cube with itself, and is the inverse Fourier transform of the PSD:

$$c = \Re \left\{ F^{-1} \left( |F(a)|^2 \right) \right\}$$

- Autocorrelations are used to identify where the function may repeat itself, e.g. to detect pure translation motion in a sequence of frames.

# NFFT: Non-equispaced Fourier Transform

- The NFFT is a generalization of the FFT so that it works on a non-uniformly sampled signal and/or returns non-uniformly sampled spatial frequency decomposition.
- Three use cases of NFFT: uniform to non-uniform (similar to DFT!), non-uniform to uniform, non-uniform to non-uniform. Write down the examples given in class.

## Variance and covariance: a reminder

- Variance of a scalar-valued random variable X:
  $\sigma_X^2 = \text{var}(X) = E[(X - E[X])^2] = E[(X - E[X]) \cdot (X - E[X])]$

- Covariance between two scalar-valued random variables $X$ and $Y$:

$$\text{cov}(X, Y) = E\left[(X - E[X])(Y - E[Y])\right] = E[XY] - E[X]E[Y]$$

- Covariance matrix of random vector $\boldsymbol{X}$:

$$\Sigma_{\boldsymbol{XX}} = \begin{bmatrix} E[(X_1 - E[X_1])(X_1 - E[X_1])] & E[(X_1 - E[X_1])(X_2 - E[X_2])] & \cdots & E[(X_1 - E[X_1])(X_n - E[X_n])] \\ E[(X_2 - E[X_2])(X_1 - E[X_1])] & E[(X_2 - E[X_2])(X_2 - E[X_2])] & \cdots & E[(X_2 - E[X_2])(X_n - E[X_n])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - E[X_n])(X_1 - E[X_1])] & E[(X_n - E[X_n])(X_2 - E[X_2])] & \cdots & E[(X_n - E[X_n])(X_n - E[X_n])] \end{bmatrix}$$
$$= E[(\boldsymbol{X} - E[\boldsymbol{X}])^\top (\boldsymbol{X} - E[\boldsymbol{X}])]$$

- Cross-covariance matrix of random vectors $\boldsymbol{X}$ and $\boldsymbol{Y}$:
  $\Sigma_{\boldsymbol{XY}} = E[(\boldsymbol{X} - E[\boldsymbol{X}])^\top (\boldsymbol{Y} - E[\boldsymbol{Y}])]$

# Inverse problems and regularization

- **Inverse problem**: recovering a wanted set of parameters from noise-corrupted data (e.g. pixel fluxes in image reconstruction)
- **Ill-posed** inverse problem: when the effective number of parameters to recover is greater than the effective number of data points, the solution may not exist or may not be unique
- **Regularization**: introduction of priors (known as regularizers) which will try to compensate for our lack of data. As usual, the prior/regularizer expresses the *a priori* probability of an object when we have no data.

## Linear inverse problems

- $y \in \mathbb{R}^m$ is our one-dimensional column vector of $m$ observations (our data)
- $x \in \mathbb{R}^n$ is the one-dimensional column vector, the object of interest to recover. For convenience, if $x$ is multi-dimensional, such as a $k \times k$ pixels image, it is sought as a column vector with $n = k^2$ pixels arranged in lexicographic order.
- $\boldsymbol{A} \in \mathbb{R}^{m \times n}$ is the observation matrix that transforms the object of interest $\boldsymbol{x}$ into data points $\boldsymbol{y}$. We have:

$$\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{n}$$

where $n$ is a noise vector. In this course we will assume Gaussian white noise.

# Linear inverse problems: examples

- In image deblurring & denoising, $y$ = noisy and blurred observed $k \times k$-pixel image, $x$ = denoised & deblurred $k \times k$-pixel object. $x$ and $y$ don't necessarily have the same size;
- In light curve inversion, $y$ is the flux received as a function of time while $x$ is a tessel map of the surface temperature.

# Direct matrix inversion

- $A$ is rectangular and in general not inversible since $A \in \mathbb{R}^{m \times n}$
- $A^\mathsf{T}A \in \mathbb{R}^{n \times n}$, $AA^\mathsf{T} \in \mathbb{R}^{m \times m}$ are invertible
- Direct linear inversion can be done in a non-Bayesian way:

$$Ax = y \quad (-n)$$
$$A^\mathsf{T}Ax = A^\mathsf{T}y$$
$$x = \left(A^\mathsf{T}A\right)^{-1} A^\mathsf{T}y = A^+ y$$

- $A^+ = \left(A^\mathsf{T}A\right)^{-1} A^\mathsf{T}$ is the Moore-Penrose pseudoinverse/generalized inverse of $A$.
- Application to denoising of unblurred images: since $A = I_n$, $x = y$. Not very helpful.

# Maximum Likelihood for inverse problems (1): the setup

- We know the distribution of the noise $\boldsymbol{n} = \boldsymbol{y} - \boldsymbol{Ax}$ and covariances $\{\sigma_{ij}\}_{i,j=1\ldots m}$ of our measurements (uncertainties).
- If the data points are independent then the inverse covariance matrix is diagonal: $\boldsymbol{\Sigma} = \begin{bmatrix} \frac{1}{\sigma_{11}^2} & & \\ & \ddots & \\ & & \frac{1}{\sigma_{mm}^2} \end{bmatrix}$, with $\sigma_{ij} = 0, \quad \forall i \neq j$.
- For a given $\boldsymbol{x}$, the realization of the noise $\boldsymbol{n} = \boldsymbol{y} - \boldsymbol{Ax}$ depends on $\boldsymbol{y}$ and is normally distributed; we want to maximize the likelihood, which is:

$$\Pr(\boldsymbol{y}|\boldsymbol{x}) = \mathcal{L}(\boldsymbol{x}|\boldsymbol{y}) \propto \prod_{i=1}^{m}\prod_{j=1}^{m} \exp\left[-\frac{(y_i - (\boldsymbol{Ax})_i)(y_j - (\boldsymbol{Ax})_j)}{2\sigma_{ij}^2}\right]$$

- We want to minimize the negative log-likelihood, which is:

$$-\log\Pr(\boldsymbol{y}|\boldsymbol{x}) = \mathrm{cst} + \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\frac{(y_i - (Ax)_i)(y_j - (Ax)_j)}{\sigma_{ij}^2} = \mathrm{cst} + \frac{1}{2}\chi^2(\boldsymbol{x})$$

# Maximum Likelihood for inverse problems (2): $\chi^2$

- Let's express the $\chi^2$ in matricial form:

$$\chi^2 = \sum_{i=1}^{m} \sum_{j=1}^{m} \frac{(y_i - (Ax)_i)(y_j - (Ax)_j)}{\sigma_{ij}^2}$$
$$= (\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x})^{\mathsf{T}} \boldsymbol{\Sigma} (\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}) = \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_{2,\boldsymbol{\Sigma}}^2$$

- Note: the squared $\ell_2$ norm of vector $\boldsymbol{\alpha}$ is:

$$\|\boldsymbol{\alpha}\|_2^2 = \boldsymbol{\alpha}^{\mathsf{T}} \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \ldots \alpha_m \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \end{bmatrix} = \sum_{i=1}^{m} \alpha_i^2$$

- The weighted $\ell_2$ norm of $\boldsymbol{\alpha}$ is $\|\boldsymbol{\alpha}\|_{2,\boldsymbol{M}}^2 = \boldsymbol{\alpha}^{\mathsf{T}} \boldsymbol{M} \boldsymbol{\alpha}$, where $\boldsymbol{M}$ is the weight matrix. For $M = \begin{bmatrix} M_{11} & & \\ & \ddots & \\ & & M_{mm} \end{bmatrix}$, $\|\boldsymbol{\alpha}\|_{2,\boldsymbol{M}}^2 = \sum_{i=1}^{m} M_{ii}\alpha_i^2$

# Maximum Likelihood for inverse problems (3): minimization

- The most likely solution is:

$$\widetilde{x} = \operatorname*{argmin}_{x \in \mathbb{R}^n} \chi^2(x) = \operatorname*{argmin}_{x \in \mathbb{R}^n} (y - Ax)^\mathsf{T} \Sigma (y - Ax)$$

- To minimize $\chi^2$, we set its gradient with respect to $x$ to 0:

$$\frac{\partial \chi^2}{\partial x} = \begin{bmatrix} \frac{\partial \chi^2}{\partial x_1} \\ \vdots \\ \frac{\partial \chi^2}{\partial x_n} \end{bmatrix} = \underbrace{\frac{\partial \left[ (y - Ax)^\mathsf{T} \Sigma (y - Ax) \right]}{\partial x}}_{-2A^\mathsf{T} \Sigma (y - Ax)} = 0$$

$$A^\mathsf{T} \Sigma A x = A^\mathsf{T} \Sigma y$$
$$\widetilde{x} = (A^\mathsf{T} \Sigma A)^{-1} A^\mathsf{T} \Sigma y$$

- $\widetilde{x}$ is the linear least squares estimator
- As with direct inversion, we may be **overfitting**, i.e. minimizing $\chi^2$ too much.

# Maximum a Posteriori: regularization and constraints

- We want to find the most probable image $\tilde{x}$ given data $y$, i.e. we want:

$$\tilde{x} = \underset{x \in \mathbb{R}^n}{\mathrm{argmax}}\, \Pr(x|y) \propto \underset{x \in \mathbb{R}^n}{\mathrm{argmax}}\, \Pr(y|x)\Pr(x)$$

- We add prior under the form of a **regularizer**: $R(x) = -\log \Pr(x)$

$$x = \underset{x \in \mathbb{U}}{\mathrm{argmin}}\, (y - Ax)^{\mathsf{T}}\Sigma(y - Ax) + \lambda R(x)$$

- Additional regularization is provided by restricting $x$ to $\mathbb{U}$. This allows to include **support constraints** ($\mathbb{U} = $ mask), positivity ($\mathbb{U} = \mathbb{R}^n_+$), or bound constraints ($\mathbb{U} = [a, b]^n$).

- The **objective function** $J(x) = (y - Ax)^{\mathsf{T}}\Sigma(y - Ax) + \lambda R(x)$ is composed of the $\chi^2$ term plus the **regularization function** $R$.

- $\lambda \in \mathbb{R}$ is a scalar **hyperparameter** governing the regularization weight

- $R$ and $\lambda$ should be chosen to prevent overfitting of the data.

# Maximum a Posteriori: generalized Tikhonov regularization

- The Tikhonov regularization is a prior with a simple $\ell_2$ norm on a linear transform of the $\boldsymbol{x}$, i.e. $R(\boldsymbol{x}) \propto \|\boldsymbol{\Gamma x}\|_2^2$

- The Tikhonov regularized maximum likelihood problem:

$$\widetilde{\boldsymbol{x}} = \underset{\boldsymbol{x} \in \mathbb{R}^n}{\operatorname{argmin}} \, (\boldsymbol{y} - \boldsymbol{Ax})^{\mathsf{T}} \boldsymbol{\Sigma} (\boldsymbol{y} - \boldsymbol{Ax}) + \lambda \|\boldsymbol{\Gamma x}\|_2^2$$

- The choice $\boldsymbol{\Gamma} = \boldsymbol{I}_n$ is called ridge regression regularization, i.e. $R(\boldsymbol{x}) = \lambda \|\boldsymbol{x}\|_2^2$. Ridge regression strongly penalizes pixels with large flux values, and weakly pixels with smaller flux values.

- The choice $\boldsymbol{\Gamma} = \boldsymbol{\nabla}$, where $\boldsymbol{\nabla}$ is the spatial gradient operator, is called Total Squared Variation. It penalizes strongly images with quickly varying fluxes (large spatial gradient), and weakly images with patches of uniform flux.

- Solving for $\widetilde{x}$, the most likely image/object of interest given the data:

$$\widetilde{x} = \underset{x \in \mathbb{R}^n}{\mathrm{argmin}} \, (y - Ax)^T \Sigma (y - Ax) + \lambda \|\Gamma x\|_2^2$$

$$\implies -2 A^T \Sigma (y - A\widetilde{x}) + 2\lambda \Gamma^T \Gamma \widetilde{x} = 0$$

$$(A^T \Sigma A + \lambda \Gamma^T \Gamma)\widetilde{x} = A^T \Sigma y$$

$$\boxed{\widetilde{x} = (A^T \Sigma A + \lambda \Gamma^T \Gamma)^{-1} A^T \Sigma y}$$

- The choice $\Gamma = I_n \implies \Gamma^T \Gamma = I_{n \times n}$.
- The choice $\Gamma = \nabla \implies \Gamma^T \Gamma = \nabla^T \nabla$.

# Choosing $\lambda$, under and overfitting, L-Curve

- Play with the Tikhonov code on our github repository

## Separable regularizers and $\ell_p$ norm

- For a set of independent parameters with identical priors (e.g. independent pixel fluxes with Poisson priors), the global regularizer is a **separable function**:

$$\Pr(\boldsymbol{x}) = \Pr(x_1) \times \Pr(x_2) \ldots \times \Pr(x_N) \implies R(\boldsymbol{x}) = \sum_i R(x_i)$$

- There are many classic regularizers. Amongst them, the $\ell_p$, $p \geq 1$ **norm** of vector $\boldsymbol{x}$ penalizes higher fluxes using a power law. The squared $\ell_2$ norm is used in Tikhonov regularization:

$$\ell_p(x) = \left[ \sum_{i=1}^{N} |x_i|^p \right]^{\frac{1}{p}}$$

- The **pseudo-norm** $\ell_0$ is even more strongly regularizing, since it counts all non-zero elements of $\boldsymbol{x}$ equally:

$$\ell_0(\boldsymbol{x}) = \sum_{i=1}^{N} \mathbb{1}(x_i > 0) \quad \text{where } \mathbb{1}() \text{ is the indicator function}$$

# Separable regularizers and $\ell_p$ norm: examples

- In practice for e.g. $\boldsymbol{x} = [0, 3, 0, 4]$
- $\ell_0$ is the number of non-zero $x_i$, $\ell_0(\boldsymbol{X}) = 2$.
- $\ell_1 = \sum_i |x_i|$ is the sum of moduli, $\ell_1(\boldsymbol{x}) = 7$.
- $\ell_2 = (\sum_i |x_i|^2)^{\frac{1}{2}}$ is the square root of the sum of square moduli, $\ell_2(\theta) = \sqrt{3^2 + 4^2} = 5$.
- Different values of $p$ penalize more or less higher vs lower fluxes.
- In MAP image reconstruction, when using $\ell_p$ regularizers in conjunction with $\chi^2$, the higher $p$, the smoother and less noisy the resulting image will be. E.g. $\ell_2$ give soft images, $\ell_1$ sharp images, $\ell_0$ very spiky images.
- $\ell_p(\theta - \gamma)$, where $\gamma$ is a default expected level, expressing defaults values for parameter set $\theta$, i.e. values that we would expect to be "normal" in absence of data. Example: the default expected flux level in a star field is zero in the absence of data, so the default image $\gamma$ is an image with zero everywhere. But the default image for an image of the Sun is a mostly uniform disc, so $\gamma$ would be this expected disc.

# Sparsity & dictionaries

- An object (image, array, vector) is said **sparse** in a basis called a **dictionary** if it can be represented in this basis by a small number of non-zero coefficients

- The **impulsion/image dictionary** is the conventional grid of square pixels. A stellar field image, with only pointlike stars and zeros elsewhere is sparse in the impulsion basis.

- JPEG stores he coefficients of images expressed in the DCT (discrete cosine transform) that work well with most natural pictures. Since fewer coefficients can be used than in image dictionary, this results in compression of the image data.

- JPEG2000 stores the coefficients of images expressed in **wavelet dictionaries** (CDF 9/7 or 5/3)

# Compressed sensing, analysis form

- **Compressed sensing theory** is a recent mathematical paradigm (from the 2000s) that asserts that the optimal regularization can be obtained by imposing sparsity in the basis where the solution is the most sparse. The MAP problem takes the **analysis form**:

$$\widetilde{x} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ - \log \mathcal{L}(x) + \lambda \, \ell_0(\mathbf{\Gamma} x) \right\}$$

- Compressed sensing literature recommends using $\ell_0$ on the sparsity basis as regularizers, in order to enforce sparsity. Since minimizing $\ell_0$ is a NP-hard problem, $\ell_1$ is often used instead.

- **Total variation** is $\ell_1(\nabla x)$, i.e. the $\ell_1$ nor of the **spatial gradient** of $x$. The dictionary in this case is the vectorial space of spatial gradients. In image reconstruction this regularizer favors patches of uniform flux with sharp transitions.

## Compressed sensing, synthesis form

- The **synthesis** approach is to look directly in the dictionary space, i.e. to look for a sparse vector $\boldsymbol{\theta} = \Gamma \boldsymbol{x}$ that represents the object of interest in the dictionary:

$$\widetilde{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \mathbb{R}^k}{\operatorname{argmin}} - \log \mathcal{L}(\Gamma^{-1}\boldsymbol{\theta}) + \lambda \, \ell_0(\boldsymbol{\theta})$$

then do $\boldsymbol{x} = \Gamma^{-1}\boldsymbol{\theta}$ to find the wanted solution.

- Doing image reconstruction by directly recovering CDF wavelets coefficients instead of pixels (impulsion dictionary).

- Analysis and synthesis are not equivalent. The dimensions of $\boldsymbol{x}$ and $\boldsymbol{\theta}$ may be very different. At this point in this course we cannot solve $\ell_0$ or $\ell_1$ regularization problems.