

Lecture Notes for EE263

Stephen Boyd

Introduction to Linear Dynamical Systems

Autumn 2008-09

Copyright Stephen Boyd. Limited copying or use for educational purposes is fine, but please acknowledge source, e.g., “taken from *Lecture Notes for EE263, Stephen Boyd, Stanford 2008.*”

Contents

Lecture 1 – Overview
Lecture 2 – Linear functions and examples
Lecture 3 – Linear algebra review
Lecture 4 – Orthonormal sets of vectors and QR factorization
Lecture 5 – Least-squares
Lecture 6 – Least-squares applications
Lecture 7 – Regularized least-squares and Gauss-Newton method
Lecture 8 – Least-norm solutions of underdetermined equations
Lecture 9 – Autonomous linear dynamical systems
Lecture 10 – Solution via Laplace transform and matrix exponential
Lecture 11 – Eigenvectors and diagonalization
Lecture 12 – Jordan canonical form
Lecture 13 – Linear dynamical systems with inputs and outputs
Lecture 14 – Example: Aircraft dynamics
Lecture 15 – Symmetric matrices, quadratic forms, matrix norm, and SVD
Lecture 16 – SVD applications
Lecture 17 – Example: Quantum mechanics
Lecture 18 – Controllability and state transfer
Lecture 19 – Observability and state estimation
Lecture 20 – Some final comments

Basic notation

Matrix primer

Crimes against matrices

Least-squares and least-norm solutions using Matlab

Solving general linear equations using Matlab

Low rank approximation and extremal gain problems

Exercises

Lecture 1

Overview

- course mechanics
- outline & topics
- what is a linear dynamical system?
- why study linear systems?
- some examples

1-1

Course mechanics

- all class info, lectures, homeworks, announcements on class web page:

`www.stanford.edu/class/ee263`

course requirements:

- weekly homework
- takehome midterm exam (date TBD)
- takehome final exam (date TBD)

Prerequisites

- exposure to linear algebra (*e.g.*, Math 103)
- exposure to Laplace transform, differential equations

not needed, but might increase appreciation:

- control systems
- circuits & systems
- dynamics

Major topics & outline

- linear algebra & applications
- autonomous linear dynamical systems
- linear dynamical systems with inputs & outputs
- basic quadratic control & estimation

Linear dynamical system

continuous-time linear dynamical system (CT LDS) has the form

$$\frac{dx}{dt} = A(t)x(t) + B(t)u(t), \quad y(t) = C(t)x(t) + D(t)u(t)$$

where:

- $t \in \mathbf{R}$ denotes *time*
- $x(t) \in \mathbf{R}^n$ is the *state* (vector)
- $u(t) \in \mathbf{R}^m$ is the *input* or *control*
- $y(t) \in \mathbf{R}^p$ is the *output*

Overview

1-5

- $A(t) \in \mathbf{R}^{n \times n}$ is the *dynamics matrix*
- $B(t) \in \mathbf{R}^{n \times m}$ is the *input matrix*
- $C(t) \in \mathbf{R}^{p \times n}$ is the *output* or *sensor matrix*
- $D(t) \in \mathbf{R}^{p \times m}$ is the *feedthrough matrix*

for lighter appearance, equations are often written

$$\dot{x} = Ax + Bu, \quad y = Cx + Du$$

- CT LDS is a first order vector *differential equation*
- also called *state equations*, or '*m*-input, *n*-state, *p*-output' LDS

Overview

1-6

Some LDS terminology

- most linear systems encountered are *time-invariant*: A, B, C, D are constant, *i.e.*, don't depend on t
- when there is no input u (hence, no B or D) system is called *autonomous*
- very often there is no feedthrough, *i.e.*, $D = 0$
- when $u(t)$ and $y(t)$ are scalar, system is called *single-input, single-output* (SISO); when input & output signal dimensions are more than one, MIMO

Discrete-time linear dynamical system

discrete-time linear dynamical system (DT LDS) has the form

$$x(t+1) = A(t)x(t) + B(t)u(t), \quad y(t) = C(t)x(t) + D(t)u(t)$$

where

- $t \in \mathbf{Z} = \{0, \pm 1, \pm 2, \dots\}$
- (vector) signals x, u, y are *sequences*

DT LDS is a first order vector *recursion*

Why study linear systems?

applications arise in **many** areas, *e.g.*

- automatic control systems
- signal processing
- communications
- economics, finance
- circuit analysis, simulation, design
- mechanical and civil engineering
- aeronautics
- navigation, guidance

Usefulness of LDS

- depends on availability of **computing power**, which is large & increasing exponentially
- used for
 - analysis & design
 - implementation, embedded in real-time systems
- like DSP, was a specialized topic & technology 30 years ago

Origins and history

- parts of LDS theory can be traced to 19th century
- builds on classical circuits & systems (1920s on) (transfer functions . . .) but with more emphasis on linear algebra
- first engineering application: aerospace, 1960s
- transitioned from specialized topic to ubiquitous in 1980s (just like digital signal processing, information theory, . . .)

Nonlinear dynamical systems

many dynamical systems are **nonlinear** (a fascinating topic) so why study **linear** systems?

- most techniques for nonlinear systems are based on linear methods
- methods for linear systems often work unreasonably well, in practice, for nonlinear systems
- if you don't understand linear dynamical systems you certainly can't understand nonlinear dynamical systems

Examples (ideas only, no details)

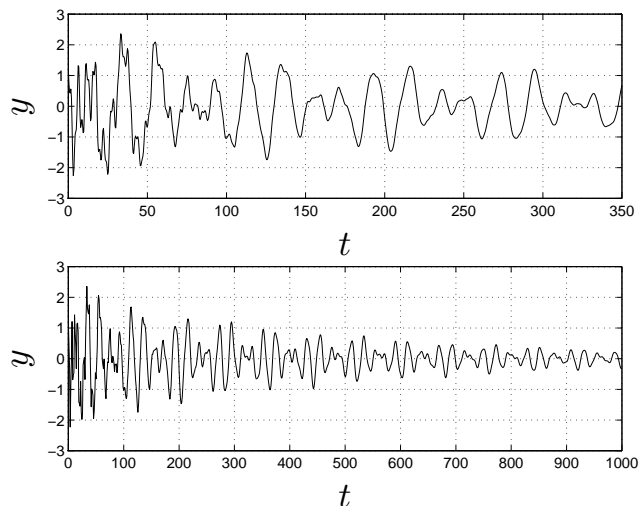
- let's consider a specific system

$$\dot{x} = Ax, \quad y = Cx$$

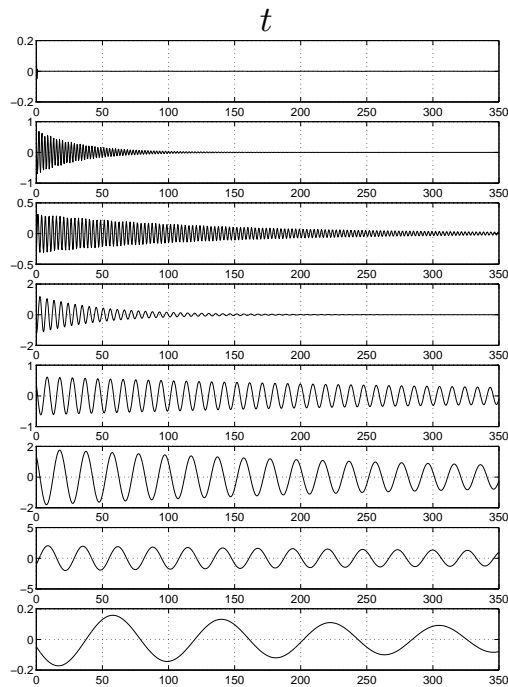
with $x(t) \in \mathbf{R}^{16}$, $y(t) \in \mathbf{R}$ (a '16-state single-output system')

- model of a lightly damped mechanical system, but it doesn't matter

typical output:



- output waveform is very complicated; looks almost random and unpredictable
- we'll see that such a solution can be decomposed into much simpler (modal) components



(idea probably familiar from ‘poles’)

Input design

add two inputs, two outputs to system:

$$\dot{x} = Ax + Bu, \quad y = Cx, \quad x(0) = 0$$

where $B \in \mathbf{R}^{16 \times 2}$, $C \in \mathbf{R}^{2 \times 16}$ (same A as before)

problem: find appropriate $u : \mathbf{R}_+ \rightarrow \mathbf{R}^2$ so that $y(t) \rightarrow y_{\text{des}} = (1, -2)$

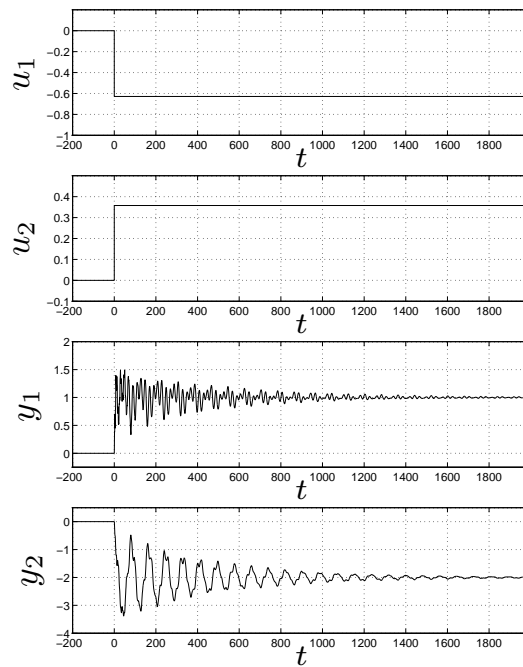
simple approach: consider static conditions (u , x , y constant):

$$\dot{x} = 0 = Ax + Bu_{\text{static}}, \quad y = y_{\text{des}} = Cx$$

solve for u to get:

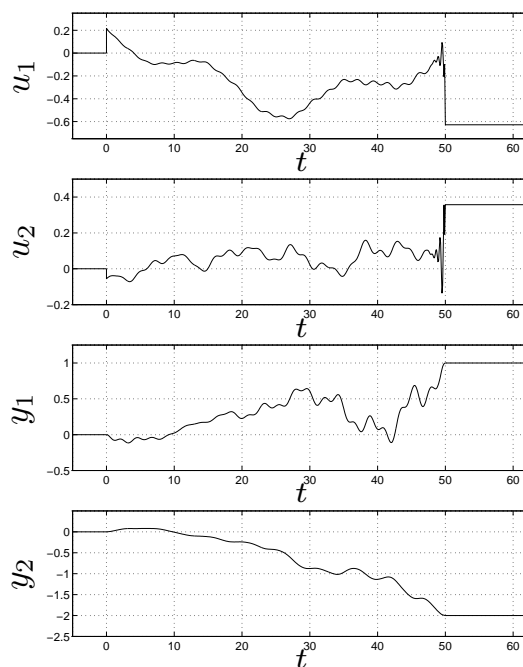
$$u_{\text{static}} = (-CA^{-1}B)^{-1} y_{\text{des}} = \begin{bmatrix} -0.63 \\ 0.36 \end{bmatrix}$$

let's apply $u = u_{\text{static}}$ and just wait for things to settle:



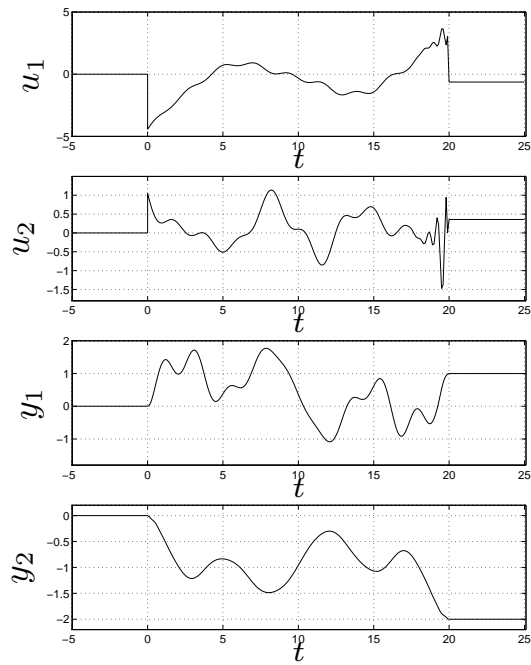
... takes about 1500 sec for $y(t)$ to converge to y_{des}

using very clever input waveforms (EE263) we can do much better, *e.g.*



... here y converges *exactly* in 50 sec

in fact by using larger inputs we do still better, *e.g.*

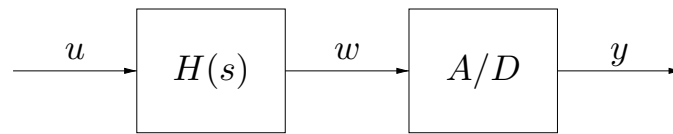


. . . here we have (exact) convergence in 20 sec

in this course we'll study

- how to synthesize or design such inputs
- the tradeoff between size of u and convergence time

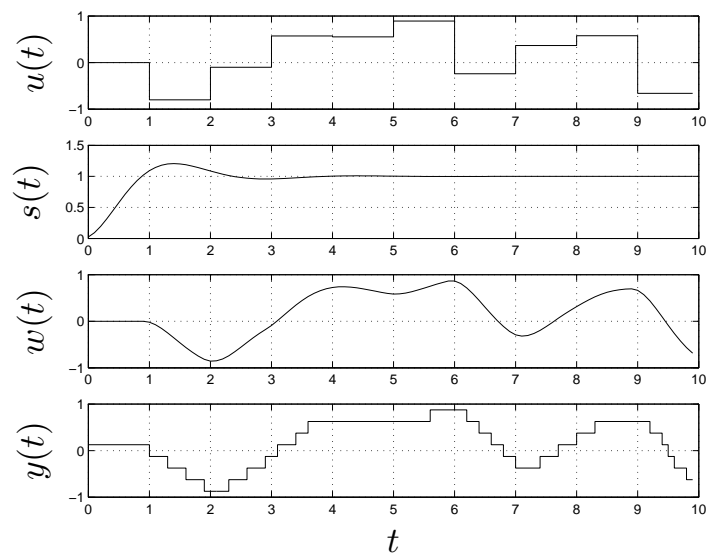
Estimation / filtering



- signal u is piecewise constant (period 1 sec)
- filtered by 2nd-order system $H(s)$, step response $s(t)$
- A/D runs at 10Hz, with 3-bit quantizer

Overview

1-21



problem: estimate original signal u , given quantized, filtered signal y

Overview

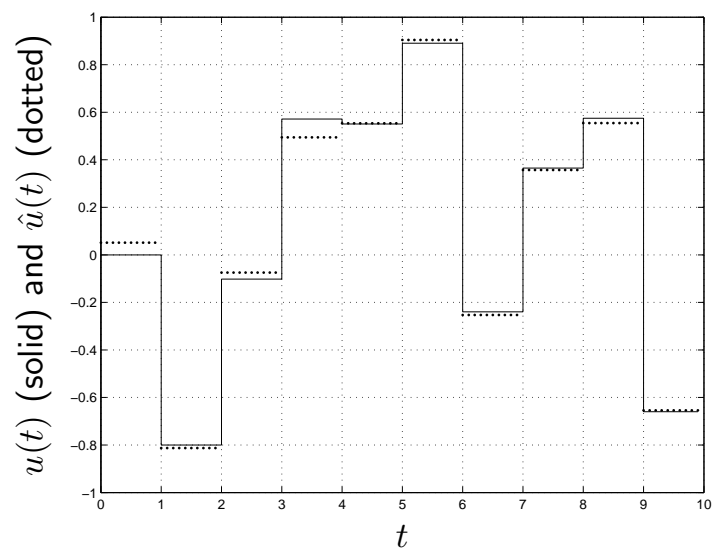
1-22

simple approach:

- ignore quantization
- design equalizer $G(s)$ for $H(s)$ (*i.e.*, $GH \approx 1$)
- approximate u as $G(s)y$

. . . yields terrible results

formulate as *estimation problem* (EE263) . . .



RMS error 0.03, well **below** quantization error (!)

Lecture 2

Linear functions and examples

- linear equations and functions
- engineering examples
- interpretations

2-1

Linear equations

consider system of linear equations

$$\begin{array}{rclclcl} y_1 & = & a_{11}x_1 & + & a_{12}x_2 & + \cdots + & a_{1n}x_n \\ y_2 & = & a_{21}x_1 & + & a_{22}x_2 & + \cdots + & a_{2n}x_n \\ & \vdots & & & & & \\ y_m & = & a_{m1}x_1 & + & a_{m2}x_2 & + \cdots + & a_{mn}x_n \end{array}$$

can be written in matrix form as $y = Ax$, where

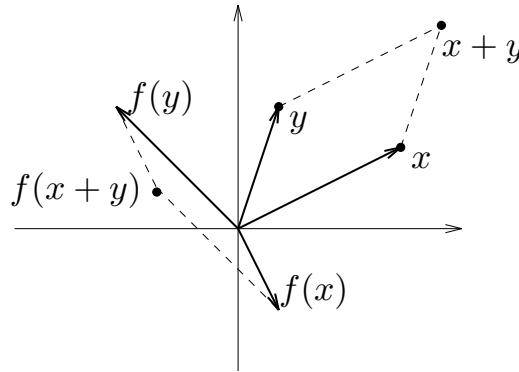
$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Linear functions

a function $f : \mathbf{R}^n \longrightarrow \mathbf{R}^m$ is *linear* if

- $f(x + y) = f(x) + f(y), \forall x, y \in \mathbf{R}^n$
- $f(\alpha x) = \alpha f(x), \forall x \in \mathbf{R}^n \forall \alpha \in \mathbf{R}$

i.e., *superposition* holds



Matrix multiplication function

- consider function $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ given by $f(x) = Ax$, where $A \in \mathbf{R}^{m \times n}$
- matrix multiplication function f is linear
- **converse** is true: **any** linear function $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ can be written as $f(x) = Ax$ for some $A \in \mathbf{R}^{m \times n}$
- representation via matrix multiplication is unique: for any linear function f there is only one matrix A for which $f(x) = Ax$ for all x
- $y = Ax$ is a concrete representation of a generic linear function

Interpretations of $y = Ax$

- y is measurement or observation; x is unknown to be determined
- x is 'input' or 'action'; y is 'output' or 'result'
- $y = Ax$ defines a function or transformation that maps $x \in \mathbf{R}^n$ into $y \in \mathbf{R}^m$

Interpretation of a_{ij}

$$y_i = \sum_{j=1}^n a_{ij}x_j$$

a_{ij} is *gain factor* from j th input (x_j) to i th output (y_i)

thus, *e.g.*,

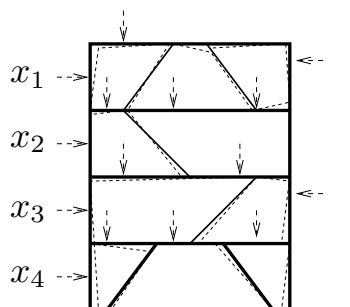
- i th *row* of A concerns i th *output*
- j th *column* of A concerns j th *input*
- $a_{27} = 0$ means 2nd output (y_2) doesn't depend on 7th input (x_7)
- $|a_{31}| \gg |a_{3j}|$ for $j \neq 1$ means y_3 depends mainly on x_1

- $|a_{52}| \gg |a_{i2}|$ for $i \neq 5$ means x_2 affects mainly y_5
- A is lower triangular, i.e., $a_{ij} = 0$ for $i < j$, means y_i only depends on x_1, \dots, x_i
- A is diagonal, i.e., $a_{ij} = 0$ for $i \neq j$, means i th output depends only on i th input

more generally, **sparsity pattern** of A , i.e., list of zero/nonzero entries of A , shows which x_j affect which y_i

Linear elastic structure

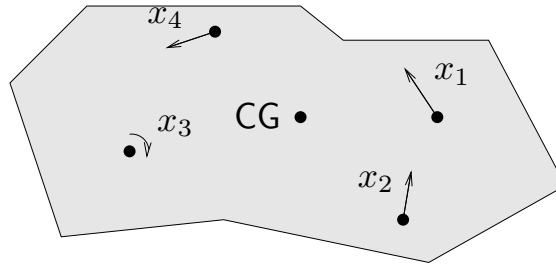
- x_j is external force applied at some node, in some fixed direction
- y_i is (small) deflection of some node, in some fixed direction



(provided x , y are small) we have $y \approx Ax$

- A is called the *compliance matrix*
- a_{ij} gives deflection i per unit force at j (in m/N)

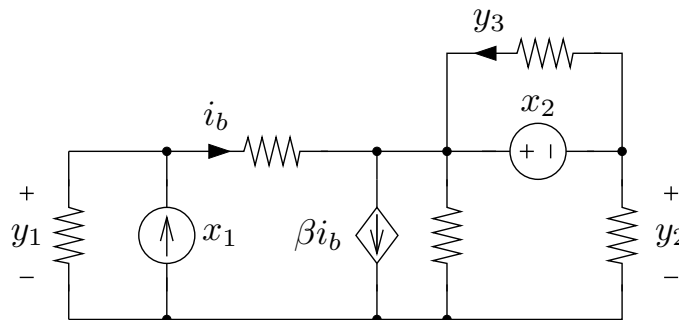
Total force/torque on rigid body



- x_j is external force/torque applied at some point/direction/axis
- $y \in \mathbf{R}^6$ is resulting total force & torque on body
(y_1, y_2, y_3 are \mathbf{x} -, \mathbf{y} -, \mathbf{z} - components of total force,
 y_4, y_5, y_6 are \mathbf{x} -, \mathbf{y} -, \mathbf{z} - components of total torque)
- we have $y = Ax$
- A depends on geometry
(of applied forces and torques with respect to center of gravity CG)
- j th column gives resulting force & torque for unit force/torque j

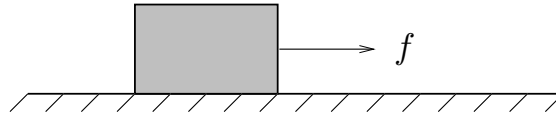
Linear static circuit

interconnection of resistors, linear dependent (controlled) sources, and independent sources



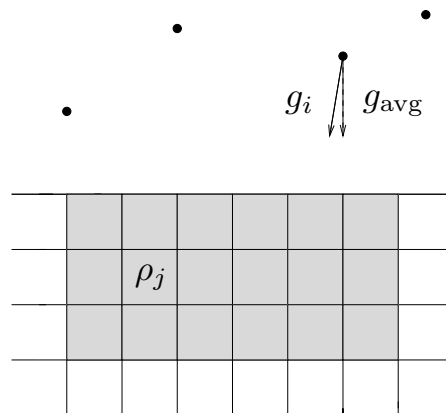
- x_j is value of independent source j
- y_i is some circuit variable (voltage, current)
- we have $y = Ax$
- if x_j are currents and y_i are voltages, A is called the *impedance* or *resistance* matrix

Final position/velocity of mass due to applied forces



- unit mass, zero position/velocity at $t = 0$, subject to force $f(t)$ for $0 \leq t \leq n$
- $f(t) = x_j$ for $j - 1 \leq t < j$, $j = 1, \dots, n$
(x is the sequence of applied forces, constant in each interval)
- y_1, y_2 are final position and velocity (*i.e.*, at $t = n$)
- we have $y = Ax$
- a_{1j} gives influence of applied force during $j - 1 \leq t < j$ on final position
- a_{2j} gives influence of applied force during $j - 1 \leq t < j$ on final velocity

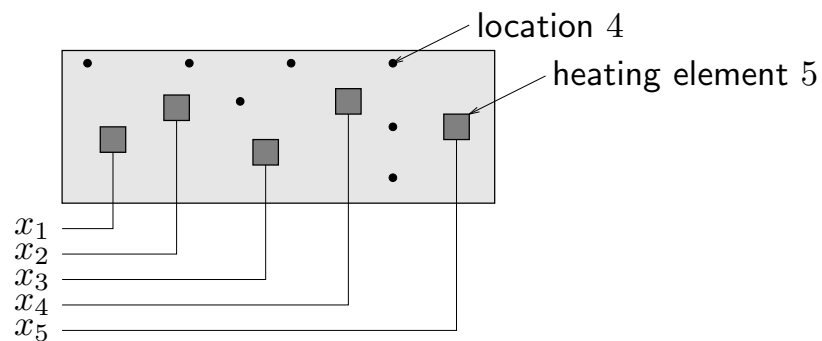
Gravimeter prospecting



- $x_j = \rho_j - \rho_{\text{avg}}$ is (excess) mass density of earth in voxel j ;
- y_i is measured *gravity anomaly* at location i , *i.e.*, some component (typically vertical) of $g_i - g_{\text{avg}}$
- $y = Ax$

- A comes from physics and geometry
- j th column of A shows sensor readings caused by unit density anomaly at voxel j
- i th row of A shows sensitivity pattern of sensor i

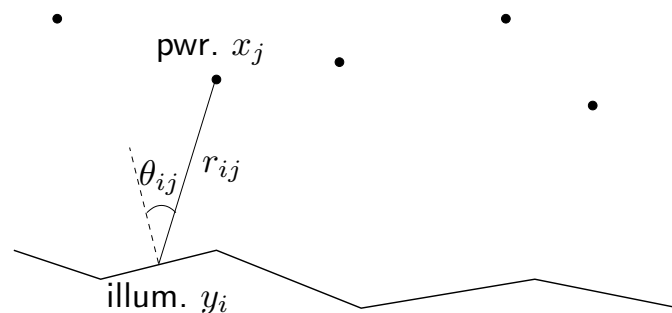
Thermal system



- x_j is power of j th heating element or heat source
- y_i is change in steady-state temperature at location i
- thermal transport via conduction
- $y = Ax$

- a_{ij} gives influence of heater j at location i (in $^{\circ}\text{C}/\text{W}$)
- j th column of A gives pattern of steady-state temperature rise due to 1W at heater j
- i th row shows how heaters affect location i

Illumination with multiple lamps



- n lamps illuminating m (small, flat) patches, no shadows
- x_j is power of j th lamp; y_i is illumination level of patch i
- $y = Ax$, where $a_{ij} = r_{ij}^{-2} \max\{\cos \theta_{ij}, 0\}$
 ($\cos \theta_{ij} < 0$ means patch i is shaded from lamp j)
- j th column of A shows illumination pattern from lamp j

Signal and interference power in wireless system

- n transmitter/receiver pairs
- transmitter j transmits to receiver j (and, inadvertently, to the other receivers)
- p_j is power of j th transmitter
- s_i is received signal power of i th receiver
- z_i is received interference power of i th receiver
- G_{ij} is path gain from transmitter j to receiver i
- we have $s = Ap$, $z = Bp$, where

$$a_{ij} = \begin{cases} G_{ii} & i = j \\ 0 & i \neq j \end{cases} \quad b_{ij} = \begin{cases} 0 & i = j \\ G_{ij} & i \neq j \end{cases}$$

- A is diagonal; B has zero diagonal (ideally, A is 'large', B is 'small')

Cost of production

production *inputs* (materials, parts, labor, . . .) are combined to make a number of *products*

- x_j is price per unit of production input j
- a_{ij} is units of production input j required to manufacture one unit of product i
- y_i is production cost per unit of product i
- we have $y = Ax$
- i th row of A is *bill of materials* for unit of product i

production inputs needed

- q_i is quantity of product i to be produced
- r_j is total quantity of production input j needed
- we have $r = A^T q$

total production cost is

$$r^T x = (A^T q)^T x = q^T A x$$

Network traffic and flows

- n flows with rates f_1, \dots, f_n pass from their source nodes to their destination nodes over fixed routes in a network
- t_i , traffic on link i , is sum of rates of flows passing through it
- flow routes given by *flow-link incidence matrix*

$$A_{ij} = \begin{cases} 1 & \text{flow } j \text{ goes over link } i \\ 0 & \text{otherwise} \end{cases}$$

- traffic and flow rates related by $t = A f$

link delays and flow latency

- let d_1, \dots, d_m be link delays, and l_1, \dots, l_n be latency (total travel time) of flows
- $l = A^T d$
- $f^T l = f^T A^T d = (Af)^T d = t^T d$, total # of packets in network

Linearization

- if $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ is differentiable at $x_0 \in \mathbf{R}^n$, then

$$x \text{ near } x_0 \implies f(x) \text{ very near } f(x_0) + Df(x_0)(x - x_0)$$

where

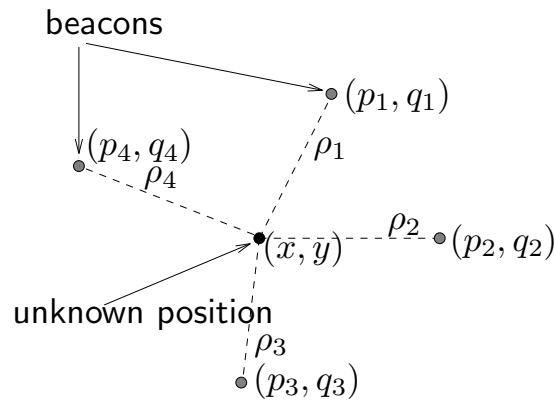
$$Df(x_0)_{ij} = \left. \frac{\partial f_i}{\partial x_j} \right|_{x_0}$$

is derivative (Jacobian) matrix

- with $y = f(x)$, $y_0 = f(x_0)$, define *input deviation* $\delta x := x - x_0$, *output deviation* $\delta y := y - y_0$
- then we have $\delta y \approx Df(x_0)\delta x$
- when deviations are small, they are (approximately) related by a linear function

Navigation by range measurement

- (x, y) unknown coordinates in plane
- (p_i, q_i) known coordinates of beacons for $i = 1, 2, 3, 4$
- ρ_i measured (known) distance or range from beacon i



- $\rho \in \mathbf{R}^4$ is a nonlinear function of $(x, y) \in \mathbf{R}^2$:

$$\rho_i(x, y) = \sqrt{(x - p_i)^2 + (y - q_i)^2}$$

- linearize around (x_0, y_0) : $\delta\rho \approx A \begin{bmatrix} \delta x \\ \delta y \end{bmatrix}$, where

$$a_{i1} = \frac{(x_0 - p_i)}{\sqrt{(x_0 - p_i)^2 + (y_0 - q_i)^2}}, \quad a_{i2} = \frac{(y_0 - q_i)}{\sqrt{(x_0 - p_i)^2 + (y_0 - q_i)^2}}$$

- i th row of A shows (approximate) change in i th range measurement for (small) shift in (x, y) from (x_0, y_0)
- first column of A shows sensitivity of range measurements to (small) change in x from x_0
- obvious application: (x_0, y_0) is last navigation fix; (x, y) is current position, a short time later

Broad categories of applications

linear model or function $y = Ax$

some broad categories of applications:

- estimation or inversion
- control or design
- mapping or transformation

(this list is not exclusive; can have combinations . . .)

Estimation or inversion

$$y = Ax$$

- y_i is i th measurement or sensor reading (which we know)
- x_j is j th parameter to be estimated or determined
- a_{ij} is sensitivity of i th sensor to j th parameter

sample problems:

- find x , given y
- find all x 's that result in y (*i.e.*, all x 's consistent with measurements)
- if there is no x such that $y = Ax$, find x s.t. $y \approx Ax$ (*i.e.*, if the sensor readings are inconsistent, find x which is almost consistent)

Control or design

$$y = Ax$$

- x is vector of design parameters or inputs (which we can choose)
- y is vector of results, or outcomes
- A describes how input choices affect results

sample problems:

- find x so that $y = y_{\text{des}}$
- find all x 's that result in $y = y_{\text{des}}$ (*i.e.*, find all designs that meet specifications)
- among x 's that satisfy $y = y_{\text{des}}$, find a small one (*i.e.*, find a small or efficient x that meets specifications)

Mapping or transformation

- x is mapped or transformed to y by linear function $y = Ax$

sample problems:

- determine if there is an x that maps to a given y
- (if possible) find *an* x that maps to y
- find *all* x 's that map to a given y
- if there is only one x that maps to y , find it (*i.e.*, decode or undo the mapping)

Matrix multiplication as mixture of columns

write $A \in \mathbf{R}^{m \times n}$ in terms of its columns:

$$A = \begin{bmatrix} a_1 & a_2 & \cdots & a_n \end{bmatrix}$$

where $a_j \in \mathbf{R}^m$

then $y = Ax$ can be written as

$$y = x_1 a_1 + x_2 a_2 + \cdots + x_n a_n$$

(x_j 's are scalars, a_j 's are m -vectors)

- y is a (linear) combination or mixture of the columns of A
- coefficients of x give coefficients of mixture

an important example: $x = e_j$, the j th *unit vector*

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots \quad e_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

then $Ae_j = a_j$, the j th column of A

(e_j corresponds to a pure mixture, giving only column j)

Matrix multiplication as inner product with rows

write A in terms of its rows:

$$A = \begin{bmatrix} \tilde{a}_1^T \\ \tilde{a}_2^T \\ \vdots \\ \tilde{a}_n^T \end{bmatrix}$$

where $\tilde{a}_i \in \mathbf{R}^n$

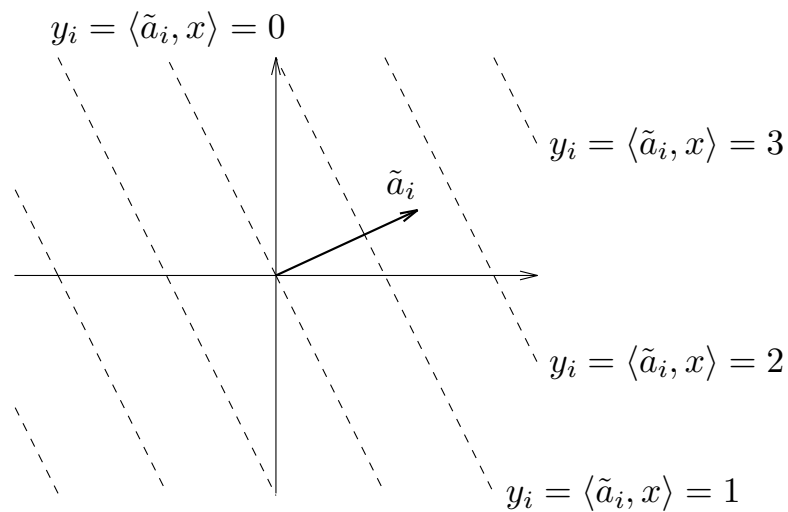
then $y = Ax$ can be written as

$$y = \begin{bmatrix} \tilde{a}_1^T x \\ \tilde{a}_2^T x \\ \vdots \\ \tilde{a}_m^T x \end{bmatrix}$$

thus $y_i = \langle \tilde{a}_i, x \rangle$, i.e., y_i is inner product of i th row of A with x

geometric interpretation:

$y_i = \tilde{a}_i^T x = \alpha$ is a hyperplane in \mathbf{R}^n (normal to \tilde{a}_i)



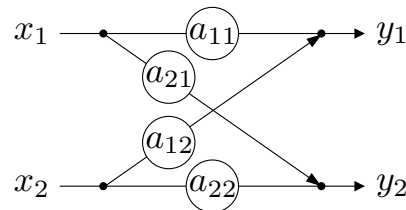
Block diagram representation

$y = Ax$ can be represented by a *signal flow graph* or *block diagram*

e.g. for $m = n = 2$, we represent

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

as



- a_{ij} is the gain along the path from j th input to i th output
- (by not drawing paths with zero gain) shows sparsity structure of A (*e.g.*, diagonal, block upper triangular, arrow . . .)

example: block upper triangular, *i.e.*,

$$A = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix}$$

where $A_{11} \in \mathbf{R}^{m_1 \times n_1}$, $A_{12} \in \mathbf{R}^{m_1 \times n_2}$, $A_{21} \in \mathbf{R}^{m_2 \times n_1}$, $A_{22} \in \mathbf{R}^{m_2 \times n_2}$

partition x and y conformably as

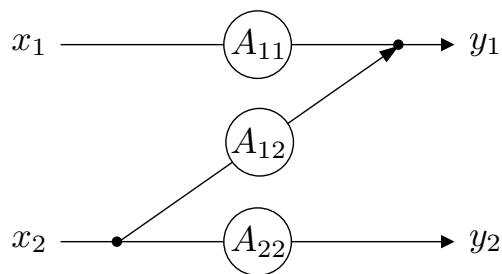
$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

($x_1 \in \mathbf{R}^{n_1}$, $x_2 \in \mathbf{R}^{n_2}$, $y_1 \in \mathbf{R}^{m_1}$, $y_2 \in \mathbf{R}^{m_2}$) so

$$y_1 = A_{11}x_1 + A_{12}x_2, \quad y_2 = A_{22}x_2,$$

i.e., y_2 doesn't depend on x_1

block diagram:



. . . no path from x_1 to y_2 , so y_2 doesn't depend on x_1

Matrix multiplication as composition

for $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{n \times p}$, $C = AB \in \mathbf{R}^{m \times p}$ where

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$

composition interpretation: $y = Cz$ represents composition of $y = Ax$ and $x = Bz$

$$z \xrightarrow[p]{\quad} \boxed{B} \xrightarrow[n]{x} \boxed{A} \xrightarrow[m]{\quad} y \quad \equiv \quad z \xrightarrow[p]{\quad} \boxed{AB} \xrightarrow[m]{\quad} y$$

(note that B is on left in block diagram)

Column and row interpretations

can write product $C = AB$ as

$$C = [c_1 \cdots c_p] = AB = [Ab_1 \cdots Ab_p]$$

i.e., i th column of C is A acting on i th column of B

similarly we can write

$$C = \begin{bmatrix} \tilde{c}_1^T \\ \vdots \\ \tilde{c}_m^T \end{bmatrix} = AB = \begin{bmatrix} \tilde{a}_1^T B \\ \vdots \\ \tilde{a}_m^T B \end{bmatrix}$$

i.e., i th row of C is i th row of A acting (on left) on B

Inner product interpretation

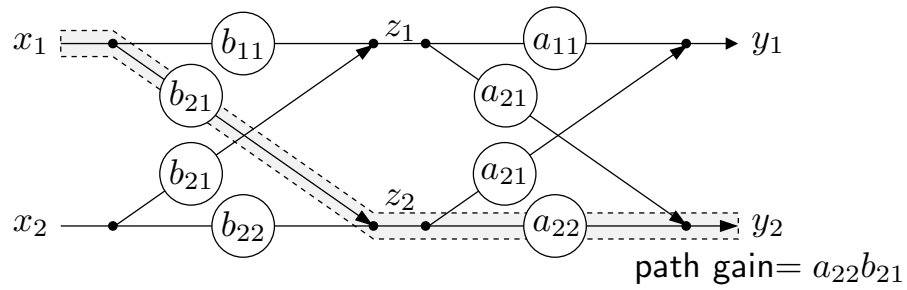
inner product interpretation:

$$c_{ij} = \tilde{a}_i^T b_j = \langle \tilde{a}_i, b_j \rangle$$

i.e., entries of C are inner products of rows of A and columns of B

- $c_{ij} = 0$ means i th row of A is orthogonal to j th column of B
- **Gram matrix** of vectors f_1, \dots, f_n defined as $G_{ij} = f_i^T f_j$
(gives inner product of each vector with the others)
- $G = [f_1 \cdots f_n]^T [f_1 \cdots f_n]$

Matrix multiplication interpretation via paths



- $a_{ik}b_{kj}$ is gain of path from input j to output i via k
- c_{ij} is sum of gains over *all* paths from input j to output i

Lecture 3

Linear algebra review

- vector space, subspaces
- independence, basis, dimension
- range, nullspace, rank
- change of coordinates
- norm, angle, inner product

3-1

Vector spaces

a *vector space* or *linear space* (over the reals) consists of

- a set \mathcal{V}
- a vector sum $+: \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{V}$
- a scalar multiplication $: \mathbf{R} \times \mathcal{V} \rightarrow \mathcal{V}$
- a distinguished element $0 \in \mathcal{V}$

which satisfy a list of properties

- $x + y = y + x, \quad \forall x, y \in \mathcal{V} \quad (+ \text{ is commutative})$
- $(x + y) + z = x + (y + z), \quad \forall x, y, z \in \mathcal{V} \quad (+ \text{ is associative})$
- $0 + x = x, \forall x \in \mathcal{V} \quad (0 \text{ is additive identity})$
- $\forall x \in \mathcal{V} \quad \exists(-x) \in \mathcal{V} \text{ s.t. } x + (-x) = 0 \quad (\text{existence of additive inverse})$
- $(\alpha\beta)x = \alpha(\beta x), \quad \forall \alpha, \beta \in \mathbf{R} \quad \forall x \in \mathcal{V} \quad (\text{scalar mult. is associative})$
- $\alpha(x + y) = \alpha x + \alpha y, \quad \forall \alpha \in \mathbf{R} \quad \forall x, y \in \mathcal{V} \quad (\text{right distributive rule})$
- $(\alpha + \beta)x = \alpha x + \beta x, \quad \forall \alpha, \beta \in \mathbf{R} \quad \forall x \in \mathcal{V} \quad (\text{left distributive rule})$
- $1x = x, \quad \forall x \in \mathcal{V}$

Examples

- $\mathcal{V}_1 = \mathbf{R}^n$, with standard (componentwise) vector addition and scalar multiplication
- $\mathcal{V}_2 = \{0\}$ (where $0 \in \mathbf{R}^n$)
- $\mathcal{V}_3 = \text{span}(v_1, v_2, \dots, v_k)$ where

$$\text{span}(v_1, v_2, \dots, v_k) = \{\alpha_1 v_1 + \dots + \alpha_k v_k \mid \alpha_i \in \mathbf{R}\}$$

and $v_1, \dots, v_k \in \mathbf{R}^n$

Subspaces

- a *subspace* of a vector space is a *subset* of a vector space which is itself a vector space
- roughly speaking, a subspace is closed under vector addition and scalar multiplication
- examples $\mathcal{V}_1, \mathcal{V}_2, \mathcal{V}_3$ above are subspaces of \mathbf{R}^n

Vector spaces of functions

- $\mathcal{V}_4 = \{x : \mathbf{R}_+ \rightarrow \mathbf{R}^n \mid x \text{ is differentiable}\}$, where vector sum is sum of functions:

$$(x + z)(t) = x(t) + z(t)$$

and scalar multiplication is defined by

$$(\alpha x)(t) = \alpha x(t)$$

(a *point* in \mathcal{V}_4 is a *trajectory* in \mathbf{R}^n)

- $\mathcal{V}_5 = \{x \in \mathcal{V}_4 \mid \dot{x} = Ax\}$
(*points* in \mathcal{V}_5 are *trajectories* of the linear system $\dot{x} = Ax$)
- \mathcal{V}_5 is a subspace of \mathcal{V}_4

Independent set of vectors

a set of vectors $\{v_1, v_2, \dots, v_k\}$ is *independent* if

$$\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_k v_k = 0 \implies \alpha_1 = \alpha_2 = \dots = 0$$

some equivalent conditions:

- coefficients of $\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_k v_k$ are uniquely determined, *i.e.*,

$$\alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_k v_k = \beta_1 v_1 + \beta_2 v_2 + \dots + \beta_k v_k$$

implies $\alpha_1 = \beta_1, \alpha_2 = \beta_2, \dots, \alpha_k = \beta_k$

- no vector v_i can be expressed as a linear combination of the other vectors $v_1, \dots, v_{i-1}, v_{i+1}, \dots, v_k$

Basis and dimension

set of vectors $\{v_1, v_2, \dots, v_k\}$ is a *basis* for a vector space \mathcal{V} if

- v_1, v_2, \dots, v_k **span** \mathcal{V} , *i.e.*, $\mathcal{V} = \text{span}(v_1, v_2, \dots, v_k)$
- $\{v_1, v_2, \dots, v_k\}$ is independent

equivalent: every $v \in \mathcal{V}$ can be uniquely expressed as

$$v = \alpha_1 v_1 + \dots + \alpha_k v_k$$

fact: for a given vector space \mathcal{V} , the number of vectors in any basis is the same

number of vectors in any basis is called the *dimension* of \mathcal{V} , denoted **dim** \mathcal{V}

(we assign **dim** $\{0\} = 0$, and **dim** $\mathcal{V} = \infty$ if there is no basis)

Nullspace of a matrix

the *nullspace* of $A \in \mathbf{R}^{m \times n}$ is defined as

$$\mathcal{N}(A) = \{ x \in \mathbf{R}^n \mid Ax = 0 \}$$

- $\mathcal{N}(A)$ is set of vectors mapped to zero by $y = Ax$
- $\mathcal{N}(A)$ is set of vectors orthogonal to all rows of A

$\mathcal{N}(A)$ gives *ambiguity* in x given $y = Ax$:

- if $y = Ax$ and $z \in \mathcal{N}(A)$, then $y = A(x + z)$
- conversely, if $y = Ax$ and $y = A\tilde{x}$, then $\tilde{x} = x + z$ for some $z \in \mathcal{N}(A)$

Zero nullspace

A is called *one-to-one* if 0 is the only element of its nullspace:

$$\mathcal{N}(A) = \{0\} \iff$$

- x can always be uniquely determined from $y = Ax$
(*i.e.*, the linear transformation $y = Ax$ doesn't 'lose' information)
- mapping from x to Ax is one-to-one: different x 's map to different y 's
- columns of A are independent (hence, a basis for their span)
- A has a *left inverse*, *i.e.*, there is a matrix $B \in \mathbf{R}^{n \times m}$ s.t. $BA = I$
- $\det(A^T A) \neq 0$

(we'll establish these later)

Interpretations of nullspace

suppose $z \in \mathcal{N}(A)$

$y = Ax$ represents **measurement** of x

- z is undetectable from sensors — get zero sensor readings
- x and $x + z$ are indistinguishable from sensors: $Ax = A(x + z)$

$\mathcal{N}(A)$ characterizes *ambiguity* in x from measurement $y = Ax$

$y = Ax$ represents **output** resulting from input x

- z is an input with no result
- x and $x + z$ have same result

$\mathcal{N}(A)$ characterizes *freedom of input choice* for given result

Range of a matrix

the *range* of $A \in \mathbf{R}^{m \times n}$ is defined as

$$\mathcal{R}(A) = \{Ax \mid x \in \mathbf{R}^n\} \subseteq \mathbf{R}^m$$

$\mathcal{R}(A)$ can be interpreted as

- the set of vectors that can be ‘hit’ by linear mapping $y = Ax$
- the span of columns of A
- the set of vectors y for which $Ax = y$ has a solution

Onto matrices

A is called *onto* if $\mathcal{R}(A) = \mathbf{R}^m \iff$

- $Ax = y$ can be solved in x for any y
- columns of A span \mathbf{R}^m
- A has a *right inverse*, i.e., there is a matrix $B \in \mathbf{R}^{n \times m}$ s.t. $AB = I$
- rows of A are independent
- $\mathcal{N}(A^T) = \{0\}$
- $\det(AA^T) \neq 0$

(some of these are not obvious; we'll establish them later)

Interpretations of range

suppose $v \in \mathcal{R}(A)$, $w \notin \mathcal{R}(A)$

$y = Ax$ represents **measurement** of x

- $y = v$ is a *possible* or *consistent* sensor signal
- $y = w$ is *impossible* or *inconsistent*; sensors have failed or model is wrong

$y = Ax$ represents **output** resulting from input x

- v is a possible result or output
- w cannot be a result or output

$\mathcal{R}(A)$ characterizes the *possible results* or *achievable outputs*

Inverse

$A \in \mathbf{R}^{n \times n}$ is *invertible* or *nonsingular* if $\det A \neq 0$

equivalent conditions:

- columns of A are a basis for \mathbf{R}^n
- rows of A are a basis for \mathbf{R}^n
- $y = Ax$ has a unique solution x for every $y \in \mathbf{R}^n$
- A has a (left and right) inverse denoted $A^{-1} \in \mathbf{R}^{n \times n}$, with $AA^{-1} = A^{-1}A = I$
- $\mathcal{N}(A) = \{0\}$
- $\mathcal{R}(A) = \mathbf{R}^n$
- $\det A^T A = \det AA^T \neq 0$

Interpretations of inverse

suppose $A \in \mathbf{R}^{n \times n}$ has inverse $B = A^{-1}$

- mapping associated with B undoes mapping associated with A (applied either before or after!)
- $x = By$ is a perfect (pre- or post-) *equalizer* for the *channel* $y = Ax$
- $x = By$ is unique solution of $Ax = y$

Dual basis interpretation

- let a_i be columns of A , and \tilde{b}_i^T be rows of $B = A^{-1}$
- from $y = x_1 a_1 + \cdots + x_n a_n$ and $x_i = \tilde{b}_i^T y$, we get

$$y = \sum_{i=1}^n (\tilde{b}_i^T y) a_i$$

thus, inner product with *rows of inverse matrix* gives the coefficients in the *expansion of a vector in the columns of the matrix*

- $\tilde{b}_1, \dots, \tilde{b}_n$ and a_1, \dots, a_n are called *dual bases*

Rank of a matrix

we define the *rank* of $A \in \mathbf{R}^{m \times n}$ as

$$\mathbf{rank}(A) = \mathbf{dim} \mathcal{R}(A)$$

(nontrivial) **facts:**

- $\mathbf{rank}(A) = \mathbf{rank}(A^T)$
- $\mathbf{rank}(A)$ is maximum number of independent columns (or rows) of A
hence $\mathbf{rank}(A) \leq \mathbf{min}(m, n)$
- $\mathbf{rank}(A) + \mathbf{dim} \mathcal{N}(A) = n$

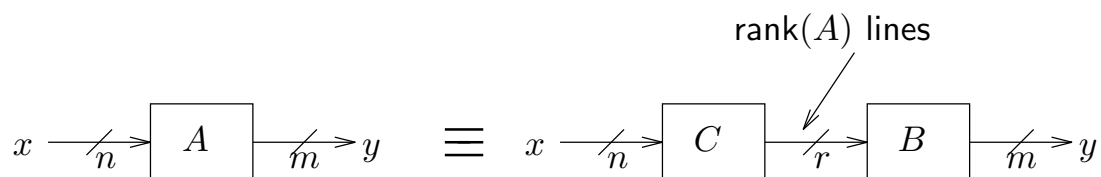
Conservation of dimension

interpretation of $\text{rank}(A) + \dim \mathcal{N}(A) = n$:

- $\text{rank}(A)$ is dimension of set 'hit' by the mapping $y = Ax$
- $\dim \mathcal{N}(A)$ is dimension of set of x 'crushed' to zero by $y = Ax$
- 'conservation of dimension': each dimension of input is either crushed to zero or ends up in output
- roughly speaking:
 - n is number of degrees of freedom in input x
 - $\dim \mathcal{N}(A)$ is number of degrees of freedom lost in the mapping from x to $y = Ax$
 - $\text{rank}(A)$ is number of degrees of freedom in output y

'Coding' interpretation of rank

- rank of product: $\text{rank}(BC) \leq \min\{\text{rank}(B), \text{rank}(C)\}$
- hence if $A = BC$ with $B \in \mathbf{R}^{m \times r}$, $C \in \mathbf{R}^{r \times n}$, then $\text{rank}(A) \leq r$
- conversely: if $\text{rank}(A) = r$ then $A \in \mathbf{R}^{m \times n}$ can be factored as $A = BC$ with $B \in \mathbf{R}^{m \times r}$, $C \in \mathbf{R}^{r \times n}$:



- $\text{rank}(A) = r$ is minimum size of vector needed to faithfully reconstruct y from x

Application: fast matrix-vector multiplication

- need to compute matrix-vector product $y = Ax$, $A \in \mathbf{R}^{m \times n}$
- A has known factorization $A = BC$, $B \in \mathbf{R}^{m \times r}$
- computing $y = Ax$ directly: mn operations
- computing $y = Ax$ as $y = B(Cx)$ (compute $z = Cx$ first, then $y = Bz$): $rn + mr = (m + n)r$ operations
- savings can be considerable if $r \ll \min\{m, n\}$

Full rank matrices

for $A \in \mathbf{R}^{m \times n}$ we always have $\text{rank}(A) \leq \min(m, n)$

we say A is *full rank* if $\text{rank}(A) = \min(m, n)$

- for **square** matrices, full rank means nonsingular
- for **skinny** matrices ($m \geq n$), full rank means columns are independent
- for **fat** matrices ($m \leq n$), full rank means rows are independent

Change of coordinates

'standard' basis vectors in \mathbf{R}^n : (e_1, e_2, \dots, e_n) where

$$e_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

(1 in i th component)

obviously we have

$$x = x_1 e_1 + x_2 e_2 + \dots + x_n e_n$$

x_i are called the coordinates of x (in the standard basis)

if (t_1, t_2, \dots, t_n) is another basis for \mathbf{R}^n , we have

$$x = \tilde{x}_1 t_1 + \tilde{x}_2 t_2 + \dots + \tilde{x}_n t_n$$

where \tilde{x}_i are the coordinates of x in the basis (t_1, t_2, \dots, t_n)

define $T = \begin{bmatrix} t_1 & t_2 & \dots & t_n \end{bmatrix}$ so $x = T\tilde{x}$, hence

$$\tilde{x} = T^{-1}x$$

(T is invertible since t_i are a basis)

T^{-1} transforms (standard basis) coordinates of x into t_i -coordinates

inner product i th row of T^{-1} with x extracts t_i -coordinate of x

consider linear transformation $y = Ax$, $A \in \mathbf{R}^{n \times n}$

express y and x in terms of t_1, t_2, \dots, t_n :

$$x = T\tilde{x}, \quad y = T\tilde{y}$$

so

$$\tilde{y} = (T^{-1}AT)\tilde{x}$$

- $A \longrightarrow T^{-1}AT$ is called *similarity transformation*
- similarity transformation by T expresses linear transformation $y = Ax$ in coordinates t_1, t_2, \dots, t_n

(Euclidean) norm

for $x \in \mathbf{R}^n$ we define the (Euclidean) norm as

$$\|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} = \sqrt{x^T x}$$

$\|x\|$ measures length of vector (from origin)

important properties:

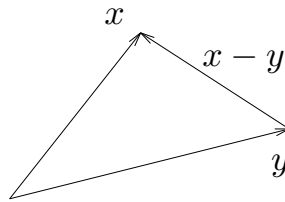
- $\|\alpha x\| = |\alpha| \|x\|$ (homogeneity)
- $\|x + y\| \leq \|x\| + \|y\|$ (triangle inequality)
- $\|x\| \geq 0$ (nonnegativity)
- $\|x\| = 0 \iff x = 0$ (definiteness)

RMS value and (Euclidean) distance

root-mean-square (RMS) value of vector $x \in \mathbf{R}^n$:

$$\mathbf{rms}(x) = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right)^{1/2} = \frac{\|x\|}{\sqrt{n}}$$

norm defines distance between vectors: **dist** $(x, y) = \|x - y\|$



Inner product

$$\langle x, y \rangle := x_1 y_1 + x_2 y_2 + \cdots + x_n y_n = x^T y$$

important properties:

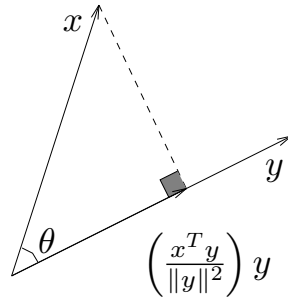
- $\langle \alpha x, y \rangle = \alpha \langle x, y \rangle$
- $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$
- $\langle x, y \rangle = \langle y, x \rangle$
- $\langle x, x \rangle \geq 0$
- $\langle x, x \rangle = 0 \iff x = 0$

$f(y) = \langle x, y \rangle$ is linear function : $\mathbf{R}^n \rightarrow \mathbf{R}$, with linear map defined by row vector x^T

Cauchy-Schwartz inequality and angle between vectors

- for any $x, y \in \mathbf{R}^n$, $|x^T y| \leq \|x\| \|y\|$
- (unsigned) angle between vectors in \mathbf{R}^n defined as

$$\theta = \angle(x, y) = \cos^{-1} \frac{x^T y}{\|x\| \|y\|}$$



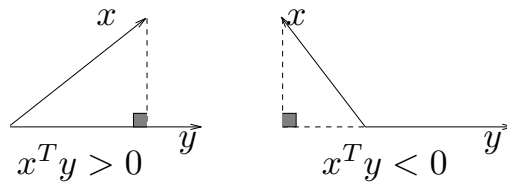
thus $x^T y = \|x\| \|y\| \cos \theta$

special cases:

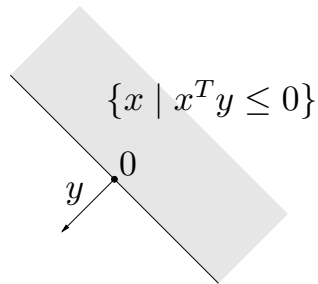
- x and y are *aligned*: $\theta = 0$; $x^T y = \|x\| \|y\|$;
(if $x \neq 0$) $y = \alpha x$ for some $\alpha \geq 0$
- x and y are *opposed*: $\theta = \pi$; $x^T y = -\|x\| \|y\|$;
(if $x \neq 0$) $y = -\alpha x$ for some $\alpha \geq 0$
- x and y are *orthogonal*: $\theta = \pi/2$ or $-\pi/2$; $x^T y = 0$
denoted $x \perp y$

interpretation of $x^T y > 0$ and $x^T y < 0$:

- $x^T y > 0$ means $\angle(x, y)$ is acute
- $x^T y < 0$ means $\angle(x, y)$ is obtuse



$\{x \mid x^T y \leq 0\}$ defines a *halfspace* with outward normal vector y , and boundary passing through 0



Lecture 4

Orthonormal sets of vectors and QR factorization

- orthonormal sets of vectors
- Gram-Schmidt procedure, QR factorization
- orthogonal decomposition induced by a matrix

4-1

Orthonormal set of vectors

set of vectors $u_1, \dots, u_k \in \mathbf{R}^n$ is

- *normalized* if $\|u_i\| = 1$, $i = 1, \dots, k$
(u_i are called *unit vectors* or *direction vectors*)
- *orthogonal* if $u_i \perp u_j$ for $i \neq j$
- *orthonormal* if both

slang: we say ' u_1, \dots, u_k are orthonormal vectors' but orthonormality (like independence) is a property of a set of vectors, not vectors individually

in terms of $U = [u_1 \cdots u_k]$, orthonormal means

$$U^T U = I_k$$

- orthonormal vectors are independent
(multiply $\alpha_1 u_1 + \alpha_2 u_2 + \cdots + \alpha_k u_k = 0$ by u_i^T)
- hence u_1, \dots, u_k is an orthonormal basis for

$$\text{span}(u_1, \dots, u_k) = \mathcal{R}(U)$$

- **warning:** if $k < n$ then $UU^T \neq I$ (since its rank is at most k)
(more on this matrix later . . .)

Geometric properties

suppose columns of $U = [u_1 \cdots u_k]$ are orthonormal

if $w = Uz$, then $\|w\| = \|z\|$

- multiplication by U does not change norm
- mapping $w = Uz$ is *isometric*: it preserves distances
- simple derivation using matrices:

$$\|w\|^2 = \|Uz\|^2 = (Uz)^T(Uz) = z^T U^T U z = z^T z = \|z\|^2$$

- *inner products* are also preserved: $\langle Uz, U\tilde{z} \rangle = \langle z, \tilde{z} \rangle$
- if $w = Uz$ and $\tilde{w} = U\tilde{z}$ then

$$\langle w, \tilde{w} \rangle = \langle Uz, U\tilde{z} \rangle = (Uz)^T (U\tilde{z}) = z^T U^T U \tilde{z} = \langle z, \tilde{z} \rangle$$

- norms and inner products preserved, so *angles* are preserved:
 $\angle(Uz, U\tilde{z}) = \angle(z, \tilde{z})$
- thus, multiplication by U preserves inner products, angles, and distances

Orthonormal basis for \mathbf{R}^n

- suppose u_1, \dots, u_n is an orthonormal *basis* for \mathbf{R}^n
- then $U = [u_1 \cdots u_n]$ is called **orthogonal**: it is square and satisfies
 $U^T U = I$
 (you'd think such matrices would be called *orthonormal*, not *orthogonal*)
- it follows that $U^{-1} = U^T$, and hence also $U U^T = I$, *i.e.*,

$$\sum_{i=1}^n u_i u_i^T = I$$

Expansion in orthonormal basis

suppose U is orthogonal, so $x = UU^T x$, i.e.,

$$x = \sum_{i=1}^n (u_i^T x) u_i$$

- $u_i^T x$ is called the *component* of x in the direction u_i
- $a = U^T x$ *resolves* x into the vector of its u_i components
- $x = Ua$ *reconstitutes* x from its u_i components
- $x = Ua = \sum_{i=1}^n a_i u_i$ is called the (u_i-) *expansion* of x

the identity $I = UU^T = \sum_{i=1}^n u_i u_i^T$ is sometimes written (in physics) as

$$I = \sum_{i=1}^n |u_i\rangle \langle u_i|$$

since

$$x = \sum_{i=1}^n |u_i\rangle \langle u_i|x\rangle$$

(but we won't use this notation)

Geometric interpretation

if U is orthogonal, then transformation $w = Uz$

- preserves *norm* of vectors, *i.e.*, $\|Uz\| = \|z\|$
- preserves *angles* between vectors, *i.e.*, $\angle(Uz, U\tilde{z}) = \angle(z, \tilde{z})$

examples:

- rotations (about some axis)
- reflections (through some plane)

Example: rotation by θ in \mathbf{R}^2 is given by

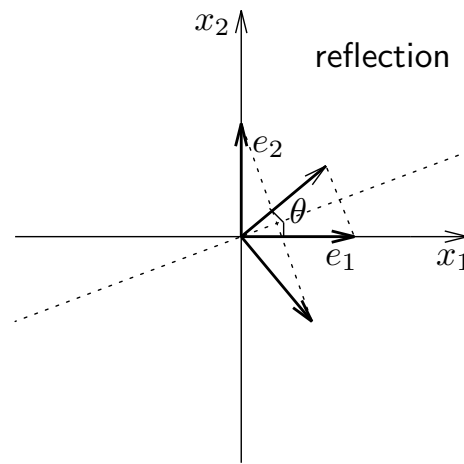
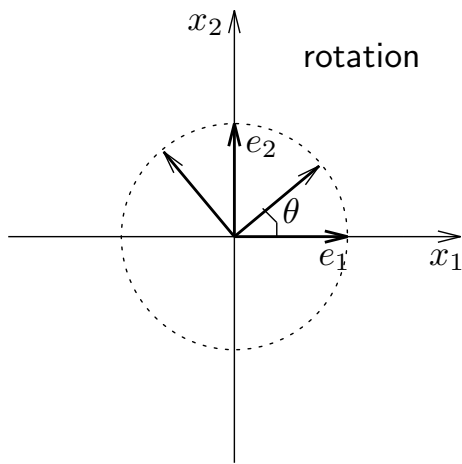
$$y = U_\theta x, \quad U_\theta = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

since $e_1 \rightarrow (\cos \theta, \sin \theta)$, $e_2 \rightarrow (-\sin \theta, \cos \theta)$

reflection across line $x_2 = x_1 \tan(\theta/2)$ is given by

$$y = R_\theta x, \quad R_\theta = \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix}$$

since $e_1 \rightarrow (\cos \theta, \sin \theta)$, $e_2 \rightarrow (\sin \theta, -\cos \theta)$



can check that U_θ and R_θ are orthogonal

Gram-Schmidt procedure

- given independent vectors $a_1, \dots, a_k \in \mathbf{R}^n$, G-S procedure finds orthonormal vectors q_1, \dots, q_k s.t.

$$\text{span}(a_1, \dots, a_r) = \text{span}(q_1, \dots, q_r) \quad \text{for } r \leq k$$

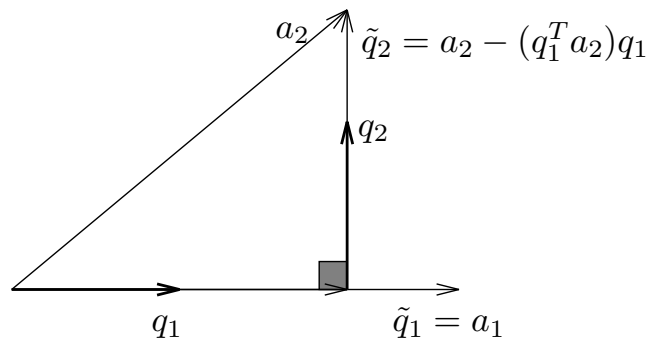
- thus, q_1, \dots, q_r is an orthonormal basis for $\text{span}(a_1, \dots, a_r)$
- rough idea of method: first *orthogonalize* each vector w.r.t. previous ones; then *normalize* result to have norm one

Gram-Schmidt procedure

- step 1a. $\tilde{q}_1 := a_1$
- step 1b. $q_1 := \tilde{q}_1 / \|\tilde{q}_1\|$ (normalize)
- step 2a. $\tilde{q}_2 := a_2 - (q_1^T a_2)q_1$ (remove q_1 component from a_2)
- step 2b. $q_2 := \tilde{q}_2 / \|\tilde{q}_2\|$ (normalize)
- step 3a. $\tilde{q}_3 := a_3 - (q_1^T a_3)q_1 - (q_2^T a_3)q_2$ (remove q_1, q_2 components)
- step 3b. $q_3 := \tilde{q}_3 / \|\tilde{q}_3\|$ (normalize)
- etc.

Orthonormal sets of vectors and QR factorization

4-13



for $i = 1, 2, \dots, k$ we have

$$\begin{aligned} a_i &= (q_1^T a_i)q_1 + (q_2^T a_i)q_2 + \cdots + (q_{i-1}^T a_i)q_{i-1} + \|\tilde{q}_i\|q_i \\ &= r_{1i}q_1 + r_{2i}q_2 + \cdots + r_{ii}q_i \end{aligned}$$

(note that the r_{ij} 's come right out of the G-S procedure, and $r_{ii} \neq 0$)

QR decomposition

written in matrix form: $A = QR$, where $A \in \mathbf{R}^{n \times k}$, $Q \in \mathbf{R}^{n \times k}$, $R \in \mathbf{R}^{k \times k}$:

$$\underbrace{\begin{bmatrix} a_1 & a_2 & \cdots & a_k \end{bmatrix}}_A = \underbrace{\begin{bmatrix} q_1 & q_2 & \cdots & q_k \end{bmatrix}}_Q \underbrace{\begin{bmatrix} r_{11} & r_{12} & \cdots & r_{1k} \\ 0 & r_{22} & \cdots & r_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & r_{kk} \end{bmatrix}}_R$$

- $Q^T Q = I_k$, and R is upper triangular & invertible
- called **QR decomposition** (or factorization) of A
- usually computed using a variation on Gram-Schmidt procedure which is less sensitive to numerical (rounding) errors
- columns of Q are orthonormal basis for $\mathcal{R}(A)$

General Gram-Schmidt procedure

- in basic G-S we assume $a_1, \dots, a_k \in \mathbf{R}^n$ are independent
- if a_1, \dots, a_k are dependent, we find $\tilde{q}_j = 0$ for some j , which means a_j is linearly dependent on a_1, \dots, a_{j-1}
- modified algorithm: when we encounter $\tilde{q}_j = 0$, skip to next vector a_{j+1} and continue:

```

r = 0;
for i = 1, ..., k
{
     $\tilde{a} = a_i - \sum_{j=1}^r q_j q_j^T a_i$ ;
    if  $\tilde{a} \neq 0$  {  $r = r + 1$ ;  $q_r = \tilde{a} / \|\tilde{a}\|$ ; }
}
    
```

on exit,

- q_1, \dots, q_r is an orthonormal basis for $\mathcal{R}(A)$ (hence $r = \mathbf{Rank}(A)$)
- each a_i is linear combination of previously generated q_j 's

in matrix notation we have $A = QR$ with $Q^T Q = I_r$ and $R \in \mathbf{R}^{r \times k}$ in *upper staircase form*:

$$\begin{bmatrix} \times & & & & & & \\ & \times & & & & & \\ & & \times & & & & \\ & & & \times & & & \\ & & & & \times & & \\ & & & & & \times & \\ & & & & & & \times \\ & & & & & & & \times \\ & & & & & & & & \times \\ & & & & & & & & & \times \end{bmatrix}$$

possibly nonzero entries

zero entries

'corner' entries (shown as \times) are nonzero

can permute columns with \times to front of matrix:

$$A = Q[\tilde{R} \ S]P$$

where:

- $Q^T Q = I_r$
- $\tilde{R} \in \mathbf{R}^{r \times r}$ is upper triangular and invertible
- $P \in \mathbf{R}^{k \times k}$ is a permutation matrix
(which moves forward the columns of a which generated a new q)

Applications

- directly yields orthonormal basis for $\mathcal{R}(A)$
- yields factorization $A = BC$ with $B \in \mathbf{R}^{n \times r}$, $C \in \mathbf{R}^{r \times k}$, $r = \mathbf{Rank}(A)$
- to check if $b \in \text{span}(a_1, \dots, a_k)$: apply Gram-Schmidt to $[a_1 \cdots a_k \ b]$
- staircase pattern in R shows which columns of A are dependent on previous ones

works incrementally: one G-S procedure yields QR factorizations of $[a_1 \cdots a_p]$ for $p = 1, \dots, k$:

$$[a_1 \cdots a_p] = [q_1 \cdots q_s] R_p$$

where $s = \mathbf{Rank}([a_1 \cdots a_p])$ and R_p is leading $s \times p$ submatrix of R

‘Full’ QR factorization

with $A = Q_1 R_1$ the QR factorization as above, write

$$A = [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$$

where $[Q_1 \ Q_2]$ is orthogonal, *i.e.*, columns of $Q_2 \in \mathbf{R}^{n \times (n-r)}$ are orthonormal, orthogonal to Q_1

to find Q_2 :

- find any matrix \tilde{A} s.t. $[A \ \tilde{A}]$ is full rank (*e.g.*, $\tilde{A} = I$)
- apply general Gram-Schmidt to $[A \ \tilde{A}]$
- Q_1 are orthonormal vectors obtained from columns of A
- Q_2 are orthonormal vectors obtained from extra columns (\tilde{A})

i.e., any set of orthonormal vectors can be *extended* to an orthonormal basis for \mathbf{R}^n

$\mathcal{R}(Q_1)$ and $\mathcal{R}(Q_2)$ are called *complementary subspaces* since

- they are orthogonal (*i.e.*, every vector in the first subspace is orthogonal to every vector in the second subspace)
- their sum is \mathbf{R}^n (*i.e.*, every vector in \mathbf{R}^n can be expressed as a sum of two vectors, one from each subspace)

this is written

- $\mathcal{R}(Q_1) \overset{\perp}{+} \mathcal{R}(Q_2) = \mathbf{R}^n$
- $\mathcal{R}(Q_2) = \mathcal{R}(Q_1)^\perp$ (and $\mathcal{R}(Q_1) = \mathcal{R}(Q_2)^\perp$)
(each subspace is the *orthogonal complement* of the other)

we know $\mathcal{R}(Q_1) = \mathcal{R}(A)$; but what is its orthogonal complement $\mathcal{R}(Q_2)$?

Orthogonal decomposition induced by A

from $A^T = \begin{bmatrix} R_1^T & 0 \end{bmatrix} \begin{bmatrix} Q_1^T \\ Q_2^T \end{bmatrix}$ we see that

$$A^T z = 0 \iff Q_1^T z = 0 \iff z \in \mathcal{R}(Q_2)$$

so $\mathcal{R}(Q_2) = \mathcal{N}(A^T)$

(in fact the columns of Q_2 are an orthonormal basis for $\mathcal{N}(A^T)$)

we conclude: $\mathcal{R}(A)$ and $\mathcal{N}(A^T)$ are *complementary subspaces*:

- $\mathcal{R}(A) \overset{\perp}{+} \mathcal{N}(A^T) = \mathbf{R}^n$ (recall $A \in \mathbf{R}^{n \times k}$)
- $\mathcal{R}(A)^\perp = \mathcal{N}(A^T)$ (and $\mathcal{N}(A^T)^\perp = \mathcal{R}(A)$)
- called *orthogonal decomposition (of \mathbf{R}^n) induced by $A \in \mathbf{R}^{n \times k}$*

- every $y \in \mathbf{R}^n$ can be written uniquely as $y = z + w$, with $z \in \mathcal{R}(A)$, $w \in \mathcal{N}(A^T)$ (we'll soon see what the vector z is . . .)
- can now prove most of the assertions from the linear algebra review lecture
- switching $A \in \mathbf{R}^{n \times k}$ to $A^T \in \mathbf{R}^{k \times n}$ gives decomposition of \mathbf{R}^k :

$$\mathcal{N}(A) \overset{\perp}{+} \mathcal{R}(A^T) = \mathbf{R}^k$$

Lecture 5

Least-squares

- least-squares (approximate) solution of overdetermined equations
- projection and orthogonality principle
- least-squares estimation
- BLUE property

5-1

Overdetermined linear equations

consider $y = Ax$ where $A \in \mathbf{R}^{m \times n}$ is (strictly) skinny, *i.e.*, $m > n$

- called *overdetermined* set of linear equations (more equations than unknowns)
- for most y , cannot solve for x

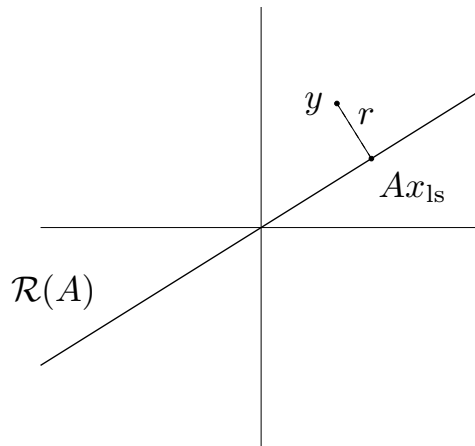
one approach to *approximately* solve $y = Ax$:

- define *residual* or error $r = Ax - y$
- find $x = x_{\text{ls}}$ that minimizes $\|r\|$

x_{ls} called *least-squares* (approximate) solution of $y = Ax$

Geometric interpretation

Ax_{ls} is point in $\mathcal{R}(A)$ closest to y (Ax_{ls} is *projection* of y onto $\mathcal{R}(A)$)



Least-squares (approximate) solution

- assume A is full rank, skinny
- to find x_{ls} , we'll minimize norm of residual squared,

$$\|r\|^2 = x^T A^T A x - 2y^T A x + y^T y$$

- set gradient w.r.t. x to zero:

$$\nabla_x \|r\|^2 = 2A^T A x - 2A^T y = 0$$

- yields the *normal equations*: $A^T A x = A^T y$
- assumptions imply $A^T A$ invertible, so we have

$$x_{\text{ls}} = (A^T A)^{-1} A^T y$$

. . . a very famous formula

- x_{ls} is linear function of y
- $x_{\text{ls}} = A^{-1}y$ if A is square
- x_{ls} solves $y = Ax_{\text{ls}}$ if $y \in \mathcal{R}(A)$
- $A^\dagger = (A^T A)^{-1} A^T$ is called the *pseudo-inverse* of A
- A^\dagger is a *left inverse* of (full rank, skinny) A :

$$A^\dagger A = (A^T A)^{-1} A^T A = I$$

Projection on $\mathcal{R}(A)$

Ax_{ls} is (by definition) the point in $\mathcal{R}(A)$ that is closest to y , *i.e.*, it is the *projection* of y onto $\mathcal{R}(A)$

$$Ax_{\text{ls}} = \mathcal{P}_{\mathcal{R}(A)}(y)$$

- the projection function $\mathcal{P}_{\mathcal{R}(A)}$ is linear, and given by

$$\mathcal{P}_{\mathcal{R}(A)}(y) = Ax_{\text{ls}} = A(A^T A)^{-1} A^T y$$

- $A(A^T A)^{-1} A^T$ is called the *projection matrix* (associated with $\mathcal{R}(A)$)

Orthogonality principle

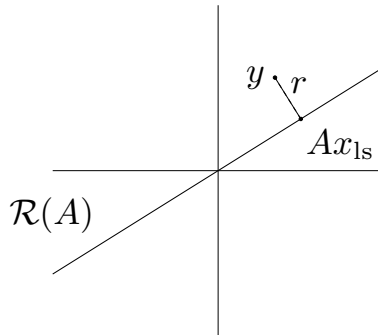
optimal residual

$$r = Ax_{\text{ls}} - y = (A(A^T A)^{-1} A^T - I)y$$

is orthogonal to $\mathcal{R}(A)$:

$$\langle r, Az \rangle = y^T (A(A^T A)^{-1} A^T - I)^T A z = 0$$

for all $z \in \mathbf{R}^n$



Least-squares

5-7

Completion of squares

since $r = Ax_{\text{ls}} - y \perp A(x - x_{\text{ls}})$ for any x , we have

$$\begin{aligned} \|Ax - y\|^2 &= \|(Ax_{\text{ls}} - y) + A(x - x_{\text{ls}})\|^2 \\ &= \|Ax_{\text{ls}} - y\|^2 + \|A(x - x_{\text{ls}})\|^2 \end{aligned}$$

this shows that for $x \neq x_{\text{ls}}$, $\|Ax - y\| > \|Ax_{\text{ls}} - y\|$

Least-squares

5-8

Least-squares via QR factorization

- $A \in \mathbf{R}^{m \times n}$ skinny, full rank
- factor as $A = QR$ with $Q^T Q = I_n$, $R \in \mathbf{R}^{n \times n}$ upper triangular, invertible
- pseudo-inverse is

$$(A^T A)^{-1} A^T = (R^T Q^T Q R)^{-1} R^T Q^T = R^{-1} Q^T$$

$$\text{so } x_{\text{ls}} = R^{-1} Q^T y$$

- projection on $\mathcal{R}(A)$ given by matrix

$$A(A^T A)^{-1} A^T = A R^{-1} Q^T = Q Q^T$$

Least-squares via full QR factorization

- full QR factorization:

$$A = [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$$

with $[Q_1 \ Q_2] \in \mathbf{R}^{m \times m}$ orthogonal, $R_1 \in \mathbf{R}^{n \times n}$ upper triangular, invertible

- multiplication by orthogonal matrix doesn't change norm, so

$$\begin{aligned} \|Ax - y\|^2 &= \left\| [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} x - y \right\|^2 \\ &= \left\| [Q_1 \ Q_2]^T [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} x - [Q_1 \ Q_2]^T y \right\|^2 \end{aligned}$$

$$\begin{aligned}
&= \left\| \begin{bmatrix} R_1 x - Q_1^T y \\ -Q_2^T y \end{bmatrix} \right\|^2 \\
&= \|R_1 x - Q_1^T y\|^2 + \|Q_2^T y\|^2
\end{aligned}$$

- this is evidently minimized by choice $x_{\text{ls}} = R_1^{-1} Q_1^T y$ (which makes first term zero)
- residual with optimal x is

$$Ax_{\text{ls}} - y = -Q_2 Q_2^T y$$

- $Q_1 Q_1^T$ gives projection onto $\mathcal{R}(A)$
- $Q_2 Q_2^T$ gives projection onto $\mathcal{R}(A)^\perp$

Least-squares estimation

many applications in inversion, estimation, and reconstruction problems have form

$$y = Ax + v$$

- x is what we want to estimate or reconstruct
- y is our sensor measurement(s)
- v is an unknown *noise* or *measurement error* (assumed small)
- i th row of A characterizes i th sensor

least-squares estimation: choose as estimate \hat{x} that minimizes

$$\|A\hat{x} - y\|$$

i.e., deviation between

- what we actually observed (y), and
- what we would observe if $x = \hat{x}$, and there were no noise ($v = 0$)

least-squares estimate is just $\hat{x} = (A^T A)^{-1} A^T y$

BLUE property

linear measurement with noise:

$$y = Ax + v$$

with A full rank, skinny

consider a *linear estimator* of form $\hat{x} = By$

- called *unbiased* if $\hat{x} = x$ whenever $v = 0$
(*i.e.*, no estimation error when there is no noise)
same as $BA = I$, *i.e.*, B is left inverse of A

- estimation error of unbiased linear estimator is

$$x - \hat{x} = x - B(Ax + v) = -Bv$$

obviously, then, we'd like B 'small' (and $BA = I$)

- **fact:** $A^\dagger = (A^T A)^{-1} A^T$ is the *smallest* left inverse of A , in the following sense:

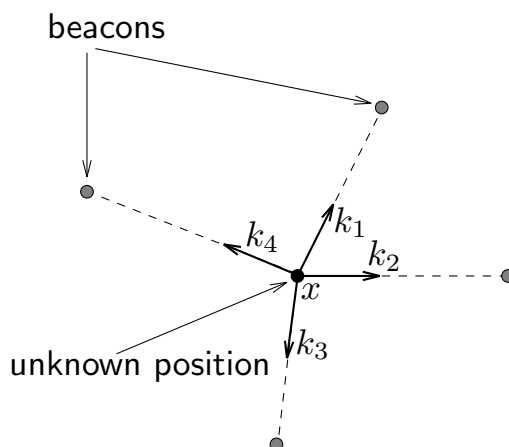
for any B with $BA = I$, we have

$$\sum_{i,j} B_{ij}^2 \geq \sum_{i,j} A_{ij}^{\dagger 2}$$

i.e., least-squares provides the *best linear unbiased estimator* (BLUE)

Navigation from range measurements

navigation using range measurements from *distant* beacons



beacons far from unknown position $x \in \mathbf{R}^2$, so linearization around $x = 0$ (say) nearly exact

ranges $y \in \mathbf{R}^4$ measured, with measurement noise v :

$$y = - \begin{bmatrix} k_1^T \\ k_2^T \\ k_3^T \\ k_4^T \end{bmatrix} x + v$$

where k_i is unit vector from 0 to beacon i

measurement errors are independent, Gaussian, with standard deviation 2
(details not important)

problem: estimate $x \in \mathbf{R}^2$, given $y \in \mathbf{R}^4$

(roughly speaking, a 2:1 measurement redundancy ratio)

actual position is $x = (5.59, 10.58)$;

measurement is $y = (-11.95, -2.84, -9.81, 2.81)$

Just enough measurements method

y_1 and y_2 suffice to find x (when $v = 0$)

compute estimate \hat{x} by inverting top (2×2) half of A :

$$\hat{x} = B_{je}y = \begin{bmatrix} 0 & -1.0 & 0 & 0 \\ -1.12 & 0.5 & 0 & 0 \end{bmatrix} y = \begin{bmatrix} 2.84 \\ 11.9 \end{bmatrix}$$

(norm of error: 3.07)

Least-squares method

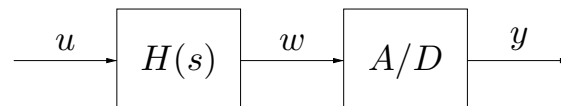
compute estimate \hat{x} by least-squares:

$$\hat{x} = A^\dagger y = \begin{bmatrix} -0.23 & -0.48 & 0.04 & 0.44 \\ -0.47 & -0.02 & -0.51 & -0.18 \end{bmatrix} y = \begin{bmatrix} 4.95 \\ 10.26 \end{bmatrix}$$

(norm of error: 0.72)

- B_{je} and A^\dagger are both left inverses of A
- larger entries in B lead to larger estimation error

Example from overview lecture



- signal u is piecewise constant, period 1 sec, $0 \leq t \leq 10$:

$$u(t) = x_j, \quad j-1 \leq t < j, \quad j = 1, \dots, 10$$

- filtered by system with impulse response $h(t)$:

$$w(t) = \int_0^t h(t-\tau)u(\tau) d\tau$$

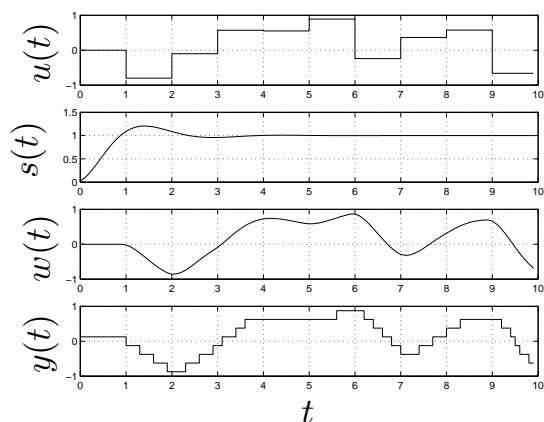
- sample at 10Hz: $\tilde{y}_i = w(0.1i)$, $i = 1, \dots, 100$

- 3-bit quantization: $y_i = Q(\tilde{y}_i)$, $i = 1, \dots, 100$, where Q is 3-bit quantizer characteristic

$$Q(a) = (1/4) (\text{round}(4a + 1/2) - 1/2)$$

- **problem:** estimate $x \in \mathbf{R}^{10}$ given $y \in \mathbf{R}^{100}$

example:



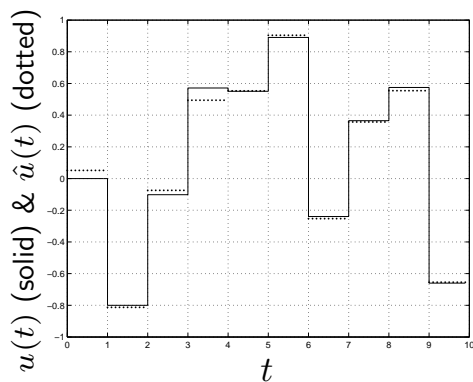
Least-squares

5-21

we have $y = Ax + v$, where

- $A \in \mathbf{R}^{100 \times 10}$ is given by $A_{ij} = \int_{j-1}^j h(0.1i - \tau) d\tau$
- $v \in \mathbf{R}^{100}$ is quantization error: $v_i = Q(\tilde{y}_i) - \tilde{y}_i$ (so $|v_i| \leq 0.125$)

least-squares estimate: $x_{\text{ls}} = (A^T A)^{-1} A^T y$



Least-squares

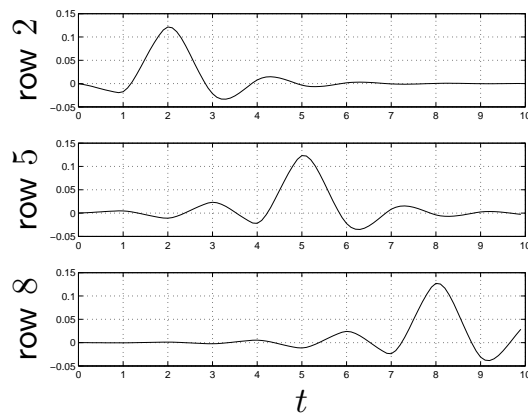
5-22

$$\text{RMS error is } \frac{\|x - x_{\text{ls}}\|}{\sqrt{10}} = 0.03$$

better than if we had no filtering! (RMS error 0.07)

more on this later . . .

some rows of $B_{\text{ls}} = (A^T A)^{-1} A^T$:



- rows show how sampled measurements of y are used to form estimate of x_i for $i = 2, 5, 8$
- to estimate x_5 , which is the original input signal for $4 \leq t < 5$, we mostly use $y(t)$ for $3 \leq t \leq 7$

Lecture 6

Least-squares applications

- least-squares data fitting
- growing sets of regressors
- system identification
- growing sets of measurements and recursive least-squares

6-1

Least-squares data fitting

we are given:

- functions $f_1, \dots, f_n : S \rightarrow \mathbf{R}$, called *regressors* or *basis functions*
- *data or measurements* (s_i, g_i) , $i = 1, \dots, m$, where $s_i \in S$ and (usually) $m \gg n$

problem: find coefficients $x_1, \dots, x_n \in \mathbf{R}$ so that

$$x_1 f_1(s_i) + \dots + x_n f_n(s_i) \approx g_i, \quad i = 1, \dots, m$$

i.e., find linear combination of functions that fits data

least-squares fit: choose x to minimize total square fitting error:

$$\sum_{i=1}^m (x_1 f_1(s_i) + \dots + x_n f_n(s_i) - g_i)^2$$

- using matrix notation, total square fitting error is $\|Ax - g\|^2$, where $A_{ij} = f_j(s_i)$
- hence, least-squares fit is given by

$$x = (A^T A)^{-1} A^T g$$

(assuming A is skinny, full rank)

- corresponding function is

$$f_{\text{lsfit}}(s) = x_1 f_1(s) + \cdots + x_n f_n(s)$$

- applications:
 - interpolation, extrapolation, smoothing of data
 - developing simple, approximate model of data

Least-squares polynomial fitting

problem: fit polynomial of degree $< n$,

$$p(t) = a_0 + a_1 t + \cdots + a_{n-1} t^{n-1},$$

to data (t_i, y_i) , $i = 1, \dots, m$

- basis functions are $f_j(t) = t^{j-1}$, $j = 1, \dots, n$
- matrix A has form $A_{ij} = t_i^{j-1}$

$$A = \begin{bmatrix} 1 & t_1 & t_1^2 & \cdots & t_1^{n-1} \\ 1 & t_2 & t_2^2 & \cdots & t_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_m & t_m^2 & \cdots & t_m^{n-1} \end{bmatrix}$$

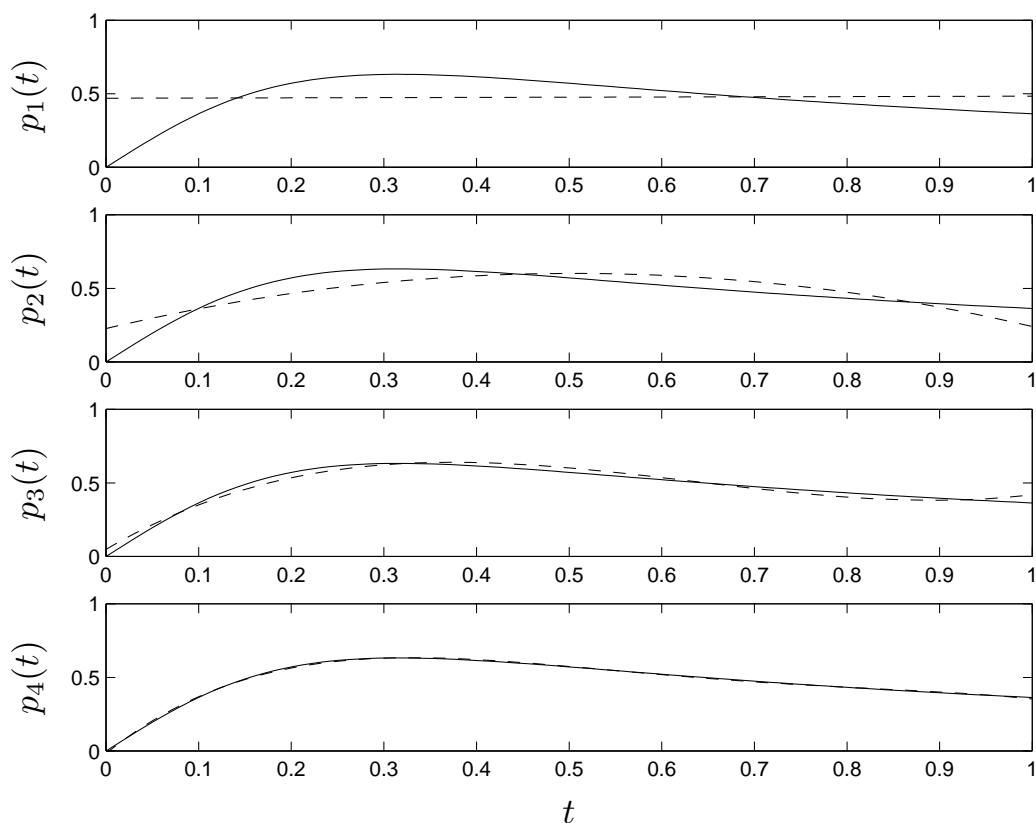
(called a *Vandermonde matrix*)

assuming $t_k \neq t_l$ for $k \neq l$ and $m \geq n$, A is full rank:

- suppose $Aa = 0$
- corresponding polynomial $p(t) = a_0 + \cdots + a_{n-1}t^{n-1}$ vanishes at m points t_1, \dots, t_m
- by fundamental theorem of algebra p can have no more than $n - 1$ zeros, so p is identically zero, and $a = 0$
- columns of A are independent, *i.e.*, A full rank

Example

- fit $g(t) = 4t/(1 + 10t^2)$ with polynomial
- $m = 100$ points between $t = 0$ & $t = 1$
- least-squares fit for degrees 1, 2, 3, 4 have RMS errors .135, .076, .025, .005, respectively



Growing sets of regressors

consider *family* of least-squares problems

$$\text{minimize } \|\sum_{i=1}^p x_i a_i - y\|$$

for $p = 1, \dots, n$

(a_1, \dots, a_p are called *regressors*)

- approximate y by linear combination of a_1, \dots, a_p
- project y onto $\text{span}\{a_1, \dots, a_p\}$
- regress y on a_1, \dots, a_p
- as p increases, get better fit, so optimal residual decreases

solution for each $p \leq n$ is given by

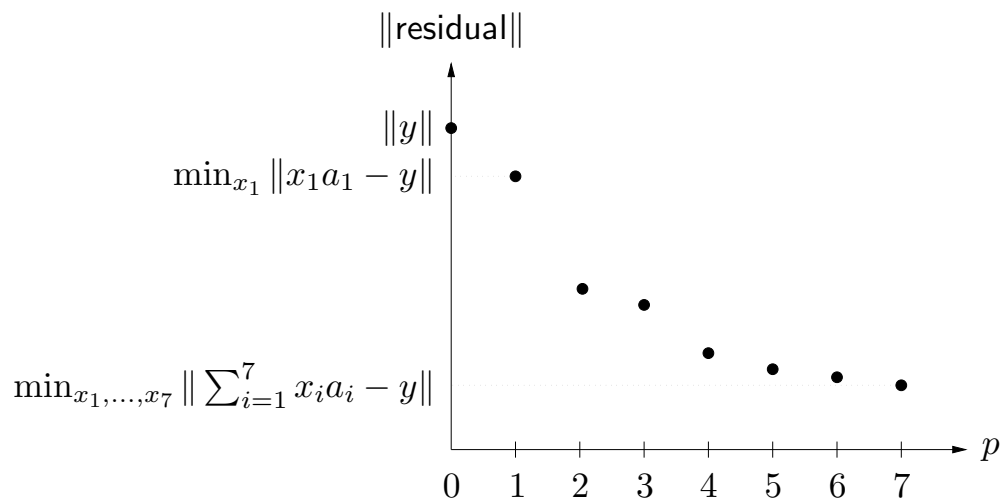
$$x_{\text{ls}}^{(p)} = (A_p^T A_p)^{-1} A_p^T y = R_p^{-1} Q_p^T y$$

where

- $A_p = [a_1 \cdots a_p] \in \mathbf{R}^{m \times p}$ is the first p columns of A
- $A_p = Q_p R_p$ is the QR factorization of A_p
- $R_p \in \mathbf{R}^{p \times p}$ is the leading $p \times p$ submatrix of R
- $Q_p = [q_1 \cdots q_p]$ is the first p columns of Q

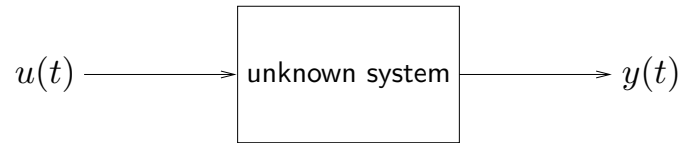
Norm of optimal residual versus p

plot of optimal residual versus p shows how well y can be matched by linear combination of a_1, \dots, a_p , as function of p



Least-squares system identification

we measure input $u(t)$ and output $y(t)$ for $t = 0, \dots, N$ of unknown system



system identification problem: find reasonable model for system based on measured I/O data u, y

example with scalar u, y (vector u, y readily handled): fit I/O data with moving-average (MA) model with n delays

$$\hat{y}(t) = h_0 u(t) + h_1 u(t-1) + \dots + h_n u(t-n)$$

where $h_0, \dots, h_n \in \mathbf{R}$

we can write model or predicted output as

$$\begin{bmatrix} \hat{y}(n) \\ \hat{y}(n+1) \\ \vdots \\ \hat{y}(N) \end{bmatrix} = \begin{bmatrix} u(n) & u(n-1) & \dots & u(0) \\ u(n+1) & u(n) & \dots & u(1) \\ \vdots & \vdots & \dots & \vdots \\ u(N) & u(N-1) & \dots & u(N-n) \end{bmatrix} \begin{bmatrix} h_0 \\ h_1 \\ \vdots \\ h_n \end{bmatrix}$$

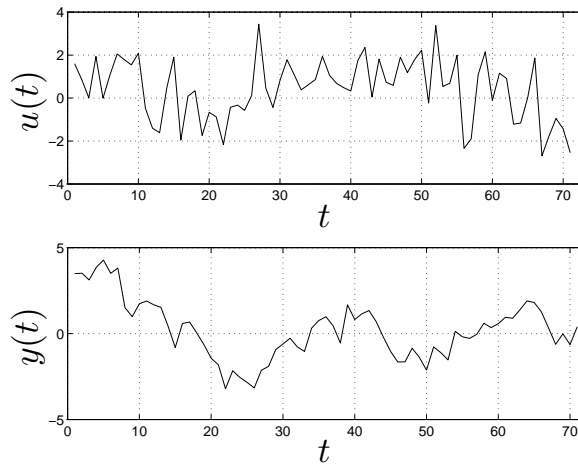
model prediction error is

$$e = (y(n) - \hat{y}(n), \dots, y(N) - \hat{y}(N))$$

least-squares identification: choose model (*i.e.*, h) that minimizes norm of model prediction error $\|e\|$

... a least-squares problem (with variables h)

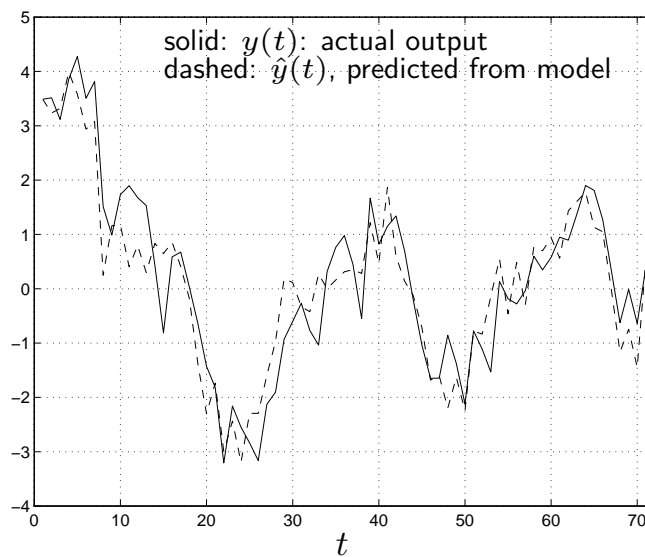
Example



for $n = 7$ we obtain MA model with

$$(h_0, \dots, h_7) = (.024, .282, .418, .354, .243, .487, .208, .441)$$

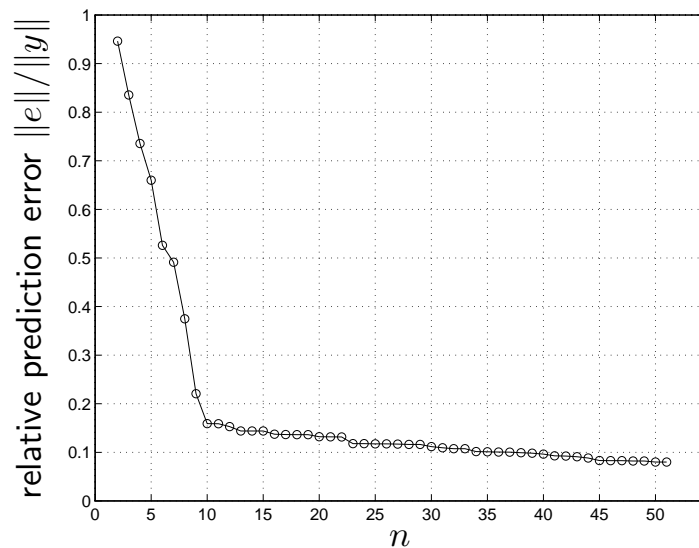
with relative prediction error $\|e\|/\|y\| = 0.37$



Model order selection

question: how large should n be?

- obviously the larger n , the smaller the prediction error *on the data used to form the model*
- suggests using largest possible model order for smallest prediction error

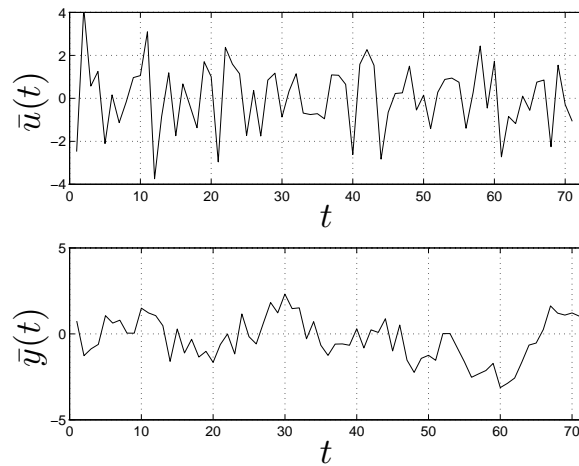


difficulty: for n too large the *predictive ability* of the model on *other I/O data* (from the same system) becomes worse

Cross-validation

evaluate model predictive performance on another I/O data set *not used to develop model*

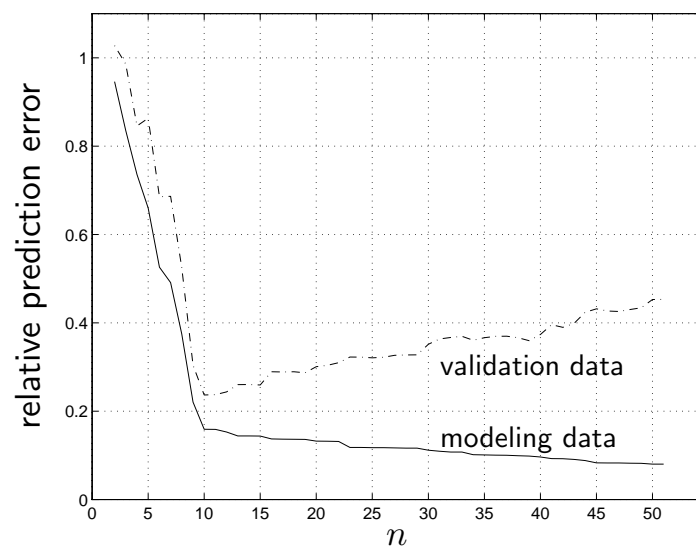
model validation data set:



Least-squares applications

6-17

now check prediction error of models (developed using *modeling data*) on *validation data*:

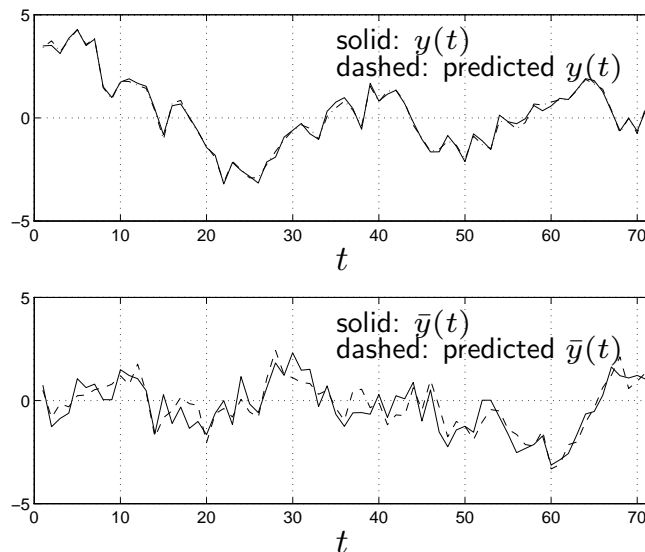


plot suggests $n = 10$ is a good choice

Least-squares applications

6-18

for $n = 50$ the actual and predicted outputs on system identification and model validation data are:



loss of predictive ability when n too large is called *model overfit* or *overmodeling*

Growing sets of measurements

least-squares problem in 'row' form:

$$\text{minimize} \quad \|Ax - y\|^2 = \sum_{i=1}^m (a_i^T x - y_i)^2$$

where a_i^T are the rows of A ($a_i \in \mathbf{R}^n$)

- $x \in \mathbf{R}^n$ is some vector to be estimated
- each pair a_i, y_i corresponds to one measurement
- solution is

$$x_{\text{ls}} = \left(\sum_{i=1}^m a_i a_i^T \right)^{-1} \sum_{i=1}^m y_i a_i$$

- suppose that a_i and y_i become available sequentially, *i.e.*, m increases with time

Recursive least-squares

we can compute $x_{\text{ls}}(m) = \left(\sum_{i=1}^m a_i a_i^T \right)^{-1} \sum_{i=1}^m y_i a_i$ recursively

- initialize $P(0) = 0 \in \mathbf{R}^{n \times n}$, $q(0) = 0 \in \mathbf{R}^n$
- for $m = 0, 1, \dots$,

$$P(m+1) = P(m) + a_{m+1} a_{m+1}^T \quad q(m+1) = q(m) + y_{m+1} a_{m+1}$$

- if $P(m)$ is invertible, we have $x_{\text{ls}}(m) = P(m)^{-1} q(m)$
- $P(m)$ is invertible $\iff a_1, \dots, a_m$ span \mathbf{R}^n
(so, once $P(m)$ becomes invertible, it stays invertible)

Fast update for recursive least-squares

we can calculate

$$P(m+1)^{-1} = (P(m) + a_{m+1} a_{m+1}^T)^{-1}$$

efficiently from $P(m)^{-1}$ using the *rank one update formula*

$$(P + aa^T)^{-1} = P^{-1} - \frac{1}{1 + a^T P^{-1} a} (P^{-1} a)(P^{-1} a)^T$$

valid when $P = P^T$, and P and $P + aa^T$ are both invertible

- gives an $O(n^2)$ method for computing $P(m+1)^{-1}$ from $P(m)^{-1}$
- standard methods for computing $P(m+1)^{-1}$ from $P(m+1)$ are $O(n^3)$

Verification of rank one update formula

$$\begin{aligned} & (P + aa^T) \left(P^{-1} - \frac{1}{1 + a^T P^{-1} a} (P^{-1} a)(P^{-1} a)^T \right) \\ &= I + aa^T P^{-1} - \frac{1}{1 + a^T P^{-1} a} P (P^{-1} a)(P^{-1} a)^T \\ &\quad - \frac{1}{1 + a^T P^{-1} a} aa^T (P^{-1} a)(P^{-1} a)^T \\ &= I + aa^T P^{-1} - \frac{1}{1 + a^T P^{-1} a} aa^T P^{-1} - \frac{a^T P^{-1} a}{1 + a^T P^{-1} a} aa^T P^{-1} \\ &= I \end{aligned}$$

Lecture 7

Regularized least-squares and Gauss-Newton method

- multi-objective least-squares
- regularized least-squares
- nonlinear least-squares
- Gauss-Newton method

7-1

Multi-objective least-squares

in many problems we have two (or more) objectives

- we want $J_1 = \|Ax - y\|^2$ small
- and also $J_2 = \|Fx - g\|^2$ small

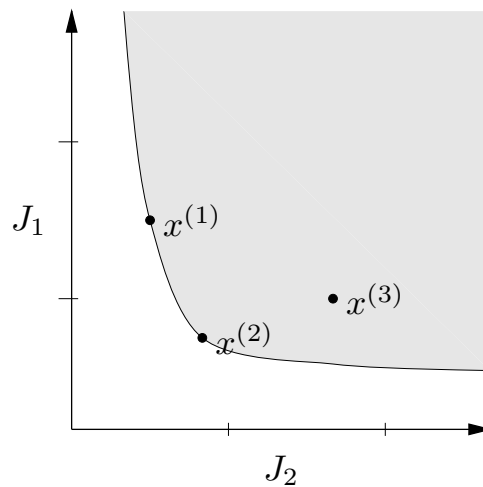
($x \in \mathbf{R}^n$ is the variable)

- usually the objectives are *competing*
- we can make one smaller, at the expense of making the other larger

common example: $F = I$, $g = 0$; we want $\|Ax - y\|$ small, with small x

Plot of achievable objective pairs

plot (J_2, J_1) for every x :



note that $x \in \mathbf{R}^n$, but this plot is in \mathbf{R}^2 ; point labeled $x^{(1)}$ is really $(J_2(x^{(1)}), J_1(x^{(1)}))$

- shaded area shows (J_2, J_1) achieved by some $x \in \mathbf{R}^n$
- clear area shows (J_2, J_1) not achieved by any $x \in \mathbf{R}^n$
- boundary of region is called *optimal trade-off curve*
- corresponding x are called *Pareto optimal*
(for the two objectives $\|Ax - y\|^2$, $\|Fx - g\|^2$)

three example choices of x : $x^{(1)}$, $x^{(2)}$, $x^{(3)}$

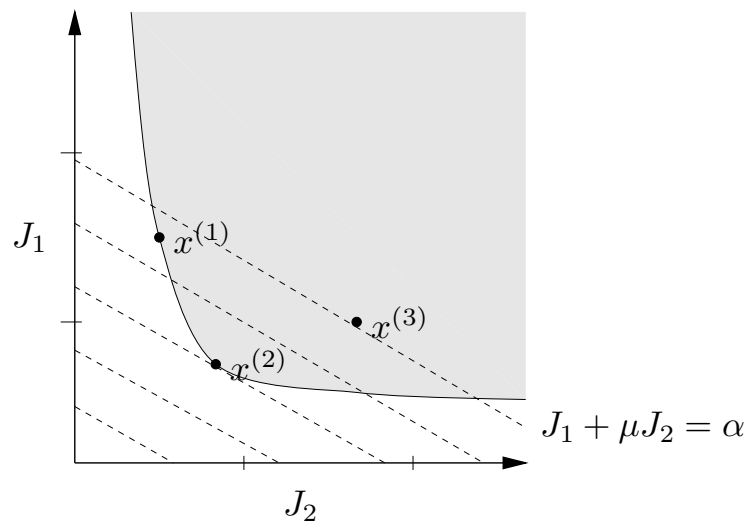
- $x^{(3)}$ is worse than $x^{(2)}$ on both counts (J_2 and J_1)
- $x^{(1)}$ is better than $x^{(2)}$ in J_2 , but worse in J_1

Weighted-sum objective

- to find Pareto optimal points, *i.e.*, x 's on optimal trade-off curve, we minimize *weighted-sum objective*

$$J_1 + \mu J_2 = \|Ax - y\|^2 + \mu \|Fx - g\|^2$$

- parameter $\mu \geq 0$ gives relative weight between J_1 and J_2
- points where weighted sum is constant, $J_1 + \mu J_2 = \alpha$, correspond to line with slope $-\mu$ on (J_2, J_1) plot



- $x^{(2)}$ minimizes weighted-sum objective for μ shown
- by varying μ from 0 to $+\infty$, can sweep out entire *optimal tradeoff curve*

Minimizing weighted-sum objective

can express weighted-sum objective as ordinary least-squares objective:

$$\begin{aligned}\|Ax - y\|^2 + \mu\|Fx - g\|^2 &= \left\| \begin{bmatrix} A \\ \sqrt{\mu}F \end{bmatrix} x - \begin{bmatrix} y \\ \sqrt{\mu}g \end{bmatrix} \right\|^2 \\ &= \left\| \tilde{A}x - \tilde{y} \right\|^2\end{aligned}$$

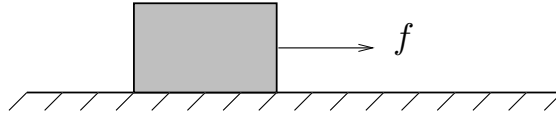
where

$$\tilde{A} = \begin{bmatrix} A \\ \sqrt{\mu}F \end{bmatrix}, \quad \tilde{y} = \begin{bmatrix} y \\ \sqrt{\mu}g \end{bmatrix}$$

hence solution is (assuming \tilde{A} full rank)

$$\begin{aligned}x &= (\tilde{A}^T \tilde{A})^{-1} \tilde{A}^T \tilde{y} \\ &= (A^T A + \mu F^T F)^{-1} (A^T y + \mu F^T g)\end{aligned}$$

Example



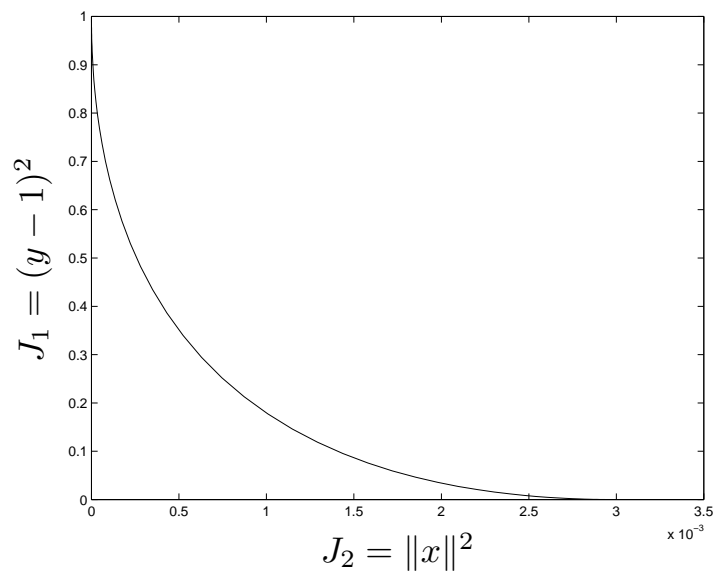
- unit mass at rest subject to forces x_i for $i - 1 < t \leq i$, $i = 1, \dots, 10$
- $y \in \mathbf{R}$ is position at $t = 10$; $y = a^T x$ where $a \in \mathbf{R}^{10}$
- $J_1 = (y - 1)^2$ (final position error squared)
- $J_2 = \|x\|^2$ (sum of squares of forces)

weighted-sum objective: $(a^T x - 1)^2 + \mu\|x\|^2$

optimal x :

$$x = (aa^T + \mu I)^{-1} a$$

optimal trade-off curve:



- upper left corner of optimal trade-off curve corresponds to $x = 0$
- bottom right corresponds to input that yields $y = 1$, *i.e.*, $J_1 = 0$

Regularized least-squares

when $F = I$, $g = 0$ the objectives are

$$J_1 = \|Ax - y\|^2, \quad J_2 = \|x\|^2$$

minimizer of weighted-sum objective,

$$x = (A^T A + \mu I)^{-1} A^T y,$$

is called *regularized* least-squares (approximate) solution of $Ax \approx y$

- also called *Tychonov regularization*
- for $\mu > 0$, works for *any* A (no restrictions on shape, rank . . .)

estimation/inversion application:

- $Ax - y$ is sensor residual
- prior information: x small
- or, model only accurate for x small
- regularized solution trades off sensor fit, size of x

Nonlinear least-squares

nonlinear least-squares (NLLS) problem: find $x \in \mathbf{R}^n$ that minimizes

$$\|r(x)\|^2 = \sum_{i=1}^m r_i(x)^2,$$

where $r : \mathbf{R}^n \rightarrow \mathbf{R}^m$

- $r(x)$ is a vector of ‘residuals’
- reduces to (linear) least-squares if $r(x) = Ax - y$

Position estimation from ranges

estimate position $x \in \mathbf{R}^2$ from approximate distances to beacons at locations $b_1, \dots, b_m \in \mathbf{R}^2$ *without* linearizing

- we measure $\rho_i = \|x - b_i\| + v_i$
(v_i is range error, unknown but assumed small)
- NLLS estimate: choose \hat{x} to minimize

$$\sum_{i=1}^m r_i(x)^2 = \sum_{i=1}^m (\rho_i - \|x - b_i\|)^2$$

Gauss-Newton method for NLLS

NLLS: find $x \in \mathbf{R}^n$ that minimizes $\|r(x)\|^2 = \sum_{i=1}^m r_i(x)^2$, where
 $r : \mathbf{R}^n \rightarrow \mathbf{R}^m$

- in general, very hard to solve exactly
- many good heuristics to compute *locally optimal* solution

Gauss-Newton method:

given starting guess for x

repeat

 linearize r near current guess

 new guess is linear LS solution, using linearized r

until convergence

Gauss-Newton method (more detail):

- linearize r near current iterate $x^{(k)}$:

$$r(x) \approx r(x^{(k)}) + Dr(x^{(k)})(x - x^{(k)})$$

where Dr is the Jacobian: $(Dr)_{ij} = \partial r_i / \partial x_j$

- write linearized approximation as

$$r(x^{(k)}) + Dr(x^{(k)})(x - x^{(k)}) = A^{(k)}x - b^{(k)}$$

$$A^{(k)} = Dr(x^{(k)}), \quad b^{(k)} = Dr(x^{(k)})x^{(k)} - r(x^{(k)})$$

- at k th iteration, we approximate NLLS problem by linear LS problem:

$$\|r(x)\|^2 \approx \|A^{(k)}x - b^{(k)}\|^2$$

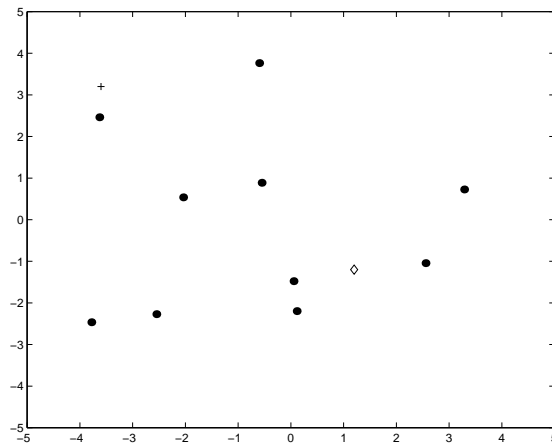
- next iterate solves this linearized LS problem:

$$x^{(k+1)} = \left(A^{(k)T} A^{(k)} \right)^{-1} A^{(k)T} b^{(k)}$$

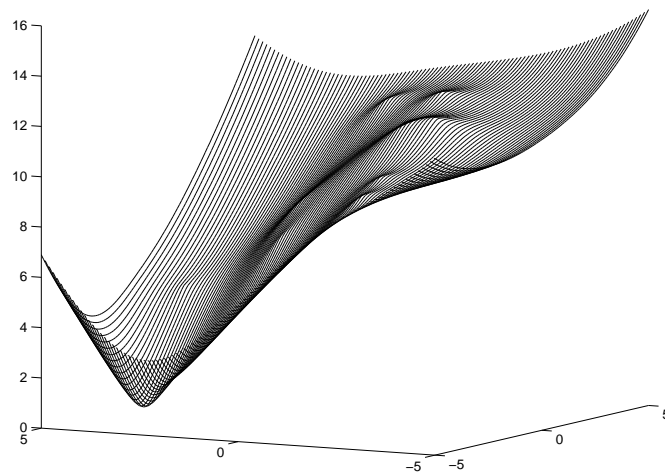
- repeat until convergence (which *isn't* guaranteed)

Gauss-Newton example

- 10 beacons
- + true position $(-3.6, 3.2)$; \diamond initial guess $(1.2, -1.2)$
- range estimates accurate to ± 0.5

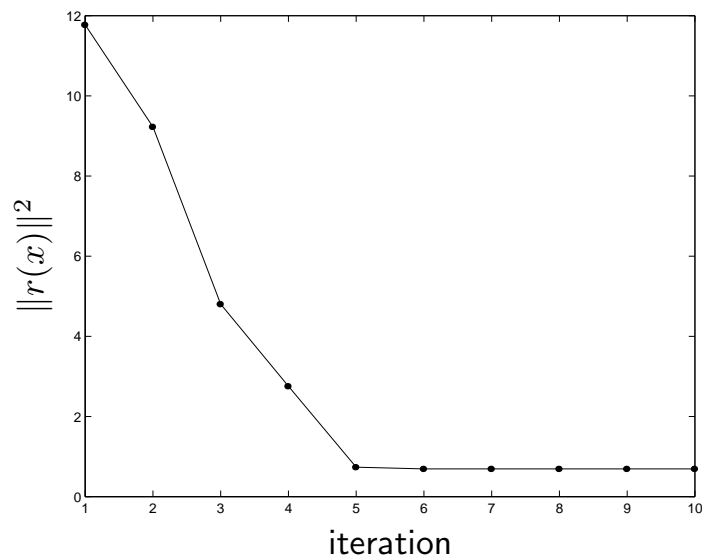


NLLS objective $\|r(x)\|^2$ versus x :



- for a linear LS problem, objective would be nice quadratic 'bowl'
- bumps in objective due to strong nonlinearity of r

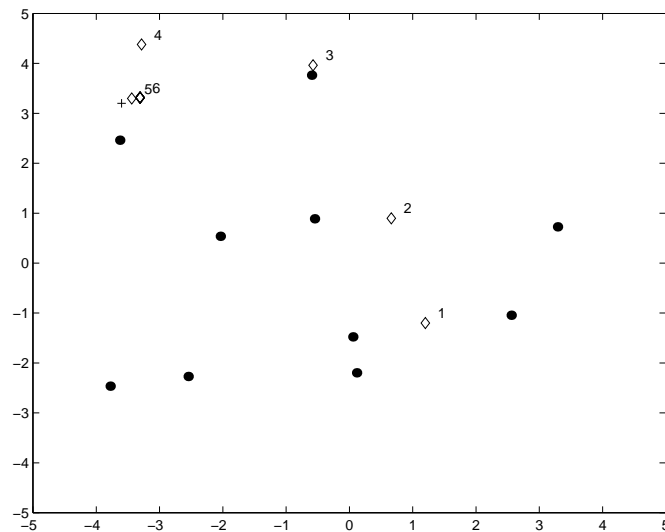
objective of Gauss-Newton iterates:



- $x^{(k)}$ converges to (in this case, global) minimum of $\|r(x)\|^2$
- convergence takes only five or so steps

- final estimate is $\hat{x} = (-3.3, 3.3)$
- estimation error is $\|\hat{x} - x\| = 0.31$
(substantially smaller than range accuracy!)

convergence of Gauss-Newton iterates:



useful variation on Gauss-Newton: add regularization term

$$\|A^{(k)}x - b^{(k)}\|^2 + \mu\|x - x^{(k)}\|^2$$

so that next iterate is not too far from previous one (hence, linearized model still pretty accurate)

Lecture 8

Least-norm solutions of undetermined equations

- least-norm solution of underdetermined equations
- minimum norm solutions via QR factorization
- derivation via Lagrange multipliers
- relation to regularized least-squares
- general norm minimization with equality constraints

8-1

Underdetermined linear equations

we consider

$$y = Ax$$

where $A \in \mathbf{R}^{m \times n}$ is fat ($m < n$), *i.e.*,

- there are more variables than equations
- x is *underspecified*, *i.e.*, many choices of x lead to the same y

we'll assume that A is full rank (m), so for each $y \in \mathbf{R}^m$, there is a solution set of all solutions has form

$$\{ x \mid Ax = y \} = \{ x_p + z \mid z \in \mathcal{N}(A) \}$$

where x_p is any ('particular') solution, *i.e.*, $Ax_p = y$

- z characterizes available choices in solution
- solution has $\dim \mathcal{N}(A) = n - m$ 'degrees of freedom'
- can choose z to satisfy other specs or optimize among solutions

Least-norm solution

one particular solution is

$$x_{\text{ln}} = A^T(AA^T)^{-1}y$$

(AA^T is invertible since A full rank)

in fact, x_{ln} is the solution of $y = Ax$ that minimizes $\|x\|$

i.e., x_{ln} is solution of optimization problem

$$\begin{array}{ll} \text{minimize} & \|x\| \\ \text{subject to} & Ax = y \end{array}$$

(with variable $x \in \mathbf{R}^n$)

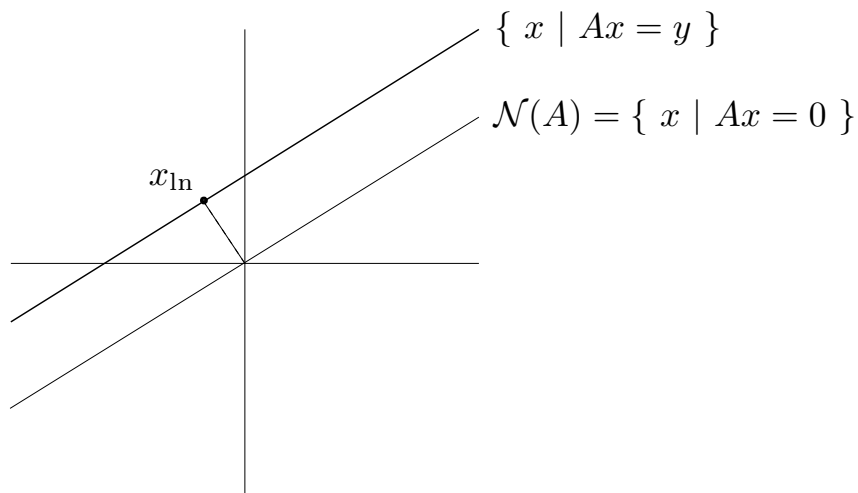
suppose $Ax = y$, so $A(x - x_{\text{ln}}) = 0$ and

$$\begin{aligned} (x - x_{\text{ln}})^T x_{\text{ln}} &= (x - x_{\text{ln}})^T A^T (AA^T)^{-1} y \\ &= (A(x - x_{\text{ln}}))^T (AA^T)^{-1} y \\ &= 0 \end{aligned}$$

i.e., $(x - x_{\text{ln}}) \perp x_{\text{ln}}$, so

$$\|x\|^2 = \|x_{\text{ln}} + x - x_{\text{ln}}\|^2 = \|x_{\text{ln}}\|^2 + \|x - x_{\text{ln}}\|^2 \geq \|x_{\text{ln}}\|^2$$

i.e., x_{ln} has smallest norm of any solution



- **orthogonality condition:** $x_{\text{ln}} \perp \mathcal{N}(A)$
- **projection interpretation:** x_{ln} is projection of 0 on solution set $\{ x \mid Ax = y \}$

- $A^\dagger = A^T(AA^T)^{-1}$ is called the *pseudo-inverse* of full rank, fat A
- $A^T(AA^T)^{-1}$ is a *right inverse* of A
- $I - A^T(AA^T)^{-1}A$ gives projection onto $\mathcal{N}(A)$

cf. analogous formulas for full rank, **skinny** matrix A :

- $A^\dagger = (A^T A)^{-1}A^T$
- $(A^T A)^{-1}A^T$ is a *left inverse* of A
- $A(A^T A)^{-1}A^T$ gives projection onto $\mathcal{R}(A)$

Least-norm solution via QR factorization

find QR factorization of A^T , i.e., $A^T = QR$, with

- $Q \in \mathbf{R}^{n \times m}$, $Q^T Q = I_m$
- $R \in \mathbf{R}^{m \times m}$ upper triangular, nonsingular

then

- $x_{\text{ln}} = A^T(AA^T)^{-1}y = QR^{-T}y$
- $\|x_{\text{ln}}\| = \|R^{-T}y\|$

Derivation via Lagrange multipliers

- least-norm solution solves optimization problem

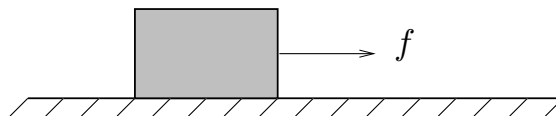
$$\begin{array}{ll} \text{minimize} & x^T x \\ \text{subject to} & Ax = y \end{array}$$

- introduce Lagrange multipliers: $L(x, \lambda) = x^T x + \lambda^T (Ax - y)$
- optimality conditions are

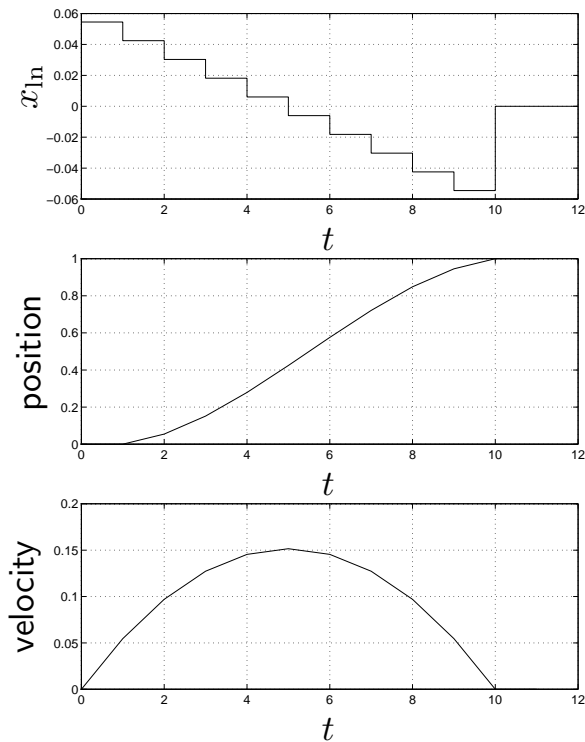
$$\nabla_x L = 2x + A^T \lambda = 0, \quad \nabla_\lambda L = Ax - y = 0$$

- from first condition, $x = -A^T \lambda / 2$
- substitute into second to get $\lambda = -2(AA^T)^{-1}y$
- hence $x = A^T(AA^T)^{-1}y$

Example: transferring mass unit distance



- unit mass at rest subject to forces x_i for $i - 1 < t \leq i$, $i = 1, \dots, 10$
- y_1 is position at $t = 10$, y_2 is velocity at $t = 10$
- $y = Ax$ where $A \in \mathbf{R}^{2 \times 10}$ (A is fat)
- find least norm force that transfers mass unit distance with zero final velocity, *i.e.*, $y = (1, 0)$



Relation to regularized least-squares

- suppose $A \in \mathbf{R}^{m \times n}$ is fat, full rank
- define $J_1 = \|Ax - y\|^2$, $J_2 = \|x\|^2$
- least-norm solution minimizes J_2 with $J_1 = 0$
- minimizer of weighted-sum objective $J_1 + \mu J_2 = \|Ax - y\|^2 + \mu \|x\|^2$ is

$$x_\mu = (A^T A + \mu I)^{-1} A^T y$$

- **fact:** $x_\mu \rightarrow x_{\text{ln}}$ as $\mu \rightarrow 0$, *i.e.*, regularized solution converges to least-norm solution as $\mu \rightarrow 0$
- in matrix terms: as $\mu \rightarrow 0$,

$$(A^T A + \mu I)^{-1} A^T \rightarrow A^T (A A^T)^{-1}$$

(for full rank, fat A)

General norm minimization with equality constraints

consider problem

$$\begin{array}{ll}\text{minimize} & \|Ax - b\| \\ \text{subject to} & Cx = d\end{array}$$

with variable x

- includes least-squares and least-norm problems as special cases

- equivalent to

$$\begin{array}{ll}\text{minimize} & (1/2)\|Ax - b\|^2 \\ \text{subject to} & Cx = d\end{array}$$

- Lagrangian is

$$\begin{aligned}L(x, \lambda) &= (1/2)\|Ax - b\|^2 + \lambda^T(Cx - d) \\ &= (1/2)x^T A^T A x - b^T A x + (1/2)b^T b + \lambda^T C x - \lambda^T d\end{aligned}$$

Least-norm solutions of undetermined equations

8-13

- optimality conditions are

$$\nabla_x L = A^T A x - A^T b + C^T \lambda = 0, \quad \nabla_\lambda L = Cx - d = 0$$

- write in block matrix form as

$$\begin{bmatrix} A^T A & C^T \\ C & 0 \end{bmatrix} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} A^T b \\ d \end{bmatrix}$$

- if the block matrix is invertible, we have

$$\begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} A^T A & C^T \\ C & 0 \end{bmatrix}^{-1} \begin{bmatrix} A^T b \\ d \end{bmatrix}$$

if $A^T A$ is invertible, we can derive a more explicit (and complicated) formula for x

- from first block equation we get

$$x = (A^T A)^{-1}(A^T b - C^T \lambda)$$

- substitute into $Cx = d$ to get

$$C(A^T A)^{-1}(A^T b - C^T \lambda) = d$$

so

$$\lambda = (C(A^T A)^{-1}C^T)^{-1} (C(A^T A)^{-1}A^T b - d)$$

- recover x from equation above (not pretty)

$$x = (A^T A)^{-1} \left(A^T b - C^T (C(A^T A)^{-1}C^T)^{-1} (C(A^T A)^{-1}A^T b - d) \right)$$

Lecture 9

Autonomous linear dynamical systems

- autonomous linear dynamical systems
- examples
- higher order systems
- linearization near equilibrium point
- linearization along trajectory

9-1

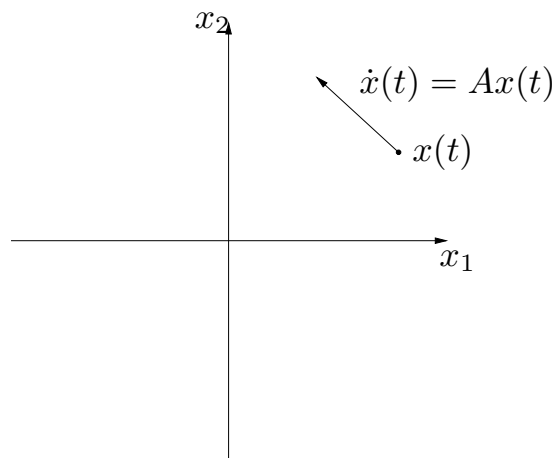
Autonomous linear dynamical systems

continuous-time autonomous LDS has form

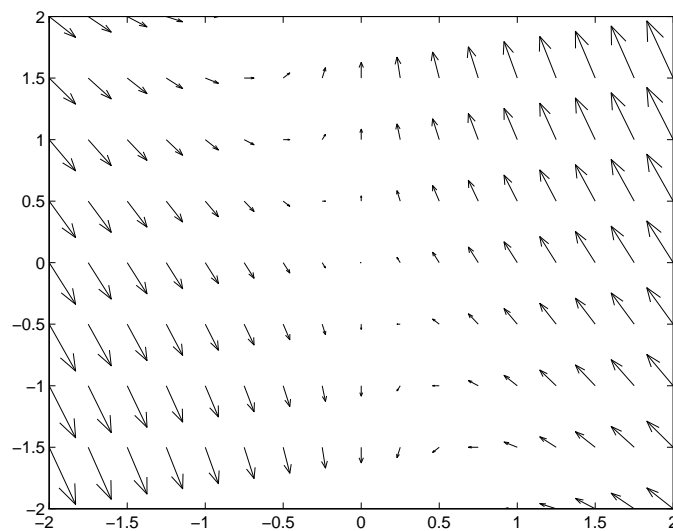
$$\dot{x} = Ax$$

- $x(t) \in \mathbf{R}^n$ is called the state
- n is the *state dimension* or (informally) the *number of states*
- A is the *dynamics matrix*
(system is *time-invariant* if A doesn't depend on t)

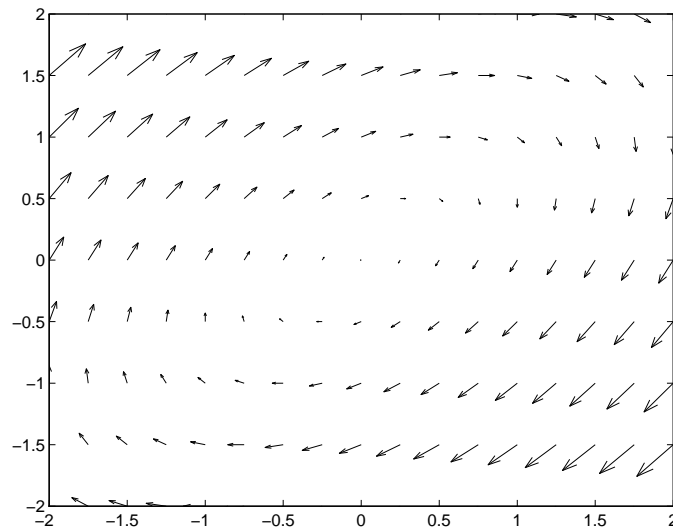
picture (*phase plane*):



example 1: $\dot{x} = \begin{bmatrix} -1 & 0 \\ 2 & 1 \end{bmatrix} x$

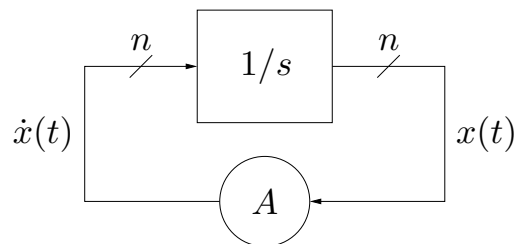


example 2: $\dot{x} = \begin{bmatrix} -0.5 & 1 \\ -1 & 0.5 \end{bmatrix} x$



Block diagram

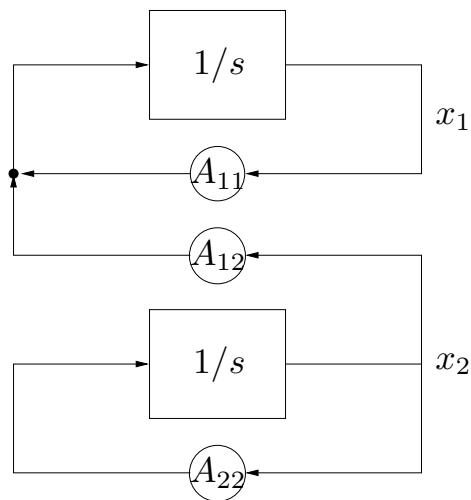
block diagram representation of $\dot{x} = Ax$:



- $1/s$ block represents n parallel scalar integrators
- coupling comes from dynamics matrix A

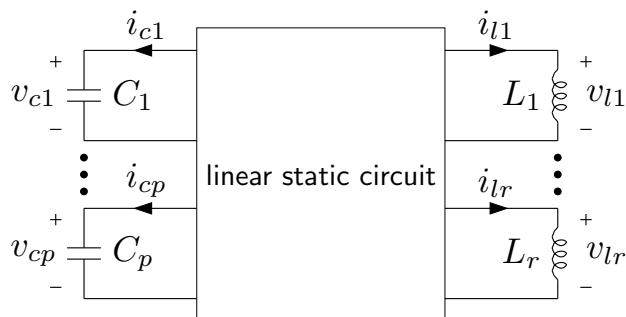
useful when A has structure, *e.g.*, block upper triangular:

$$\dot{x} = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix} x$$



here x_1 doesn't affect x_2 at all

Linear circuit



circuit equations are

$$C \frac{dv_c}{dt} = i_c, \quad L \frac{di_l}{dt} = v_l, \quad \begin{bmatrix} i_c \\ v_l \end{bmatrix} = F \begin{bmatrix} v_c \\ i_l \end{bmatrix}$$

$$C = \mathbf{diag}(C_1, \dots, C_p), \quad L = \mathbf{diag}(L_1, \dots, L_r)$$

with state $x = \begin{bmatrix} v_c \\ i_l \end{bmatrix}$, we have

$$\dot{x} = \begin{bmatrix} C^{-1} & 0 \\ 0 & L^{-1} \end{bmatrix} Fx$$

Chemical reactions

- reaction involving n chemicals; x_i is concentration of chemical i
- linear model of reaction kinetics

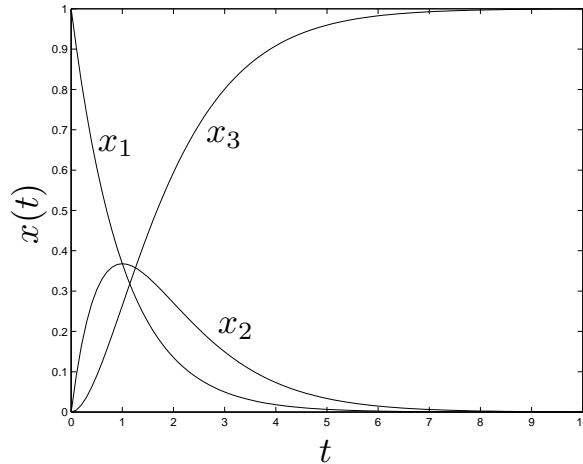
$$\frac{dx_i}{dt} = a_{i1}x_1 + \cdots + a_{in}x_n$$

- good model for some reactions; A is usually sparse

Example: series reaction $A \xrightarrow{k_1} B \xrightarrow{k_2} C$ with linear dynamics

$$\dot{x} = \begin{bmatrix} -k_1 & 0 & 0 \\ k_1 & -k_2 & 0 \\ 0 & k_2 & 0 \end{bmatrix} x$$

plot for $k_1 = k_2 = 1$, initial $x(0) = (1, 0, 0)$



Finite-state discrete-time Markov chain

$z(t) \in \{1, \dots, n\}$ is a random sequence with

$$\mathbf{Prob}(z(t+1) = i \mid z(t) = j) = P_{ij}$$

where $P \in \mathbf{R}^{n \times n}$ is the matrix of *transition probabilities*

can represent probability distribution of $z(t)$ as n -vector

$$p(t) = \begin{bmatrix} \mathbf{Prob}(z(t) = 1) \\ \vdots \\ \mathbf{Prob}(z(t) = n) \end{bmatrix}$$

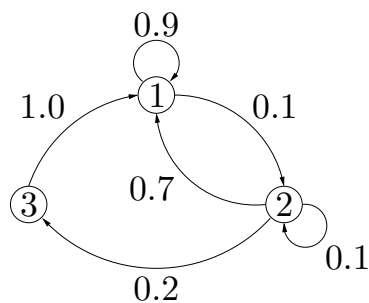
(so, e.g., $\mathbf{Prob}(z(t) = 1, 2, \text{ or } 3) = [1 \ 1 \ 1 \ 0 \cdots 0]p(t)$)

then we have $p(t+1) = Pp(t)$

P is often sparse; Markov chain is depicted graphically

- nodes are states
- edges show transition probabilities

example:



- state 1 is 'system OK'
- state 2 is 'system down'
- state 3 is 'system being repaired'

$$p(t+1) = \begin{bmatrix} 0.9 & 0.7 & 1.0 \\ 0.1 & 0.1 & 0 \\ 0 & 0.2 & 0 \end{bmatrix} p(t)$$

Numerical integration of continuous system

compute approximate solution of $\dot{x} = Ax$, $x(0) = x_0$

suppose h is small time step (x doesn't change much in h seconds)

simple ('forward Euler') approximation:

$$x(t+h) \approx x(t) + h\dot{x}(t) = (I + hA)x(t)$$

by carrying out this recursion (discrete-time LDS), starting at $x(0) = x_0$, we get approximation

$$x(kh) \approx (I + hA)^k x(0)$$

(forward Euler is never used in practice)

Higher order linear dynamical systems

$$x^{(k)} = A_{k-1}x^{(k-1)} + \dots + A_1x^{(1)} + A_0x, \quad x(t) \in \mathbf{R}^n$$

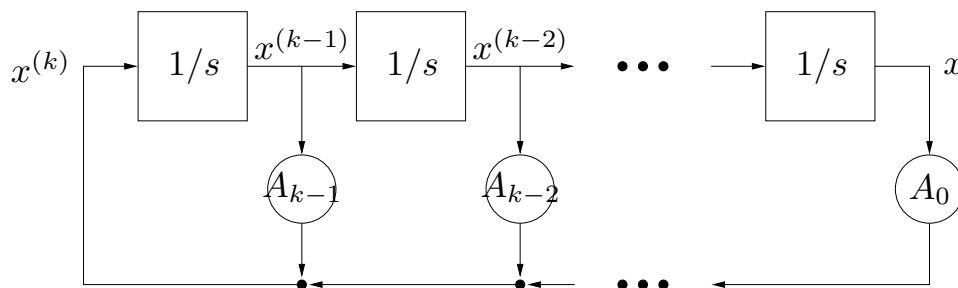
where $x^{(m)}$ denotes m th derivative

define new variable $z = \begin{bmatrix} x \\ x^{(1)} \\ \vdots \\ x^{(k-1)} \end{bmatrix} \in \mathbf{R}^{nk}$, so

$$\dot{z} = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(k)} \end{bmatrix} = \begin{bmatrix} 0 & I & 0 & \dots & 0 \\ 0 & 0 & I & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \dots & I \\ A_0 & A_1 & A_2 & \dots & A_{k-1} \end{bmatrix} z$$

a (first order) LDS (with bigger state)

block diagram:



Mechanical systems

mechanical system with k degrees of freedom undergoing small motions:

$$M\ddot{q} + D\dot{q} + Kq = 0$$

- $q(t) \in \mathbf{R}^k$ is the vector of generalized displacements
- M is the *mass matrix*
- K is the *stiffness matrix*
- D is the *damping matrix*

with state $x = \begin{bmatrix} q \\ \dot{q} \end{bmatrix}$ we have

$$\dot{x} = \begin{bmatrix} \dot{q} \\ \ddot{q} \end{bmatrix} = \begin{bmatrix} 0 & I \\ -M^{-1}K & -M^{-1}D \end{bmatrix} x$$

Linearization near equilibrium point

nonlinear, time-invariant differential equation (DE):

$$\dot{x} = f(x)$$

where $f : \mathbf{R}^n \rightarrow \mathbf{R}^n$

suppose x_e is an *equilibrium point*, i.e., $f(x_e) = 0$

(so $x(t) = x_e$ satisfies DE)

now suppose $x(t)$ is near x_e , so

$$\dot{x}(t) = f(x(t)) \approx f(x_e) + Df(x_e)(x(t) - x_e)$$

with $\delta x(t) = x(t) - x_e$, rewrite as

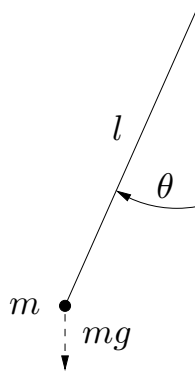
$$\dot{\delta x}(t) \approx Df(x_e)\delta x(t)$$

replacing \approx with $=$ yields *linearized approximation* of DE near x_e

we *hope* solution of $\dot{\delta x} = Df(x_e)\delta x$ is a good approximation of $x - x_e$

(more later)

example: pendulum



2nd order nonlinear DE $ml^2\ddot{\theta} = -lmg \sin \theta$

rewrite as first order DE with state $x = \begin{bmatrix} \theta \\ \dot{\theta} \end{bmatrix}$:

$$\dot{x} = \begin{bmatrix} x_2 \\ -(g/l) \sin x_1 \end{bmatrix}$$

equilibrium point (pendulum down): $x = 0$

linearized system near $x_e = 0$:

$$\dot{\delta x} = \begin{bmatrix} 0 & 1 \\ -g/l & 0 \end{bmatrix} \delta x$$

Does linearization 'work' ?

the linearized system usually, but not always, gives a good idea of the system behavior near x_e

example 1: $\dot{x} = -x^3$ near $x_e = 0$

for $x(0) > 0$ solutions have form $x(t) = (x(0)^{-2} + 2t)^{-1/2}$

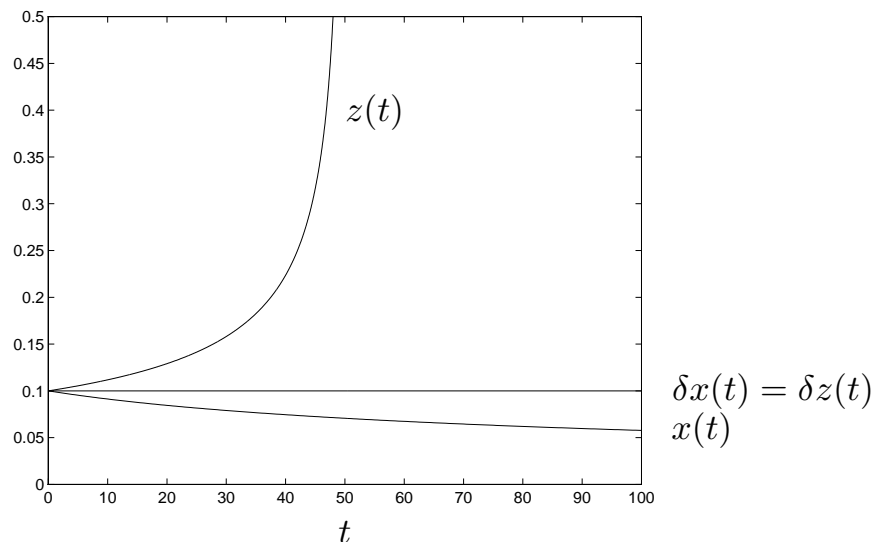
linearized system is $\dot{\delta x} = 0$; solutions are constant

example 2: $\dot{z} = z^3$ near $z_e = 0$

for $z(0) > 0$ solutions have form $z(t) = (z(0)^{-2} - 2t)^{-1/2}$

(finite escape time at $t = z(0)^{-2}/2$)

linearized system is $\dot{\delta z} = 0$; solutions are constant



- systems with very different behavior have same linearized system
- linearized systems do not predict qualitative behavior of either system

Linearization along trajectory

- suppose $x_{\text{traj}} : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ satisfies $\dot{x}_{\text{traj}}(t) = f(x_{\text{traj}}(t), t)$
- suppose $x(t)$ is another trajectory, *i.e.*, $\dot{x}(t) = f(x(t), t)$, and is near $x_{\text{traj}}(t)$
- then

$$\frac{d}{dt}(x - x_{\text{traj}}) = f(x, t) - f(x_{\text{traj}}, t) \approx D_x f(x_{\text{traj}}, t)(x - x_{\text{traj}})$$

- (time-varying) LDS

$$\dot{\delta x} = D_x f(x_{\text{traj}}, t) \delta x$$

is called *linearized* or *variational system* along trajectory x_{traj}

example: linearized oscillator

suppose $x_{\text{traj}}(t)$ is T -periodic solution of nonlinear DE:

$$\dot{x}_{\text{traj}}(t) = f(x_{\text{traj}}(t)), \quad x_{\text{traj}}(t + T) = x_{\text{traj}}(t)$$

linearized system is

$$\dot{\delta x} = A(t) \delta x$$

where $A(t) = Df(x_{\text{traj}}(t))$

$A(t)$ is T -periodic, so linearized system is called *T -periodic linear system*.

used to study:

- startup dynamics of clock and oscillator circuits
- effects of power supply and other disturbances on clock behavior

Lecture 10

Solution via Laplace transform and matrix exponential

- Laplace transform
- solving $\dot{x} = Ax$ via Laplace transform
- state transition matrix
- matrix exponential
- qualitative behavior and stability

10-1

Laplace transform of matrix valued function

suppose $z : \mathbf{R}_+ \rightarrow \mathbf{R}^{p \times q}$

Laplace transform: $Z = \mathcal{L}(z)$, where $Z : D \subseteq \mathbf{C} \rightarrow \mathbf{C}^{p \times q}$ is defined by

$$Z(s) = \int_0^\infty e^{-st} z(t) dt$$

- integral of matrix is done term-by-term
- convention: upper case denotes Laplace transform
- D is the *domain* or *region of convergence* of Z
- D includes at least $\{s \mid \Re s > a\}$, where a satisfies $|z_{ij}(t)| \leq \alpha e^{at}$ for $t \geq 0$, $i = 1, \dots, p$, $j = 1, \dots, q$

Derivative property

$$\mathcal{L}(\dot{z}) = sZ(s) - z(0)$$

to derive, integrate by parts:

$$\begin{aligned}\mathcal{L}(\dot{z})(s) &= \int_0^{\infty} e^{-st} \dot{z}(t) dt \\ &= e^{-st} z(t) \Big|_{t=0}^{t \rightarrow \infty} + s \int_0^{\infty} e^{-st} z(t) dt \\ &= sZ(s) - z(0)\end{aligned}$$

Laplace transform solution of $\dot{x} = Ax$

consider continuous-time time-invariant (TI) LDS

$$\dot{x} = Ax$$

for $t \geq 0$, where $x(t) \in \mathbf{R}^n$

- take Laplace transform: $sX(s) - x(0) = AX(s)$
- rewrite as $(sI - A)X(s) = x(0)$
- hence $X(s) = (sI - A)^{-1}x(0)$
- take inverse transform

$$x(t) = \mathcal{L}^{-1} \left((sI - A)^{-1} \right) x(0)$$

Resolvent and state transition matrix

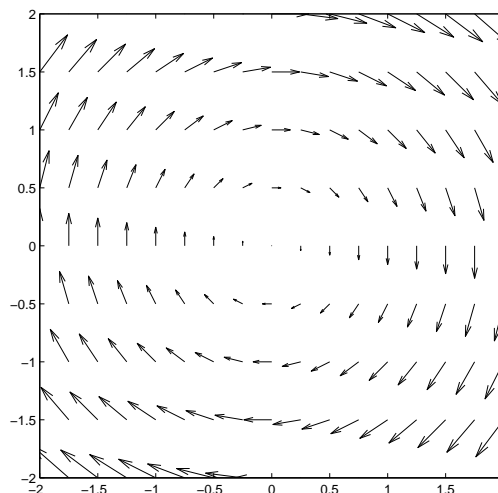
- $(sI - A)^{-1}$ is called the *resolvent* of A
- resolvent defined for $s \in \mathbf{C}$ except eigenvalues of A , *i.e.*, s such that $\det(sI - A) = 0$
- $\Phi(t) = \mathcal{L}^{-1}((sI - A)^{-1})$ is called the *state-transition matrix*; it maps the initial state to the state at time t :

$$x(t) = \Phi(t)x(0)$$

(in particular, state $x(t)$ is a linear function of initial state $x(0)$)

Example 1: Harmonic oscillator

$$\dot{x} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} x$$



$$sI - A = \begin{bmatrix} s & -1 \\ 1 & s \end{bmatrix}, \text{ so resolvent is}$$

$$(sI - A)^{-1} = \begin{bmatrix} \frac{s}{s^2+1} & \frac{1}{s^2+1} \\ \frac{-1}{s^2+1} & \frac{s}{s^2+1} \end{bmatrix}$$

(eigenvalues are $\pm j$)

state transition matrix is

$$\Phi(t) = \mathcal{L}^{-1} \left(\begin{bmatrix} \frac{s}{s^2+1} & \frac{1}{s^2+1} \\ \frac{-1}{s^2+1} & \frac{s}{s^2+1} \end{bmatrix} \right) = \begin{bmatrix} \cos t & \sin t \\ -\sin t & \cos t \end{bmatrix}$$

a rotation matrix ($-t$ radians)

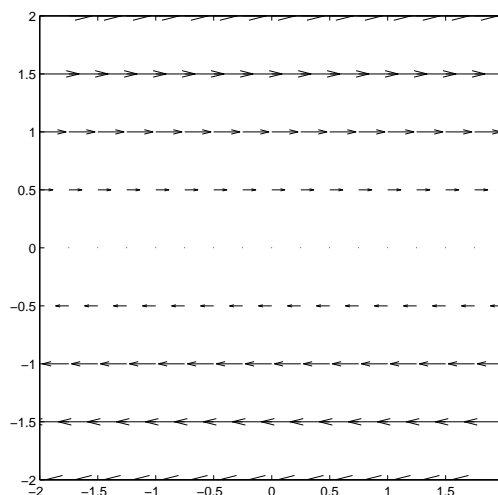
$$\text{so we have } x(t) = \begin{bmatrix} \cos t & \sin t \\ -\sin t & \cos t \end{bmatrix} x(0)$$

Solution via Laplace transform and matrix exponential

10-7

Example 2: Double integrator

$$\dot{x} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} x$$



Solution via Laplace transform and matrix exponential

10-8

$sI - A = \begin{bmatrix} s & -1 \\ 0 & s \end{bmatrix}$, so resolvent is

$$(sI - A)^{-1} = \begin{bmatrix} \frac{1}{s} & \frac{1}{s^2} \\ 0 & \frac{1}{s} \end{bmatrix}$$

(eigenvalues are 0, 0)

state transition matrix is

$$\Phi(t) = \mathcal{L}^{-1} \left(\begin{bmatrix} \frac{1}{s} & \frac{1}{s^2} \\ 0 & \frac{1}{s} \end{bmatrix} \right) = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix}$$

so we have $x(t) = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix} x(0)$

Characteristic polynomial

$\mathcal{X}(s) = \det(sI - A)$ is called the *characteristic polynomial* of A

- $\mathcal{X}(s)$ is a polynomial of degree n , with leading (*i.e.*, s^n) coefficient one
- roots of \mathcal{X} are the eigenvalues of A
- \mathcal{X} has real coefficients, so eigenvalues are either real or occur in conjugate pairs
- there are n eigenvalues (if we count multiplicity as roots of \mathcal{X})

Eigenvalues of A and poles of resolvent

i, j entry of resolvent can be expressed via Cramer's rule as

$$(-1)^{i+j} \frac{\det \Delta_{ij}}{\det(sI - A)}$$

where Δ_{ij} is $sI - A$ with j th row and i th column deleted

- $\det \Delta_{ij}$ is a polynomial of degree less than n , so i, j entry of resolvent has form $f_{ij}(s)/\mathcal{X}(s)$ where f_{ij} is polynomial with degree less than n
- poles of entries of resolvent must be eigenvalues of A
- but not all eigenvalues of A show up as poles of each entry (when there are cancellations between $\det \Delta_{ij}$ and $\mathcal{X}(s)$)

Solution via Laplace transform and matrix exponential

10-11

Matrix exponential

$$(I - C)^{-1} = I + C + C^2 + C^3 + \dots \text{ (if series converges)}$$

- series expansion of resolvent:

$$(sI - A)^{-1} = (1/s)(I - A/s)^{-1} = \frac{I}{s} + \frac{A}{s^2} + \frac{A^2}{s^3} + \dots$$

(valid for $|s|$ large enough) so

$$\Phi(t) = \mathcal{L}^{-1}((sI - A)^{-1}) = I + tA + \frac{(tA)^2}{2!} + \dots$$

Solution via Laplace transform and matrix exponential

10-12

- looks like ordinary power series

$$e^{at} = 1 + ta + \frac{(ta)^2}{2!} + \dots$$

with square matrices instead of scalars . . .

- define **matrix exponential** as

$$e^M = I + M + \frac{M^2}{2!} + \dots$$

for $M \in \mathbf{R}^{n \times n}$ (which in fact converges for all M)

- with this definition, state-transition matrix is

$$\Phi(t) = \mathcal{L}^{-1}((sI - A)^{-1}) = e^{tA}$$

Matrix exponential solution of autonomous LDS

solution of $\dot{x} = Ax$, with $A \in \mathbf{R}^{n \times n}$ and constant, is

$$x(t) = e^{tA}x(0)$$

generalizes scalar case: solution of $\dot{x} = ax$, with $a \in \mathbf{R}$ and constant, is

$$x(t) = e^{ta}x(0)$$

- matrix exponential is *meant* to look like scalar exponential
- some things you'd guess hold for the matrix exponential (by analogy with the scalar exponential) do in fact hold
- but **many things you'd guess are wrong**

example: you might guess that $e^{A+B} = e^A e^B$, but it's false (in general)

$$A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$e^A = \begin{bmatrix} 0.54 & 0.84 \\ -0.84 & 0.54 \end{bmatrix}, \quad e^B = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

$$e^{A+B} = \begin{bmatrix} 0.16 & 1.40 \\ -0.70 & 0.16 \end{bmatrix} \neq e^A e^B = \begin{bmatrix} 0.54 & 1.38 \\ -0.84 & -0.30 \end{bmatrix}$$

Solution via Laplace transform and matrix exponential

10-15

however, we do have $e^{A+B} = e^A e^B$ if $AB = BA$, *i.e.*, A and B commute

thus for $t, s \in \mathbf{R}$, $e^{(tA+sA)} = e^{tA} e^{sA}$

with $s = -t$ we get

$$e^{tA} e^{-tA} = e^{tA-tA} = e^0 = I$$

so e^{tA} is nonsingular, with inverse

$$(e^{tA})^{-1} = e^{-tA}$$

Solution via Laplace transform and matrix exponential

10-16

example: let's find e^A , where $A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$

we already found

$$e^{tA} = \mathcal{L}^{-1}(sI - A)^{-1} = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix}$$

so, plugging in $t = 1$, we get $e^A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$

let's check power series:

$$e^A = I + A + \frac{A^2}{2!} + \dots = I + A$$

since $A^2 = A^3 = \dots = 0$

Time transfer property

for $\dot{x} = Ax$ we know

$$x(t) = \Phi(t)x(0) = e^{tA}x(0)$$

interpretation: the matrix e^{tA} propagates initial condition into state at time t

more generally we have, for *any* t and τ ,

$$x(\tau + t) = e^{tA}x(\tau)$$

(to see this, apply result above to $z(t) = x(t + \tau)$)

interpretation: the matrix e^{tA} propagates state t seconds forward in time (backward if $t < 0$)

- recall first order (forward Euler) *approximate* state update, for small t :

$$x(\tau + t) \approx x(\tau) + t\dot{x}(\tau) = (I + tA)x(\tau)$$

- *exact* solution is

$$x(\tau + t) = e^{tA}x(\tau) = (I + tA + (tA)^2/2! + \dots)x(\tau)$$

- forward Euler is just first two terms in series

Sampling a continuous-time system

suppose $\dot{x} = Ax$

sample x at times $t_1 \leq t_2 \leq \dots$: define $z(k) = x(t_k)$

then $z(k+1) = e^{(t_{k+1}-t_k)A}z(k)$

for uniform sampling $t_{k+1} - t_k = h$, so

$$z(k+1) = e^{hA}z(k),$$

a discrete-time LDS (called *discretized version* of continuous-time system)

Piecewise constant system

consider *time-varying* LDS $\dot{x} = A(t)x$, with

$$A(t) = \begin{cases} A_0 & 0 \leq t < t_1 \\ A_1 & t_1 \leq t < t_2 \\ \vdots & \end{cases}$$

where $0 < t_1 < t_2 < \dots$ (sometimes called jump linear system)

for $t \in [t_i, t_{i+1}]$ we have

$$x(t) = e^{(t-t_i)A_i} \dots e^{(t_3-t_2)A_2} e^{(t_2-t_1)A_1} e^{t_1 A_0} x(0)$$

(matrix on righthand side is called state transition matrix for system, and denoted $\Phi(t)$)

Qualitative behavior of $x(t)$

suppose $\dot{x} = Ax$, $x(t) \in \mathbf{R}^n$

then $x(t) = e^{tA}x(0)$; $X(s) = (sI - A)^{-1}x(0)$

i th component $X_i(s)$ has form

$$X_i(s) = \frac{a_i(s)}{\mathcal{X}(s)}$$

where a_i is a polynomial of degree $< n$

thus the poles of X_i are all eigenvalues of A (but not necessarily the other way around)

first assume eigenvalues λ_i are distinct, so $X_i(s)$ cannot have repeated poles

then $x_i(t)$ has form

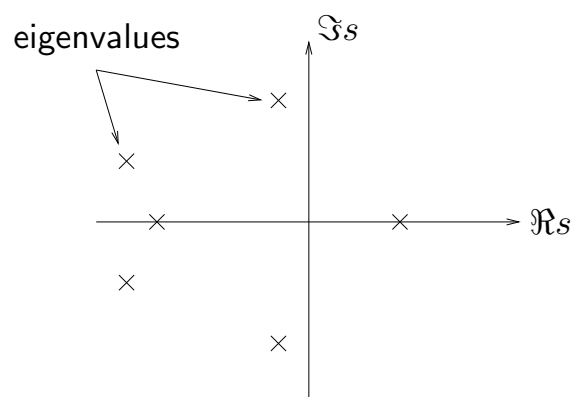
$$x_i(t) = \sum_{j=1}^n \beta_{ij} e^{\lambda_j t}$$

where β_{ij} depend on $x(0)$ (linearly)

eigenvalues determine (possible) qualitative behavior of x :

- eigenvalues give exponents that can occur in exponentials
- real eigenvalue λ corresponds to an exponentially decaying or growing term $e^{\lambda t}$ in solution
- complex eigenvalue $\lambda = \sigma + j\omega$ corresponds to decaying or growing sinusoidal term $e^{\sigma t} \cos(\omega t + \phi)$ in solution

- $\Re \lambda_j$ gives exponential growth rate (if > 0), or exponential decay rate (if < 0) of term
- $\Im \lambda_j$ gives frequency of oscillatory term (if $\neq 0$)



now suppose A has repeated eigenvalues, so X_i can have repeated poles

express eigenvalues as $\lambda_1, \dots, \lambda_r$ (distinct) with multiplicities n_1, \dots, n_r , respectively ($n_1 + \dots + n_r = n$)

then $x_i(t)$ has form

$$x_i(t) = \sum_{j=1}^r p_{ij}(t) e^{\lambda_j t}$$

where $p_{ij}(t)$ is a polynomial of degree $< n_j$ (that depends linearly on $x(0)$)

Stability

we say system $\dot{x} = Ax$ is *stable* if $e^{tA} \rightarrow 0$ as $t \rightarrow \infty$

meaning:

- state $x(t)$ converges to 0, as $t \rightarrow \infty$, no matter what $x(0)$ is
- all trajectories of $\dot{x} = Ax$ converge to 0 as $t \rightarrow \infty$

fact: $\dot{x} = Ax$ is stable if and only if all eigenvalues of A have negative real part:

$$\Re \lambda_i < 0, \quad i = 1, \dots, n$$

the 'if' part is clear since

$$\lim_{t \rightarrow \infty} p(t)e^{\lambda t} = 0$$

for any polynomial, if $\Re \lambda < 0$

we'll see the 'only if' part next lecture

more generally, $\max_i \Re \lambda_i$ determines the maximum asymptotic logarithmic growth rate of $x(t)$ (or decay, if < 0)

Lecture 11

Eigenvectors and diagonalization

- eigenvectors
- dynamic interpretation: invariant sets
- complex eigenvectors & invariant planes
- left eigenvectors
- diagonalization
- modal form
- discrete-time stability

11-1

Eigenvectors and eigenvalues

$\lambda \in \mathbf{C}$ is an *eigenvalue* of $A \in \mathbf{C}^{n \times n}$ if

$$\mathcal{X}(\lambda) = \det(\lambda I - A) = 0$$

equivalent to:

- there exists nonzero $v \in \mathbf{C}^n$ s.t. $(\lambda I - A)v = 0$, *i.e.*,

$$Av = \lambda v$$

any such v is called an *eigenvector* of A (associated with eigenvalue λ)

- there exists nonzero $w \in \mathbf{C}^n$ s.t. $w^T(\lambda I - A) = 0$, *i.e.*,

$$w^T A = \lambda w^T$$

any such w is called a *left eigenvector* of A

- if v is an eigenvector of A with eigenvalue λ , then so is αv , for any $\alpha \in \mathbf{C}$, $\alpha \neq 0$
- even when A is real, eigenvalue λ and eigenvector v can be complex
- when A and λ are real, we can always find a real eigenvector v associated with λ : if $Av = \lambda v$, with $A \in \mathbf{R}^{n \times n}$, $\lambda \in \mathbf{R}$, and $v \in \mathbf{C}^n$, then

$$A\Re v = \lambda\Re v, \quad A\Im v = \lambda\Im v$$

so $\Re v$ and $\Im v$ are real eigenvectors, if they are nonzero (and at least one is)

- *conjugate symmetry*: if A is real and $v \in \mathbf{C}^n$ is an eigenvector associated with $\lambda \in \mathbf{C}$, then \bar{v} is an eigenvector associated with $\bar{\lambda}$: taking conjugate of $Av = \lambda v$ we get $\overline{Av} = \overline{\lambda v}$, so

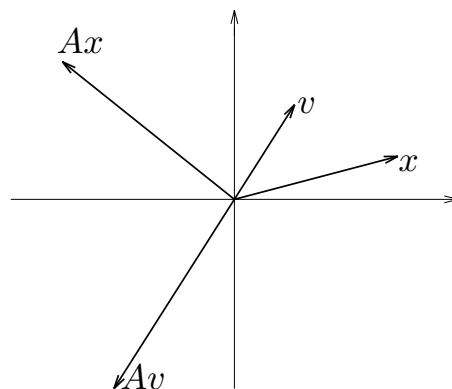
$$A\bar{v} = \bar{\lambda}\bar{v}$$

we'll assume A is real from now on . . .

Scaling interpretation

(assume $\lambda \in \mathbf{R}$ for now; we'll consider $\lambda \in \mathbf{C}$ later)

if v is an eigenvector, effect of A on v is very simple: scaling by λ



(what is λ here?)

- $\lambda \in \mathbf{R}, \lambda > 0$: v and Av point in same direction
- $\lambda \in \mathbf{R}, \lambda < 0$: v and Av point in opposite directions
- $\lambda \in \mathbf{R}, |\lambda| < 1$: Av smaller than v
- $\lambda \in \mathbf{R}, |\lambda| > 1$: Av larger than v

(we'll see later how this relates to stability of continuous- and discrete-time systems. . .)

Dynamic interpretation

suppose $Av = \lambda v, v \neq 0$

if $\dot{x} = Ax$ and $x(0) = v$, then $x(t) = e^{\lambda t}v$

several ways to see this, *e.g.*,

$$\begin{aligned}
 x(t) = e^{tA}v &= \left(I + tA + \frac{(tA)^2}{2!} + \cdots \right) v \\
 &= v + \lambda t v + \frac{(\lambda t)^2}{2!} v + \cdots \\
 &= e^{\lambda t} v
 \end{aligned}$$

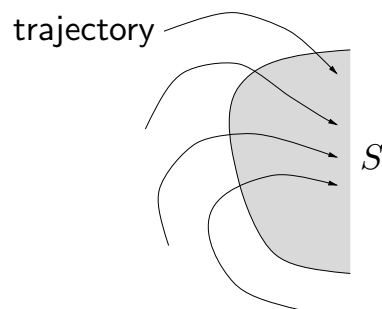
(since $(tA)^k v = (\lambda t)^k v$)

- for $\lambda \in \mathbf{C}$, solution is complex (we'll interpret later); for now, assume $\lambda \in \mathbf{R}$
 - if initial state is an eigenvector v , resulting motion is very simple — always on the line spanned by v
 - solution $x(t) = e^{\lambda t}v$ is called *mode* of system $\dot{x} = Ax$ (associated with eigenvalue λ)
-
- for $\lambda \in \mathbf{R}$, $\lambda < 0$, mode contracts or shrinks as $t \uparrow$
 - for $\lambda \in \mathbf{R}$, $\lambda > 0$, mode expands or grows as $t \uparrow$

Invariant sets

a set $S \subseteq \mathbf{R}^n$ is *invariant* under $\dot{x} = Ax$ if whenever $x(t) \in S$, then $x(\tau) \in S$ for all $\tau \geq t$

i.e.: once trajectory enters S , it stays in S



vector field interpretation: trajectories only cut *into* S , never out

suppose $Av = \lambda v$, $v \neq 0$, $\lambda \in \mathbf{R}$

- line $\{ tv \mid t \in \mathbf{R} \}$ is invariant
(in fact, ray $\{ tv \mid t > 0 \}$ is invariant)
- if $\lambda < 0$, line segment $\{ tv \mid 0 \leq t \leq a \}$ is invariant

Complex eigenvectors

suppose $Av = \lambda v$, $v \neq 0$, λ is complex

for $a \in \mathbf{C}$, (complex) trajectory $ae^{\lambda t}v$ satisfies $\dot{x} = Ax$

hence so does (real) trajectory

$$\begin{aligned} x(t) &= \Re(ae^{\lambda t}v) \\ &= e^{\sigma t} \begin{bmatrix} v_{\text{re}} & v_{\text{im}} \end{bmatrix} \begin{bmatrix} \cos \omega t & \sin \omega t \\ -\sin \omega t & \cos \omega t \end{bmatrix} \begin{bmatrix} \alpha \\ -\beta \end{bmatrix} \end{aligned}$$

where

$$v = v_{\text{re}} + jv_{\text{im}}, \quad \lambda = \sigma + j\omega, \quad a = \alpha + j\beta$$

- trajectory stays in *invariant plane* $\text{span}\{v_{\text{re}}, v_{\text{im}}\}$
- σ gives logarithmic growth/decay factor
- ω gives angular velocity of rotation in plane

Dynamic interpretation: left eigenvectors

suppose $w^T A = \lambda w^T$, $w \neq 0$

then

$$\frac{d}{dt}(w^T x) = w^T \dot{x} = w^T A x = \lambda(w^T x)$$

i.e., $w^T x$ satisfies the DE $d(w^T x)/dt = \lambda(w^T x)$

hence $w^T x(t) = e^{\lambda t} w^T x(0)$

- even if trajectory x is complicated, $w^T x$ is simple
- if, e.g., $\lambda \in \mathbf{R}$, $\lambda < 0$, halfspace $\{ z \mid w^T z \leq a \}$ is invariant (for $a \geq 0$)
- for $\lambda = \sigma + j\omega \in \mathbf{C}$, $(\Re w)^T x$ and $(\Im w)^T x$ both have form

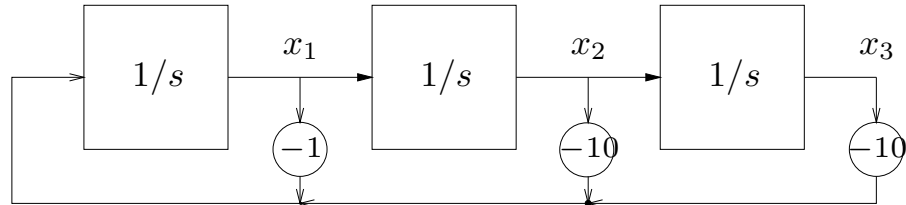
$$e^{\sigma t} (\alpha \cos(\omega t) + \beta \sin(\omega t))$$

Summary

- *right eigenvectors* are initial conditions from which resulting motion is simple (i.e., remains on line or in plane)
- *left eigenvectors* give linear functions of state that are simple, for any initial condition

example 1: $\dot{x} = \begin{bmatrix} -1 & -10 & -10 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} x$

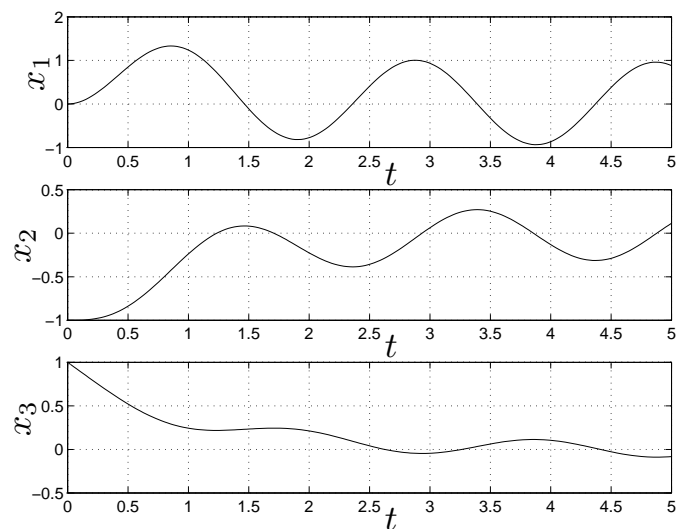
block diagram:



$$\mathcal{X}(s) = s^3 + s^2 + 10s + 10 = (s + 1)(s^2 + 10)$$

eigenvalues are $-1, \pm j\sqrt{10}$

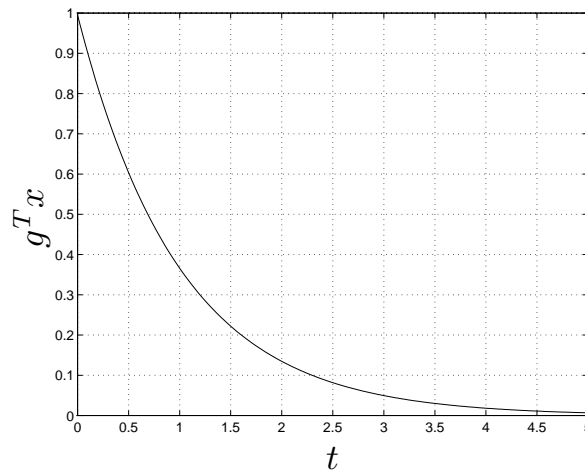
trajectory with $x(0) = (0, -1, 1)$:



left eigenvector associated with eigenvalue -1 is

$$g = \begin{bmatrix} 0.1 \\ 0 \\ 1 \end{bmatrix}$$

let's check $g^T x(t)$ when $x(0) = (0, -1, 1)$ (as above):



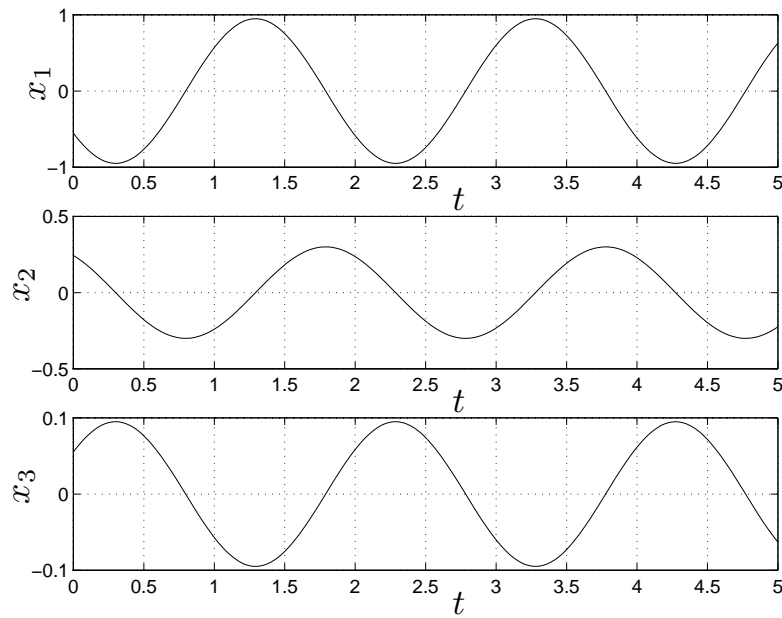
eigenvector associated with eigenvalue $j\sqrt{10}$ is

$$v = \begin{bmatrix} -0.554 + j0.771 \\ 0.244 + j0.175 \\ 0.055 - j0.077 \end{bmatrix}$$

so an invariant plane is spanned by

$$v_{\text{re}} = \begin{bmatrix} -0.554 \\ 0.244 \\ 0.055 \end{bmatrix}, \quad v_{\text{im}} = \begin{bmatrix} 0.771 \\ 0.175 \\ -0.077 \end{bmatrix}$$

for example, with $x(0) = v_{\text{re}}$ we have



Example 2: Markov chain

probability distribution satisfies $p(t+1) = Pp(t)$

$p_i(t) = \mathbf{Prob}(z(t) = i)$ so $\sum_{i=1}^n p_i(t) = 1$

$P_{ij} = \mathbf{Prob}(z(t+1) = i \mid z(t) = j)$, so $\sum_{i=1}^n P_{ij} = 1$
(such matrices are called *stochastic*)

rewrite as:

$$[1 \ 1 \ \cdots \ 1]P = [1 \ 1 \ \cdots \ 1]$$

i.e., $[1 \ 1 \ \cdots \ 1]$ is a left eigenvector of P with e.v. 1

hence $\det(I - P) = 0$, so there is a right eigenvector $v \neq 0$ with $Pv = v$

it can be shown that v can be chosen so that $v_i \geq 0$, hence we can
normalize v so that $\sum_{i=1}^n v_i = 1$

interpretation: v is an *equilibrium distribution*; *i.e.*, if $p(0) = v$ then
 $p(t) = v$ for all $t \geq 0$

(if v is unique it is called the *steady-state distribution* of the Markov chain)

Diagonalization

suppose v_1, \dots, v_n is a *linearly independent* set of eigenvectors of $A \in \mathbf{R}^{n \times n}$:

$$Av_i = \lambda_i v_i, \quad i = 1, \dots, n$$

express as

$$A \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

define $T = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix}$ and $\Lambda = \mathbf{diag}(\lambda_1, \dots, \lambda_n)$, so

$$AT = T\Lambda$$

and finally

$$T^{-1}AT = \Lambda$$

- T invertible since v_1, \dots, v_n linearly independent
- similarity transformation by T diagonalizes A

conversely if there is a $T = [v_1 \cdots v_n]$ s.t.

$$T^{-1}AT = \Lambda = \mathbf{diag}(\lambda_1, \dots, \lambda_n)$$

then $AT = T\Lambda$, i.e.,

$$Av_i = \lambda_i v_i, \quad i = 1, \dots, n$$

so v_1, \dots, v_n is a linearly independent set of eigenvectors of A

we say A is *diagonalizable* if

- there exists T s.t. $T^{-1}AT = \Lambda$ is diagonal
- A has a set of linearly independent eigenvectors

(if A is not diagonalizable, it is sometimes called *defective*)

Not all matrices are diagonalizable

example: $A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$

characteristic polynomial is $\mathcal{X}(s) = s^2$, so $\lambda = 0$ is only eigenvalue

eigenvectors satisfy $Av = 0v = 0$, *i.e.*

$$\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

so all eigenvectors have form $v = \begin{bmatrix} v_1 \\ 0 \end{bmatrix}$ where $v_1 \neq 0$

thus, A cannot have two independent eigenvectors

Distinct eigenvalues

fact: if A has distinct eigenvalues, *i.e.*, $\lambda_i \neq \lambda_j$ for $i \neq j$, then A is diagonalizable

(the converse is false — A can have repeated eigenvalues but still be diagonalizable)

Diagonalization and left eigenvectors

rewrite $T^{-1}AT = \Lambda$ as $T^{-1}A = \Lambda T^{-1}$, or

$$\begin{bmatrix} w_1^T \\ \vdots \\ w_n^T \end{bmatrix} A = \Lambda \begin{bmatrix} w_1^T \\ \vdots \\ w_n^T \end{bmatrix}$$

where w_1^T, \dots, w_n^T are the rows of T^{-1}

thus

$$w_i^T A = \lambda_i w_i^T$$

i.e., the rows of T^{-1} are (lin. indep.) left eigenvectors, normalized so that

$$w_i^T v_j = \delta_{ij}$$

(*i.e.*, left & right eigenvectors chosen this way are *dual bases*)

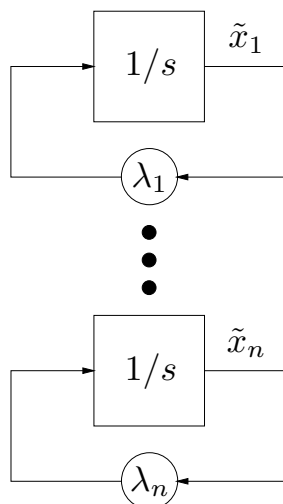
Modal form

suppose A is diagonalizable by T

define new coordinates by $x = T\tilde{x}$, so

$$T\dot{\tilde{x}} = AT\tilde{x} \quad \Leftrightarrow \quad \dot{\tilde{x}} = T^{-1}AT\tilde{x} \quad \Leftrightarrow \quad \dot{\tilde{x}} = \Lambda\tilde{x}$$

in new coordinate system, system is diagonal (decoupled):



trajectories consist of n independent modes, *i.e.*,

$$\tilde{x}_i(t) = e^{\lambda_i t} \tilde{x}_i(0)$$

hence the name *modal form*

Real modal form

when eigenvalues (hence T) are complex, system can be put in *real modal form*:

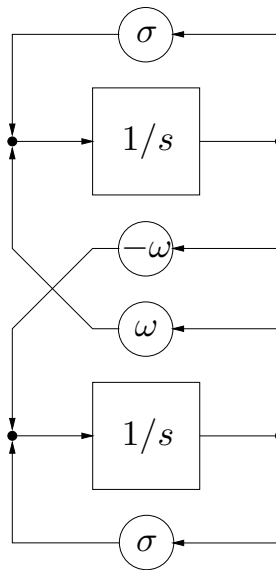
$$S^{-1}AS = \mathbf{diag}(\Lambda_r, M_{r+1}, M_{r+3}, \dots, M_{n-1})$$

where $\Lambda_r = \mathbf{diag}(\lambda_1, \dots, \lambda_r)$ are the real eigenvalues, and

$$M_i = \begin{bmatrix} \sigma_i & \omega_i \\ -\omega_i & \sigma_i \end{bmatrix}, \quad \lambda_i = \sigma_i + j\omega_i, \quad i = r+1, r+3, \dots, n$$

where λ_i are the complex eigenvalues (one from each conjugate pair)

block diagram of 'complex mode':



diagonalization simplifies many matrix expressions

e.g., resolvent:

$$\begin{aligned}
 (sI - A)^{-1} &= (sTT^{-1} - T\Lambda T^{-1})^{-1} \\
 &= (T(sI - \Lambda)T^{-1})^{-1} \\
 &= T(sI - \Lambda)^{-1}T^{-1} \\
 &= T \mathbf{diag} \left(\frac{1}{s - \lambda_1}, \dots, \frac{1}{s - \lambda_n} \right) T^{-1}
 \end{aligned}$$

powers (*i.e.*, discrete-time solution):

$$\begin{aligned}
 A^k &= (T\Lambda T^{-1})^k \\
 &= (T\Lambda T^{-1}) \dots (T\Lambda T^{-1}) \\
 &= T\Lambda^k T^{-1} \\
 &= T \mathbf{diag}(\lambda_1^k, \dots, \lambda_n^k) T^{-1}
 \end{aligned}$$

(for $k < 0$ only if A invertible, *i.e.*, all $\lambda_i \neq 0$)

exponential (*i.e.*, continuous-time solution):

$$\begin{aligned}
 e^A &= I + A + A^2/2! + \dots \\
 &= I + T\Lambda T^{-1} + (T\Lambda T^{-1})^2/2! + \dots \\
 &= T(I + \Lambda + \Lambda^2/2! + \dots)T^{-1} \\
 &= Te^{\Lambda}T^{-1} \\
 &= T \mathbf{diag}(e^{\lambda_1}, \dots, e^{\lambda_n})T^{-1}
 \end{aligned}$$

Analytic function of a matrix

for any analytic function $f : \mathbf{R} \rightarrow \mathbf{R}$, *i.e.*, given by power series

$$f(a) = \beta_0 + \beta_1 a + \beta_2 a^2 + \beta_3 a^3 + \dots$$

we can define $f(A)$ for $A \in \mathbf{R}^{n \times n}$ (*i.e.*, overload f) as

$$f(A) = \beta_0 I + \beta_1 A + \beta_2 A^2 + \beta_3 A^3 + \dots$$

substituting $A = T\Lambda T^{-1}$, we have

$$\begin{aligned}
 f(A) &= \beta_0 I + \beta_1 A + \beta_2 A^2 + \beta_3 A^3 + \dots \\
 &= \beta_0 T T^{-1} + \beta_1 T \Lambda T^{-1} + \beta_2 (T \Lambda T^{-1})^2 + \dots \\
 &= T (\beta_0 I + \beta_1 \Lambda + \beta_2 \Lambda^2 + \dots) T^{-1} \\
 &= T \mathbf{diag}(f(\lambda_1), \dots, f(\lambda_n)) T^{-1}
 \end{aligned}$$

Solution via diagonalization

assume A is diagonalizable

consider LDS $\dot{x} = Ax$, with $T^{-1}AT = \Lambda$

then

$$\begin{aligned}x(t) &= e^{tA}x(0) \\&= Te^{\Lambda t}T^{-1}x(0) \\&= \sum_{i=1}^n e^{\lambda_i t} (w_i^T x(0)) v_i\end{aligned}$$

thus: any trajectory can be expressed as linear combination of modes

interpretation:

- (left eigenvectors) decompose initial state $x(0)$ into modal components $w_i^T x(0)$
- $e^{\lambda_i t}$ term propagates i th mode forward t seconds
- reconstruct state as linear combination of (right) eigenvectors

application: for what $x(0)$ do we have $x(t) \rightarrow 0$ as $t \rightarrow \infty$?

divide eigenvalues into those with negative real parts

$$\Re \lambda_1 < 0, \dots, \Re \lambda_s < 0,$$

and the others,

$$\Re \lambda_{s+1} \geq 0, \dots, \Re \lambda_n \geq 0$$

from

$$x(t) = \sum_{i=1}^n e^{\lambda_i t} (w_i^T x(0)) v_i$$

condition for $x(t) \rightarrow 0$ is:

$$x(0) \in \text{span}\{v_1, \dots, v_s\},$$

or equivalently,

$$w_i^T x(0) = 0, \quad i = s+1, \dots, n$$

(can you prove this?)

Stability of discrete-time systems

suppose A diagonalizable

consider discrete-time LDS $x(t+1) = Ax(t)$

if $A = T\Lambda T^{-1}$, then $A^k = T\Lambda^k T^{-1}$

then

$$x(t) = A^t x(0) = \sum_{i=1}^n \lambda_i^t (w_i^T x(0)) v_i \rightarrow 0 \quad \text{as } t \rightarrow \infty$$

for all $x(0)$ if and only if

$$|\lambda_i| < 1, \quad i = 1, \dots, n.$$

we will see later that this is true even when A is not diagonalizable, so we have

fact: $x(t+1) = Ax(t)$ is stable if and only if all eigenvalues of A have magnitude less than one

Lecture 12

Jordan canonical form

- Jordan canonical form
- generalized modes
- Cayley-Hamilton theorem

12-1

Jordan canonical form

what if A cannot be diagonalized?

any matrix $A \in \mathbf{R}^{n \times n}$ can be put in *Jordan canonical form* by a similarity transformation, *i.e.*

$$T^{-1}AT = J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_q \end{bmatrix}$$

where

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix} \in \mathbf{C}^{n_i \times n_i}$$

is called a *Jordan block* of size n_i with eigenvalue λ_i (so $n = \sum_{i=1}^q n_i$)

- J is upper bidiagonal
- J diagonal is the special case of n Jordan blocks of size $n_i = 1$
- Jordan form is unique (up to permutations of the blocks)
- can have multiple blocks with same eigenvalue

note: JCF is a *conceptual tool*, never used in numerical computations!

$$\chi(s) = \det(sI - A) = (s - \lambda_1)^{n_1} \cdots (s - \lambda_q)^{n_q}$$

hence distinct eigenvalues $\Rightarrow n_i = 1 \Rightarrow A$ diagonalizable

$\dim \mathcal{N}(\lambda I - A)$ is the number of Jordan blocks with eigenvalue λ

more generally,

$$\dim \mathcal{N}(\lambda I - A)^k = \sum_{\lambda_i = \lambda} \min\{k, n_i\}$$

so from $\dim \mathcal{N}(\lambda I - A)^k$ for $k = 1, 2, \dots$ we can determine the sizes of the Jordan blocks associated with λ

- factor out T and T^{-1} , $\lambda I - A = T(\lambda I - J)T^{-1}$

- for, say, a block of size 3:

$$\lambda_i I - J_i = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix} \quad (\lambda_i I - J_i)^2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (\lambda_i I - J_i)^3 = 0$$

- for other blocks (say, size 3, for $k \geq 2$)

$$(\lambda_i I - J_j)^k = \begin{bmatrix} (\lambda_i - \lambda_j)^k & -k(\lambda_i - \lambda_j)^{k-1} & (k(k-1)/2)(\lambda_i - \lambda_j)^{k-2} \\ 0 & (\lambda_j - \lambda_i)^k & -k(\lambda_j - \lambda_i)^{k-1} \\ 0 & 0 & (\lambda_j - \lambda_i)^k \end{bmatrix}$$

Generalized eigenvectors

suppose $T^{-1}AT = J = \text{diag}(J_1, \dots, J_q)$

express T as

$$T = [T_1 \ T_2 \ \cdots \ T_q]$$

where $T_i \in \mathbf{C}^{n \times n_i}$ are the columns of T associated with i th Jordan block J_i

we have $AT_i = T_i J_i$

let $T_i = [v_{i1} \ v_{i2} \ \cdots \ v_{in_i}]$

then we have:

$$Av_{i1} = \lambda_i v_{i1},$$

i.e., the first column of each T_i is an eigenvector associated with e.v. λ_i

for $j = 2, \dots, n_i$,

$$Av_{ij} = v_{i,j-1} + \lambda_i v_{ij}$$

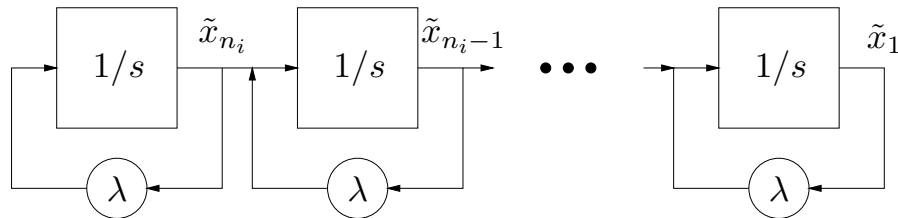
the vectors v_{i1}, \dots, v_{in_i} are sometimes called *generalized eigenvectors*

Jordan form LDS

consider LDS $\dot{x} = Ax$

by change of coordinates $x = T\tilde{x}$, can put into form $\dot{\tilde{x}} = J\tilde{x}$

system is decomposed into independent 'Jordan block systems' $\dot{\tilde{x}}_i = J_i \tilde{x}_i$



Jordan blocks are sometimes called Jordan chains

(block diagram shows why)

Resolvent, exponential of Jordan block

resolvent of $k \times k$ Jordan block with eigenvalue λ :

$$\begin{aligned}
 (sI - J_\lambda)^{-1} &= \begin{bmatrix} s - \lambda & -1 & & \\ & s - \lambda & \ddots & \\ & & \ddots & -1 \\ & & & s - \lambda \end{bmatrix}^{-1} \\
 &= \begin{bmatrix} (s - \lambda)^{-1} & (s - \lambda)^{-2} & \cdots & (s - \lambda)^{-k} \\ & (s - \lambda)^{-1} & \cdots & (s - \lambda)^{-k+1} \\ & & \ddots & \vdots \\ & & & (s - \lambda)^{-1} \end{bmatrix} \\
 &= (s - \lambda)^{-1}I + (s - \lambda)^{-2}F_1 + \cdots + (s - \lambda)^{-k}F_{k-1}
 \end{aligned}$$

where F_i is the matrix with ones on the i th upper diagonal

by inverse Laplace transform, exponential is:

$$\begin{aligned}
 e^{tJ\lambda} &= e^{t\lambda} (I + tF_1 + \cdots + (t^{k-1}/(k-1)!)F_{k-1}) \\
 &= e^{t\lambda} \begin{bmatrix} 1 & t & \cdots & t^{k-1}/(k-1)! \\ & 1 & \cdots & t^{k-2}/(k-2)! \\ & & \ddots & \vdots \\ & & & 1 \end{bmatrix}
 \end{aligned}$$

Jordan blocks yield:

- repeated poles in resolvent
- terms of form $t^p e^{t\lambda}$ in e^{tA}

Generalized modes

consider $\dot{x} = Ax$, with

$$x(0) = a_1 v_{i1} + \cdots + a_{n_i} v_{in_i} = T_i a$$

then $x(t) = T e^{Jt} \tilde{x}(0) = T_i e^{J_i t} a$

- trajectory stays in span of generalized eigenvectors
- coefficients have form $p(t)e^{\lambda t}$, where p is polynomial
- such solutions are called *generalized modes* of the system

with general $x(0)$ we can write

$$x(t) = e^{tA}x(0) = Te^{tJ}T^{-1}x(0) = \sum_{i=1}^q T_i e^{tJ_i} (S_i^T x(0))$$

where

$$T^{-1} = \begin{bmatrix} S_1^T \\ \vdots \\ S_q^T \end{bmatrix}$$

hence: all solutions of $\dot{x} = Ax$ are linear combinations of (generalized) modes

Cayley-Hamilton theorem

if $p(s) = a_0 + a_1s + \cdots + a_k s^k$ is a polynomial and $A \in \mathbf{R}^{n \times n}$, we define

$$p(A) = a_0 I + a_1 A + \cdots + a_k A^k$$

Cayley-Hamilton theorem: for any $A \in \mathbf{R}^{n \times n}$ we have $\mathcal{X}(A) = 0$, where $\mathcal{X}(s) = \det(sI - A)$

example: with $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ we have $\mathcal{X}(s) = s^2 - 5s - 2$, so

$$\begin{aligned} \mathcal{X}(A) &= A^2 - 5A - 2I \\ &= \begin{bmatrix} 7 & 10 \\ 15 & 22 \end{bmatrix} - 5 \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} - 2I \\ &= 0 \end{aligned}$$

corollary: for every $p \in \mathbf{Z}_+$, we have

$$A^p \in \text{span} \{ I, A, A^2, \dots, A^{n-1} \}$$

(and if A is invertible, also for $p \in \mathbf{Z}$)

i.e., every power of A can be expressed as linear combination of I, A, \dots, A^{n-1}

proof: divide $\mathcal{X}(s)$ into s^p to get $s^p = q(s)\mathcal{X}(s) + r(s)$

$r = \alpha_0 + \alpha_1 s + \dots + \alpha_{n-1} s^{n-1}$ is remainder polynomial

then

$$A^p = q(A)\mathcal{X}(A) + r(A) = r(A) = \alpha_0 I + \alpha_1 A + \dots + \alpha_{n-1} A^{n-1}$$

for $p = -1$: rewrite C-H theorem

$$\mathcal{X}(A) = A^n + a_{n-1}A^{n-1} + \dots + a_0 I = 0$$

as

$$I = A \left(-(a_1/a_0)I - (a_2/a_0)A - \dots - (1/a_0)A^{n-1} \right)$$

(A is invertible $\Leftrightarrow a_0 \neq 0$) so

$$A^{-1} = -(a_1/a_0)I - (a_2/a_0)A - \dots - (1/a_0)A^{n-1}$$

i.e., inverse is linear combination of A^k , $k = 0, \dots, n-1$

Proof of C-H theorem

first assume A is diagonalizable: $T^{-1}AT = \Lambda$

$$\mathcal{X}(s) = (s - \lambda_1) \cdots (s - \lambda_n)$$

since

$$\mathcal{X}(A) = \mathcal{X}(T\Lambda T^{-1}) = T\mathcal{X}(\Lambda)T^{-1}$$

it suffices to show $\mathcal{X}(\Lambda) = 0$

$$\begin{aligned}\mathcal{X}(\Lambda) &= (\Lambda - \lambda_1 I) \cdots (\Lambda - \lambda_n I) \\ &= \mathbf{diag}(0, \lambda_2 - \lambda_1, \dots, \lambda_n - \lambda_1) \cdots \mathbf{diag}(\lambda_1 - \lambda_n, \dots, \lambda_{n-1} - \lambda_n, 0) \\ &= 0\end{aligned}$$

now let's do general case: $T^{-1}AT = J$

$$\mathcal{X}(s) = (s - \lambda_1)^{n_1} \cdots (s - \lambda_q)^{n_q}$$

suffices to show $\mathcal{X}(J_i) = 0$

$$\mathcal{X}(J_i) = (J_i - \lambda_1 I)^{n_1} \cdots \underbrace{\begin{bmatrix} 0 & 1 & 0 & \cdots \\ 0 & 0 & 1 & \cdots \\ & & \ddots & \ddots \end{bmatrix}}_{(J_i - \lambda_i I)^{n_i}} \cdots (J_i - \lambda_q I)^{n_q} = 0$$

Lecture 13

Linear dynamical systems with inputs & outputs

- inputs & outputs: interpretations
- transfer matrix
- impulse and step matrices
- examples

13-1

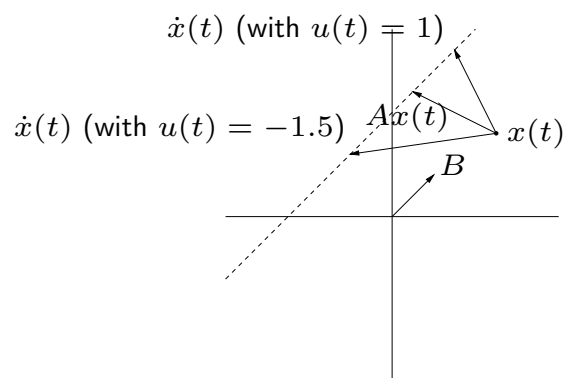
Inputs & outputs

recall continuous-time time-invariant LDS has form

$$\dot{x} = Ax + Bu, \quad y = Cx + Du$$

- Ax is called the *drift term* (of \dot{x})
- Bu is called the *input term* (of \dot{x})

picture, with $B \in \mathbf{R}^{2 \times 1}$:



Interpretations

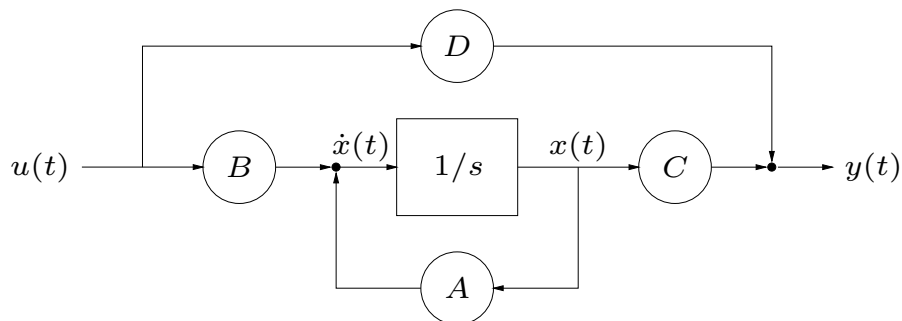
write $\dot{x} = Ax + b_1u_1 + \cdots + b_mu_m$, where $B = [b_1 \cdots b_m]$

- state derivative is sum of autonomous term (Ax) and one term per input (b_iu_i)
- each input u_i gives another degree of freedom for \dot{x} (assuming columns of B independent)

write $\dot{x} = Ax + Bu$ as $\dot{x}_i = \tilde{a}_i^T x + \tilde{b}_i^T u$, where $\tilde{a}_i^T, \tilde{b}_i^T$ are the rows of A, B

- i th state derivative is linear function of state x and input u

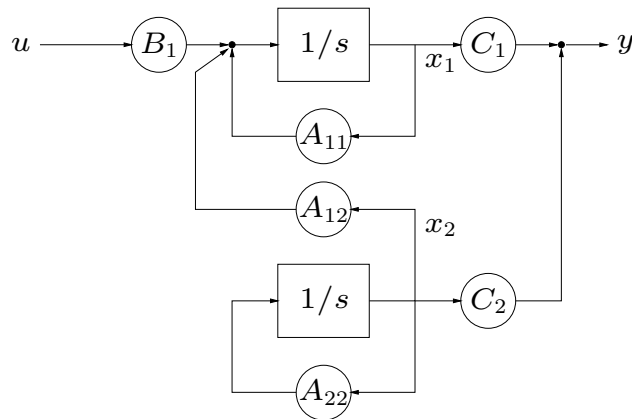
Block diagram



- A_{ij} is gain factor from state x_j into integrator i
- B_{ij} is gain factor from input u_j into integrator i
- C_{ij} is gain factor from state x_j into output y_i
- D_{ij} is gain factor from input u_j into output y_i

interesting when there is structure, *e.g.*, with $x_1 \in \mathbf{R}^{n_1}$, $x_2 \in \mathbf{R}^{n_2}$:

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} B_1 \\ 0 \end{bmatrix} u, \quad y = \begin{bmatrix} C_1 & C_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



- x_2 is not affected by input u , *i.e.*, x_2 propagates autonomously
- x_2 affects y directly and through x_1

Transfer matrix

take Laplace transform of $\dot{x} = Ax + Bu$:

$$sX(s) - x(0) = AX(s) + BU(s)$$

hence

$$X(s) = (sI - A)^{-1}x(0) + (sI - A)^{-1}BU(s)$$

so

$$x(t) = e^{tA}x(0) + \int_0^t e^{(t-\tau)A}Bu(\tau) \, d\tau$$

- $e^{tA}x(0)$ is the unforced or autonomous response
- $e^{tA}B$ is called the input-to-state impulse matrix
- $(sI - A)^{-1}B$ is called the *input-to-state transfer matrix* or *transfer function*

with $y = Cx + Du$ we have:

$$Y(s) = C(sI - A)^{-1}x(0) + (C(sI - A)^{-1}B + D)U(s)$$

so

$$y(t) = Ce^{tA}x(0) + \int_0^t Ce^{(t-\tau)A}Bu(\tau) d\tau + Du(t)$$

- output term $Ce^{tA}x(0)$ due to initial condition
- $H(s) = C(sI - A)^{-1}B + D$ is called the *transfer function* or *transfer matrix*
- $h(t) = Ce^{tA}B + D\delta(t)$ is called the *impulse matrix* or *impulse response* (δ is the Dirac delta function)

with zero initial condition we have:

$$Y(s) = H(s)U(s), \quad y = h * u$$

where $*$ is convolution (of matrix valued functions)

intepretation:

- H_{ij} is transfer function from input u_j to output y_i

Impulse matrix

impulse matrix $h(t) = Ce^{tA}B + D\delta(t)$

with $x(0) = 0$, $y = h * u$, i.e.,

$$y_i(t) = \sum_{j=1}^m \int_0^t h_{ij}(t - \tau) u_j(\tau) d\tau$$

interpretations:

- $h_{ij}(t)$ is impulse response from j th input to i th output
- $h_{ij}(t)$ gives y_i when $u(t) = e_j\delta$
- $h_{ij}(\tau)$ shows how dependent output i is, on what input j was, τ seconds ago
- i indexes output; j indexes input; τ indexes time lag

Step matrix

the *step matrix* or *step response matrix* is given by

$$s(t) = \int_0^t h(\tau) d\tau$$

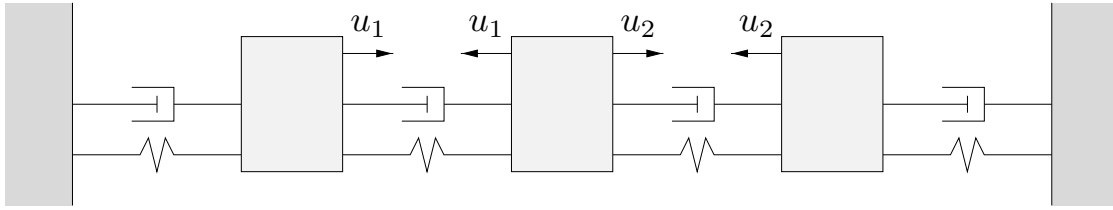
interpretations:

- $s_{ij}(t)$ is step response from j th input to i th output
- $s_{ij}(t)$ gives y_i when $u = e_j$ for $t \geq 0$

for invertible A , we have

$$s(t) = CA^{-1} (e^{tA} - I) B + D$$

Example 1



- unit masses, springs, dampers
- u_1 is tension between 1st & 2nd masses
- u_2 is tension between 2nd & 3rd masses
- $y \in \mathbf{R}^3$ is displacement of masses 1,2,3
- $x = \begin{bmatrix} y \\ \dot{y} \end{bmatrix}$

Linear dynamical systems with inputs & outputs

13-11

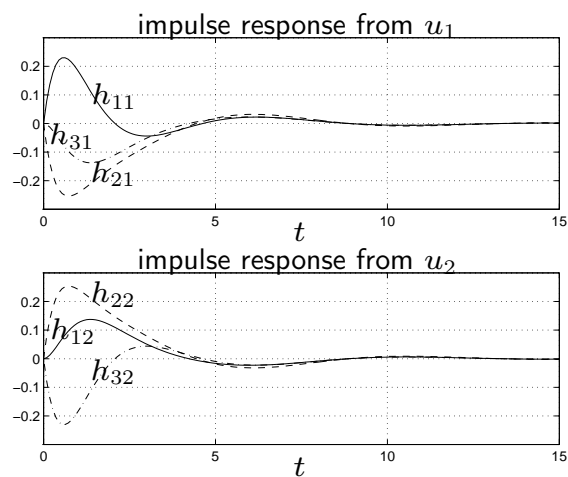
system is:

$$\dot{x} = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ -2 & 1 & 0 & -2 & 1 & 0 \\ 1 & -2 & 1 & 1 & -2 & 1 \\ 0 & 1 & -2 & 0 & 1 & -2 \end{bmatrix} x + \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ -1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

eigenvalues of A are

$$-1.71 \pm j0.71, \quad -1.00 \pm j1.00, \quad -0.29 \pm j0.71$$

impulse matrix:

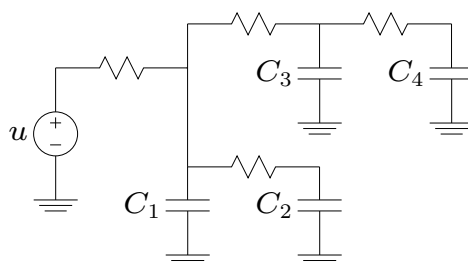


roughly speaking:

- impulse at u_1 affects third mass less than other two
- impulse at u_2 affects first mass later than other two

Example 2

interconnect circuit:



- $u(t) \in \mathbf{R}$ is input (drive) voltage
- x_i is voltage across C_i
- output is state: $y = x$
- unit resistors, unit capacitors
- step response matrix shows delay to each node

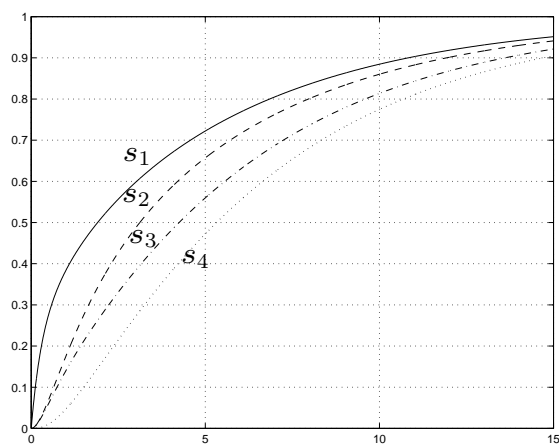
system is

$$\dot{x} = \begin{bmatrix} -3 & 1 & 1 & 0 \\ 1 & -1 & 0 & 0 \\ 1 & 0 & -2 & 1 \\ 0 & 0 & 1 & -1 \end{bmatrix} x + \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} u, \quad y = x$$

eigenvalues of A are

$$-0.17, \quad -0.66, \quad -2.21, \quad -3.96$$

step response matrix $s(t) \in \mathbf{R}^{4 \times 1}$:



- shortest delay to x_1 ; longest delay to x_4
- delays ≈ 10 , consistent with slowest (*i.e.*, dominant) eigenvalue -0.17

DC or static gain matrix

- transfer matrix at $s = 0$ is $H(0) = -CA^{-1}B + D \in \mathbf{R}^{m \times p}$
- DC transfer matrix describes system under *static* conditions, *i.e.*, x , u , y constant:

$$0 = \dot{x} = Ax + Bu, \quad y = Cx + Du$$

eliminate x to get $y = H(0)u$

- if system is stable,

$$H(0) = \int_0^\infty h(t) dt = \lim_{t \rightarrow \infty} s(t)$$

$$(\text{recall: } H(s) = \int_0^\infty e^{-st} h(t) dt, \quad s(t) = \int_0^t h(\tau) d\tau)$$

if $u(t) \rightarrow u_\infty \in \mathbf{R}^m$, then $y(t) \rightarrow y_\infty \in \mathbf{R}^p$ where $y_\infty = H(0)u_\infty$

DC gain matrix for example 1 (springs):

$$H(0) = \begin{bmatrix} 1/4 & 1/4 \\ -1/2 & 1/2 \\ -1/4 & -1/4 \end{bmatrix}$$

DC gain matrix for example 2 (RC circuit):

$$H(0) = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

(do these make sense?)

Discretization with piecewise constant inputs

linear system $\dot{x} = Ax + Bu$, $y = Cx + Du$

suppose $u_d : \mathbf{Z}_+ \rightarrow \mathbf{R}^m$ is a sequence, and

$$u(t) = u_d(k) \quad \text{for } kh \leq t < (k+1)h, \quad k = 0, 1, \dots$$

define sequences

$$x_d(k) = x(kh), \quad y_d(k) = y(kh), \quad k = 0, 1, \dots$$

- $h > 0$ is called the *sample interval* (for x and y) or *update interval* (for u)
- u is piecewise constant (called *zero-order-hold*)
- x_d, y_d are sampled versions of x, y

$$\begin{aligned} x_d(k+1) &= x((k+1)h) \\ &= e^{hA}x(kh) + \int_0^h e^{\tau A} Bu((k+1)h - \tau) d\tau \\ &= e^{hA}x_d(k) + \left(\int_0^h e^{\tau A} d\tau \right) B u_d(k) \end{aligned}$$

x_d, u_d , and y_d satisfy discrete-time LDS equations

$$x_d(k+1) = A_d x_d(k) + B_d u_d(k), \quad y_d(k) = C_d x_d(k) + D_d u_d(k)$$

where

$$A_d = e^{hA}, \quad B_d = \left(\int_0^h e^{\tau A} d\tau \right) B, \quad C_d = C, \quad D_d = D$$

called *discretized system*

if A is invertible, we can express integral as

$$\int_0^h e^{\tau A} d\tau = A^{-1} (e^{hA} - I)$$

stability: if eigenvalues of A are $\lambda_1, \dots, \lambda_n$, then eigenvalues of A_d are $e^{h\lambda_1}, \dots, e^{h\lambda_n}$

discretization preserves stability properties since

$$\Re \lambda_i < 0 \quad \Leftrightarrow \quad |e^{h\lambda_i}| < 1$$

for $h > 0$

extensions/variations:

- *offsets:* updates for u and sampling of x, y are offset in time
- *multirate:* u_i updated, y_i sampled at different intervals
(usually integer multiples of a common interval h)

both very common in practice

Dual system

the *dual system* associated with system

$$\dot{x} = Ax + Bu, \quad y = Cx + Du$$

is given by

$$\dot{z} = A^T z + C^T v, \quad w = B^T z + D^T v$$

- all matrices are transposed
- role of B and C are swapped

transfer function of dual system:

$$(B^T)(sI - A^T)^{-1}(C^T) + D^T = H(s)^T$$

where $H(s) = C(sI - A)^{-1}B + D$

(for SISO case, TF of dual is same as original)

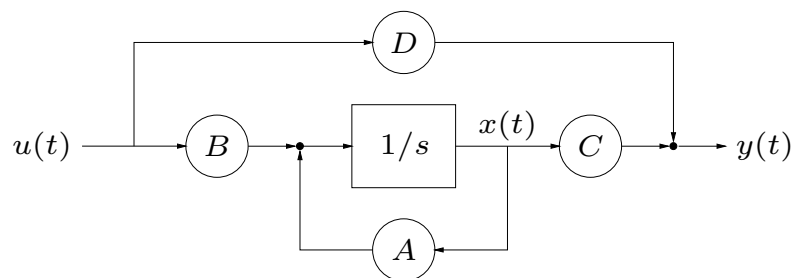
eigenvalues (hence stability properties) are the same

Dual via block diagram

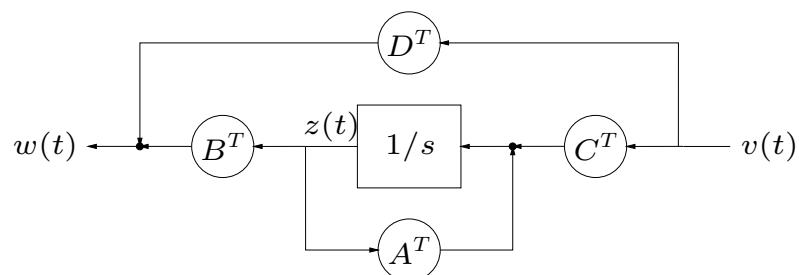
in terms of block diagrams, dual is formed by:

- transpose all matrices
- swap inputs and outputs on all boxes
- reverse directions of signal flow arrows
- swap solder joints and summing junctions

original system:



dual system:



Causality

interpretation of

$$\begin{aligned}x(t) &= e^{tA}x(0) + \int_0^t e^{(t-\tau)A}Bu(\tau) d\tau \\y(t) &= Ce^{tA}x(0) + \int_0^t Ce^{(t-\tau)A}Bu(\tau) d\tau + Du(t)\end{aligned}$$

for $t \geq 0$:

current state ($x(t)$) and output ($y(t)$) depend on *past* input ($u(\tau)$ for $\tau \leq t$)

i.e., mapping from input to state and output is *causal* (with fixed *initial* state)

now consider fixed *final* state $x(T)$: for $t \leq T$,

$$x(t) = e^{(t-T)A}x(T) + \int_T^t e^{(t-\tau)A}Bu(\tau) d\tau,$$

i.e., current state (and output) depend on future input!

so for fixed final condition, same system is anti-causal

Idea of state

$x(t)$ is called *state* of system at time t since:

- future output depends only on current state and future input
- future output depends on past input only through current state
- state summarizes effect of past inputs on future output
- state is bridge between past inputs and future outputs

Change of coordinates

start with LDS $\dot{x} = Ax + Bu, y = Cx + Du$

change coordinates in \mathbf{R}^n to \tilde{x} , with $x = T\tilde{x}$

then

$$\dot{\tilde{x}} = T^{-1}\dot{x} = T^{-1}(Ax + Bu) = T^{-1}AT\tilde{x} + T^{-1}Bu$$

hence LDS can be expressed as

$$\dot{\tilde{x}} = \tilde{A}\tilde{x} + \tilde{B}u, \quad y = \tilde{C}\tilde{x} + \tilde{D}u$$

where

$$\tilde{A} = T^{-1}AT, \quad \tilde{B} = T^{-1}B, \quad \tilde{C} = CT, \quad \tilde{D} = D$$

TF is same (since u, y aren't affected):

$$\tilde{C}(sI - \tilde{A})^{-1}\tilde{B} + \tilde{D} = C(sI - A)^{-1}B + D$$

Standard forms for LDS

can change coordinates to put A in various forms (diagonal, real modal, Jordan . . .)

e.g., to put LDS in *diagonal form*, find T s.t.

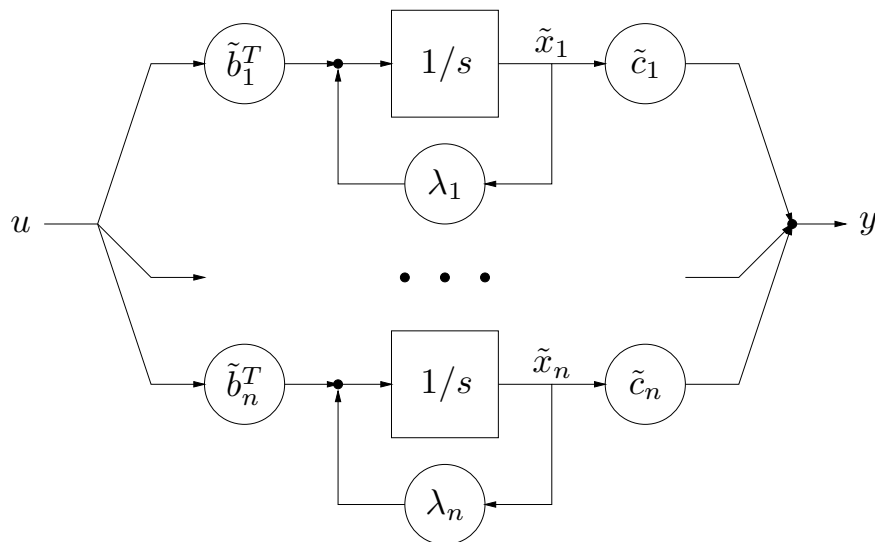
$$T^{-1}AT = \mathbf{diag}(\lambda_1, \dots, \lambda_n)$$

write

$$T^{-1}B = \begin{bmatrix} \tilde{b}_1^T \\ \vdots \\ \tilde{b}_n^T \end{bmatrix}, \quad CT = \begin{bmatrix} \tilde{c}_1 & \cdots & \tilde{c}_n \end{bmatrix}$$

so

$$\dot{\tilde{x}}_i = \lambda_i \tilde{x}_i + \tilde{b}_i^T u, \quad y = \sum_{i=1}^n \tilde{c}_i \tilde{x}_i$$

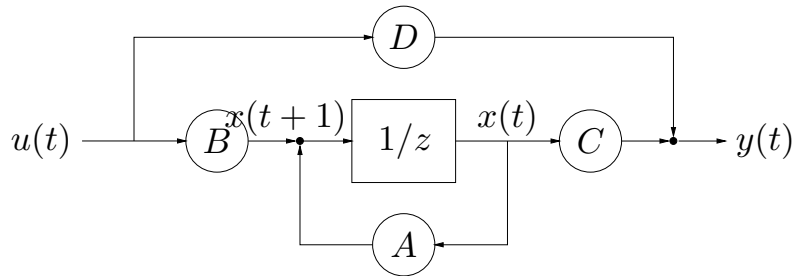


(here we assume $D = 0$)

Discrete-time systems

discrete-time LDS:

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t)$$



- only difference w/cts-time: z instead of s
- interpretation of z^{-1} block:
 - unit delayor (shifts sequence back in time one epoch)
 - latch (plus small delay to avoid race condition)

we have:

$$x(1) = Ax(0) + Bu(0),$$

$$\begin{aligned} x(2) &= Ax(1) + Bu(1) \\ &= A^2x(0) + ABu(0) + Bu(1), \end{aligned}$$

and in general, for $t \in \mathbf{Z}_+$,

$$x(t) = A^t x(0) + \sum_{\tau=0}^{t-1} A^{(t-1-\tau)} Bu(\tau)$$

hence

$$y(t) = CA^t x(0) + h * u$$

where $*$ is discrete-time convolution and

$$h(t) = \begin{cases} D, & t = 0 \\ CA^{t-1}B, & t > 0 \end{cases}$$

is the impulse response

\mathcal{Z} -transform

suppose $w \in \mathbf{R}^{p \times q}$ is a sequence (discrete-time signal), *i.e.*,

$$w : \mathbf{Z}_+ \rightarrow \mathbf{R}^{p \times q}$$

recall \mathcal{Z} -transform $W = \mathcal{Z}(w)$:

$$W(z) = \sum_{t=0}^{\infty} z^{-t} w(t)$$

where $W : D \subseteq \mathbf{C} \rightarrow \mathbf{C}^{p \times q}$ (D is domain of W)

time-advanced or shifted signal v :

$$v(t) = w(t+1), \quad t = 0, 1, \dots$$

\mathcal{Z} -transform of time-advanced signal:

$$\begin{aligned} V(z) &= \sum_{t=0}^{\infty} z^{-t} w(t+1) \\ &= z \sum_{t=1}^{\infty} z^{-t} w(t) \\ &= zW(z) - zw(0) \end{aligned}$$

Discrete-time transfer function

take \mathcal{Z} -transform of system equations

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t)$$

yields

$$zX(z) - zx(0) = AX(z) + BU(z), \quad Y(z) = CX(z) + DU(z)$$

solve for $X(z)$ to get

$$X(z) = (zI - A)^{-1}zx(0) + (zI - A)^{-1}BU(z)$$

(note extra z in first term!)

hence

$$Y(z) = H(z)U(z) + C(zI - A)^{-1}zx(0)$$

where $H(z) = C(zI - A)^{-1}B + D$ is the *discrete-time transfer function*

note power series expansion of resolvent:

$$(zI - A)^{-1} = z^{-1}I + z^{-2}A + z^{-3}A^2 + \dots$$

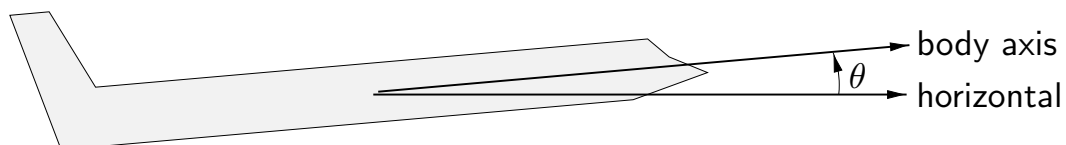
Lecture 14

Example: Aircraft dynamics

- longitudinal aircraft dynamics
- wind gust & control inputs
- linearized dynamics
- steady-state analysis
- eigenvalues & modes
- impulse matrices

14-1

Longitudinal aircraft dynamics



variables are (small) deviations from operating point or *trim conditions*
state (components):

- u : velocity of aircraft along body axis
- v : velocity of aircraft perpendicular to body axis
(down is positive)
- θ : angle between body axis and horizontal
(up is positive)
- $q = \dot{\theta}$: angular velocity of aircraft (pitch rate)

Inputs

disturbance inputs:

- u_w : velocity of wind along body axis
- v_w : velocity of wind perpendicular to body axis

control or actuator inputs:

- δ_e : elevator angle ($\delta_e > 0$ is down)
- δ_t : thrust

Example: Aircraft dynamics

14-3

Linearized dynamics

for 747, level flight, 40000 ft, 774 ft/sec,

$$\begin{bmatrix} \dot{u} \\ \dot{v} \\ \dot{q} \\ \dot{\theta} \end{bmatrix} = \begin{bmatrix} -.003 & .039 & 0 & -.322 \\ -.065 & -.319 & 7.74 & 0 \\ .020 & -.101 & -.429 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} u - u_w \\ v - v_w \\ q \\ \theta \end{bmatrix} + \begin{bmatrix} .01 & 1 \\ -.18 & -.04 \\ -1.16 & .598 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \delta_e \\ \delta_t \end{bmatrix}$$

- units: ft, sec, crad ($= 0.01\text{rad} \approx 0.57^\circ$)
- matrix coefficients are called *stability derivatives*

Example: Aircraft dynamics

14-4

outputs of interest:

- aircraft speed u (deviation from trim)
- climb rate $\dot{h} = -v + 7.74\theta$

Steady-state analysis

DC gain from $(u_w, v_w, \delta_e, \delta_t)$ to (u, \dot{h}) :

$$H(0) = -CA^{-1}B + D = \begin{bmatrix} 1 & 0 & 27.2 & -15.0 \\ 0 & -1 & -1.34 & 24.9 \end{bmatrix}$$

gives steady-state change in speed & climb rate due to wind, elevator & thrust changes

solve for control variables in terms of wind velocities, desired speed & climb rate

$$\begin{bmatrix} \delta_e \\ \delta_t \end{bmatrix} = \begin{bmatrix} .0379 & .0229 \\ .0020 & .0413 \end{bmatrix} \begin{bmatrix} u - u_w \\ \dot{h} + v_w \end{bmatrix}$$

- level flight, increase in speed is obtained mostly by increasing elevator (*i.e.*, downwards)
- constant speed, increase in climb rate is obtained by increasing thrust and increasing elevator (*i.e.*, downwards)

(thrust on 747 gives strong pitch up torque)

Eigenvalues and modes

eigenvalues are

$$-0.3750 \pm 0.8818j, \quad -0.0005 \pm 0.0674j$$

- two complex modes, called *short-period* and *phugoid*, respectively
- system is stable (but lightly damped)
- hence step responses converge (eventually) to DC gain matrix

eigenvectors are

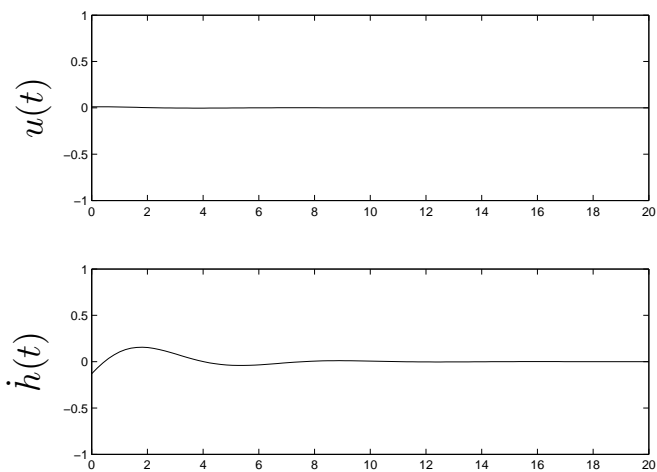
$$\begin{aligned}x_{\text{short}} &= \begin{bmatrix} 0.0005 \\ -0.5433 \\ -0.0899 \\ -0.0283 \end{bmatrix} \pm j \begin{bmatrix} 0.0135 \\ 0.8235 \\ -0.0677 \\ 0.1140 \end{bmatrix}, \\x_{\text{phug}} &= \begin{bmatrix} -0.7510 \\ -0.0962 \\ -0.0111 \\ 0.1225 \end{bmatrix} \pm j \begin{bmatrix} 0.6130 \\ 0.0941 \\ 0.0082 \\ 0.1637 \end{bmatrix}\end{aligned}$$

Example: Aircraft dynamics

14-9

Short-period mode

$y(t) = Ce^{tA}(\Re x_{\text{short}})$ (pure short-period mode motion)



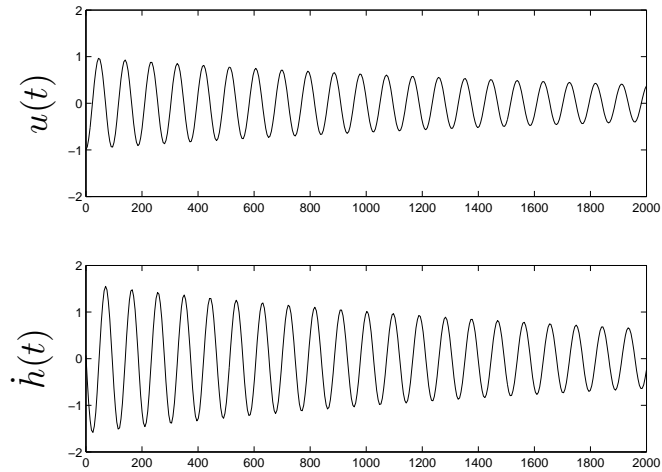
- only small effect on speed u
- period ≈ 7 sec, decays in ≈ 10 sec

Example: Aircraft dynamics

14-10

Phugoid mode

$$y(t) = Ce^{tA}(\Re x_{\text{phug}}) \text{ (pure phugoid mode motion)}$$



- affects both speed and climb rate
- period ≈ 100 sec; decays in ≈ 5000 sec

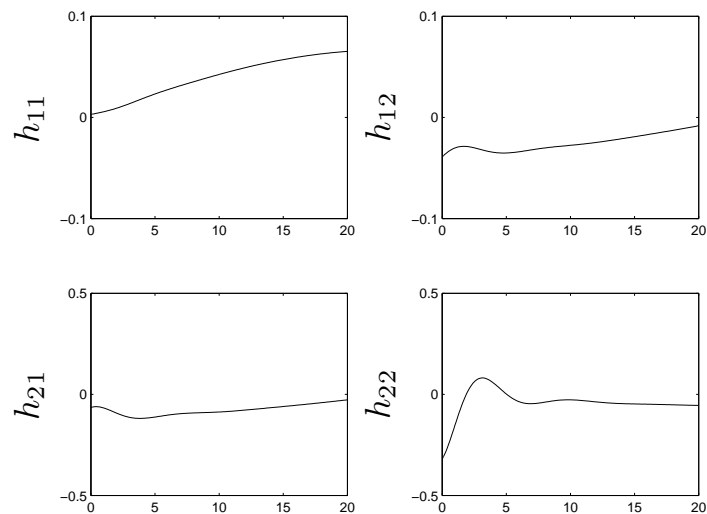
Example: Aircraft dynamics

14–11

Dynamic response to wind gusts

impulse response matrix from (u_w, v_w) to (u, \dot{h}) (gives response to short wind bursts)

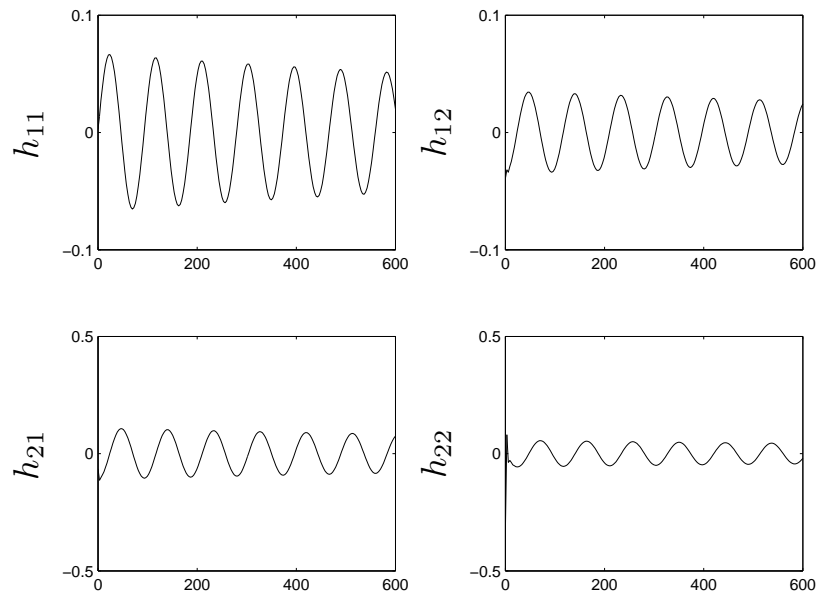
over time period $[0, 20]$:



Example: Aircraft dynamics

14–12

over time period $[0, 600]$:



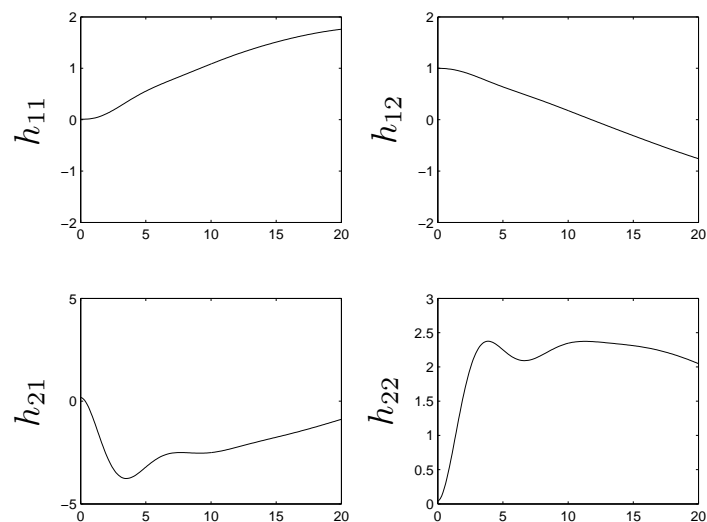
Example: Aircraft dynamics

14-13

Dynamic response to actuators

impulse response matrix from (δ_e, δ_t) to (u, \dot{h})

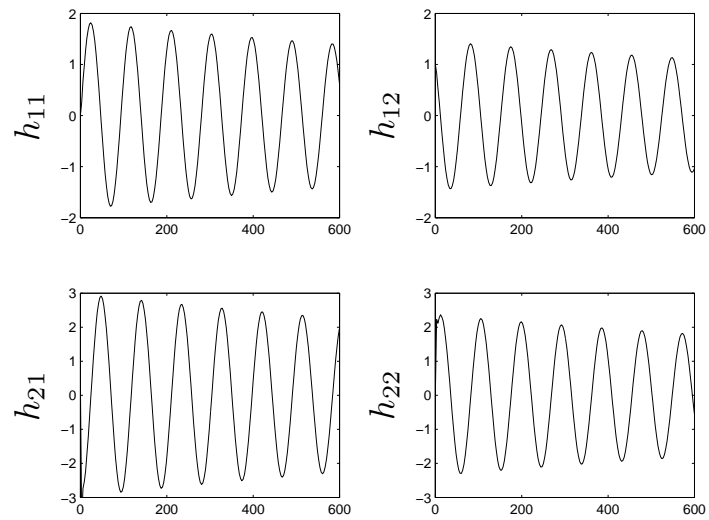
over time period $[0, 20]$:



Example: Aircraft dynamics

14-14

over time period $[0, 600]$:



Example: Aircraft dynamics

14–15

Lecture 15

Symmetric matrices, quadratic forms, matrix norm, and SVD

- eigenvectors of symmetric matrices
- quadratic forms
- inequalities for quadratic forms
- positive semidefinite matrices
- norm of a matrix
- singular value decomposition

15-1

Eigenvalues of symmetric matrices

suppose $A \in \mathbf{R}^{n \times n}$ is symmetric, *i.e.*, $A = A^T$

fact: the eigenvalues of A are real

to see this, suppose $Av = \lambda v$, $v \neq 0$, $v \in \mathbf{C}^n$

then

$$\bar{v}^T Av = \bar{v}^T (Av) = \lambda \bar{v}^T v = \lambda \sum_{i=1}^n |v_i|^2$$

but also

$$\bar{v}^T Av = \overline{(Av)}^T v = \overline{(\lambda v)}^T v = \bar{\lambda} \sum_{i=1}^n |v_i|^2$$

so we have $\lambda = \bar{\lambda}$, *i.e.*, $\lambda \in \mathbf{R}$ (hence, can assume $v \in \mathbf{R}^n$)

Eigenvectors of symmetric matrices

fact: there is a set of orthonormal eigenvectors of A , i.e., q_1, \dots, q_n s.t.
 $Aq_i = \lambda_i q_i$, $q_i^T q_j = \delta_{ij}$

in matrix form: there is an orthogonal Q s.t.

$$Q^{-1}AQ = Q^T AQ = \Lambda$$

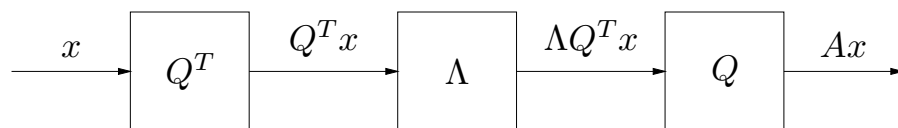
hence we can express A as

$$A = Q\Lambda Q^T = \sum_{i=1}^n \lambda_i q_i q_i^T$$

in particular, q_i are both left and right eigenvectors

Interpretations

$$A = Q\Lambda Q^T$$



linear mapping $y = Ax$ can be decomposed as

- resolve into q_i coordinates
- scale coordinates by λ_i
- reconstitute with basis q_i

or, geometrically,

- rotate by Q^T
- diagonal real scale ('dilation') by Λ
- rotate back by Q

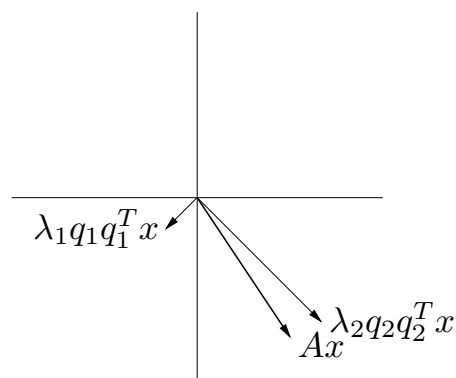
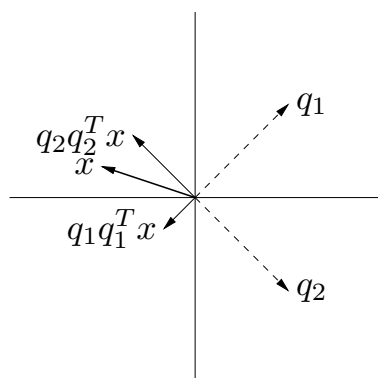
decomposition

$$A = \sum_{i=1}^n \lambda_i q_i q_i^T$$

expresses A as linear combination of 1-dimensional projections

example:

$$\begin{aligned} A &= \begin{bmatrix} -1/2 & 3/2 \\ 3/2 & -1/2 \end{bmatrix} \\ &= \left(\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \right) \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix} \left(\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \right)^T \end{aligned}$$



proof (case of λ_i distinct)

since λ_i distinct, can find v_1, \dots, v_n , a set of linearly independent eigenvectors of A :

$$Av_i = \lambda_i v_i, \quad \|v_i\| = 1$$

then we have

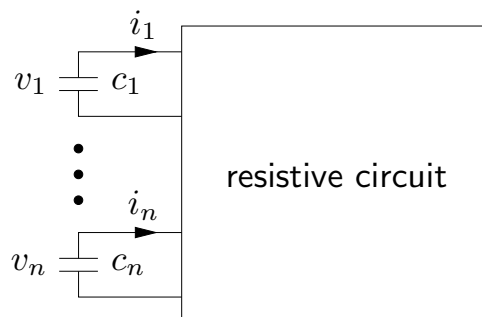
$$v_i^T (Av_j) = \lambda_j v_i^T v_j = (Av_i)^T v_j = \lambda_i v_i^T v_j$$

$$\text{so } (\lambda_i - \lambda_j) v_i^T v_j = 0$$

for $i \neq j$, $\lambda_i \neq \lambda_j$, hence $v_i^T v_j = 0$

- in this case we can say: eigenvectors *are* orthogonal
- in general case (λ_i not distinct) we must say: eigenvectors *can be chosen* to be orthogonal

Example: RC circuit



$$c_k \dot{v}_k = -i_k, \quad i = Gv$$

$G = G^T \in \mathbf{R}^{n \times n}$ is conductance matrix of resistive circuit

thus $\dot{v} = -C^{-1}Gv$ where $C = \mathbf{diag}(c_1, \dots, c_n)$

note $-C^{-1}G$ is not symmetric

use state $x_i = \sqrt{c_i}v_i$, so

$$\dot{x} = C^{1/2}\dot{v} = -C^{-1/2}GC^{-1/2}x$$

where $C^{1/2} = \mathbf{diag}(\sqrt{c_1}, \dots, \sqrt{c_n})$

we conclude:

- eigenvalues $\lambda_1, \dots, \lambda_n$ of $-C^{-1/2}GC^{-1/2}$ (hence, $-C^{-1}G$) are real
- eigenvectors q_i (in x_i coordinates) can be chosen orthogonal
- eigenvectors in voltage coordinates, $s_i = C^{-1/2}q_i$, satisfy

$$-C^{-1}Gs_i = \lambda_i s_i, \quad s_i^T C s_i = \delta_{ij}$$

Quadratic forms

a function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ of the form

$$f(x) = x^T A x = \sum_{i,j=1}^n A_{ij} x_i x_j$$

is called a *quadratic form*

in a quadratic form we may as well assume $A = A^T$ since

$$x^T A x = x^T ((A + A^T)/2) x$$

$((A + A^T)/2)$ is called the *symmetric part* of A)

uniqueness: if $x^T A x = x^T B x$ for all $x \in \mathbf{R}^n$ and $A = A^T$, $B = B^T$, then $A = B$

Examples

- $\|Bx\|^2 = x^T B^T B x$
- $\sum_{i=1}^{n-1} (x_{i+1} - x_i)^2$
- $\|Fx\|^2 - \|Gx\|^2$

sets defined by quadratic forms:

- $\{ x \mid f(x) = a \}$ is called a *quadratic surface*
- $\{ x \mid f(x) \leq a \}$ is called a *quadratic region*

Inequalities for quadratic forms

suppose $A = A^T$, $A = Q\Lambda Q^T$ with eigenvalues sorted so $\lambda_1 \geq \dots \geq \lambda_n$

$$\begin{aligned} x^T A x &= x^T Q \Lambda Q^T x \\ &= (Q^T x)^T \Lambda (Q^T x) \\ &= \sum_{i=1}^n \lambda_i (q_i^T x)^2 \\ &\leq \lambda_1 \sum_{i=1}^n (q_i^T x)^2 \\ &= \lambda_1 \|x\|^2 \end{aligned}$$

i.e., we have $x^T A x \leq \lambda_1 x^T x$

similar argument shows $x^T Ax \geq \lambda_n \|x\|^2$, so we have

$$\lambda_n x^T x \leq x^T Ax \leq \lambda_1 x^T x$$

sometimes λ_1 is called λ_{\max} , λ_n is called λ_{\min}

note also that

$$q_1^T A q_1 = \lambda_1 \|q_1\|^2, \quad q_n^T A q_n = \lambda_n \|q_n\|^2,$$

so the inequalities are tight

Positive semidefinite and positive definite matrices

suppose $A = A^T \in \mathbf{R}^{n \times n}$

we say A is *positive semidefinite* if $x^T Ax \geq 0$ for all x

- denoted $A \geq 0$ (and sometimes $A \succeq 0$)
- $A \geq 0$ if and only if $\lambda_{\min}(A) \geq 0$, *i.e.*, all eigenvalues are nonnegative
- **not** the same as $A_{ij} \geq 0$ for all i, j

we say A is *positive definite* if $x^T Ax > 0$ for all $x \neq 0$

- denoted $A > 0$
- $A > 0$ if and only if $\lambda_{\min}(A) > 0$, *i.e.*, all eigenvalues are positive

Matrix inequalities

- we say A is *negative semidefinite* if $-A \geq 0$
- we say A is *negative definite* if $-A > 0$
- otherwise, we say A is *indefinite*

matrix inequality: if $B = B^T \in \mathbf{R}^n$ we say $A \geq B$ if $A - B \geq 0$, $A < B$ if $B - A > 0$, etc.

for example:

- $A \geq 0$ means A is positive semidefinite
- $A > B$ means $x^T A x > x^T B x$ for all $x \neq 0$

many properties that you'd guess hold actually do, *e.g.*,

- if $A \geq B$ and $C \geq D$, then $A + C \geq B + D$
- if $B \leq 0$ then $A + B \leq A$
- if $A \geq 0$ and $\alpha \geq 0$, then $\alpha A \geq 0$
- $A^2 \geq 0$
- if $A > 0$, then $A^{-1} > 0$

matrix inequality is only a *partial order*: we can have

$$A \not\geq B, \quad B \not\geq A$$

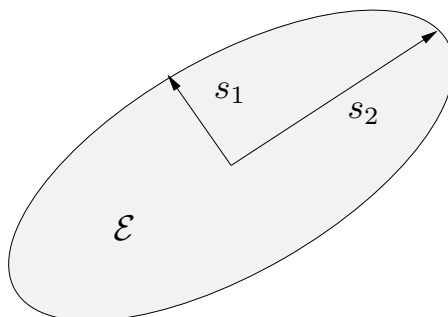
(such matrices are called *incomparable*)

Ellipsoids

if $A = A^T > 0$, the set

$$\mathcal{E} = \{ x \mid x^T A x \leq 1 \}$$

is an *ellipsoid* in \mathbf{R}^n , centered at 0



semi-axes are given by $s_i = \lambda_i^{-1/2} q_i$, *i.e.*:

- eigenvectors determine directions of semiaxes
- eigenvalues determine lengths of semiaxes

note:

- in direction q_1 , $x^T A x$ is *large*, hence ellipsoid is *thin* in direction q_1
- in direction q_n , $x^T A x$ is *small*, hence ellipsoid is *fat* in direction q_n
- $\sqrt{\lambda_{\max}/\lambda_{\min}}$ gives maximum *eccentricity*

if $\tilde{\mathcal{E}} = \{ x \mid x^T B x \leq 1 \}$, where $B > 0$, then $\mathcal{E} \subseteq \tilde{\mathcal{E}} \iff A \geq B$

Gain of a matrix in a direction

suppose $A \in \mathbf{R}^{m \times n}$ (not necessarily square or symmetric)

for $x \in \mathbf{R}^n$, $\|Ax\|/\|x\|$ gives the *amplification factor* or *gain* of A in the direction x

obviously, gain varies with direction of input x

questions:

- what is maximum gain of A
(and corresponding maximum gain direction)?
- what is minimum gain of A
(and corresponding minimum gain direction)?
- how does gain of A vary with direction?

Matrix norm

the maximum gain

$$\max_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

is called the *matrix norm* or *spectral norm* of A and is denoted $\|A\|$

$$\max_{x \neq 0} \frac{\|Ax\|^2}{\|x\|^2} = \max_{x \neq 0} \frac{x^T A^T A x}{\|x\|^2} = \lambda_{\max}(A^T A)$$

so we have $\|A\| = \sqrt{\lambda_{\max}(A^T A)}$

similarly the minimum gain is given by

$$\min_{x \neq 0} \|Ax\|/\|x\| = \sqrt{\lambda_{\min}(A^T A)}$$

note that

- $A^T A \in \mathbf{R}^{n \times n}$ is symmetric and $A^T A \geq 0$ so $\lambda_{\min}, \lambda_{\max} \geq 0$
- ‘max gain’ input direction is $x = q_1$, eigenvector of $A^T A$ associated with λ_{\max}
- ‘min gain’ input direction is $x = q_n$, eigenvector of $A^T A$ associated with λ_{\min}

example: $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$

$$\begin{aligned} A^T A &= \begin{bmatrix} 35 & 44 \\ 44 & 56 \end{bmatrix} \\ &= \begin{bmatrix} 0.620 & 0.785 \\ 0.785 & -0.620 \end{bmatrix} \begin{bmatrix} 90.7 & 0 \\ 0 & 0.265 \end{bmatrix} \begin{bmatrix} 0.620 & 0.785 \\ 0.785 & -0.620 \end{bmatrix}^T \end{aligned}$$

then $\|A\| = \sqrt{\lambda_{\max}(A^T A)} = 9.53$:

$$\left\| \begin{bmatrix} 0.620 \\ 0.785 \end{bmatrix} \right\| = 1, \quad \left\| A \begin{bmatrix} 0.620 \\ 0.785 \end{bmatrix} \right\| = \left\| \begin{bmatrix} 2.18 \\ 4.99 \\ 7.78 \end{bmatrix} \right\| = 9.53$$

min gain is $\sqrt{\lambda_{\min}(A^T A)} = 0.514$:

$$\left\| \begin{bmatrix} 0.785 \\ -0.620 \end{bmatrix} \right\| = 1, \quad \left\| A \begin{bmatrix} 0.785 \\ -0.620 \end{bmatrix} \right\| = \left\| \begin{bmatrix} 0.46 \\ 0.14 \\ -0.18 \end{bmatrix} \right\| = 0.514$$

for all $x \neq 0$, we have

$$0.514 \leq \frac{\|Ax\|}{\|x\|} \leq 9.53$$

Properties of matrix norm

- consistent with vector norm: matrix norm of $a \in \mathbf{R}^{n \times 1}$ is $\sqrt{\lambda_{\max}(a^T a)} = \sqrt{a^T a}$
- for any x , $\|Ax\| \leq \|A\|\|x\|$
- scaling: $\|aA\| = |a|\|A\|$
- triangle inequality: $\|A + B\| \leq \|A\| + \|B\|$
- definiteness: $\|A\| = 0 \iff A = 0$
- norm of product: $\|AB\| \leq \|A\|\|B\|$

Singular value decomposition

more complete picture of gain properties of A given by *singular value decomposition* (SVD) of A :

$$A = U\Sigma V^T$$

where

- $A \in \mathbf{R}^{m \times n}$, $\mathbf{Rank}(A) = r$
- $U \in \mathbf{R}^{m \times r}$, $U^T U = I$
- $V \in \mathbf{R}^{n \times r}$, $V^T V = I$
- $\Sigma = \mathbf{diag}(\sigma_1, \dots, \sigma_r)$, where $\sigma_1 \geq \dots \geq \sigma_r > 0$

with $U = [u_1 \cdots u_r]$, $V = [v_1 \cdots v_r]$,

$$A = U\Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T$$

- σ_i are the (nonzero) *singular values* of A
- v_i are the *right* or *input singular vectors* of A
- u_i are the *left* or *output singular vectors* of A

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^2 V^T$$

hence:

- v_i are eigenvectors of $A^T A$ (corresponding to nonzero eigenvalues)
- $\sigma_i = \sqrt{\lambda_i(A^T A)}$ (and $\lambda_i(A^T A) = 0$ for $i > r$)
- $\|A\| = \sigma_1$

similarly,

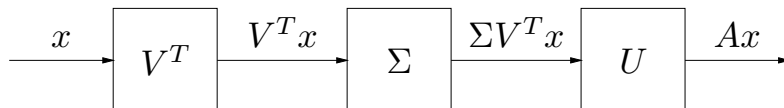
$$A A^T = (U \Sigma V^T)(U \Sigma V^T)^T = U \Sigma^2 U^T$$

hence:

- u_i are eigenvectors of $A A^T$ (corresponding to nonzero eigenvalues)
- $\sigma_i = \sqrt{\lambda_i(A A^T)}$ (and $\lambda_i(A A^T) = 0$ for $i > r$)
- u_1, \dots, u_r are orthonormal basis for $\text{range}(A)$
- v_1, \dots, v_r are orthonormal basis for $\mathcal{N}(A)^\perp$

Interpretations

$$A = U\Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T$$



linear mapping $y = Ax$ can be decomposed as

- compute coefficients of x along input directions v_1, \dots, v_r
- scale coefficients by σ_i
- reconstitute along output directions u_1, \dots, u_r

difference with eigenvalue decomposition for symmetric A : input and output directions are *different*

- v_1 is most sensitive (highest gain) input direction
- u_1 is highest gain output direction
- $Av_1 = \sigma_1 u_1$

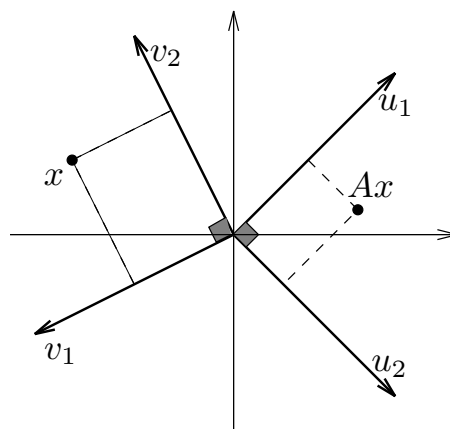
SVD gives clearer picture of gain as function of input/output directions

example: consider $A \in \mathbf{R}^{4 \times 4}$ with $\Sigma = \mathbf{diag}(10, 7, 0.1, 0.05)$

- input components along directions v_1 and v_2 are amplified (by about 10) and come out mostly along plane spanned by u_1, u_2
- input components along directions v_3 and v_4 are attenuated (by about 10)
- $\|Ax\|/\|x\|$ can range between 10 and 0.05
- A is nonsingular
- for some applications you might say A is *effectively* rank 2

example: $A \in \mathbf{R}^{2 \times 2}$, with $\sigma_1 = 1, \sigma_2 = 0.5$

- resolve x along v_1, v_2 : $v_1^T x = 0.5, v_2^T x = 0.6$, *i.e.*, $x = 0.5v_1 + 0.6v_2$
- now form $Ax = (v_1^T x)\sigma_1 u_1 + (v_2^T x)\sigma_2 u_2 = (0.5)(1)u_1 + (0.6)(0.5)u_2$



Lecture 16

SVD Applications

- general pseudo-inverse
- full SVD
- image of unit ball under linear transformation
- SVD in estimation/inversion
- sensitivity of linear equations to data error
- low rank approximation via SVD

16-1

General pseudo-inverse

if $A \neq 0$ has SVD $A = U\Sigma V^T$,

$$A^\dagger = V\Sigma^{-1}U^T$$

is the *pseudo-inverse* or *Moore-Penrose inverse* of A

if A is skinny and full rank,

$$A^\dagger = (A^T A)^{-1} A^T$$

gives the least-squares approximate solution $x_{\text{ls}} = A^\dagger y$

if A is fat and full rank,

$$A^\dagger = A^T (A A^T)^{-1}$$

gives the least-norm solution $x_{\text{ln}} = A^\dagger y$

in general case:

$$X_{\text{ls}} = \{ z \mid \|Az - y\| = \min_w \|Aw - y\| \}$$

is set of least-squares approximate solutions

$x_{\text{pinv}} = A^\dagger y \in X_{\text{ls}}$ has minimum norm on X_{ls} , *i.e.*, x_{pinv} is the minimum-norm, least-squares approximate solution

Pseudo-inverse via regularization

for $\mu > 0$, let x_μ be (unique) minimizer of

$$\|Ax - y\|^2 + \mu\|x\|^2$$

i.e.,

$$x_\mu = (A^T A + \mu I)^{-1} A^T y$$

here, $A^T A + \mu I > 0$ and so is invertible

then we have $\lim_{\mu \rightarrow 0} x_\mu = A^\dagger y$

in fact, we have $\lim_{\mu \rightarrow 0} (A^T A + \mu I)^{-1} A^T = A^\dagger$

(check this!)

Full SVD

SVD of $A \in \mathbf{R}^{m \times n}$ with $\text{Rank}(A) = r$:

$$A = U_1 \Sigma_1 V_1^T = \begin{bmatrix} u_1 & \cdots & u_r \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_r^T \end{bmatrix}$$

- find $U_2 \in \mathbf{R}^{m \times (m-r)}$, $V_2 \in \mathbf{R}^{n \times (n-r)}$ s.t. $U = [U_1 \ U_2] \in \mathbf{R}^{m \times m}$ and $V = [V_1 \ V_2] \in \mathbf{R}^{n \times n}$ are orthogonal
- add zero rows/cols to Σ_1 to form $\Sigma \in \mathbf{R}^{m \times n}$:

$$\Sigma = \left[\begin{array}{c|c} \Sigma_1 & 0_{r \times (n-r)} \\ \hline 0_{(m-r) \times r} & 0_{(m-r) \times (n-r)} \end{array} \right]$$

SVD Applications

16-5

then we have

$$A = U_1 \Sigma_1 V_1^T = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \left[\begin{array}{c|c} \Sigma_1 & 0_{r \times (n-r)} \\ \hline 0_{(m-r) \times r} & 0_{(m-r) \times (n-r)} \end{array} \right] \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix}$$

i.e.:

$$A = U \Sigma V^T$$

called *full SVD* of A

(SVD with positive singular values only called *compact SVD*)

Image of unit ball under linear transformation

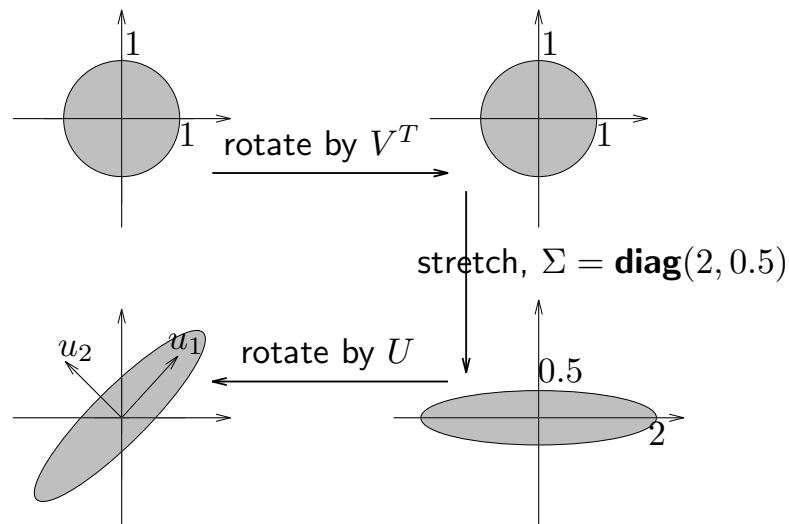
full SVD:

$$A = U\Sigma V^T$$

gives interpretation of $y = Ax$:

- rotate (by V^T)
- stretch along axes by σ_i ($\sigma_i = 0$ for $i > r$)
- zero-pad (if $m > n$) or truncate (if $m < n$) to get m -vector
- rotate (by U)

Image of unit ball under A



$\{Ax \mid \|x\| \leq 1\}$ is *ellipsoid* with principal axes $\sigma_i u_i$.

SVD in estimation/inversion

suppose $y = Ax + v$, where

- $y \in \mathbf{R}^m$ is measurement
- $x \in \mathbf{R}^n$ is vector to be estimated
- v is a measurement noise or error

'norm-bound' model of noise: we assume $\|v\| \leq \alpha$ but otherwise know nothing about v (α gives max norm of noise)

- consider estimator $\hat{x} = By$, with $BA = I$ (i.e., unbiased)
- estimation or inversion error is $\tilde{x} = \hat{x} - x = Bv$
- set of possible estimation errors is ellipsoid

$$\tilde{x} \in \mathcal{E}_{\text{unc}} = \{ Bv \mid \|v\| \leq \alpha \}$$

- $x = \hat{x} - \tilde{x} \in \hat{x} - \mathcal{E}_{\text{unc}} = \hat{x} + \mathcal{E}_{\text{unc}}$, i.e.:
true x lies in *uncertainty ellipsoid* \mathcal{E}_{unc} , centered at estimate \hat{x}
- 'good' estimator has 'small' \mathcal{E}_{unc} (with $BA = I$, of course)

semiaxes of \mathcal{E}_{unc} are $\alpha\sigma_i u_i$ (singular values & vectors of B)

e.g., maximum norm of error is $\alpha\|B\|$, *i.e.*, $\|\hat{x} - x\| \leq \alpha\|B\|$

optimality of least-squares: suppose $BA = I$ is any estimator, and $B_{\text{ls}} = A^\dagger$ is the least-squares estimator

then:

- $B_{\text{ls}}B_{\text{ls}}^T \leq BB^T$
- $\mathcal{E}_{\text{ls}} \subseteq \mathcal{E}$
- in particular $\|B_{\text{ls}}\| \leq \|B\|$

i.e., the least-squares estimator gives the *smallest* uncertainty ellipsoid

Example: navigation using range measurements (lect. 4)

we have

$$y = - \begin{bmatrix} k_1^T \\ k_2^T \\ k_3^T \\ k_4^T \end{bmatrix} x$$

where $k_i \in \mathbf{R}^2$

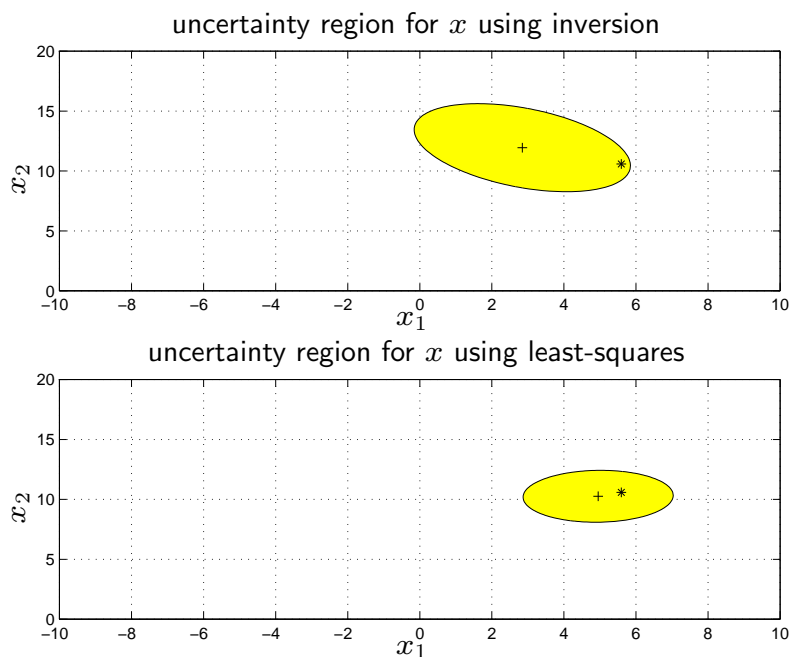
using first two measurements and inverting:

$$\hat{x} = - \begin{bmatrix} \begin{bmatrix} k_1^T \\ k_2^T \end{bmatrix}^{-1} & 0_{2 \times 2} \end{bmatrix} y$$

using all four measurements and least-squares:

$$\hat{x} = A^\dagger y$$

uncertainty regions (with $\alpha = 1$):



Proof of optimality property

suppose $A \in \mathbf{R}^{m \times n}$, $m > n$, is full rank

SVD: $A = U\Sigma V^T$, with V orthogonal

$B_{\text{ls}} = A^\dagger = V\Sigma^{-1}U^T$, and B satisfies $BA = I$

define $Z = B - B_{\text{ls}}$, so $B = B_{\text{ls}} + Z$

then $ZA = ZU\Sigma V^T = 0$, so $ZU = 0$ (multiply by $V\Sigma^{-1}$ on right)

therefore

$$\begin{aligned}
 BB^T &= (B_{\text{ls}} + Z)(B_{\text{ls}} + Z)^T \\
 &= B_{\text{ls}}B_{\text{ls}}^T + B_{\text{ls}}Z^T + ZB_{\text{ls}}^T + ZZ^T \\
 &= B_{\text{ls}}B_{\text{ls}}^T + ZZ^T \\
 &\geq B_{\text{ls}}B_{\text{ls}}^T
 \end{aligned}$$

using $ZB_{\text{ls}}^T = (ZU)\Sigma^{-1}V^T = 0$

Sensitivity of linear equations to data error

consider $y = Ax$, $A \in \mathbf{R}^{n \times n}$ invertible; of course $x = A^{-1}y$

suppose we have an error or noise in y , *i.e.*, y becomes $y + \delta y$

then x becomes $x + \delta x$ with $\delta x = A^{-1}\delta y$

hence we have $\|\delta x\| = \|A^{-1}\delta y\| \leq \|A^{-1}\|\|\delta y\|$

if $\|A^{-1}\|$ is large,

- small errors in y can lead to large errors in x
- can't solve for x given y (with small errors)
- hence, A can be considered singular in practice

a more refined analysis uses *relative* instead of *absolute* errors in x and y

since $y = Ax$, we also have $\|y\| \leq \|A\|\|x\|$, hence

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\|\|A^{-1}\| \frac{\|\delta y\|}{\|y\|}$$

$$\kappa(A) = \|A\|\|A^{-1}\| = \sigma_{\max}(A)/\sigma_{\min}(A)$$

is called the *condition number* of A

we have:

relative error in solution $x \leq$ condition number \cdot relative error in data y

or, in terms of # bits of guaranteed accuracy:

$$\# \text{ bits accuracy in solution} \approx \# \text{ bits accuracy in data} - \log_2 \kappa$$

we say

- A is well conditioned if κ is small
- A is poorly conditioned if κ is large

(definition of 'small' and 'large' depend on application)

same analysis holds for least-squares approximate solutions with A nonsquare, $\kappa = \sigma_{\max}(A)/\sigma_{\min}(A)$

Low rank approximations

suppose $A \in \mathbf{R}^{m \times n}$, $\mathbf{Rank}(A) = r$, with SVD $A = U\Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T$

we seek matrix \hat{A} , $\mathbf{Rank}(\hat{A}) \leq p < r$, s.t. $\hat{A} \approx A$ in the sense that $\|A - \hat{A}\|$ is minimized

solution: optimal rank p approximator is

$$\hat{A} = \sum_{i=1}^p \sigma_i u_i v_i^T$$

- hence $\|A - \hat{A}\| = \left\| \sum_{i=p+1}^r \sigma_i u_i v_i^T \right\| = \sigma_{p+1}$
- interpretation: SVD dyads $u_i v_i^T$ are ranked in order of 'importance'; take p to get rank p approximant

proof: suppose $\mathbf{Rank}(B) \leq p$

then $\dim \mathcal{N}(B) \geq n - p$

also, $\dim \text{span}\{v_1, \dots, v_{p+1}\} = p + 1$

hence, the two subspaces intersect, *i.e.*, there is a unit vector $z \in \mathbf{R}^n$ s.t.

$$Bz = 0, \quad z \in \text{span}\{v_1, \dots, v_{p+1}\}$$

$$(A - B)z = Az = \sum_{i=1}^{p+1} \sigma_i u_i v_i^T z$$

$$\|(A - B)z\|^2 = \sum_{i=1}^{p+1} \sigma_i^2 (v_i^T z)^2 \geq \sigma_{p+1}^2 \|z\|^2$$

hence $\|A - B\| \geq \sigma_{p+1} = \|A - \hat{A}\|$

Distance to singularity

another interpretation of σ_i :

$$\sigma_i = \min\{ \|A - B\| \mid \mathbf{Rank}(B) \leq i - 1 \}$$

i.e., the distance (measured by matrix norm) to the nearest rank $i - 1$ matrix

for example, if $A \in \mathbf{R}^{n \times n}$, $\sigma_n = \sigma_{\min}$ is distance to nearest singular matrix

hence, small σ_{\min} means A is near to a singular matrix

application: model simplification

suppose $y = Ax + v$, where

- $A \in \mathbf{R}^{100 \times 30}$ has SVs

10, 7, 2, 0.5, 0.01, ..., 0.0001

- $\|x\|$ is on the order of 1
- unknown error or noise v has norm on the order of 0.1

then the terms $\sigma_i u_i v_i^T x$, for $i = 5, \dots, 30$, are substantially smaller than the noise term v

simplified model:

$$y = \sum_{i=1}^4 \sigma_i u_i v_i^T x + v$$

Lecture 17

Example: Quantum mechanics

- wave function and Schrodinger equation
- discretization
- preservation of probability
- eigenvalues & eigenstates
- example

17-1

Quantum mechanics

- single particle in interval $[0, 1]$, mass m
- potential $V : [0, 1] \rightarrow \mathbf{R}$

$\Psi : [0, 1] \times \mathbf{R}_+ \rightarrow \mathbf{C}$ is (complex-valued) *wave function*

interpretation: $|\Psi(x, t)|^2$ is probability density of particle at position x , time t

(so $\int_0^1 |\Psi(x, t)|^2 dx = 1$ for all t)

evolution of Ψ governed by *Schrodinger* equation:

$$i\hbar \dot{\Psi} = \left(V - \frac{\hbar^2}{2m} \nabla_x^2 \right) \Psi = H\Psi$$

where H is *Hamiltonian* operator, $i = \sqrt{-1}$

Discretization

let's discretize position x into N discrete points, k/N , $k = 1, \dots, N$

wave function is approximated as *vector* $\Psi(t) \in \mathbf{C}^N$

∇_x^2 operator is approximated as *matrix*

$$\nabla^2 = N^2 \begin{bmatrix} -2 & 1 & & & \\ & 1 & -2 & 1 & \\ & & 1 & -2 & 1 \\ & & & \ddots & \ddots & \ddots \\ & & & & 1 & -2 \end{bmatrix}$$

so $w = \nabla^2 v$ means

$$w_k = \frac{(v_{k+1} - v_k)/(1/N) - (v_k - v_{k-1})(1/N)}{1/N}$$

(which approximates $w = \partial^2 v / \partial x^2$)

Example: Quantum mechanics

17-3

discretized Schrodinger equation is (complex) linear dynamical system

$$\dot{\Psi} = (-i/\hbar)(V - (\hbar/2m)\nabla^2)\Psi = (-i/\hbar)H\Psi$$

where V is a diagonal matrix with $V_{kk} = V(k/N)$

hence we analyze using linear dynamical system theory (with complex vectors & matrices):

$$\dot{\Psi} = (-i/\hbar)H\Psi$$

solution of Shrodinger equation: $\Psi(t) = e^{(-i/\hbar)tH}\Psi(0)$

matrix $e^{(-i/\hbar)tH}$ propogates wave function forward in time t seconds (backward if $t < 0$)

Example: Quantum mechanics

17-4

Preservation of probability

$$\begin{aligned}
 \frac{d}{dt} \|\Psi\|^2 &= \frac{d}{dt} \Psi^* \Psi \\
 &= \dot{\Psi}^* \Psi + \Psi^* \dot{\Psi} \\
 &= ((-i/\hbar) H \Psi)^* \Psi + \Psi^* ((-i/\hbar) H \Psi) \\
 &= (i/\hbar) \Psi^* H \Psi + (-i/\hbar) \Psi^* H \Psi \\
 &= 0
 \end{aligned}$$

(using $H = H^T \in \mathbf{R}^{N \times N}$)

hence, $\|\Psi(t)\|^2$ is constant; our discretization preserves probability *exactly*

Example: Quantum mechanics

17-5

$U = e^{-(i/\hbar)tH}$ is *unitary*, meaning $U^*U = I$

unitary is extension of *orthogonal* for complex matrix: if $U \in \mathbf{C}^{N \times N}$ is unitary and $z \in \mathbf{C}^N$, then

$$\|Uz\|^2 = (Uz)^*(Uz) = z^* U^* U z = z^* z = \|z\|^2$$

Example: Quantum mechanics

17-6

Eigenvalues & eigenstates

H is symmetric, so

- its eigenvalues $\lambda_1, \dots, \lambda_N$ are real ($\lambda_1 \leq \dots \leq \lambda_N$)
- its eigenvectors v_1, \dots, v_N can be chosen to be orthogonal (and real)

from $Hv = \lambda v \Leftrightarrow (-i/\hbar)Hv = (-i/\hbar)\lambda v$ we see:

- eigenvectors of $(-i/\hbar)H$ are same as eigenvectors of H , *i.e.*, v_1, \dots, v_N
- eigenvalues of $(-i/\hbar)H$ are $(-i/\hbar)\lambda_1, \dots, (-i/\hbar)\lambda_N$ (which are pure imaginary)

Example: Quantum mechanics

17-7

- eigenvectors v_k are called *eigenstates* of system
- eigenvalue λ_k is *energy* of eigenstate v_k
- for mode $\Psi(t) = e^{(-i/\hbar)\lambda_k t} v_k$, probability density

$$|\Psi_m(t)|^2 = \left| e^{(-i/\hbar)\lambda_k t} v_k \right|^2 = |v_{mk}|^2$$

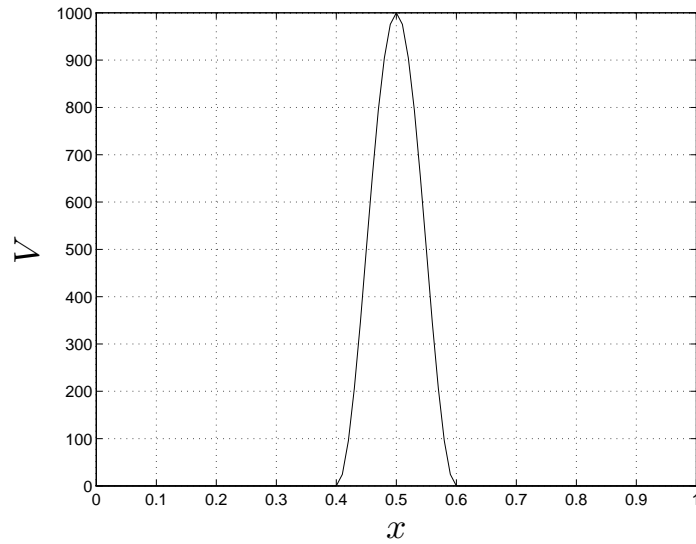
doesn't change with time (v_{mk} is m th entry of v_k)

Example: Quantum mechanics

17-8

Example

Potential Function $V(x)$

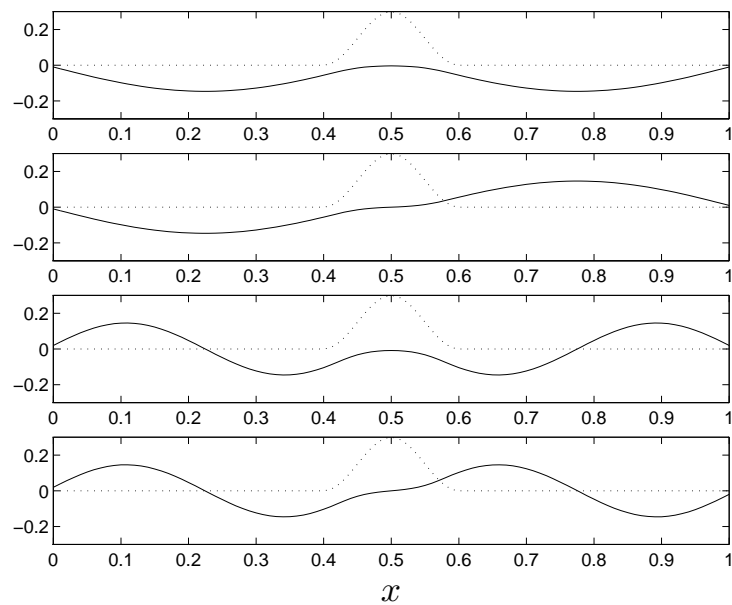


- potential bump in middle of infinite potential well
- (for this example, we set $\hbar = 1$, $m = 1 \dots$)

Example: Quantum mechanics

17-9

lowest energy eigenfunctions



- potential V shown as dotted line (scaled to fit plot)
- four eigenstates with lowest energy shown (*i.e.*, v_1, v_2, v_3, v_4)

Example: Quantum mechanics

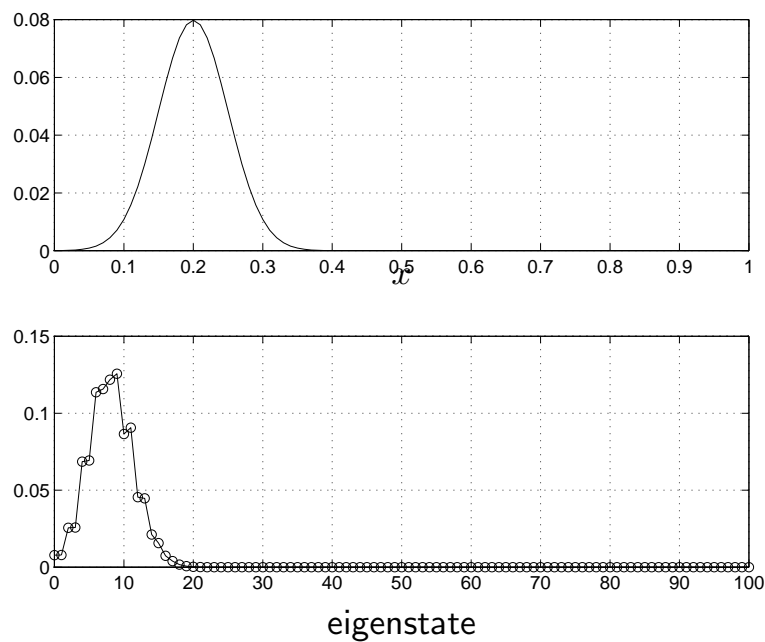
17-10

now let's look at a trajectory of Ψ , with initial wave function $\Psi(0)$

- particle near $x = 0.2$
- with momentum to right (can't see in plot of $|\Psi|^2$)
- (expected) kinetic energy half potential bump height

Example: Quantum mechanics

17-11



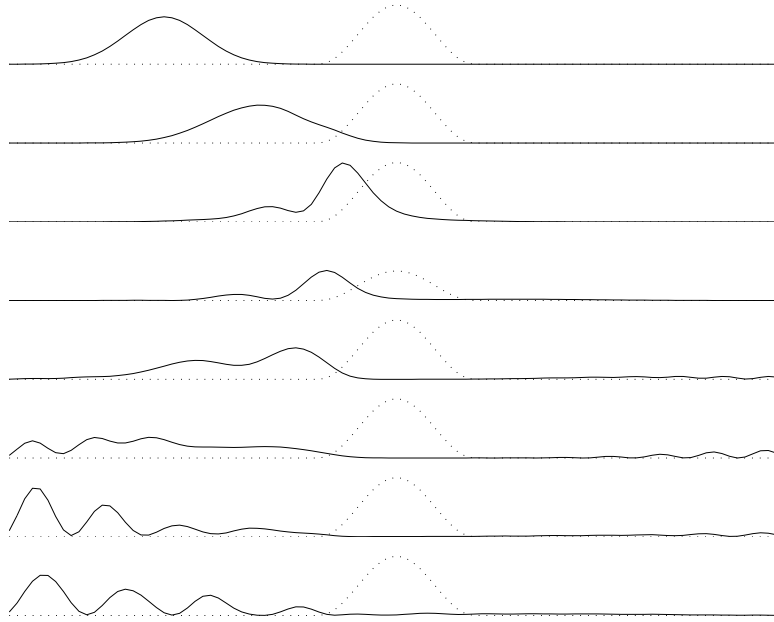
- top plot shows initial probability density $|\Psi(0)|^2$
- bottom plot shows $|v_k^* \Psi(0)|^2$, *i.e.*, resolution of $\Psi(0)$ into eigenstates

Example: Quantum mechanics

17-12

time evolution, for $t = 0, 40, 80, \dots, 320$:

$$|\Psi(t)|^2$$



Example: Quantum mechanics

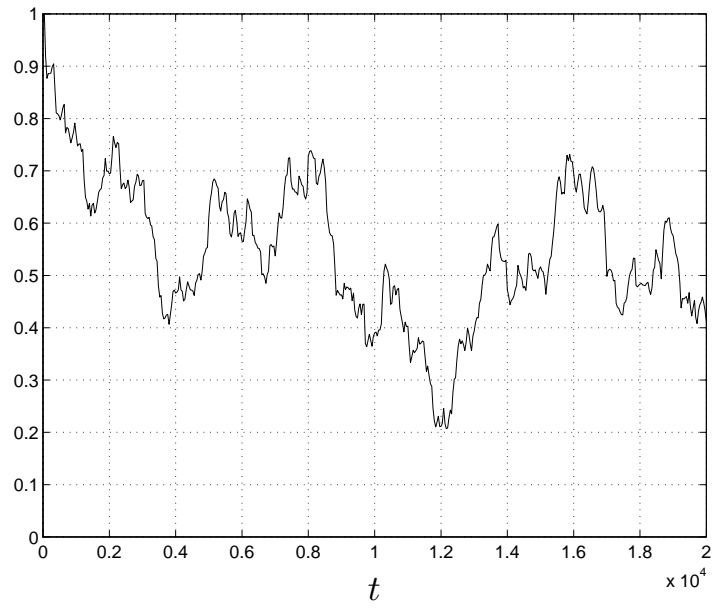
17-13

cf. classical solution:

- particle rolls half way up potential bump, stops, then rolls back down
- reverses velocity when it hits the wall at left (perfectly elastic collision)
- then repeats

Example: Quantum mechanics

17-14



plot shows probability that particle is in left half of well, *i.e.*, $\sum_{k=1}^{N/2} |\Psi_k(t)|^2$,
versus time t

Lecture 18

Controllability and state transfer

- state transfer
- reachable set, controllability matrix
- minimum norm inputs
- infinite-horizon minimum norm transfer

18-1

State transfer

consider $\dot{x} = Ax + Bu$ (or $x(t+1) = Ax(t) + Bu(t)$) over time interval $[t_i, t_f]$

we say input $u : [t_i, t_f] \rightarrow \mathbf{R}^m$ *steers* or *transfers* state from $x(t_i)$ to $x(t_f)$ (over time interval $[t_i, t_f]$)

(subscripts stand for *initial* and *final*)

questions:

- where can $x(t_i)$ be transferred to at $t = t_f$?
- how quickly can $x(t_i)$ be transferred to some x_{target} ?
- how do we find a u that transfers $x(t_i)$ to $x(t_f)$?
- how do we find a ‘small’ or ‘efficient’ u that transfers $x(t_i)$ to $x(t_f)$?

Reachability

consider state transfer from $x(0) = 0$ to $x(t)$

we say $x(t)$ is *reachable* (in t seconds or epochs)

we define $\mathcal{R}_t \subseteq \mathbf{R}^n$ as the set of points reachable in t seconds or epochs

for CT system $\dot{x} = Ax + Bu$,

$$\mathcal{R}_t = \left\{ \int_0^t e^{(t-\tau)A} Bu(\tau) d\tau \mid u : [0, t] \rightarrow \mathbf{R}^m \right\}$$

and for DT system $x(t+1) = Ax(t) + Bu(t)$,

$$\mathcal{R}_t = \left\{ \sum_{\tau=0}^{t-1} A^{t-1-\tau} Bu(\tau) \mid u(t) \in \mathbf{R}^m \right\}$$

- \mathcal{R}_t is a subspace of \mathbf{R}^n
- $\mathcal{R}_t \subseteq \mathcal{R}_s$ if $t \leq s$
(*i.e.*, can reach more points given more time)

we define the *reachable set* \mathcal{R} as the set of points reachable for some t :

$$\mathcal{R} = \bigcup_{t \geq 0} \mathcal{R}_t$$

Reachability for discrete-time LDS

DT system $x(t+1) = Ax(t) + Bu(t)$, $x(t) \in \mathbf{R}^n$

$$x(t) = \mathcal{C}_t \begin{bmatrix} u(t-1) \\ \vdots \\ u(0) \end{bmatrix}$$

where $\mathcal{C}_t = [B \quad AB \quad \dots \quad A^{t-1}B]$

so reachable set at t is $\mathcal{R}_t = \text{range}(\mathcal{C}_t)$

by C-H theorem, we can express each A^k for $k \geq n$ as linear combination of A^0, \dots, A^{n-1}

hence for $t \geq n$, $\text{range}(\mathcal{C}_t) = \text{range}(\mathcal{C}_n)$

thus we have

$$\mathcal{R}_t = \begin{cases} \text{range}(\mathcal{C}_t) & t < n \\ \text{range}(\mathcal{C}) & t \geq n \end{cases}$$

where $\mathcal{C} = \mathcal{C}_n$ is called the *controllability matrix*

- any state that can be reached can be reached by $t = n$
- the reachable set is $\mathcal{R} = \text{range}(\mathcal{C})$

Controllable system

system is called *reachable* or *controllable* if all states are reachable (*i.e.*, $\mathcal{R} = \mathbf{R}^n$)

system is reachable if and only if $\text{Rank}(\mathcal{C}) = n$

example: $x(t+1) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} x(t) + \begin{bmatrix} 1 \\ 1 \end{bmatrix} u(t)$

controllability matrix is $\mathcal{C} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$

hence system is not controllable; reachable set is

$$\mathcal{R} = \text{range}(\mathcal{C}) = \{ x \mid x_1 = x_2 \}$$

General state transfer

with $t_f > t_i$,

$$x(t_f) = A^{t_f-t_i}x(t_i) + \mathcal{C}_{t_f-t_i} \begin{bmatrix} u(t_f-1) \\ \vdots \\ u(t_i) \end{bmatrix}$$

hence can transfer $x(t_i)$ to $x(t_f) = x_{\text{des}}$

$$\Leftrightarrow x_{\text{des}} - A^{t_f-t_i}x(t_i) \in \mathcal{R}_{t_f-t_i}$$

- general state transfer reduces to reachability problem
- if system is controllable any state transfer can be achieved in $\leq n$ steps
- important special case: driving state to zero (sometimes called regulating or controlling state)

Least-norm input for reachability

assume system is reachable, $\text{Rank}(\mathcal{C}_t) = n$

to steer $x(0) = 0$ to $x(t) = x_{\text{des}}$, inputs $u(0), \dots, u(t-1)$ must satisfy

$$x_{\text{des}} = \mathcal{C}_t \begin{bmatrix} u(t-1) \\ \vdots \\ u(0) \end{bmatrix}$$

among all u that steer $x(0) = 0$ to $x(t) = x_{\text{des}}$, the one that minimizes

$$\sum_{\tau=0}^{t-1} \|u(\tau)\|^2$$

is given by

$$\begin{bmatrix} u_{\text{ln}}(t-1) \\ \vdots \\ u_{\text{ln}}(0) \end{bmatrix} = \mathcal{C}_t^T (\mathcal{C}_t \mathcal{C}_t^T)^{-1} x_{\text{des}}$$

u_{ln} is called *least-norm* or *minimum energy* input that effects state transfer

can express as

$$u_{\text{ln}}(\tau) = B^T (A^T)^{(t-1-\tau)} \left(\sum_{s=0}^{t-1} A^s B B^T (A^T)^s \right)^{-1} x_{\text{des}},$$

for $\tau = 0, \dots, t-1$

\mathcal{E}_{\min} , the minimum value of $\sum_{\tau=0}^{t-1} \|u(\tau)\|^2$ required to reach $x(t) = x_{\text{des}}$, is sometimes called *minimum energy* required to reach $x(t) = x_{\text{des}}$

$$\begin{aligned}
 \mathcal{E}_{\min} &= \sum_{\tau=0}^{t-1} \|u_{\text{ln}}(\tau)\|^2 \\
 &= \left(\mathcal{C}_t^T (\mathcal{C}_t \mathcal{C}_t^T)^{-1} x_{\text{des}} \right)^T \mathcal{C}_t^T (\mathcal{C}_t \mathcal{C}_t^T)^{-1} x_{\text{des}} \\
 &= x_{\text{des}}^T (\mathcal{C}_t \mathcal{C}_t^T)^{-1} x_{\text{des}} \\
 &= x_{\text{des}}^T \left(\sum_{\tau=0}^{t-1} A^\tau B B^T (A^T)^\tau \right)^{-1} x_{\text{des}}
 \end{aligned}$$

- $\mathcal{E}_{\min}(x_{\text{des}}, t)$ gives measure of how hard it is to reach $x(t) = x_{\text{des}}$ from $x(0) = 0$ (*i.e.*, how large a u is required)
- $\mathcal{E}_{\min}(x_{\text{des}}, t)$ gives practical measure of controllability/reachability (as function of x_{des}, t)
- ellipsoid $\{ z \mid \mathcal{E}_{\min}(z, t) \leq 1 \}$ shows points in state space reachable at t with one unit of energy
(shows directions that can be reached with small inputs, and directions that can be reached only with large inputs)

\mathcal{E}_{\min} as function of t :

if $t \geq s$ then

$$\sum_{\tau=0}^{t-1} A^\tau B B^T (A^T)^\tau \geq \sum_{\tau=0}^{s-1} A^\tau B B^T (A^T)^\tau$$

hence

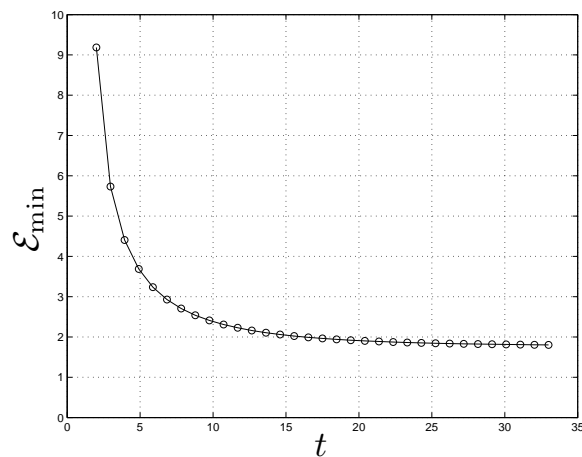
$$\left(\sum_{\tau=0}^{t-1} A^\tau B B^T (A^T)^\tau \right)^{-1} \leq \left(\sum_{\tau=0}^{s-1} A^\tau B B^T (A^T)^\tau \right)^{-1}$$

so $\mathcal{E}_{\min}(x_{\text{des}}, t) \leq \mathcal{E}_{\min}(x_{\text{des}}, s)$

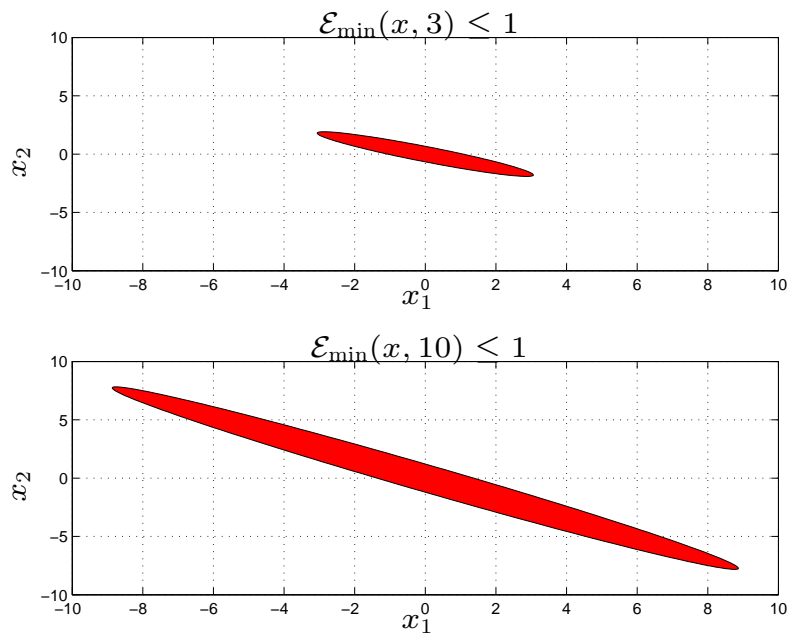
i.e.: takes less energy to get somewhere more leisurely

example: $x(t+1) = \begin{bmatrix} 1.75 & 0.8 \\ -0.95 & 0 \end{bmatrix} x(t) + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u(t)$

$\mathcal{E}_{\min}(z, t)$ for $z = [1 \ 1]^T$:



ellipsoids $\mathcal{E}_{\min} \leq 1$ for $t = 3$ and $t = 10$:



Minimum energy over infinite horizon

the matrix

$$P = \lim_{t \rightarrow \infty} \left(\sum_{\tau=0}^{t-1} A^\tau B B^T (A^T)^\tau \right)^{-1}$$

always exists, and gives the minimum energy required to reach a point x_{des} (with no limit on t):

$$\min \left\{ \sum_{\tau=0}^{t-1} \|u(\tau)\|^2 \mid x(0) = 0, x(t) = x_{\text{des}} \right\} = x_{\text{des}}^T P x_{\text{des}}$$

if A is stable, $P > 0$ (*i.e.*, can't get anywhere for free)

if A is not stable, then P can have nonzero nullspace

- $Pz = 0, z \neq 0$ means can get to z using u 's with energy as small as you like
(u just gives a little kick to the state; the instability carries it out to z efficiently)
- basis of highly maneuverable, unstable aircraft

Continuous-time reachability

consider now $\dot{x} = Ax + Bu$ with $x(t) \in \mathbf{R}^n$

reachable set at time t is

$$\mathcal{R}_t = \left\{ \int_0^t e^{(t-\tau)A} Bu(\tau) d\tau \mid u : [0, t] \rightarrow \mathbf{R}^m \right\}$$

fact: for $t > 0$, $\mathcal{R}_t = \mathcal{R} = \text{range}(\mathcal{C})$, where

$$\mathcal{C} = [B \quad AB \quad \cdots \quad A^{n-1}B]$$

is the controllability matrix of (A, B)

- same \mathcal{R} as discrete-time system
- for continuous-time system, any reachable point can be reached as fast as you like (with large enough u)

first let's show for *any* u (and $x(0) = 0$) we have $x(t) \in \text{range}(\mathcal{C})$

write e^{tA} as power series:

$$e^{tA} = I + \frac{t}{1!}A + \frac{t^2}{2!}A^2 + \dots$$

by C-H, express A^n, A^{n+1}, \dots in terms of A^0, \dots, A^{n-1} and collect powers of A :

$$e^{tA} = \alpha_0(t)I + \alpha_1(t)A + \dots + \alpha_{n-1}(t)A^{n-1}$$

therefore

$$\begin{aligned} x(t) &= \int_0^t e^{\tau A} B u(t - \tau) d\tau \\ &= \int_0^t \left(\sum_{i=0}^{n-1} \alpha_i(\tau) A^i \right) B u(t - \tau) d\tau \end{aligned}$$

$$\begin{aligned} &= \sum_{i=0}^{n-1} A^i B \int_0^t \alpha_i(\tau) u(t - \tau) d\tau \\ &= \mathcal{C} z \end{aligned}$$

where $z_i = \int_0^t \alpha_i(\tau) u(t - \tau) d\tau$

hence, $x(t)$ is always in $\text{range}(\mathcal{C})$

need to show converse: every point in $\text{range}(\mathcal{C})$ can be reached

Impulsive inputs

suppose $x(0_-) = 0$ and we apply input $u(t) = \delta^{(k)}(t)f$, where $\delta^{(k)}$ denotes k th derivative of δ and $f \in \mathbf{R}^m$

then $U(s) = s^k f$, so

$$\begin{aligned} X(s) &= (sI - A)^{-1} B s^k f \\ &= (s^{-1}I + s^{-2}A + \dots) B s^k f \\ &= (\underbrace{s^{k-1} + \dots + sA^{k-2} + A^{k-1}}_{\text{impulsive terms}} + s^{-1}A^k + \dots) B f \end{aligned}$$

hence

$$x(t) = \text{impulsive terms} + A^k B f + A^{k+1} B f \frac{t}{1!} + A^{k+2} B f \frac{t^2}{2!} + \dots$$

in particular, $x(0_+) = A^k B f$

thus, input $u = \delta^{(k)}f$ transfers state from $x(0_-) = 0$ to $x(0_+) = A^k B f$

now consider input of form

$$u(t) = \delta(t)f_0 + \dots + \delta^{(n-1)}(t)f_{n-1}$$

where $f_i \in \mathbf{R}^m$

by linearity we have

$$x(0_+) = B f_0 + \dots + A^{n-1} B f_{n-1} = \mathcal{C} \begin{bmatrix} f_0 \\ \vdots \\ f_{n-1} \end{bmatrix}$$

hence we can reach any point in $\text{range}(\mathcal{C})$

(at least, using impulse inputs)

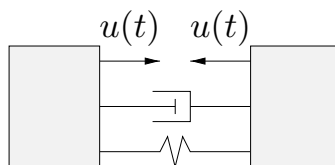
can also be shown that any point in $\text{range}(\mathcal{C})$ can be reached for any $t > 0$ using *nonimpulsive* inputs

fact: if $x(0) \in \mathcal{R}$, then $x(t) \in \mathcal{R}$ for all t (no matter what u is)

to show this, need to show $e^{tA}x(0) \in \mathcal{R}$ if $x(0) \in \mathcal{R} \dots$

Example

- unit masses at y_1, y_2 , connected by unit springs, dampers
- input is tension between masses
- state is $x = [y^T \dot{y}^T]^T$



system is

$$\dot{x} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -1 & 1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix} x + \begin{bmatrix} 0 \\ 0 \\ 1 \\ -1 \end{bmatrix} u$$

- can we maneuver state anywhere, starting from $x(0) = 0$?
- if not, where can we maneuver state?

controllability matrix is

$$C = [B \quad AB \quad A^2B \quad A^3B] = \begin{bmatrix} 0 & 1 & -2 & 2 \\ 0 & -1 & 2 & -2 \\ 1 & -2 & 2 & 0 \\ -1 & 2 & -2 & 0 \end{bmatrix}$$

hence reachable set is

$$\mathcal{R} = \text{span} \left\{ \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ -1 \end{bmatrix} \right\}$$

we can reach states with $y_1 = -y_2$, $\dot{y}_1 = -\dot{y}_2$, *i.e.*, precisely the differential motions

it's obvious — internal force does not affect center of mass position or total momentum!

Least-norm input for reachability

(also called *minimum energy input*)

assume that $\dot{x} = Ax + Bu$ is reachable

we seek u that steers $x(0) = 0$ to $x(t) = x_{\text{des}}$ and minimizes

$$\int_0^t \|u(\tau)\|^2 d\tau$$

let's discretize system with interval $h = t/N$

(we'll let $N \rightarrow \infty$ later)

thus u is piecewise constant:

$$u(\tau) = u_d(k) \quad \text{for } kh \leq \tau < (k+1)h, \quad k = 0, \dots, N-1$$

so

$$x(t) = \begin{bmatrix} B_d & A_d B_d & \cdots & A_d^{N-1} B_d \end{bmatrix} \begin{bmatrix} u_d(N-1) \\ \vdots \\ u_d(0) \end{bmatrix}$$

where

$$A_d = e^{hA}, \quad B_d = \int_0^h e^{\tau A} d\tau B$$

least-norm u_d that yields $x(t) = x_{\text{des}}$ is

$$u_{\text{dln}}(k) = B_d^T (A_d^T)^{(N-1-k)} \left(\sum_{i=0}^{N-1} A_d^i B_d B_d^T (A_d^T)^i \right)^{-1} x_{\text{des}}$$

let's express in terms of A :

$$B_d^T (A_d^T)^{(N-1-k)} = B_d^T e^{(t-\tau)A^T}$$

where $\tau = t(k+1)/N$

for N large, $B_d \approx (t/N)B$, so this is approximately

$$(t/N)B^T e^{(t-\tau)A^T}$$

similarly

$$\begin{aligned} \sum_{i=0}^{N-1} A_d^i B_d B_d^T (A_d^T)^i &= \sum_{i=0}^{N-1} e^{(ti/N)A} B_d B_d^T e^{(ti/N)A^T} \\ &\approx (t/N) \int_0^t e^{\bar{t}A} B B^T e^{\bar{t}A^T} d\bar{t} \end{aligned}$$

for large N

hence least-norm discretized input is approximately

$$u_{\text{ln}}(\tau) = B^T e^{(t-\tau)A^T} \left(\int_0^t e^{\bar{t}A} B B^T e^{\bar{t}A^T} d\bar{t} \right)^{-1} x_{\text{des}}, \quad 0 \leq \tau \leq t$$

for large N

hence, this is the least-norm continuous input

- can make t small, but get larger u
- cf. DT solution: sum becomes integral

min energy is

$$\int_0^t \|u_{\text{ln}}(\tau)\|^2 d\tau = x_{\text{des}}^T Q(t)^{-1} x_{\text{des}}$$

where

$$Q(t) = \int_0^t e^{\tau A} B B^T e^{\tau A^T} d\tau$$

can show

$$\begin{aligned} (A, B) \text{ controllable} &\Leftrightarrow Q(t) > 0 \text{ for all } t > 0 \\ &\Leftrightarrow Q(s) > 0 \text{ for some } s > 0 \end{aligned}$$

in fact, $\text{range}(Q(t)) = \mathcal{R}$ for any $t > 0$

Minimum energy over infinite horizon

the matrix

$$P = \lim_{t \rightarrow \infty} \left(\int_0^t e^{\tau A} B B^T e^{\tau A^T} d\tau \right)^{-1}$$

always exists, and gives minimum energy required to reach a point x_{des} (with no limit on t):

$$\min \left\{ \int_0^t \|u(\tau)\|^2 d\tau \mid x(0) = 0, x(t) = x_{\text{des}} \right\} = x_{\text{des}}^T P x_{\text{des}}$$

- if A is stable, $P > 0$ (*i.e.*, can't get anywhere for free)
- if A is not stable, then P can have nonzero nullspace
- $Pz = 0, z \neq 0$ means can get to z using u 's with energy as small as you like (u just gives a little kick to the state; the instability carries it out to z efficiently)

General state transfer

consider state transfer from $x(t_i)$ to $x(t_f) = x_{\text{des}}, t_f > t_i$

since

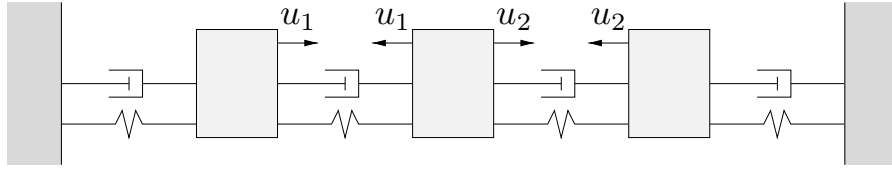
$$x(t_f) = e^{(t_f - t_i)A} x(t_i) + \int_{t_i}^{t_f} e^{(t_f - \tau)A} B u(\tau) d\tau$$

u steers $x(t_i)$ to $x(t_f) = x_{\text{des}} \Leftrightarrow$

u (shifted by t_i) steers $x(0) = 0$ to $x(t_f - t_i) = x_{\text{des}} - e^{(t_f - t_i)A} x(t_i)$

- general state transfer reduces to reachability problem
- if system is controllable, any state transfer can be effected
 - in 'zero' time with impulsive inputs
 - in any positive time with non-impulsive inputs

Example



- unit masses, springs, dampers
- u_1 is force between 1st & 2nd masses
- u_2 is force between 2nd & 3rd masses
- $y \in \mathbf{R}^3$ is displacement of masses 1,2,3
- $x = \begin{bmatrix} y \\ \dot{y} \end{bmatrix}$

Controllability and state transfer

18-33

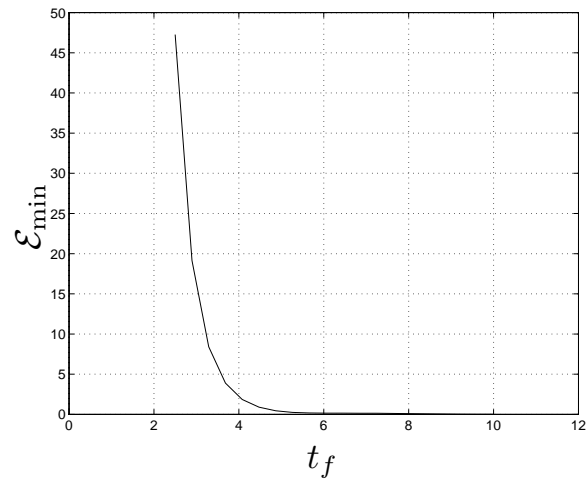
system is:

$$\dot{x} = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ -2 & 1 & 0 & -2 & 1 & 0 \\ 1 & -2 & 1 & 1 & -2 & 1 \\ 0 & 1 & -2 & 0 & 1 & -2 \end{bmatrix} x + \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ -1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

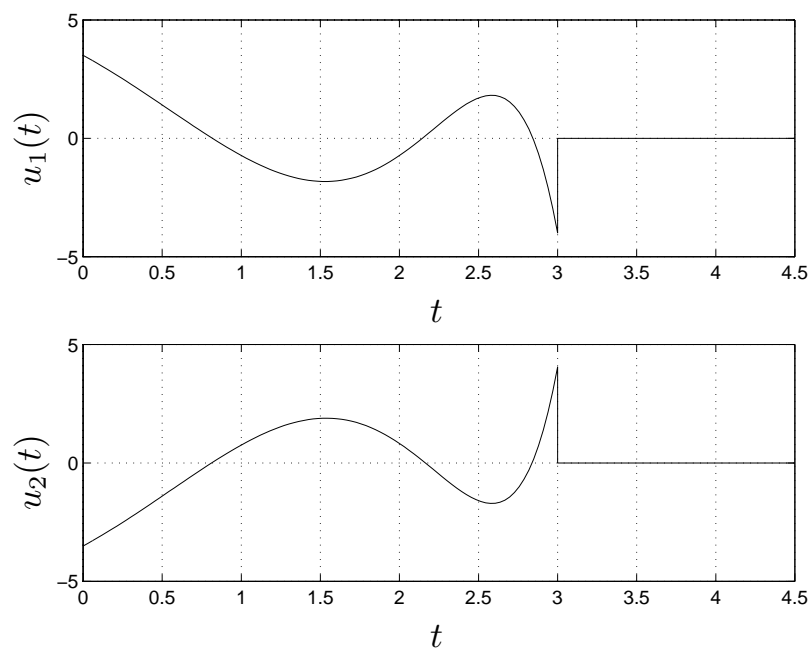
steer state from $x(0) = e_1$ to $x(t_f) = 0$

i.e., control initial state e_1 to zero at $t = t_f$

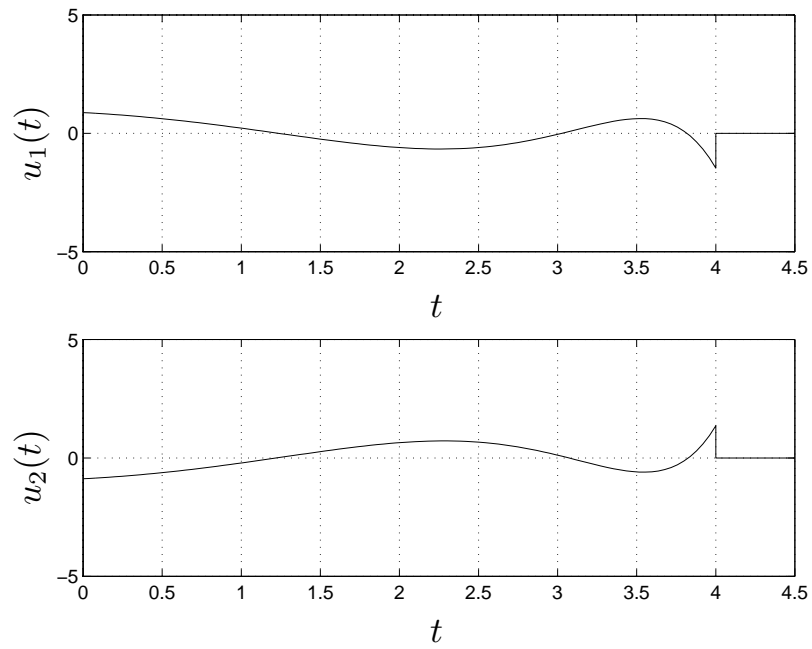
$$\mathcal{E}_{\min} = \int_0^{t_f} \|u_{\text{ln}}(\tau)\|^2 d\tau \text{ vs. } t_f:$$



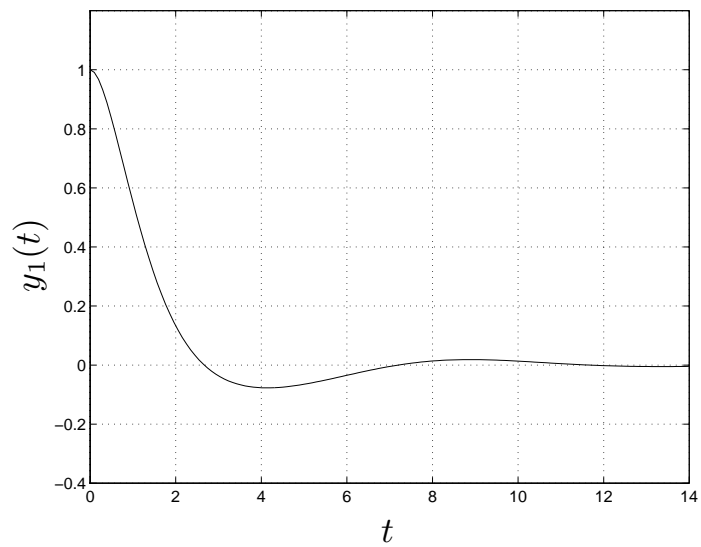
for $t_f = 3$, $u = u_{\text{ln}}$ is:



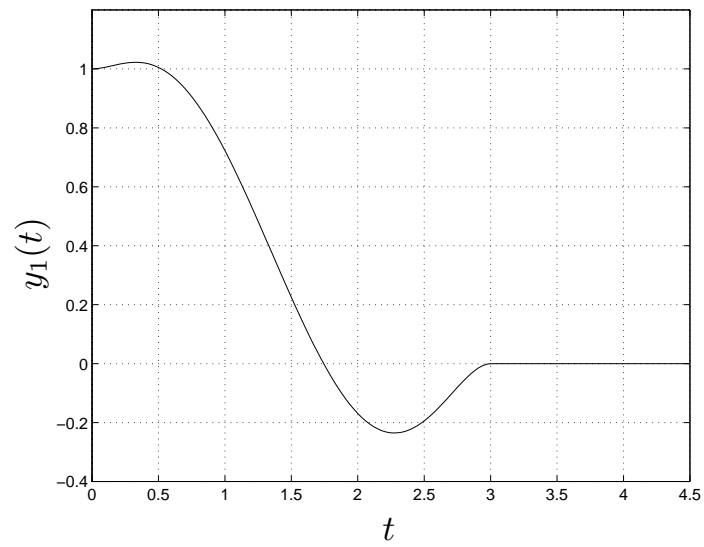
and for $t_f = 4$:



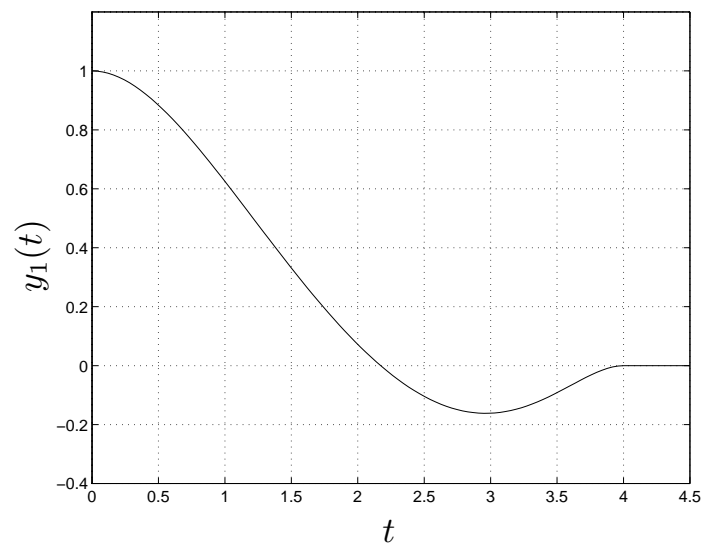
output y_1 for $u = 0$:



output y_1 for $u = u_{\text{in}}$ with $t_f = 3$:



output y_1 for $u = u_{\text{in}}$ with $t_f = 4$:



Lecture 19

Observability and state estimation

- state estimation
- discrete-time observability
- observability – controllability duality
- observers for noiseless case
- continuous-time observability
- least-squares observers
- example

19-1

State estimation set up

we consider the discrete-time system

$$x(t+1) = Ax(t) + Bu(t) + w(t), \quad y(t) = Cx(t) + Du(t) + v(t)$$

- w is state *disturbance* or *noise*
- v is sensor *noise* or *error*
- A , B , C , and D are known
- u and y are observed over time interval $[0, t-1]$
- w and v are not known, but can be described statistically, or assumed small (*e.g.*, in RMS value)

State estimation problem

state estimation problem: estimate $x(s)$ from

$$u(0), \dots, u(t-1), y(0), \dots, y(t-1)$$

- $s = 0$: estimate initial state
- $s = t - 1$: estimate current state
- $s = t$: estimate (*i.e.*, predict) next state

an algorithm or system that yields an estimate $\hat{x}(s)$ is called an *observer* or *state estimator*

$\hat{x}(s)$ is denoted $\hat{x}(s|t-1)$ to show what information estimate is based on (read, “ $\hat{x}(s)$ given $t-1$ ”)

Noiseless case

let's look at finding $x(0)$, with no state or measurement noise:

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t)$$

with $x(t) \in \mathbf{R}^n$, $u(t) \in \mathbf{R}^m$, $y(t) \in \mathbf{R}^p$

then we have

$$\begin{bmatrix} y(0) \\ \vdots \\ y(t-1) \end{bmatrix} = \mathcal{O}_t x(0) + \mathcal{I}_t \begin{bmatrix} u(0) \\ \vdots \\ u(t-1) \end{bmatrix}$$

where

$$\mathcal{O}_t = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{t-1} \end{bmatrix}, \quad \mathcal{T}_t = \begin{bmatrix} D & 0 & \cdots & \cdots \\ CB & D & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ CA^{t-2}B & CA^{t-3}B & \cdots & CB & D \end{bmatrix}$$

- \mathcal{O}_t maps initial state into resulting output over $[0, t-1]$
- \mathcal{T}_t maps input to output over $[0, t-1]$

hence we have

$$\mathcal{O}_t x(0) = \begin{bmatrix} y(0) \\ \vdots \\ y(t-1) \end{bmatrix} - \mathcal{T}_t \begin{bmatrix} u(0) \\ \vdots \\ u(t-1) \end{bmatrix}$$

RHS is known, $x(0)$ is to be determined

hence:

- can uniquely determine $x(0)$ if and only if $\mathcal{N}(\mathcal{O}_t) = \{0\}$
- $\mathcal{N}(\mathcal{O}_t)$ gives ambiguity in determining $x(0)$
- if $x(0) \in \mathcal{N}(\mathcal{O}_t)$ and $u = 0$, output is zero over interval $[0, t-1]$
- input u does not affect ability to determine $x(0)$;
its effect can be subtracted out

Observability matrix

by C-H theorem, each A^k is linear combination of A^0, \dots, A^{n-1}

hence for $t \geq n$, $\mathcal{N}(\mathcal{O}_t) = \mathcal{N}(\mathcal{O})$ where

$$\mathcal{O} = \mathcal{O}_n = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix}$$

is called the *observability matrix*

if $x(0)$ can be deduced from u and y over $[0, t-1]$ for any t , then $x(0)$ can be deduced from u and y over $[0, n-1]$

$\mathcal{N}(\mathcal{O})$ is called *unobservable subspace*; describes ambiguity in determining state from input and output

system is called *observable* if $\mathcal{N}(\mathcal{O}) = \{0\}$, i.e., $\mathbf{Rank}(\mathcal{O}) = n$

Observability – controllability duality

let $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ be dual of system (A, B, C, D) , i.e.,

$$\tilde{A} = A^T, \quad \tilde{B} = C^T, \quad \tilde{C} = B^T, \quad \tilde{D} = D^T$$

controllability matrix of dual system is

$$\begin{aligned} \tilde{\mathcal{C}} &= [\tilde{B} \ \tilde{A}\tilde{B} \ \dots \ \tilde{A}^{n-1}\tilde{B}] \\ &= [C^T \ A^T C^T \ \dots \ (A^T)^{n-1} C^T] \\ &= \mathcal{O}^T, \end{aligned}$$

transpose of observability matrix

similarly we have $\tilde{\mathcal{O}} = \mathcal{C}^T$

thus, system is observable (controllable) if and only if dual system is controllable (observable)

in fact,

$$\mathcal{N}(\mathcal{O}) = \text{range}(\mathcal{O}^T)^\perp = \text{range}(\tilde{\mathcal{C}})^\perp$$

i.e., unobservable subspace is orthogonal complement of controllable subspace of dual

Observers for noiseless case

suppose $\text{Rank}(\mathcal{O}_t) = n$ (*i.e.*, system is observable) and let F be any left inverse of \mathcal{O}_t , *i.e.*, $F\mathcal{O}_t = I$

then we have the observer

$$x(0) = F \left(\begin{bmatrix} y(0) \\ \vdots \\ y(t-1) \end{bmatrix} - \mathcal{T}_t \begin{bmatrix} u(0) \\ \vdots \\ u(t-1) \end{bmatrix} \right)$$

which deduces $x(0)$ (exactly) from u, y over $[0, t-1]$

in fact we have

$$x(\tau - t + 1) = F \left(\begin{bmatrix} y(\tau - t + 1) \\ \vdots \\ y(\tau) \end{bmatrix} - \mathcal{T}_t \begin{bmatrix} u(\tau - t + 1) \\ \vdots \\ u(\tau) \end{bmatrix} \right)$$

i.e., our observer estimates what state was $t - 1$ epochs ago, given past $t - 1$ inputs & outputs

observer is (multi-input, multi-output) *finite impulse response* (FIR) filter, with inputs u and y , and output \hat{x}

Invariance of unobservable set

fact: the unobservable subspace $\mathcal{N}(\mathcal{O})$ is invariant, *i.e.*, if $z \in \mathcal{N}(\mathcal{O})$, then $Az \in \mathcal{N}(\mathcal{O})$

proof: suppose $z \in \mathcal{N}(\mathcal{O})$, *i.e.*, $CA^k z = 0$ for $k = 0, \dots, n - 1$

evidently $CA^k(Az) = 0$ for $k = 0, \dots, n - 2$;

$$CA^{n-1}(Az) = CA^n z = - \sum_{i=0}^{n-1} \alpha_i CA^i z = 0$$

(by C-H) where

$$\det(sI - A) = s^n + \alpha_{n-1}s^{n-1} + \dots + \alpha_0$$

Continuous-time observability

continuous-time system with no sensor or state noise:

$$\dot{x} = Ax + Bu, \quad y = Cx + Du$$

can we deduce state x from u and y ?

let's look at derivatives of y :

$$y = Cx + Du$$

$$\dot{y} = C\dot{x} + D\dot{u} = CAx + CBu + D\dot{u}$$

$$\ddot{y} = CA^2x + CABu + CB\dot{u} + D\ddot{u}$$

and so on

hence we have

$$\begin{bmatrix} y \\ \dot{y} \\ \vdots \\ y^{(n-1)} \end{bmatrix} = \mathcal{O}x + \mathcal{T} \begin{bmatrix} u \\ \dot{u} \\ \vdots \\ u^{(n-1)} \end{bmatrix}$$

where \mathcal{O} is the observability matrix and

$$\mathcal{T} = \begin{bmatrix} D & 0 & \cdots & \cdots \\ CB & D & 0 & \cdots \\ \vdots & & & \\ CA^{n-2}B & CA^{n-3}B & \cdots & CB & D \end{bmatrix}$$

(same matrices we encountered in discrete-time case!)

rewrite as

$$\mathcal{O}x = \begin{bmatrix} y \\ \dot{y} \\ \vdots \\ y^{(n-1)} \end{bmatrix} - \mathcal{T} \begin{bmatrix} u \\ \dot{u} \\ \vdots \\ u^{(n-1)} \end{bmatrix}$$

RHS is known; x is to be determined

hence if $\mathcal{N}(\mathcal{O}) = \{0\}$ we can deduce $x(t)$ from derivatives of $u(t)$, $y(t)$ up to order $n - 1$

in this case we say system is observable

can construct an observer using any left inverse F of \mathcal{O} :

$$x = F \left(\begin{bmatrix} y \\ \dot{y} \\ \vdots \\ y^{(n-1)} \end{bmatrix} - \mathcal{T} \begin{bmatrix} u \\ \dot{u} \\ \vdots \\ u^{(n-1)} \end{bmatrix} \right)$$

- reconstructs $x(t)$ (exactly and instantaneously) from

$$u(t), \dots, u^{(n-1)}(t), y(t), \dots, y^{(n-1)}(t)$$

- derivative-based state reconstruction is dual of state transfer using impulsive inputs

A converse

suppose $z \in \mathcal{N}(\mathcal{O})$ (the unobservable subspace), and u is any input, with x, y the corresponding state and output, *i.e.*,

$$\dot{x} = Ax + Bu, \quad y = Cx + Du$$

then state trajectory $\tilde{x} = x + e^{tA}z$ satisfies

$$\dot{\tilde{x}} = A\tilde{x} + Bu, \quad y = C\tilde{x} + Du$$

i.e., input/output signals u, y consistent with both state trajectories x, \tilde{x}

hence if system is unobservable, no signal processing of any kind applied to u and y can deduce x

unobservable subspace $\mathcal{N}(\mathcal{O})$ gives fundamental ambiguity in deducing x from u, y

Least-squares observers

discrete-time system, with sensor noise:

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t) + v(t)$$

we assume $\mathbf{Rank}(\mathcal{O}_t) = n$ (hence, system is observable)

least-squares observer uses pseudo-inverse:

$$\hat{x}(0) = \mathcal{O}_t^\dagger \left(\begin{bmatrix} y(0) \\ \vdots \\ y(t-1) \end{bmatrix} - \mathcal{T}_t \begin{bmatrix} u(0) \\ \vdots \\ u(t-1) \end{bmatrix} \right)$$

where $\mathcal{O}_t^\dagger = (\mathcal{O}_t^T \mathcal{O}_t)^{-1} \mathcal{O}_t^T$

interpretation: $\hat{x}_{ls}(0)$ minimizes discrepancy between

- output \hat{y} that *would be* observed, with input u and initial state $x(0)$ (and no sensor noise), and
- output y that *was* observed,

measured as $\sum_{\tau=0}^{t-1} \|\hat{y}(\tau) - y(\tau)\|^2$

can express least-squares initial state estimate as

$$\hat{x}_{ls}(0) = \left(\sum_{\tau=0}^{t-1} (A^T)^\tau C^T C A^\tau \right)^{-1} \sum_{\tau=0}^{t-1} (A^T)^\tau C^T \tilde{y}(\tau)$$

where \tilde{y} is observed output with portion due to input subtracted:
 $\tilde{y} = y - h * u$ where h is impulse response

Least-squares observer uncertainty ellipsoid

since $\mathcal{O}_t^\dagger \mathcal{O}_t = I$, we have

$$\tilde{x}(0) = \hat{x}_{ls}(0) - x(0) = \mathcal{O}_t^\dagger \begin{bmatrix} v(0) \\ \vdots \\ v(t-1) \end{bmatrix}$$

where $\tilde{x}(0)$ is the estimation error of the initial state

in particular, $\hat{x}_{ls}(0) = x(0)$ if sensor noise is zero
(i.e., observer recovers exact state in noiseless case)

now assume sensor noise is unknown, but has RMS value $\leq \alpha$,

$$\frac{1}{t} \sum_{\tau=0}^{t-1} \|v(\tau)\|^2 \leq \alpha^2$$

set of possible estimation errors is ellipsoid

$$\tilde{x}(0) \in \mathcal{E}_{\text{unc}} = \left\{ \mathcal{O}_t^\dagger \begin{bmatrix} v(0) \\ \vdots \\ v(t-1) \end{bmatrix} \mid \frac{1}{t} \sum_{\tau=0}^{t-1} \|v(\tau)\|^2 \leq \alpha^2 \right\}$$

\mathcal{E}_{unc} is 'uncertainty ellipsoid' for $x(0)$ (least-square gives best \mathcal{E}_{unc})

shape of uncertainty ellipsoid determined by matrix

$$(\mathcal{O}_t^T \mathcal{O}_t)^{-1} = \left(\sum_{\tau=0}^{t-1} (A^T)^\tau C^T C A^\tau \right)^{-1}$$

maximum norm of error is

$$\|\hat{x}_{\text{ls}}(0) - x(0)\| \leq \alpha \sqrt{t} \|\mathcal{O}_t^\dagger\|$$

Infinite horizon uncertainty ellipsoid

the matrix

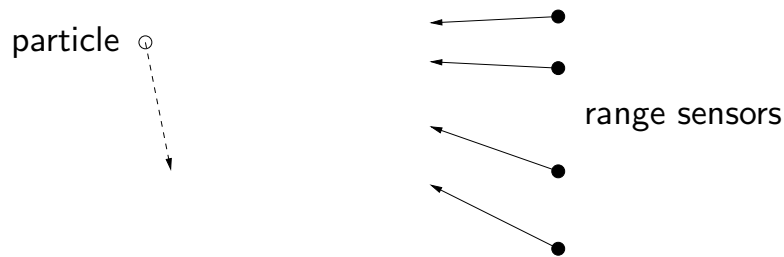
$$P = \lim_{t \rightarrow \infty} \left(\sum_{\tau=0}^{t-1} (A^T)^\tau C^T C A^\tau \right)^{-1}$$

always exists, and gives the limiting uncertainty in estimating $x(0)$ from u , y over longer and longer periods:

- if A is stable, $P > 0$
i.e., can't estimate initial state perfectly even with infinite number of measurements $u(t)$, $y(t)$, $t = 0, \dots$ (since memory of $x(0)$ fades . . .)
- if A is not stable, then P can have nonzero nullspace
i.e., initial state estimation error gets arbitrarily small (at least in some directions) as more and more of signals u and y are observed

Example

- particle in \mathbf{R}^2 moves with uniform velocity
- (linear, noisy) range measurements from directions $-15^\circ, 0^\circ, 20^\circ, 30^\circ$, once per second
- range noises IID $\mathcal{N}(0, 1)$; can assume RMS value of v is not much more than 2
- no assumptions about initial position & velocity



problem: estimate initial position & velocity from range measurements

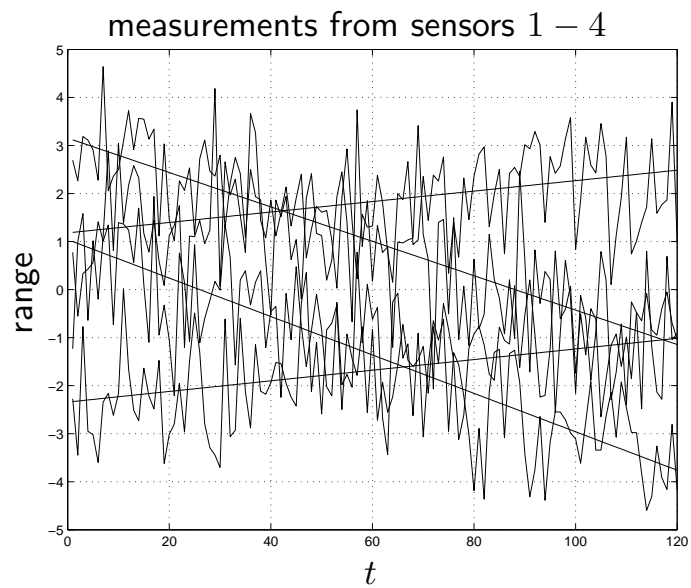
express as linear system

$$x(t+1) = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} x(t), \quad y(t) = \begin{bmatrix} k_1^T \\ \vdots \\ k_4^T \end{bmatrix} x(t) + v(t)$$

- $(x_1(t), x_2(t))$ is position of particle
- $(x_3(t), x_4(t))$ is velocity of particle
- can assume RMS value of v is around 2
- k_i is unit vector from sensor i to origin

true initial position & velocities: $x(0) = (1 \quad -3 \quad -0.04 \quad 0.03)$

range measurements (& noiseless versions):

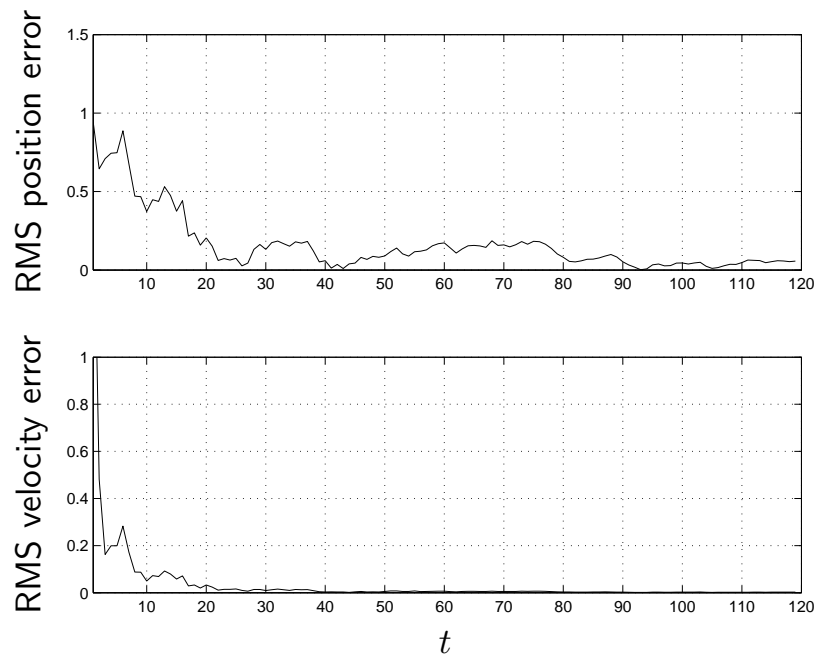


- estimate based on $(y(0), \dots, y(t))$ is $\hat{x}(0|t)$

- actual RMS position error is

$$\sqrt{(\hat{x}_1(0|t) - x_1(0))^2 + (\hat{x}_2(0|t) - x_2(0))^2}$$

(similarly for actual RMS velocity error)



Continuous-time least-squares state estimation

assume $\dot{x} = Ax + Bu$, $y = Cx + Du + v$ is observable

least-squares estimate of initial state $x(0)$, given $u(\tau)$, $y(\tau)$, $0 \leq \tau \leq t$:
choose $\hat{x}_{ls}(0)$ to minimize integral square residual

$$J = \int_0^t \|\tilde{y}(\tau) - Ce^{\tau A}x(0)\|^2 d\tau$$

where $\tilde{y} = y - h * u$ is observed output minus part due to input

let's expand as $J = x(0)^T Q x(0) + 2r^T x(0) + s$,

$$Q = \int_0^t e^{\tau A^T} C^T C e^{\tau A} d\tau, \quad r = \int_0^t e^{\tau A^T} C^T \tilde{y}(\tau) d\tau,$$

$$q = \int_0^t \tilde{y}(\tau)^T \tilde{y}(\tau) d\tau$$

setting $\nabla_{x(0)} J$ to zero, we obtain the least-squares observer

$$\hat{x}_{\text{ls}}(0) = Q^{-1} r = \left(\int_0^t e^{\tau A^T} C^T C e^{\tau A} d\tau \right)^{-1} \int_0^t e^{A^T \tau} C^T \tilde{y}(\tau) d\tau$$

estimation error is

$$\tilde{x}(0) = \hat{x}_{\text{ls}}(0) - x(0) = \left(\int_0^t e^{\tau A^T} C^T C e^{\tau A} d\tau \right)^{-1} \int_0^t e^{\tau A^T} C^T v(\tau) d\tau$$

therefore if $v = 0$ then $\hat{x}_{\text{ls}}(0) = x(0)$

Lecture 20

Some parting thoughts . . .

- linear algebra
- levels of understanding
- what's next?

20-1

Linear algebra

- comes up in *many* practical contexts (EE, ME, CE, AA, OR, Econ, . . .)
- nowadays is readily *done*
cf. 10 yrs ago (when it was mostly *talked about*)
- Matlab or equiv for fooling around
- real codes (*e.g.*, LAPACK) widely available
- current level of linear algebra technology:
 - 500 – 1000 vbles: easy with general purpose codes
 - much more possible with special structure, special codes (*e.g.*, sparse, convolution, banded, . . .)

Levels of understanding

Simple, intuitive view:

- 17 vbles, 17 eqns: usually has unique solution
- 80 vbles, 60 eqns: 20 extra degrees of freedom

Platonic view:

- singular, rank, range, nullspace, Jordan form, controllability
- everything is precise & unambiguous
- gives insight & deeper understanding
- sometimes misleading in practice

Some parting thoughts . . .

20-3

Quantitative view:

- based on ideas like least-squares, SVD
- gives numerical measures for ideas like singularity, rank, etc.
- interpretation depends on (practical) context
- very useful in practice

Some parting thoughts . . .

20-4

- must have understanding at one level before moving to next
- **never forget** which level you are operating in

What's next?

- EE363 — linear dynamical systems (Win 08-09)
- EE364a — convex optimization I (Spr 08-09)

(plus lots of other EE, CS, ICME, MS&E, Stat, ME, AA courses on signal processing, control, graphics & vision, adaptive systems, machine learning, computational geometry, numerical linear algebra, . . .)

Basic Notation

Basic set notation

$\{a_1, \dots, a_r\}$	the set with elements a_1, \dots, a_r .
$a \in S$	a is in the set S .
$S = T$	the sets S and T are equal, <i>i.e.</i> , every element of S is in T and every element of T is in S .
$S \subseteq T$	the set S is a subset of the set T , <i>i.e.</i> , every element of S is also an element of T .
$\exists a \in S \mathcal{P}(a)$	there exists an a in S for which the property \mathcal{P} holds.
$\forall x \in S \mathcal{P}(x)$	property \mathcal{P} holds for every element in S .
$\{a \in S \mid \mathcal{P}(a)\}$	the set of all a in S for which \mathcal{P} holds (the set S is sometimes omitted if it can be determined from context).
$A \cup B$	union of sets, $A \cup B = \{x \mid x \in A \text{ or } x \in B\}$.
$A \cap B$	intersection of sets, $A \cap B = \{x \mid x \in A \text{ and } x \in B\}$.
$A \times B$	Cartesian product of two sets, $A \times B = \{(a, b) \mid a \in A, b \in B\}$.

Some specific sets

\mathbf{R}	the set of real numbers.
\mathbf{R}^n	the set of real n -vectors ($n \times 1$ matrices).
$\mathbf{R}^{1 \times n}$	the set of real n -row-vectors ($1 \times n$ matrices).
$\mathbf{R}^{m \times n}$	the set of real $m \times n$ matrices.
j	can mean $\sqrt{-1}$, in the company of electrical engineers.
i	can mean $\sqrt{-1}$, for normal people; i is the polite term in mixed company (<i>i.e.</i> , when non-electrical engineers are present).
$\mathbf{C}, \mathbf{C}^n, \mathbf{C}^{m \times n}$	the set of complex numbers, complex n -vectors, complex $m \times n$ matrices.
\mathbf{Z}	the set of integers: $\mathbf{Z} = \{\dots, -1, 0, 1, \dots\}$.
\mathbf{R}_+	the nonnegative real numbers, <i>i.e.</i> , $\mathbf{R}_+ = \{x \in \mathbf{R} \mid x \geq 0\}$.
$[a, b], (a, b), [a, b), (a, b)$	the real intervals $\{x \mid a \leq x \leq b\}$, $\{x \mid a < x \leq b\}$, $\{x \mid a \leq x < b\}$, and $\{x \mid a < x < b\}$, respectively.

Vectors and matrices

We use square brackets [and] to construct matrices and vectors, with white space delineating the entries in a row, and a new line indicating a new row. For example [1 2] is a row vector in $\mathbf{R}^{1 \times 2}$, and $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$ is matrix in $\mathbf{R}^{2 \times 3}$. $[1 \ 2]^T$ denotes a column vector, *i.e.*, an element of $\mathbf{R}^{2 \times 1}$, which we abbreviate as \mathbf{R}^2 .

We use curved brackets (and) surrounding lists of entries, delineated by commas, as an alternative method to construct (column) vectors. Thus, we have three ways to denote a column vector:

$$(1, 2) = [1 \ 2]^T = \begin{bmatrix} 1 \\ 2 \end{bmatrix}.$$

Note that in our notation scheme (which is fairly standard), $[1, 2, 3]$ and $(1 \ 2 \ 3)$ aren't used. We also use square and curved brackets to construct block matrices and vectors. For example if $x, y, z \in \mathbf{R}^n$, we have

$$\begin{bmatrix} x & y & z \end{bmatrix} \in \mathbf{R}^{n \times 3},$$

a matrix with columns x, y , and z . We can construct a block vector as

$$(x, y, z) = \begin{bmatrix} x \\ y \\ z \end{bmatrix} \in \mathbf{R}^{3n}.$$

Functions

The notation $f : A \rightarrow B$ means that f is a function on the set A into the set B . The notation $b = f(a)$ means b is the value of the function f at the point a , where $a \in A$ and $b \in B$. The set A is called the *domain* of the function f ; it can thought of as the set of legal parameter values that can be passed to the function f . The set B is called the *codomain* (or sometimes range) of the function f ; it can thought of as a set that contains all possible returned values of the function f .

There are several ways to think of a function. The formal definition is that f is a subset of $A \times B$, with the property that for every $a \in A$, there is exactly one $b \in B$ such that $(a, b) \in f$. We denote this as $b = f(a)$.

Perhaps a better way to think of a function is as a *black box* or (software) *function* or *subroutine*. The domain is the set of all legal values (or data types or structures) that can be passed to f . The codomain of f gives the data type or data structure of the values returned by f .

Thus $f(a)$ is *meaningless* if $a \notin A$. If $a \in A$, then $b = f(a)$ is an element of B . Also note that the *function* is denoted f ; it is *wrong* to say 'the function $f(a)$ ' (since $f(a)$ is an element

of B , not a function). Having said that, we do sometimes use sloppy notation such as ‘the function $f(t) = t^3$ ’. To say this more clearly you could say ‘the function $f : \mathbf{R} \rightarrow \mathbf{R}$ defined by $f(t) = t^3$ for $t \in \mathbf{R}$ ’.

Examples

- $-0.1 \in \mathbf{R}$, $\sqrt{2} \in \mathbf{R}_+$, $1 - 2j \in \mathbf{C}$ (with $j = \sqrt{-1}$).
- The matrix

$$A = \begin{bmatrix} 0.3 & 6.1 & -0.12 \\ 7.2 & 0 & 0.01 \end{bmatrix}$$

is an element in $\mathbf{R}^{2 \times 3}$. We can define a function $f : \mathbf{R}^3 \rightarrow \mathbf{R}^2$ as $f(x) = Ax$ for any $x \in \mathbf{R}^3$. If $x \in \mathbf{R}^3$, then $f(x)$ is a particular vector in \mathbf{R}^2 . We can say ‘the function f is linear’. To say ‘the function $f(x)$ is linear’ is technically wrong since $f(x)$ is a vector, not a function. Similarly we can’t say ‘ A is linear’; it is just a matrix.

- We can define a function $f : \{a \in \mathbf{R} \mid a \neq 0\} \times \mathbf{R}^n \rightarrow \mathbf{R}^n$ by $f(a, x) = (1/a)x$, for any $a \in \mathbf{R}$, $a \neq 0$, and any $x \in \mathbf{R}^n$. The function f could be informally described as division of a vector by a nonzero scalar.
- Consider the set $A = \{0, -1, 3.2\}$. The elements of A are 0, -1 and 3.2. Therefore, for example, $-1 \in A$ and $\{0, 3.2\} \subseteq A$. Also, we can say that $\forall x \in A, -1 \leq x \leq 4$ or $\exists x \in A, x > 3$.
- Suppose $A = \{1, -1\}$. Another representation for A is $A = \{x \in \mathbf{R} \mid x^2 = 1\}$.
- Suppose $A = \{1, -2, 0\}$ and $B = \{3, -2\}$. Then

$$A \cup B = \{1, -2, 0, 3\}, \quad A \cap B = \{-2\}.$$

- Suppose $A = \{1, -2, 0\}$ and $B = \{1, 3\}$. Then

$$A \times B = \{(1, 1), (1, 3), (-2, 1), (-2, 3), (0, 1), (0, 3)\}.$$

- $f : \mathbf{R} \rightarrow \mathbf{R}$ with $f(x) = x^2 - x$ defines a function from \mathbf{R} to \mathbf{R} while $u : \mathbf{R}_+ \rightarrow \mathbf{R}^2$ with

$$u(t) = \begin{bmatrix} t \cos t \\ 2e^{-t} \end{bmatrix}.$$

defines a function from \mathbf{R}_+ to \mathbf{R}^2 .

A primer on matrices

Stephen Boyd

September 16, 2008

These notes describe the notation of matrices, the mechanics of matrix manipulation, and how to use matrices to formulate and solve sets of simultaneous linear equations.

We *won't* cover

- linear algebra, *i.e.*, the underlying mathematics of matrices
- numerical linear algebra, *i.e.*, the algorithms used to manipulate matrices and solve linear equations
- software for forming and manipulating matrices, *e.g.*, Matlab, Mathematica, or Octave
- how to represent and manipulate matrices, or solve linear equations, in computer languages such as C/C++ or Java
- applications, for example in statistics, mechanics, economics, circuit analysis, or graph theory

1 Matrix terminology and notation

Matrices

A *matrix* is a rectangular array of numbers (also called *scalars*), written between square brackets, as in

$$A = \begin{bmatrix} 0 & 1 & -2.3 & 0.1 \\ 1.3 & 4 & -0.1 & 0 \\ 4.1 & -1 & 0 & 1.7 \end{bmatrix}.$$

An important attribute of a matrix is its *size* or *dimensions*, *i.e.*, the numbers of *rows* and *columns*. The matrix A above, for example, has 3 rows and 4 columns, so its size is 3×4 . (Size is always given as rows \times columns.) A matrix with m rows and n columns is called an $m \times n$ matrix.

An $m \times n$ matrix is called *square* if $m = n$, *i.e.*, if it has an equal number of rows and columns. Some authors refer to an $m \times n$ matrix as *fat* if $m < n$ (fewer rows than columns), or *skinny* if $m > n$ (more rows than columns). The matrix A above is fat.

The *entries* or *coefficients* of a matrix are the values in the array. The i, j entry is the value in the i th row and j th column, denoted by double subscripts: the i, j entry of a matrix C is denoted C_{ij} (which is a number). The positive integers i and j are called the (row and column, respectively) *indices*. For our example above, $A_{13} = -2.3$, $A_{32} = -1$. The row index of the bottom left entry (which has value 4.1) is 3; its column index is 1.

Two matrices are *equal* if they are the same size and all the corresponding entries (which are numbers) are equal.

Vectors and scalars

A matrix with only one column, *i.e.*, with size $n \times 1$, is called a *column vector* or just a *vector*. Sometimes the size is specified by calling it an n -*vector*. The entries of a vector are denoted with just one subscript (since the other is 1), as in a_3 . The entries are sometimes called the *components* of the vector, and the number of rows of a vector is sometimes called its *dimension*. As an example,

$$v = \begin{bmatrix} 1 \\ -2 \\ 3.3 \\ 0.3 \end{bmatrix}$$

is a 4-vector (or 4×1 matrix, or vector of dimension 4); its third component is $v_3 = 3.3$.

Similarly, a matrix with only one row, *i.e.*, with size $1 \times n$, is called a *row vector*. As an example,

$$w = \begin{bmatrix} -2.1 & -3 & 0 \end{bmatrix}$$

is a row vector (or 1×3 matrix); its third component is $w_3 = 0$.

Sometimes a 1×1 matrix is considered to be the same as an ordinary number. In the context of matrices and scalars, ordinary numbers are often called *scalars*.

Notational conventions for matrices, vectors, and scalars

Some authors try to use notation that helps the reader distinguish between matrices, vectors, and scalars (numbers). For example, Greek letters (α, β, \dots) might be used for numbers, lower-case letters (a, x, f, \dots) for vectors, and capital letters (A, F, \dots) for matrices. Other notational conventions include matrices given in bold font (\mathbf{G}), or vectors written with arrows above them (\vec{a}).

Unfortunately, there are about as many notational conventions as authors, so you should be prepared to figure out what things are (*i.e.*, scalars, vectors, matrices) despite the author's notational scheme (if any exists).

Zero and identity matrices

The zero matrix (of size $m \times n$) is the matrix with all entries equal to zero. Sometimes the zero matrix is written as $0_{m \times n}$, where the subscript denotes the size. But often, a zero matrix is denoted just 0, the same symbol used to denote the number 0. In this case you'll

have to figure out the size of the zero matrix from the context. (More on this later.) When a zero matrix is a (row or column) vector, we call it a zero (row or column) vector.

Note that zero matrices of different sizes are different matrices, even though we use the same symbol to denote them (*i.e.*, 0). In programming we call this *overloading*: we say the symbol 0 is overloaded because it can mean different things depending on its context (*i.e.*, the equation it appears in).

An identity matrix is another common matrix. It is always square, *i.e.*, has the same number of rows as columns. Its *diagonal* entries, *i.e.*, those with equal row and column index, are all equal to one, and its off-diagonal entries (those with unequal row and column indices) are zero. Identity matrices are denoted by the letter I . Sometimes a subscript denotes the size, as in I_4 or $I_{2 \times 2}$. But more often the size must be determined from context (just like zero matrices). Formally, the identity matrix of size n is defined by

$$I_{ij} = \begin{cases} 1 & i = j, \\ 0 & i \neq j. \end{cases}$$

Perhaps more illuminating are the examples

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

which are the 2×2 and 4×4 identity matrices. (Remember that both are denoted with the same symbol, namely, I .) The importance of the identity matrix will become clear later.

Unit and ones vectors

A vector with one component one and all others zero is called a *unit vector*. The i th unit vector, whose i th component is 1 and all others are zero, is usually denoted e_i . As with zero or identity matrices, you usually have to figure out the dimension of a unit vector from context. The three unit 3-vectors are:

$$e_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad e_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

Note that the n columns of the $n \times n$ identity matrix are the n unit n -vectors. Another term for e_i is *ith standard basis vector*. Also, you should watch out, because some authors use the term ‘unit vector’ to mean a vector of length one. (We’ll explain that later.)

Another common vector is the one with all components one, sometimes called the *ones vector*, and denoted $\mathbf{1}$ (by some authors) or e (by others). For example, the 4-dimensional ones vector is

$$\mathbf{1} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}.$$

2 Matrix operations

Matrices can be combined using various operations to form other matrices.

Matrix transpose

If A is an $m \times n$ matrix, its *transpose*, denoted A^T (or sometimes A'), is the $n \times m$ matrix given by $(A^T)_{ij} = A_{ji}$. In words, the rows and columns of A are transposed in A^T . For example,

$$\begin{bmatrix} 0 & 4 \\ 7 & 0 \\ 3 & 1 \end{bmatrix}^T = \begin{bmatrix} 0 & 7 & 3 \\ 4 & 0 & 1 \end{bmatrix}.$$

Transposition converts row vectors into column vectors, and vice versa. If we transpose a matrix twice, we get back the original matrix: $(A^T)^T = A$.

Matrix addition

Two matrices *of the same size* can be added together, to form another matrix (of the same size), by adding the corresponding entries (which are numbers). Matrix addition is denoted by the symbol $+$. (Thus the symbol $+$ is overloaded to mean scalar addition when scalars appear on its left and right hand side, and matrix addition when matrices appear on its left and right hand sides.) For example,

$$\begin{bmatrix} 0 & 4 \\ 7 & 0 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 1 & 6 \\ 9 & 3 \\ 3 & 5 \end{bmatrix}$$

A pair of row or column vectors of the same size can be added, but you cannot add together a row vector and a column vector (except when they are both scalars!).

Matrix subtraction is similar. As an example,

$$\begin{bmatrix} 1 & 6 \\ 9 & 3 \end{bmatrix} - I = \begin{bmatrix} 0 & 6 \\ 9 & 2 \end{bmatrix}.$$

Note that this gives an example where we have to figure out what size the identity matrix is. Since you can only add (or subtract) matrices of the same size, we conclude that I must refer to a 2×2 identity matrix.

Matrix addition is commutative, *i.e.*, if A and B are matrices of the same size, then $A + B = B + A$. It's also associative, *i.e.*, $(A + B) + C = A + (B + C)$, so we write both as $A + B + C$. We always have $A + 0 = 0 + A = A$, *i.e.*, adding the zero matrix to a matrix has no effect. (This is another example where you have to figure out the exact dimensions of the zero matrix from context. Here, the zero matrix must have the same dimensions as A ; otherwise they could not be added.)

Scalar multiplication

Another operation is *scalar multiplication*: multiplying a matrix by a scalar (*i.e.*, number), which is done by multiplying every entry of the matrix by the scalar. Scalar multiplication is usually denoted by juxtaposition, with the scalar on the left, as in

$$(-2) \begin{bmatrix} 1 & 6 \\ 9 & 3 \\ 6 & 0 \end{bmatrix} = \begin{bmatrix} -2 & -12 \\ -18 & -6 \\ -12 & 0 \end{bmatrix}.$$

Sometimes you see scalar multiplication with the scalar on the right, or even scalar division with the scalar shown in the denominator (which just means scalar multiplication by one over the scalar), as in

$$\begin{bmatrix} 1 & 6 \\ 9 & 3 \\ 6 & 0 \end{bmatrix} \cdot 2 = \begin{bmatrix} 2 & 12 \\ 18 & 6 \\ 12 & 0 \end{bmatrix}, \quad \frac{\begin{bmatrix} 9 & 6 & 9 \\ 6 & 0 & 3 \end{bmatrix}}{3} = \begin{bmatrix} 3 & 2 & 3 \\ 2 & 0 & 1 \end{bmatrix},$$

but I think these look pretty ugly.

Scalar multiplication obeys several laws you can figure out for yourself, *e.g.*, if A is any matrix and α, β are any scalars, then

$$(\alpha + \beta)A = \alpha A + \beta A.$$

It's a useful exercise to identify the symbols appearing in this formula. The $+$ symbol on the left is addition of scalars, while the $+$ symbol on the right denotes matrix addition.

Another simple property is $(\alpha\beta)A = (\alpha)(\beta A)$, where α and β are scalars and A is a matrix. On the left hand side we see scalar-scalar multiplication $(\alpha\beta)$ and scalar-matrix multiplication; on the right we see two cases of scalar-matrix multiplication.

Note that $0 \cdot A = 0$ (where the lefthand zero is the scalar zero, and the righthand zero is a matrix zero of the same size as A).

Matrix multiplication

It's also possible to multiply two matrices using *matrix multiplication*. You can multiply two matrices A and B provided their dimensions are *compatible*, which means the number of columns of A (*i.e.*, its *width*) equals the number of rows of B (*i.e.*, its *height*). Suppose A and B are compatible, *i.e.*, A has size $m \times p$ and B has size $p \times n$. The product matrix $C = AB$, which has size $m \times n$, is defined by

$$C_{ij} = \sum_{k=1}^p a_{ik}b_{kj} = a_{i1}b_{1j} + \cdots + a_{ip}b_{pj}, \quad i = 1, \dots, m, \quad j = 1, \dots, n.$$

This rule looks complicated, but there are several ways to remember it. To find the i, j entry of the product $C = AB$, you need to know the i th row of A and the j th column of B . The

summation above can be interpreted as ‘moving left to right along the i th row of A ’ while moving ‘top to bottom’ down the j th column of B . As you go, you keep a running sum of the product of the corresponding entries from A and B .

As an example, let’s find the product $C = AB$, where

$$A = \begin{bmatrix} 1 & 2 & 3 \\ -1 & 0 & 4 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & -3 \\ 2 & 1 \\ -1 & 0 \end{bmatrix}.$$

First, we check that they are compatible: A has three columns, and B has three rows, so they’re compatible. The product matrix C will have two rows (the number of rows of A) and two columns (the number of columns of B). Now let’s find the entries of the product C . To find the 1, 1 entry, we move across the first row of A and down the first column of B , summing the products of corresponding entries:

$$C_{11} = (1)(0) + (2)(2) + (3)(-1) = 1.$$

To find the 1, 2 entry, we move across the first row of A and down the second column of B :

$$C_{12} = (1)(-3) + (2)(1) + (3)(0) = -1.$$

In each product term here, the lefthand number comes from the first row of A , and the righthand number comes from the first column of B . Two more similar calculations give us the remaining entries C_{21} and C_{22} :

$$\begin{bmatrix} 1 & 2 & 3 \\ -1 & 0 & 4 \end{bmatrix} \begin{bmatrix} 0 & -3 \\ 2 & 1 \\ -1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -4 & 3 \end{bmatrix}.$$

At this point, matrix multiplication probably looks very complicated to you. It is, but once you see all the uses for it, you’ll get used to it.

Some properties of matrix multiplication

Now we can explain why the identity has its name: if A is any $m \times n$ matrix, then $AI = A$ and $IA = A$, *i.e.*, when you multiply a matrix by an identity matrix, it has no effect. (The identity matrices in the formulas $AI = A$ and $IA = A$ have different sizes — what are they?)

One very important fact about matrix multiplication is that it is (in general) *not commutative*: we *don’t* (in general) have $AB = BA$. In fact, BA may not even make sense, or, if it makes sense, may be a different size than AB (so that equality in $AB = BA$ is meaningless). For example, if A is 2×3 and B is 3×4 , then AB makes sense (the dimensions are compatible) but BA doesn’t even make sense (much less equal AB). Even when AB and BA both make sense and are the same size, *i.e.*, when A and B are square, we don’t (in general) have $AB = BA$. As a simple example, consider:

$$\begin{bmatrix} 1 & 6 \\ 9 & 3 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} -6 & 11 \\ -3 & -3 \end{bmatrix}, \quad \begin{bmatrix} 0 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 1 & 6 \\ 9 & 3 \end{bmatrix} = \begin{bmatrix} -9 & -3 \\ 17 & 0 \end{bmatrix}.$$

Matrix multiplication *is* associative, *i.e.*, $(AB)C = A(BC)$ (provided the products make sense). Therefore we write the product simply as ABC . Matrix multiplication is also associative with scalar multiplication, *i.e.*, $\alpha(AB) = (\alpha A)B$, where α is a scalar and A and B are matrices (that can be multiplied). Matrix multiplication distributes across matrix addition: $A(B + C) = AB + AC$ and $(A + B)C = AC + BC$.

Matrix-vector product

A very important and common case of matrix multiplication is $y = Ax$, where A is an $m \times n$ matrix, x is an n -vector, and y is an m -vector. We can think of matrix vector multiplication (with an $m \times n$ matrix) as a function that transforms n -vectors into m -vectors. The formula is

$$y_i = A_{i1}x_1 + \cdots + A_{in}x_n, \quad i = 1, \dots, m$$

Inner product

Another important special case of matrix multiplication occurs when v is a row n -vector and w is a column n vector. Then the product vw makes sense, and has size 1×1 , *i.e.*, is a scalar:

$$vw = v_1w_1 + \cdots + v_nw_n.$$

This occurs often in the form x^Ty where x and y are both n -vectors. In this case the product (which is a number) is called the *inner product* or *dot product* of the vectors x and y . Other notation for the inner product is $\langle x, y \rangle$ or $x \cdot y$. If x and y are n -vectors, then their inner product is

$$\langle x, y \rangle = x^Ty = x_1y_1 + \cdots + x_ny_n.$$

But remember that the matrix product xy doesn't make sense (unless they are both scalars).

Matrix powers

When a matrix A is square, then it makes sense to multiply A by itself, *i.e.*, to form $A \cdot A$. We refer to this matrix as A^2 . Similarly, k copies of A multiplied together is denoted A^k .

(Non-integer powers, such as $A^{1/2}$ (the matrix squareroot), are pretty tricky — they might not make sense, or be ambiguous, unless certain conditions on A hold. This is an advanced topic in linear algebra.)

By convention we set $A^0 = I$ (usually only when A is invertible — see below).

Matrix inverse

If A is square, and there is a matrix F such that $FA = I$, then we say that A is *invertible* or *nonsingular*. We call F the *inverse* of A , and denote it A^{-1} . We can then also define $A^{-k} = (A^{-1})^k$. If a matrix is not invertible, we say it is *singular* or *noninvertible*.

It's important to understand that not all square matrices are invertible, *i.e.*, have inverses. (For example, a zero matrix never has an inverse.) As a less obvious example, you might try to show that the matrix

$$\begin{bmatrix} 1 & -1 \\ -2 & 2 \end{bmatrix}$$

does not have an inverse.

As an example of the matrix inverse, we have

$$\begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}^{-1} = \frac{1}{3} \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix}$$

(you should check this!).

When a matrix is invertible, the inverse of the inverse is the original matrix, *i.e.*, $(A^{-1})^{-1} = A$. A basic result of linear algebra is that $AA^{-1} = I$. In other words, if you multiply a matrix by its inverse on the *right* (as well as the left), you get the identity.

It's very useful to know the general formula for a 2×2 matrix inverse:

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

provided $ad - bc \neq 0$. (If $ad - bc = 0$, the matrix is not invertible.) There are similar, but much more complicated, formulas for the inverse of larger (invertible) square matrices, but they are not used in practice.

The importance of the matrix inverse will become clear when we study linear equations.

Useful identities

We've already mentioned a handful of matrix identities, that you could figure out yourself, *e.g.*, $A + 0 = A$. Here we list a few others, that are not hard to derive, and quite useful. (We're making no claims that our list is complete!)

- transpose of product: $(AB)^T = B^T A^T$
- transpose of sum: $(A + B)^T = A^T + B^T$
- inverse of product: $(AB)^{-1} = B^{-1} A^{-1}$ provided A and B are square (of the same size) and invertible
- products of powers: $A^k A^l = A^{k+l}$ (for $k, l \geq 1$ in general, and for all k, l if A is invertible)

Block matrices and submatrices

In some applications it's useful to form matrices whose entries are themselves matrices, as in

$$\begin{bmatrix} A & B & C \end{bmatrix}, \quad \begin{bmatrix} F & I \\ 0 & G \end{bmatrix},$$

where A , B , C , F , and G are matrices (as are 0 and I). Such matrices are called *block matrices*; the entries A , B , etc. are called 'blocks' and are sometimes named by indices. Thus, F is called the 1, 1 block of the second matrix.

Of course the block matrices must have the right dimensions to be able to fit together: matrices in the same (block) row must have the same number of rows (*i.e.*, the same 'height'); matrices in the same (block) column must have the same number of columns (*i.e.*, the same 'width'). Thus in the examples above, A , B and C must have the same number of rows (*e.g.*, they could be 2×3 , 2×2 , and 2×1). The second example is more interesting. Suppose that F is $m \times n$. Then the identity matrix in the 1, 2 position must have size $m \times m$ (since it must have the same number of rows as F). We also see that G must have m columns, say, dimensions $p \times m$. That fixes the dimensions of the 0 matrix in the 2, 1 block — it must be $p \times n$.

As a specific example, suppose that

$$C = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}, \quad D = \begin{bmatrix} 0 & 2 & 3 \\ 5 & 4 & 7 \end{bmatrix}.$$

Then we have

$$\begin{bmatrix} D & C \end{bmatrix} = \begin{bmatrix} 0 & 2 & 3 & 2 & 2 \\ 5 & 4 & 7 & 1 & 3 \end{bmatrix}.$$

Continuing this example, the expression

$$\begin{bmatrix} C \\ D \end{bmatrix}$$

doesn't make sense, because the top block has two columns and the bottom block has three. But the block expression

$$\begin{bmatrix} C \\ D^T \end{bmatrix}$$

does make sense, because now the bottom block has two columns, just like the top block.

You can also divide a larger matrix (or vector) into 'blocks'. In this context the blocks are sometimes called *submatrices* of the big matrix. For example, it's often useful to write an $m \times n$ matrix as a $1 \times n$ block matrix of m -vectors (which are just its columns), or as an $m \times 1$ block matrix of n -row-vectors (which are its rows).

Block matrices can be added and multiplied as if the entries were numbers, provided the corresponding entries have the right sizes (*i.e.*, 'conform') and you're careful about the order of multiplication. Thus we have

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} AX + BY \\ CX + DY \end{bmatrix}$$

provided the products AX , BY , CX , and DY makes sense.

3 Linear equations and matrices

Linear functions

Suppose that f is a function that takes as argument (input) n -vectors and returns (as output) m -vectors. We say f is *linear* if it satisfies two properties:

- scaling: for any n -vector x and any scalar α , $f(\alpha x) = \alpha f(x)$
- superposition: for any n -vectors u and v , $f(u + v) = f(u) + f(v)$

It's not hard to show that such a function can always be represented as matrix-vector multiplication: there is an $m \times n$ matrix A such that $f(x) = Ax$ for all n -vectors x . (Conversely, functions defined by matrix-vector multiplication are linear.)

We can also write out the linear function in explicit form, *i.e.*, $f(x) = y$ where

$$y_i = \sum_{j=1}^n A_{ij}x_j = A_{i1}x_1 + \cdots + A_{in}x_n, \quad i = 1, \dots, m$$

This gives a simple interpretation of A_{ij} : it gives the coefficient by which y_i depends on x_j .

Suppose an m -vector y is a linear function of the n -vector x , *i.e.*, $y = Ax$ where A is $m \times n$. Suppose also that a p -vector z is a linear function of y , *i.e.*, $z = By$ where B is $p \times m$. Then z is a linear function of x , which we can express in the simple form $z = By = (BA)x$. So matrix multiplication corresponds to composition of linear functions (*i.e.*, linear functions of linear functions of some variables).

Linear equations

Any set of m linear equations in (scalar) variables x_1, \dots, x_n can be represented by the compact matrix equation $Ax = b$, where x is a vector made from the variables, A is an $m \times n$ matrix and b is an m -vector. Let's start with a simple example of two equations in three variables:

$$1 + x_2 = x_3 - 2x_1, \quad x_3 = x_2 - 2.$$

The first thing to do is to rewrite the equations with the variables lined up in columns, and the constants on the righthand side:

$$\begin{array}{rrrr} 2x_1 & +x_2 & -x_3 & = & -1 \\ 0x_1 & -x_2 & +x_3 & = & -2 \end{array}$$

Now it's easy to rewrite the equations as a single matrix equation:

$$\begin{bmatrix} 2 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ -2 \end{bmatrix},$$

so we can express the two equations in the three variables as $Ax = b$ where

$$A = \begin{bmatrix} 2 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad b = \begin{bmatrix} -1 \\ -2 \end{bmatrix}$$

Solving linear equations

Now suppose we have n linear equations in n variables x_1, \dots, x_n , written as the compact matrix equation $Ax = b$, where A is an $n \times n$ matrix, and b is an n -vector. Suppose that A is invertible, *i.e.*, the inverse A^{-1} exists. Let's multiply both sides of the matrix equation $Ax = b$ on the left by A^{-1} :

$$A^{-1}(Ax) = A^{-1}b.$$

The lefthand side simplifies to $A^{-1}Ax = Ix = x$, so we have actually solved the simultaneous linear equations: $x = A^{-1}b$.

Now you can see the importance of the matrix inverse: it can be used to solve simultaneous linear equations. Here we should make a comment about matrix notation. The power of matrix notation is that just a few symbols (*e.g.*, A^{-1}) can express a lot. Another (perhaps more pessimistic) way to put this is, a lot of work can be hidden behind just a few symbols (*e.g.*, A^{-1}).

Of course, you can't *always* solve n linear equations for n variables. One or more of the equations might be redundant (*i.e.*, can be obtained from the others), or the equations could be inconsistent as in $x_1 = 1, x_1 = 2$. When these pathologies occur, the matrix A is singular, *i.e.*, noninvertible. Conversely, when a matrix A is singular, it turns out the simultaneous linear equations $Ax = b$ are redundant or inconsistent. (These facts are studied in linear algebra.)

From a practical point of view, then, A is singular means that the equations in $Ax = b$ are redundant or inconsistent — a sign that you have set up the wrong equations (or don't have enough of them). Otherwise, A^{-1} exists, and the equations can be solved as $x = A^{-1}b$.

Solving linear equations in practice

When we solve linear equations in practice, (*i.e.*, by computer) we do not first compute the matrix A^{-1} , and then multiply it on the right by the vector b , to get the solution $x = A^{-1}b$ (although that procedure would, of course, work). Practical methods compute the solution $x = A^{-1}b$ directly.

The most common methods for computing $x = A^{-1}b$ (*i.e.*, solving a system of n simultaneous linear equations) require on the order of n^3 basic arithmetic operations. But modern computers are very fast, so solving say a set of 1000 equations on a small PC class computer takes only a second or so. (A 1000×1000 matrix requires storage for 10^6 doubles, around 10MB.) But solving larger sets of equations, for example, 5000, will take much ($125\times$) longer (on a PC class computer). (A 5000×5000 matrix requires around 250MB to store.)

In many applications the matrix A has many, or almost all, of its entries equal to zero, in which case we say it is *sparse*. In terms of the associated linear equations, this means each

equation involves only some (often just a few) of the variables. It turns out that such sparse equations can be solved by computer very efficiently, using *sparse matrix techniques*. It's not uncommon to solve for hundreds of thousands of variables, with hundreds of thousands of (sparse) equations. Even on a PC class computer, solving a system of 10000 simultaneous sparse linear equations is feasible, and might take only a few seconds (but it depends on how sparse the equations are).

Crimes Against Matrices

In this note we list some matrix crimes that we have, sadly, witnessed too often. Be very careful to avoid committing any of these crimes; in EE263 we have a *zero-tolerance* policy for *crimes against matrices*, at least on things you hand in to us. (What you do with matrices in your spare time, or on scratch paper, is of course your own business. But we recommend you avoid these crimes at all times, in order to not build bad habits.)

Check your work — don't become just another sad statistic!

Syntax crimes

In a syntax crime, the perpetrator attempts to combine matrices (or other mathematical objects) in ways that violate basic syntax rules. These are serious crimes of negligence, since it is so easy to check your work for potential violations. We list some typical examples below.

- Adding, subtracting, or equating, matrices (or vectors) of different dimensions.
Example: writing $A + B$, when $A \in \mathbf{R}^{2 \times 3}$ and $B \in \mathbf{R}^{3 \times 3}$.
- Violating the rules of constructing block matrices (*e.g.*, the submatrices in any row of a block matrix must have the same number of rows).
Example: forming the block matrix $[A \ B]$, when $A \in \mathbf{R}^{2 \times 3}$ and $B \in \mathbf{R}^{3 \times 3}$.
- Multiplying matrices with incompatible dimensions (*i.e.*, forming AB , when the number of columns of A does not equal the number of rows of B).
Example: forming $A^T B$, when $A \in \mathbf{R}^{2 \times 3}$ and $B \in \mathbf{R}^{3 \times 3}$.
- Taking the inverse, determinant, trace, or powers of a nonsquare matrix.
Example: forming A^{-1} , when $A \in \mathbf{R}^{2 \times 3}$.

Semantic crimes

In a semantic crime, the perpetrator forms an expression or makes an assertion that does not break any syntax rules, but is wrong because of the meaning. These crimes are a bit harder to detect than syntax crimes, so you need to be more vigilant to avoid committing them.

- Taking the inverse of a square, but singular matrix. (Taking the inverse of a nonsquare matrix is a syntax crime—see above.)

Example: forming $(ww^T)^{-1}$, where $w \in \mathbf{R}^2$.

Note: writing $(ww^T)^{-1} = (w^T)^{-1}w^{-1}$, when $w \in \mathbf{R}^2$, involves both a syntax and semantic crime.

- Referring to a left inverse of a strictly fat matrix or a right inverse of a strictly skinny matrix.

Example: writing $QQ^T = I$, when $Q \in \mathbf{R}^{5 \times 3}$.

- Cancelling matrices on the left or right in inappropriate circumstances, *e.g.*, concluding that $B = C$ from $AB = AC$, when A is not known to be one-to-one (*i.e.*, have independent columns).

Example: concluding $x = y$ from $a^T x = a^T y$, when $a, x, y \in \mathbf{R}^4$.

- *Dimension crimes.* Alleging that a set of m vectors in \mathbf{R}^n is independent, when $m > n$. Alleging that a set of m vectors in \mathbf{R}^n span \mathbf{R}^n , when $m < n$.

Miscellaneous crimes

Some crimes are hard to classify, or involve both syntax and semantic elements. Incorrect use of a matrix identity often falls in this class.

- Using $(AB)^T = A^T B^T$ (instead of the correct formula $(AB)^T = B^T A^T$).

Note: this also violates syntax rules, if $A^T B^T$ is not a valid product.

- Using $(AB)^{-1} = A^{-1} B^{-1}$ (instead of the correct formula $(AB)^{-1} = B^{-1} A^{-1}$).

Note: $(AB)^{-1} = A^{-1} B^{-1}$ violates syntax rules, if A or B is not square; it violates semantic rules if A or B is not invertible.

- Using $(A + B)^2 = A^2 + 2AB + B^2$. This (false) identity relies on the very useful, but unfortunately false, identity $AB = BA$.

An example

Let's consider the expression $(A^T B)^{-1}$, where $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{k \times p}$. Here's how you might check for various crimes you might commit in forming this expression.

- We multiply A^T , which is $n \times m$, and B , which is $k \times p$, so we better have $m = k$ to avoid a syntax violation.

Note: if A is a scalar, then $A^T B$ might be a strange thing to write, but can be argued to not violate syntax, even though $m \neq k$. In a similar way, when B is scalar, you can write $A^T B$, and argue that syntax is not violated.

- The product $A^T B$ is $n \times p$, so we better have $n = p$ in order to (attempt to) invert it. At this point, we know that the dimensions of A and B must be the same (ignoring the case where one or the other is interpreted as a scalar).
- If $A^T B$ is a strictly skinny–strictly fat product (*i.e.*, A and B are strictly fat), then $A^T B$ cannot possibly be invertible, so we have a semantic violation. To avoid this, we must have A and B square or skinny, *i.e.*, $m \geq n$.

Summary: to write $(A^T B)^{-1}$ (assuming neither A nor B is interpreted as a scalar), A and B must have the same dimensions, and be skinny or square.

Of course, even if A and B have the same dimensions, and are skinny or square, the matrix $A^T B$ can be singular, in which case $(A^T B)^{-1}$ is meaningless. The point of our analysis above is that if A and B don't have the same dimension, or if A and B are strictly fat, then $(A^T B)^{-1}$ is *guaranteed* to be meaningless, no matter what values A and B might have in your application or argument.

Least squares and least norm in Matlab

Least squares approximate solution

Suppose $A \in \mathbf{R}^{m \times n}$ is skinny (or square), *i.e.*, $m \geq n$, and full rank, which means that $\text{Rank}(A) = n$. The least-squares approximate solution of $Ax = y$ is given by

$$x_{\text{ls}} = (A^T A)^{-1} A^T y.$$

This is the unique $x \in \mathbf{R}^n$ that minimizes $\|Ax - y\|$.

There are several ways to compute x_{ls} in Matlab. The simplest method is to use the *backslash operator*:

```
xls=A\y;
```

If A is square (and invertible), the backslash operator just solves the linear equations, *i.e.*, it computes $A^{-1}y$. If A is not full rank, then $A \backslash b$ will generate an error message, and then a least-squares solution will be returned.

You can also use the formula to compute the least squares approximate solution:

```
xls=inv(A'*A)*A'*y;
```

We encourage you to verify that you get the same answer (except possibly for some small difference due to roundoff error) as you do using the backslash operator. Here you'll get an error if A is not full rank, or fat. Compared with the backslash operator, this method is a bit less efficient, and a bit more sensitive to roundoff errors, but you won't notice it unless m and n are a thousand or so.

You can also use the pseudo-inverse function `pinv()`, which computes the pseudo-inverse, which is

$$A^\dagger = (A^T A)^{-1} A^T$$

when A is full rank and skinny (or square). (The pseudo-inverse is also defined when A is not full rank, but it's not given by the formula above.) To find the least-squares approximate solution using the pseudo-inverse, you can use

```
xls=pinv(A)*y;
```

You better be sure here that A is skinny (or square) and full rank; otherwise you'll compute something (with no warning messages) that isn't the least-squares approximate solution.

Yet another method is via a QR decomposition. In Matlab, `[Q,R]=qr(A)` returns the 'full' QR decomposition, with square, orthogonal $Q \in \mathbf{R}^{m \times m}$, and $R \in \mathbf{R}^{m \times n}$ upper triangular (with zeros in its bottom part). The 'economy' QR decomposition, in which $Q \in \mathbf{R}^{m \times n}$ (with orthonormal columns) and invertible R , is obtained using `[Q,R]=qr(A,0)`. You can compute the least-squares approximate solution using the economy QR decomposition using, for example,

```
[Q,R]=qr(A,0); % compute economy QR decomposition
xls=R\ (Q'*y);
```

To be sure you've really computed the least-squares approximate solution, we encourage you to check that the residual is orthogonal to the columns of A , for example with the commands

```
r=A*x-y; % compute residual
A'*r % compute inner product of columns of A and r
```

and checking that the result is very small.

Least norm solution

Now suppose $A \in \mathbf{R}^{m \times n}$ and is fat (or square), *i.e.*, $m \leq n$, and full rank, which means that $\mathbf{Rank}(A) = m$. The least-norm solution of $Ax = y$ is given by

$$x_{\text{ln}} = A^T(AA^T)^{-1}y.$$

Among all solutions of $Ax = y$, x_{ln} has the smallest norm.

We can compute x_{ln} in Matlab in several ways. You can use the formula to compute x_{ln} :

```
xln=A'*inv(A*A')*y;
```

You'll get an error here if A is not full rank, or if A is skinny.

You can also use the pseudo-inverse function `pinv()`, which computes the pseudo-inverse, which is

$$A^\dagger = A^T(AA^T)^{-1}$$

when A is full rank and fat or square. (The pseudo-inverse is also defined when A is not full rank, but it's not given by the formula above.)

```
xln=pinv(A)*y;
```

You better be sure here that A is fat (or square) and full rank; otherwise you'll compute something (with no warning messages) that isn't the least-norm solution.

You can find the least-norm solution via QR as well:

```
[Q,R]=qr(A',0);
xln=Q*(R'\y);
```

Warning! We should also mention a method that *doesn't work*: the *backslash operator*. If A is fat and full rank, then $A \backslash y$ gives a solution of the $Ax = y$, not necessarily the least-norm solution.

Solving general linear equations using Matlab

In this note we consider the following problem: Determine whether there is a solution $x \in \mathbf{R}^n$ of the (set of) m linear equations $Ax = b$, and if so, find one. To check existence of a solution is the same as checking if $b \in \mathcal{R}(A)$. We consider here the general case, with $A \in \mathbf{R}^{m \times n}$, with $\text{Rank}(A) = r$. In particular, we do not assume that A is full rank.

Existence of solution via rank

A simple way to check if $b \in \mathcal{R}(A)$ is to check the rank of $[A \ b]$, which is either r (*i.e.*, the rank of A) if $b \in \mathcal{R}(A)$, or $r + 1$, if $b \notin \mathcal{R}(A)$. This can be checked in Matlab using

```
rank([A b]) == rank(A)
```

(or evaluating the two ranks separately and comparing them). If the two ranks above are equal, then $Ax = b$ has a solution. But this method does not give us a solution, when one exists. This method also has a hidden catch: Matlab uses a numerical tolerance to decide on the rank of a matrix, and this tolerance might not be appropriate for your particular application.

Using the backslash and pseudo-inverse operator

In Matlab, the easiest way to determine whether $Ax = b$ has a solution, and to find such a solution when it does, is to use the backslash operator. Exactly what $A \backslash b$ returns is a bit complicated to describe in the most general case, but *if there is a solution to $Ax = b$* , then $A \backslash b$ *returns one*. A couple of warnings: First, $A \backslash b$ returns a result in many cases when there is no solution to $Ax = b$. For example, when A is skinny and full rank (*i.e.*, $m > n = r$), $A \backslash b$ returns the least-squares approximate solution, which in general is not a solution of $Ax = b$ (unless we happen to have $b \in \mathcal{R}(A)$). Second, $A \backslash b$ sometimes causes a warning to be issued, even when it returns a solution of $Ax = b$. This means that you can't just use the backslash operator: you have to *check* that what it returns is a solution. (In any case, it's just good common sense to check numerical computations as you do them.) In Matlab this can be done as follows:

```
x = A\b; % possibly a solution to Ax=b
norm(A*x-b) % if this is zero or very small, we have a solution
```

If the second line yields a result that is not very small, we conclude that $Ax = b$ does not have a solution. Note that executing the first line might cause a warning to be issued.

In contrast to the rank method described above, *you* decide on the numerical tolerance you'll accept (*i.e.*, how small $\|Ax - b\|$ has to be before you accept x as a solution of $Ax = b$). A common test that works well in many applications is $\|Ax - b\| \leq 10^{-5}\|b\|$.

You can also use the pseudo-inverse operator: $\mathbf{x} = \text{pinv}(\mathbf{A}) * \mathbf{b}$ is also guaranteed to solve $Ax = b$, if $Ax = b$ has a solution. As with the backslash operator, you have to check that the result satisfies $Ax = b$, since in general, it doesn't have to.

Using the QR factorization

While the backslash operator is a convenient way to check if $Ax = b$ has a solution, it's a bit opaque. Here we describe a method that is transparent, and can be fully explained and understood using material we've seen in the course.

We start with the full QR factorization of A with column permutations:

$$AP = QR = \begin{bmatrix} Q_1 & Q_2 \end{bmatrix} \begin{bmatrix} R_1 & R_2 \\ 0 & 0 \end{bmatrix}.$$

Here $Q \in \mathbf{R}^{m \times m}$ is orthogonal, $R \in \mathbf{R}^{m \times n}$ is upper triangular, and $P \in \mathbf{R}^{n \times n}$ is a permutation matrix. The submatrices have the following dimensions: $Q_1 \in \mathbf{R}^{m \times r}$, $Q_2 \in \mathbf{R}^{m \times (m-r)}$, $R_1 \in \mathbf{R}^{r \times r}$ is upper triangular with nonzero elements along its main diagonal, and $R_2 \in \mathbf{R}^{r \times (n-r)}$. The zero submatrices in the bottom (block) row of R have $m - r$ rows.

Using $A = QRP^T$ we can write $Ax = b$ as

$$QRP^T x = QRz = b,$$

where $z = P^T x$. Multiplying both sides of this equation by Q^T gives the equivalent set of m equations $Rz = Q^T b$. Expanding this into subcomponents gives

$$Rz = \begin{bmatrix} R_1 & R_2 \\ 0 & 0 \end{bmatrix} z = \begin{bmatrix} Q_1^T b \\ Q_2^T b \end{bmatrix}.$$

We see immediately that there is no solution of $Ax = b$, unless we have $Q_2^T b = 0$, because the bottom component of Rz is always zero.

Now let's assume that we do have $Q_2^T b = 0$. Then the equations reduce to

$$R_1 z_1 + R_2 z_2 = Q_1^T b,$$

a set r linear equations in n variables. We can find a solution of these equations by setting $z_2 = 0$. With this form for z , the equation above becomes $R_1 z_1 = Q_1^T b$, from which we get $z_1 = R_1^{-1} Q_1^T b$. Now we have a z that satisfies $Rz = Q^T b$: $z = [z_1^T \ 0]^T$. We get the corresponding x from $x = Pz$:

$$x = P \begin{bmatrix} R_1^{-1} Q_1^T b \\ 0 \end{bmatrix}.$$

This x satisfies $Ax = b$, provided we have $Q_2^T b = 0$. Whew.

Actually, the construction outlined above is pretty much what $\mathbf{A} \backslash \mathbf{b}$ does.

In Matlab, we can carry out this construction as follows:

```
[m,n]=size(A);
[Q,R,P]=qr(A); % full QR factorization
r=rank(A); % could also get rank directly from QR factorization ...

% construct the submatrices
Q1=Q(:,1:r);
Q2=Q(:,r+1:m);
R1=R(1:r,1:r);

% check if b is in range(A)
norm(Q2'*b) % if this is zero or very small, b is in range(A)

% construct a solution
x=P*[R1\'(Q1'*b); zeros(n-r,1)]; % satisfies Ax=b, if b is in range(A)

% check alleged solution (just to be sure)
norm(A*x-b)
```

Low Rank Approximation and Extremal Gain Problems

These notes pull together some similar results that depend on partial or truncated SVD or eigenvector expansions.

1 Low rank approximation

In lecture 15 we considered the following problem. We are given a matrix $A \in \mathbf{R}^{m \times n}$ with rank r , and we want to find the nearest matrix $\hat{A} \in \mathbf{R}^{m \times n}$ with rank p (with $p \leq r$), where ‘nearest’ is measured by the matrix norm, *i.e.*, $\|A - \hat{A}\|$. We found that a solution is

$$\hat{A} = \sum_{i=1}^p \sigma_i u_i v_i^T,$$

where

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T$$

is the SVD of A . The matrix \hat{A} need not be the only rank p matrix that is closest to A ; there can be other matrices, also of rank p , that satisfy $\|A - \tilde{A}\| = \|A - \hat{A}\| = \sigma_{p+1}$.

It turns out that the same matrix \hat{A} is also the nearest rank p matrix to A , as measured in the Frobenius norm, *i.e.*,

$$\|A - \hat{A}\|_F = \left(\text{Tr}(A - \hat{A})^T (A - \hat{A}) \right)^{1/2} = \left(\sum_{i=1}^m \sum_{j=1}^n (A_{ij} - \hat{A}_{ij})^2 \right)^{1/2}.$$

(The Frobenius norm is just the Euclidean norm of the matrix, written out as a long column vector.) In this case, however, \hat{A} is the unique rank p closest matrix to A , as measured in the Frobenius norm.

2 Nearest positive semidefinite matrix

Suppose that $A = A^T \in \mathbf{R}^{n \times n}$, with eigenvalue decomposition

$$A = \sum_{i=1}^n \lambda_i q_i q_i^T,$$

where $\{q_1, \dots, q_n\}$ are orthonormal, and $\lambda_1 \geq \dots \geq \lambda_n$. Consider the problem of finding a nearest positive semidefinite matrix, *i.e.*, a matrix $\hat{A} = \hat{A}^T \succeq 0$ that minimizes $\|A - \hat{A}\|$. A

solution to this problem is

$$\hat{A} = A = \sum_{i=1}^n \max\{\lambda_i, 0\} q_i q_i^T.$$

Thus, to get a nearest positive semidefinite matrix, you simply remove the terms in the eigenvector expansion that correspond to negative eigenvalues. The matrix \hat{A} is sometimes called the *positive semidefinite part* of A .

As you might guess, the matrix \hat{A} is also the nearest positive semidefinite matrix to A , as measured in the Frobenius norm.

3 Extremal gain problems

Suppose $A \in \mathbf{R}^{m \times n}$ has SVD

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T.$$

You already know that $v = v_1$ maximizes $\|Ax\|$ over all x with norm one. In other words, v_1 defines a direction of maximum gain for A . We can also find a direction of minimum gain. If $r < n$, then any unit vector x in $\mathcal{N}(A)$ minimizes $\|Ax\|$. If $r = n$, then the vector v_n minimizes $\|Ax\|$ among all vectors of norm one.

These results can be extended to finding *subspaces* on which A has large or small gain. Let \mathcal{V} be a subspace of \mathbf{R}^n . We define the *minimum gain* of $A \in \mathbf{R}^{m \times n}$ on \mathcal{V} as $\min\{\|Ax\| \mid x \in \mathcal{V}, \|x\| = 1\}$. We can then pose the question: find a subspace of dimension p , on which A has the largest possible minimum gain. The solution is what you'd guess, provided $p \leq r$:

$$\mathcal{V} = \text{span}\{v_1, \dots, v_p\},$$

the span of the right singular vectors associated with the p largest singular values. The minimum gain of A on this subspace is σ_p .

If $p > r$, then any subspace of dimension p intersects the nullspace of A , and therefore has minimum gain zero. So when $p > r$ you can take \mathcal{V} as any subspace of dimension p ; they all have the same minimum gain, namely, zero.

We can also find a subspace \mathcal{V} of dimension p that has the smallest *maximum gain* of A , defined as $\max\{\|Ax\| \mid x \in \mathcal{V}, \|x\| = 1\}$. Assuming $r = n$ (*i.e.*, A has nullspace $\{0\}$), one such subspace is

$$\mathcal{V} = \text{span}\{v_{r-p+1}, \dots, v_r\},$$

the span of the right singular vectors associated with the p smallest singular values.

We can put state these results in a more concrete form using matrices. To define a subspace of dimension p we use an orthonormal basis, $\mathcal{V} = \text{span}\{q_1, \dots, q_p\}$. Defining $Q = [q_1 \cdots q_p]$, we have $Q^T Q = I_p$, where I_p is the $p \times p$ identity matrix. We can express the minimum gain of A on \mathcal{V} as

$$\sigma_{\min}(AQ).$$

The problem of finding a subspace of dimension p that maximizes the minimum gain of A can be stated as

$$\begin{aligned} & \text{maximize} && \sigma_{\min}(AQ) \\ & \text{subject to} && Q^T Q = I_p. \end{aligned}$$

One solution to this problem is $Q = [v_1 \cdots v_p]$.

4 Extremal trace problems

Let $A \in \mathbf{R}^{n \times n}$ be symmetric, with eigenvalue decomposition $A = \sum_{i=1}^n \lambda_i q_i q_i^T$, with $\lambda_1 \geq \cdots \geq \lambda_n$, and $\{q_1, \dots, q_n\}$ orthonormal. You know that a solution of the problem

$$\begin{aligned} & \text{minimize} && x^T A x \\ & \text{subject to} && x^T x = 1, \end{aligned}$$

where the variable is $x \in \mathbf{R}^n$, is $x = q_n$. The related maximization problem is

$$\begin{aligned} & \text{maximize} && x^T A x \\ & \text{subject to} && x^T x = 1, \end{aligned}$$

with variable $x \in \mathbf{R}^n$. A solution to this problem is $x = q_1$.

Now consider the following generalization of the first problem:

$$\begin{aligned} & \text{minimize} && \mathbf{Tr}(X^T A X) \\ & \text{subject to} && X^T X = I_k, \end{aligned}$$

where the variable is $X \in \mathbf{R}^{n \times k}$, and I_k denotes the $k \times k$ identity matrix, and we assume $k \leq n$. (The constraint means that the columns of X are orthonormal.) A solution of this problem is $X = [q_{n-k+1} \cdots q_n]$. Note that when $k = 1$, this reduces to the first problem above.

The related maximization problem is

$$\begin{aligned} & \text{maximize} && \mathbf{Tr}(X^T A X) \\ & \text{subject to} && X^T X = I_k, \end{aligned}$$

with variable $X \in \mathbf{R}^{n \times k}$. A solution of this problem is $X = [q_1 \cdots q_k]$.

EE263 homework problems

Lecture 2 – Linear functions and examples

2.1 *A simple power control algorithm for a wireless network.* First some background. We consider a network of n transmitter/receiver pairs. Transmitter i transmits at power level p_i (which is positive). The path gain from transmitter j to receiver i is G_{ij} (which are all nonnegative, and G_{ii} are positive). The signal power at receiver i is given by $s_i = G_{ii}p_i$. The noise plus interference power at receiver i is given by

$$q_i = \sigma + \sum_{j \neq i} G_{ij}p_j$$

where $\sigma > 0$ is the self-noise power of the receivers (assumed to be the same for all receivers). The *signal to interference plus noise ratio* (SINR) at receiver i is defined as $S_i = s_i/q_i$. For signal reception to occur, the SINR must exceed some threshold value γ (which is often in the range 3 – 10). Various *power control algorithms* are used to adjust the powers p_i to ensure that $S_i \geq \gamma$ (so that each receiver can receive the signal transmitted by its associated transmitter). In this problem, we consider a simple power control update algorithm. The powers are all updated synchronously at a fixed time interval, denoted by $t = 0, 1, 2, \dots$. Thus the quantities p , q , and S are discrete-time signals, so for example $p_3(5)$ denotes the transmit power of transmitter 3 at time epoch $t = 5$. What we'd like is

$$S_i(t) = s_i(t)/q_i(t) = \alpha\gamma$$

where $\alpha > 1$ is an SINR safety margin (of, for example, one or two dB). Note that increasing $p_i(t)$ (power of the i th transmitter) increases S_i but decreases all other S_j . A very simple power update algorithm is given by

$$p_i(t+1) = p_i(t)(\alpha\gamma/S_i(t)). \quad (1)$$

This scales the power at the next time step to be the power that would achieve $S_i = \alpha\gamma$, if the interference plus noise term were to stay the same. But unfortunately, changing the transmit powers also changes the interference powers, so it's not that simple! Finally, we get to the problem.

- (a) Show that the power control algorithm (1) can be expressed as a linear dynamical system with constant input, *i.e.*, in the form

$$p(t+1) = Ap(t) + b,$$

where $A \in \mathbf{R}^{n \times n}$ and $b \in \mathbf{R}^n$ are constant. Describe A and b explicitly in terms of σ, γ, α and the components of G .

- (b) *Matlab simulation.* Use matlab to simulate the power control algorithm (1), starting from various initial (positive) power levels. Use the problem data

$$G = \begin{bmatrix} 1 & .2 & .1 \\ .1 & 2 & .1 \\ .3 & .1 & 3 \end{bmatrix}, \quad \gamma = 3, \quad \alpha = 1.2, \quad \sigma = 0.01.$$

Plot S_i and p as a function of t , and compare it to the target value $\alpha\gamma$. Repeat for $\gamma = 5$. Comment briefly on what you observe. *Comment:* You'll soon understand what you see.

2.2 *State equations for a linear mechanical system.* The equations of motion of a lumped mechanical system undergoing small motions can be expressed as

$$M\ddot{q} + D\dot{q} + Kq = f$$

where $q(t) \in \mathbf{R}^k$ is the vector of deflections, M , D , and K are the *mass*, *damping*, and *stiffness* matrices, respectively, and $f(t) \in \mathbf{R}^k$ is the vector of externally applied forces. Assuming M is invertible, write linear system equations for the mechanical system, with state $x = [q^T \dot{q}^T]^T$, input $u = f$, and output $y = q$.

- 2.3 *Some standard time-series models.* A time series is just a discrete-time signal, *i.e.*, a function from \mathbf{Z}_+ into \mathbf{R} . We think of $u(k)$ as the value of the signal or quantity u at time (or *epoch*) k . The study of time series predates the extensive study of state-space linear systems, and is used in many fields (*e.g.*, econometrics). Let u and y be two time series (input and output, respectively). The relation (or *time series model*)

$$y(k) = a_0 u(k) + a_1 u(k-1) + \cdots + a_r u(k-r)$$

is called a *moving average (MA) model*, since the output at time k is a weighted average of the previous r inputs, and the set of variables over which we average ‘slides along’ with time. Another model is given by

$$y(k) = u(k) + b_1 y(k-1) + \cdots + b_p y(k-p).$$

This model is called an *autoregressive (AR) model*, since the current output is a linear combination of (*i.e.*, regression on) the current input and some previous values of the output. Another widely used model is the *autoregressive moving average (ARMA) model*, which combines the MA and AR models:

$$y(k) = b_1 y(k-1) + \cdots + b_p y(k-p) + a_0 u(k) + \cdots + a_r u(k-r).$$

Finally, the problem: Express each of these models as a linear dynamical system with input u and output y . For the MA model, use state

$$x(k) = \begin{bmatrix} u(k-1) \\ \vdots \\ u(k-r) \end{bmatrix},$$

and for the AR model, use state

$$x(k) = \begin{bmatrix} y(k-1) \\ \vdots \\ y(k-p) \end{bmatrix}.$$

You decide on an appropriate state vector for the ARMA model. (There are many possible choices for the state here, even with different dimensions. We recommend you choose a state for the ARMA model that makes it easy for you to derive the state equations.) **Remark:** multi-input, multi-output time-series models (*i.e.*, $u(k) \in \mathbf{R}^m$, $y(k) \in \mathbf{R}^p$) are readily handled by allowing the coefficients a_i , b_i to be matrices.

- 2.4 *Representing linear functions as matrix multiplication.* Suppose that $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ is linear. Show that there is a matrix $A \in \mathbf{R}^{m \times n}$ such that for all $x \in \mathbf{R}^n$, $f(x) = Ax$. (Explicitly describe how you get the coefficients A_{ij} from f , and then verify that $f(x) = Ax$ for any $x \in \mathbf{R}^n$.) Is the matrix A that represents f unique? In other words, if $\tilde{A} \in \mathbf{R}^{m \times n}$ is another matrix such that $f(x) = \tilde{A}x$ for all $x \in \mathbf{R}^n$, then do we have $\tilde{A} = A$? Either show that this is so, or give an explicit counterexample.
- 2.5 *Some linear functions associated with a convolution system.* Suppose that u and y are scalar-valued discrete-time signals (*i.e.*, sequences) related via convolution:

$$y(k) = \sum_j h_j u(k-j), \quad k \in \mathbf{Z},$$

where $h_k \in \mathbf{R}$. You can assume that the convolution is *causal*, *i.e.*, $h_j = 0$ when $j < 0$.

- (a) *The input/output (Toeplitz) matrix.* Assume that $u(k) = 0$ for $k < 0$, and define

$$U = \begin{bmatrix} u(0) \\ u(1) \\ \vdots \\ u(N) \end{bmatrix}, \quad Y = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(N) \end{bmatrix}.$$

Thus U and Y are vectors that give the first $N + 1$ values of the input and output signals, respectively. Find the matrix T such that $Y = TU$. The matrix T describes the linear mapping from (a chunk of) the input to (a chunk of) the output. T is called the input/output or Toeplitz matrix (of size $N + 1$) associated with the convolution system.

- (b) *The Hankel matrix.* Now assume that $u(k) = 0$ for $k > 0$ or $k < -N$ and let

$$U = \begin{bmatrix} u(0) \\ u(-1) \\ \vdots \\ u(-N) \end{bmatrix}, \quad Y = \begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(N) \end{bmatrix}.$$

Here U gives the *past input* to the system, and Y gives (a chunk of) the resulting future output. Find the matrix H such that $Y = HU$. H is called the Hankel matrix (of size $N + 1$) associated with the convolution system.

- 2.6 *Matrix representation of polynomial differentiation.* We can represent a polynomial of degree less than n ,

$$p(x) = a_{n-1}x^{n-1} + a_{n-2}x^{n-2} + \cdots + a_1x + a_0,$$

as the vector $[a_0 \ a_1 \ \cdots \ a_{n-1}]^T \in \mathbf{R}^n$. Consider the linear transformation \mathcal{D} that differentiates polynomials, i.e., $\mathcal{D}p = dp/dx$. Find the matrix D that represents \mathcal{D} (i.e., if the coefficients of p are given by a , then the coefficients of dp/dx are given by Da).

- 2.7 Consider the (discrete-time) linear dynamical system

$$x(t+1) = A(t)x(t) + B(t)u(t), \quad y(t) = C(t)x(t) + D(t)u(t).$$

Find a matrix G such that

$$\begin{bmatrix} y(0) \\ y(1) \\ \vdots \\ y(N) \end{bmatrix} = G \begin{bmatrix} x(0) \\ u(0) \\ \vdots \\ u(N) \end{bmatrix}.$$

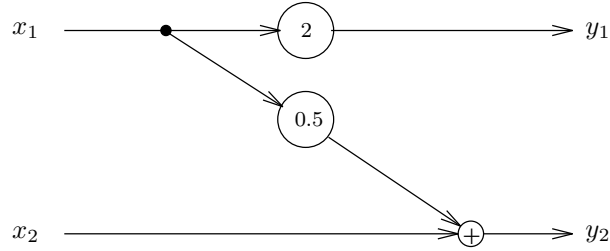
The matrix G shows how the output at $t = 0, \dots, N$ depends on the initial state $x(0)$ and the sequence of inputs $u(0), \dots, u(N)$.

- 2.8 *Some sparsity patterns.*

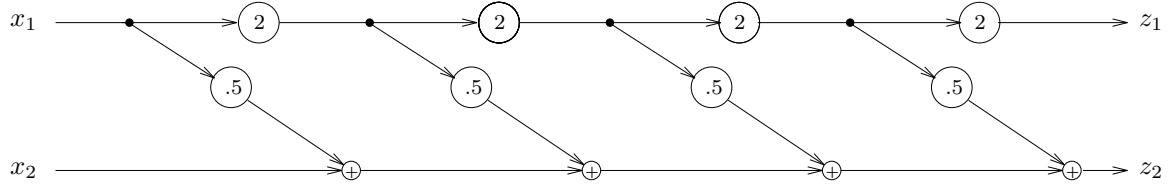
- (a) A matrix $A \in \mathbf{R}^{n \times n}$ is *tridiagonal* if $A_{ij} = 0$ for $|i - j| > 1$. Draw a block diagram of $y = Ax$ for A tridiagonal.
- (b) Consider a certain linear mapping $y = Ax$ with $A \in \mathbf{R}^{m \times n}$. For i odd, y_i depends only on x_j for j even. Similarly, for i even, y_i depends only on x_j for j odd. Describe the sparsity structure of A . Give the structure a reasonable, suggestive name.

- 2.9 *Matrices and signal flow graphs.*

- (a) Find $A \in \mathbf{R}^{2 \times 2}$ such that $y = Ax$ in the system below:



(b) Find $B \in \mathbf{R}^{2 \times 2}$ such that $z = Bx$ in the system below:



Do this two ways: first, by expressing the matrix B in terms of A from the previous part (explaining why they are related as you claim); and second, by directly evaluating all possible paths from each x_j to each z_i .

2.10 *Mass/force example.* Find the matrix A for the mass/force example in the lecture notes. For $n = 4$, find a specific input force sequence x that moves the mass to final position 1 and final velocity zero.

2.11 *Optimal force design.* We consider the mass/force example in the lecture notes, and in exercise 10, with $n = 4$, and the requirement that the final position is 1 and final velocity is 0. Roughly speaking, you have four variables and two equations, and therefore two extra degrees of freedom. In this problem you use the extra degrees of freedom to achieve other objectives, *i.e.*, minimize some cost functions that are described below.

- Find f that meets the specifications and minimizes the sum of squares of the forces, *i.e.*, $f_1^2 + f_2^2 + f_3^2 + f_4^2$.
- Find f that meets the specifications and minimizes the maximum force applied, *i.e.*, $\max\{|f_1|, |f_2|, |f_3|, |f_4|\}$.
- Find f that meets the specifications and minimizes the sum of absolute values of the forces applied, *i.e.*, $|f_1| + |f_2| + |f_3| + |f_4|$. (This would correspond to *minimum fuel usage* if the force comes from a thruster.)

There might be more than one minimizer; we just want one. If you can solve these optimization problems exactly, that's fine; but you are free to solve the problems numerically (and even approximately). We don't need an analysis verifying that your solution is indeed the best one. Make sure your solutions make sense to you. *Hints:*

- You can pick f_1 and f_2 to be anything; then f_3 and f_4 can be expressed in terms of f_1 and f_2 .
- You can plot, or just evaluate, each cost function as a function of f_1 and f_2 , over some reasonable range, to find the minimum.
- To evaluate the cost functions (which will be functions of two variables) for their minima, you might find one or more of the Matlab commands `norm`, `sum`, `abs`, `max`, `min` and `find` useful. (You can type `help norm`, for example, to see what `norm` does.) For example, if A is a matrix, `min(min(A))` returns the smallest element of the matrix. The statement `[i,j]=find(A==min(min(A)))` finds a pair of i, j indices in the matrix that correspond to the minimum value. Of course, there are many other ways to evaluate these functions; you don't have to use these commands.

- 2.12 *Counting paths in an undirected graph.* Consider an undirected graph with n nodes, and no self loops (*i.e.*, all branches connect two different nodes). Let $A \in \mathbf{R}^{n \times n}$ be the *node adjacency matrix*, defined as

$$A_{ij} = \begin{cases} 1 & \text{if there is a branch from node } i \text{ to node } j \\ 0 & \text{if there is no branch from node } i \text{ to node } j \end{cases}$$

Note that $A = A^T$, and $A_{ii} = 0$ since there are no self loops. We can interpret A_{ij} (which is either zero or one) as the number of branches that connect node i to node j . Let $B = A^k$, where $k \in \mathbf{Z}$, $k \geq 1$. Give a simple interpretation of B_{ij} in terms of the original graph. (You might need to use the concept of a *path* of length m from node p to node q .)

- 2.13 *Counting sequences in a language or code.* We consider a language or code with an alphabet of n symbols $1, 2, \dots, n$. A sentence is a finite sequence of symbols, k_1, \dots, k_m where $k_i \in \{1, \dots, n\}$. A language or code consists of a set of sequences, which we will call the *allowable sequences*. A language is called *Markov* if the allowed sequences can be described by giving the allowable transitions between consecutive symbols. For each symbol we give a set of symbols which are allowed to follow the symbol. As a simple example, consider a Markov language with three symbols $1, 2, 3$. Symbol 1 can be followed by 1 or 3; symbol 2 must be followed by 3; and symbol 3 can be followed by 1 or 2. The sentence 1132313 is allowable (*i.e.*, in the language); the sentence 1132312 is not allowable (*i.e.*, not in the language). To describe the allowed symbol transitions we can define a matrix $A \in \mathbf{R}^{n \times n}$ by

$$A_{ij} = \begin{cases} 1 & \text{if symbol } i \text{ is allowed to follow symbol } j \\ 0 & \text{if symbol } i \text{ is not allowed to follow symbol } j \end{cases}$$

- (a) Let $B = A^k$. Give an interpretation of B_{ij} in terms of the language.
- (b) Consider the Markov language with five symbols $1, 2, 3, 4, 5$, and the following transition rules:
 - 1 must be followed by 2 or 3
 - 2 must be followed by 2 or 5
 - 3 must be followed by 1
 - 4 must be followed by 4 or 2 or 5
 - 5 must be followed by 1 or 3

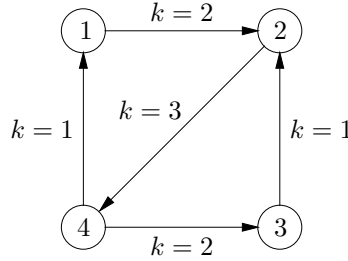
Find the total number of allowed sentences of length 10. Compare this number to the simple code that consists of all sequences from the alphabet (*i.e.*, all symbol transitions are allowed). In addition to giving the answer, you must explain how you solve the problem. Do not hesitate to use Matlab.

- 2.14 *Most common symbol in a given position.* Consider the Markov language from exercise 13, with five symbols $1, 2, 3, 4, 5$, and the following symbol transition rules:

- 1 must be followed by 2 or 3
- 2 must be followed by 2 or 5
- 3 must be followed by 1
- 4 must be followed by 4 or 2 or 5
- 5 must be followed by 1 or 3

Among all allowed sequences of length 10, find the most common value for the seventh symbol. In principle you could solve this problem by writing down all allowed sequences of length 10, and counting how many of these have symbol i as the seventh symbol, for $i = 1, \dots, 5$. (We're interested in the symbol for which this count is largest.) But we'd like you to use a smarter approach. Explain clearly how you solve the problem, as well as giving the specific answer.

2.15 *Communication over a wireless network with time-slots.* We consider a network with n nodes, labeled $1, \dots, n$. A directed graph shows which nodes can send messages (directly) to which other nodes; specifically, an edge from node j to node i means that node j can transmit a message directly to node i . Each edge is assigned to one of K *time-slots*, which are labeled $1, \dots, K$. At time period $t = 1$, only the edges assigned to time-slot 1 can transmit a message; at time period $t = 2$, only the edges assigned to time-slot 2 can transmit a message, and so on. After time period $t = K$, the pattern repeats. At time period $t = K + 1$, the edges assigned to time-slot 1 are again active; at $t = K + 2$, the edges assigned to time-slot 2 are active, etc. This cycle repeats indefinitely: when $t = mK + k$, where m is an integer, and $k \in \{1, \dots, K\}$, transmissions can occur only over edges assigned to time-slot k . Although it doesn't matter for the problem, we mention some reasons why the possible transmissions are assigned to time-slots. Two possible transmissions are assigned to different time-slots if they would interfere with each other, or if they would violate some limit (such as on the total power available at a node) if the transmissions occurred simultaneously. A message or packet can be sent from one node to another by a sequence of transmissions from node to node. At time period t , the message can be sent across any edge that is active at period t . It is also possible to store a message at a node during any time period, presumably for transmission during a later period. If a message is sent from node j to node i in period t , then in period $t + 1$ the message is at node i , and can be stored there, or transmitted across any edge emanating from node i and active at time period $t + 1$. To make sure the terminology is clear, we consider the very simple example shown below, with $n = 4$ nodes, and $K = 3$ time-slots.



In this example, we can send a message that starts in node 1 to node 3 as follows:

- During period $t = 1$ (time-slot $k = 1$), store it at node 1.
- During period $t = 2$ (time-slot $k = 2$), transmit it to node 2.
- During period $t = 3$ (time-slot $k = 3$), transmit it to node 4.
- During period $t = 4$ (time-slot $k = 1$), store it at node 4.
- During period $t = 5$ (time-slot $k = 2$), transmit it to node 3.

You can check that at each period, the transmission used is active, *i.e.*, assigned to the associated time-slot. The sequence of transmissions (and storing) described above gets the message from node 1 to node 3 in 5 periods. Finally, the problem. We consider a specific network with $n = 20$ nodes, and $K = 3$ time-slots, with edges and time-slot assignments given in `ts_data.m`. The labeled graph that specifies the possible transmissions and the associated time-slot assignments are given in a matrix $A \in \mathbf{R}^{n \times n}$, as follows:

$$A_{ij} = \begin{cases} k & \text{if transmission from node } j \text{ to node } i \text{ is allowed, and assigned to time-slot } k \\ 0 & \text{if transmission from node } j \text{ to node } i \text{ is never allowed} \\ 0 & i = j. \end{cases}$$

Note that we set $A_{ii} = 0$ for convenience. This choice has no significance; you can store a message at any node in any period. To illustrate this encoding of the graph, consider the simple example described

above. For this example, we have

$$A_{\text{example}} = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 2 & 0 & 1 & 0 \\ 0 & 0 & 0 & 2 \\ 0 & 3 & 0 & 0 \end{bmatrix}.$$

The problems below concern the network described in the mfile `ts_data.m`, and *not* the simple example given above.

- (a) *Minimum-time point-to-point routing.* Find the fastest way to get a message that starts at node 5, to node 18. Give your solution as a prescription ordered in time from $t = 1$ to $t = T$ (the last transmission), as in the example above. At each time period, give the transmission (as in ‘transmit from node 7 to node 9’) or state that the message is to be stored (as in ‘store at node 13’). Be sure that transmissions only occur during the associated time-slots. You only need to give *one* prescription for getting the message from node 5 to node 18 in minimum time.
- (b) *Minimum time flooding.* In this part of the problem, we assume that once the message reaches a node, a copy is kept there, even when the message is transmitted to another node. Thus, the message is available at the node to be transmitted along any active edge emanating from that node, at any future period. Moreover, we allow multi-cast: if during a time period there are multiple active edges emanating from a node that has (a copy of) the message, then transmission can occur during that time period across *all* (or any subset) of the active edges. In this part of the problem, we are interested in getting a message that starts at a particular node, to all others, and we attach no cost to storage or transmission, so there is no harm in assuming that at each time period, every node that has the message forwards it to all nodes it is able to transmit to. What is the minimum time it takes before all nodes have a message that starts at node 7?

For both parts of the problem, you must give the specific solution, as well as a description of your approach and method.

- 2.16 *Solving triangular linear equations.* Consider the linear equations $y = Rx$, where $R \in \mathbf{R}^{n \times n}$ is upper triangular and invertible. Suggest a simple algorithm to solve for x given R and y . *Hint:* first find x_n ; then find x_{n-1} (remembering that now you know x_n); then find x_{n-2} (remembering that now you know x_n and x_{n-1}); etc. **Remark:** the algorithm you will discover is called *back substitution*. It requires order n^2 floating point operations (flops); most methods for solving $y = Ax$ for general $A \in \mathbf{R}^{n \times n}$ require order n^3 flops.
- 2.17 *Gradient of some common functions.* Recall that the gradient of a differentiable function $f : \mathbf{R}^n \rightarrow \mathbf{R}$, at a point $x \in \mathbf{R}^n$, is defined as the vector

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix},$$

where the partial derivatives are evaluated at the point x . The first order Taylor approximation of f , near x , is given by

$$\hat{f}_{\text{tay}}(z) = f(x) + \nabla f(x)^T(z - x).$$

This function is affine, *i.e.*, a linear function plus a constant. For z near x , the Taylor approximation \hat{f}_{tay} is very near f . Find the gradient of the following functions. Express the gradients using matrix notation.

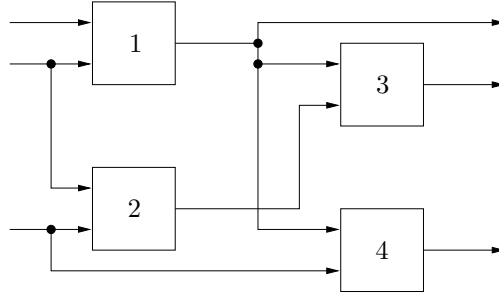
- (a) $f(x) = a^T x + b$, where $a \in \mathbf{R}^n$, $b \in \mathbf{R}$.
- (b) $f(x) = x^T A x$, for $A \in \mathbf{R}^{n \times n}$.

(c) $f(x) = x^T A x$, where $A = A^T \in \mathbf{R}^{n \times n}$. (Yes, this is a special case of the previous one.)

- 2.18 *Digital circuit gate sizing.* A digital circuit consists of a set of n (logic) gates, interconnected by wires. Each gate has one or more inputs (typically between one and four), and one output, which is connected via the wires to other gate inputs and possibly to some external circuitry. When the output of gate i is connected to an input of gate j , we say that gate i *drives* gate j , or that gate j is in the *fan-out* of gate i . We describe the topology of the circuit by the *fan-out list* for each gate, which tells us which other gates the output of a gate connects to. We denote the fan-out list of gate i as $\text{FO}(i) \subseteq \{1, \dots, n\}$. We can have $\text{FO}(i) = \emptyset$, which means that the output of gate i does not connect to the inputs of any of the gates $1, \dots, n$ (presumably the output of gate i connects to some external circuitry). It's common to order the gates in such a way that each gate only drives gates with higher indices, *i.e.*, we have $\text{FO}(i) \subseteq \{i+1, \dots, n\}$. We'll assume that's the case here. (This means that the gate interconnections form a directed acyclic graph.)

To illustrate the notation, a simple digital circuit with $n = 4$ gates, each with 2 inputs, is shown below. For this circuit we have

$$\text{FO}(1) = \{3, 4\}, \quad \text{FO}(2) = \{3\}, \quad \text{FO}(3) = \emptyset, \quad \text{FO}(4) = \emptyset.$$



The 3 input signals arriving from the left are called *primary inputs*, and the 3 output signals emerging from the right are called *primary outputs* of the circuit. (You don't need to know this, however, to solve this problem.)

Each gate has a (real) *scale factor* or *size* x_i . These scale factors are the design variables in the gate sizing problem. They must satisfy $1 \leq x_i \leq x^{\max}$, where x^{\max} is a given maximum allowed gate scale factor (typically on the order of 100). The total area of the circuit has the form

$$A = \sum_{i=1}^n a_i x_i,$$

where a_i are positive constants.

Each gate has an *input capacitance* C_i^{in} , which depends on the scale factor x_i as

$$C_i^{\text{in}} = \alpha_i x_i,$$

where α_i are positive constants.

Each gate has a *delay* d_i , which is given by

$$d_i = \beta_i + \gamma_i C_i^{\text{load}} / x_i,$$

where β_i and γ_i are positive constants, and C_i^{load} is the *load capacitance* of gate i . Note that the gate delay d_i is always larger than β_i , which can be interpreted as the minimum possible delay of gate i , achieved only in the limit as the gate scale factor becomes large.

The load capacitance of gate i is given by

$$C_i^{\text{load}} = C_i^{\text{ext}} + \sum_{j \in \text{FO}(i)} C_j^{\text{in}},$$

where C_i^{ext} is a positive constant that accounts for the capacitance of the interconnect wires and external circuitry.

We will follow a simple design method, which assigns an equal delay T to all gates in the circuit, *i.e.*, we have $d_i = T$, where $T > 0$ is given. For a given value of T , there may or may not exist a feasible design (*i.e.*, a choice of the x_i , with $1 \leq x_i \leq x^{\text{max}}$) that yields $d_i = T$ for $i = 1, \dots, n$. We can assume, of course, that $T > \max_i \beta_i$, *i.e.*, T is larger than the largest minimum delay of the gates.

Finally, we get to the problem.

- (a) Explain how to find a design $x^* \in \mathbf{R}^n$ that minimizes T , subject to a given area constraint $A \leq A^{\text{max}}$. You can assume the fanout lists, and all constants in the problem description are known; your job is to find the scale factors x_i . Be sure to explain how you determine if the design problem is feasible, *i.e.*, whether or not there is an x that gives $d_i = T$, with $1 \leq x_i \leq x^{\text{max}}$, and $A \leq A^{\text{max}}$.

Your method can involve any of the methods or concepts we have seen so far in the course. It can also involve a simple search procedure, *e.g.*, trying (many) different values of T over a range.

Note: this problem concerns the general case, and not the simple example shown above.

- (b) Carry out your method on the particular circuit with data given in the file `gate_sizing_data.m`. The fan-out lists are given as an $n \times n$ matrix \mathbf{F} , with i, j entry one if $j \in \text{FO}(i)$, and zero otherwise. In other words, the i th row of \mathbf{F} gives the fanout of gate i . The j th entry in the i th row is 1 if gate j is in the fan-out of gate i , and 0 otherwise.

Comment. You do not need to know anything about digital circuits; *everything* you need to know is stated above.

- 2.19 *Some matrices from signal processing.* We consider $x \in \mathbf{R}^n$ as a signal, with x_i the (scalar) value of the signal at (discrete) time period i , for $i = 1, \dots, n$. Below we describe several transformations of the signal x , that produce a new signal y (whose dimension varies). For each one, find a matrix A for which $y = Ax$. formula for

- (a) $2 \times$ *up-conversion with linear interpolation.* We take $y \in \mathbf{R}^{2n-1}$. For i odd, $y_i = x_{(i+1)/2}$. For i even, $y_i = (x_{i/2} + x_{i/2+1})/2$. Roughly speaking, this operation doubles the sample rate, inserting new samples in between the original ones using linear interpolation.
- (b) $2 \times$ *down-sampling.* We assume here that n is even, and take $y \in \mathbf{R}^{n/2}$, with $y_i = x_{2i}$.
- (c) $2 \times$ *down-sampling with averaging.* We assume here that n is even, and take $y \in \mathbf{R}^{n/2}$, with $y_i = (x_{2i-1} + x_{2i})/2$.

- 2.20 *Affine functions.* A function $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ is called *affine* if for any $x, y \in \mathbf{R}^n$ and any $\alpha, \beta \in \mathbf{R}$ with $\alpha + \beta = 1$, we have

$$f(\alpha x + \beta y) = \alpha f(x) + \beta f(y).$$

(Without the restriction $\alpha + \beta = 1$, this would be the definition of linearity.)

- (a) Suppose that $A \in \mathbf{R}^{m \times n}$ and $b \in \mathbf{R}^m$. Show that the function $f(x) = Ax + b$ is affine.
- (b) Now the converse: Show that any affine function f can be represented as $f(x) = Ax + b$, for some $A \in \mathbf{R}^{m \times n}$ and $b \in \mathbf{R}^m$. (This representation is unique: for a given affine function f there is only one A and one b for which $f(x) = Ax + b$ for all x .)

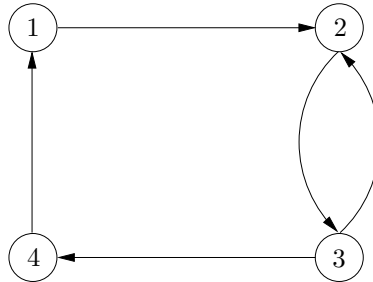
Hint. Show that the function $g(x) = f(x) - f(0)$ is linear.

You can think of an affine function as a linear function, plus an offset. In some contexts, affine functions are (mistakenly, or informally) called linear, even though in general they are not. (Example: $y = mx + b$ is described as ‘linear’ in US high schools.)

- 2.21 *Paths and cycles in a directed graph.* We consider a directed graph with n nodes. The graph is specified by its *node adjacency matrix* $A \in \mathbf{R}^{n \times n}$, defined as

$$A_{ij} = \begin{cases} 1 & \text{if there is an edge from node } j \text{ to node } i \\ 0 & \text{otherwise.} \end{cases}$$

Note that the edges are *oriented*, i.e., $A_{34} = 1$ means there is an edge from node 4 to node 3. For simplicity we do not allow self-loops, i.e., $A_{ii} = 0$ for all i , $1 \leq i \leq n$. A simple example illustrating this notation is shown below.



The node adjacency matrix for this example is

$$A = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

In this example, nodes 2 and 3 are connected in both directions, i.e., there is an edge from 2 to 3 and also an edge from 3 to 2. A *path* of length $l > 0$ from node j to node i is a sequence $s_0 = j, s_1, \dots, s_l = i$ of nodes, with $A_{s_{k+1}, s_k} = 1$ for $k = 0, 1, \dots, l-1$. For example, in the graph shown above, 1, 2, 3, 2 is a path of length 3. A *cycle* of length l is a path of length l , with the same starting and ending node, with no repeated nodes other than the endpoints. In other words, a cycle is a sequence of nodes of the form $s_0, s_1, \dots, s_{l-1}, s_0$, with

$$A_{s_1, s_0} = 1, \quad A_{s_2, s_1} = 1, \quad \dots \quad A_{s_0, s_{l-1}} = 1,$$

and

$$s_i \neq s_j \text{ for } i \neq j, \quad i, j = 0, \dots, l-1.$$

For example, in the graph shown above, 1, 2, 3, 4, 1 is a cycle of length 4. The rest of this problem concerns a specific graph, given in the file `directed_graph.m` on the course web site. For each of the following questions, you must give the answer explicitly (for example, enclosed in a box). You must also explain clearly how you arrived at your answer.

- What is the length of a shortest cycle? (Shortest means minimum length.)
- What is the length of a shortest path from node 13 to node 17? (If there are no paths from node 13 to node 17, you can give the answer as ‘infinity’.)
- What is the length of a shortest path from node 13 to node 17, that *does not* pass through node 3?

- (d) What is the length of a shortest path from node 13 to node 17, that *does* pass through node 9?
- (e) Among all paths of length 10 that start at node 5, find the most common ending node.
- (f) Among all paths of length 10 that end at node 8, find the most common starting node.
- (g) Among all paths of length 10, find the most common pair of starting and ending nodes. In other words, find i, j which maximize the number of paths of length 10 from i to j .

Lecture 3 – Linear algebra review

3.1 *Price elasticity of demand.* The demand for n different goods as a function of their prices is described by a function $f(\cdot)$ from \mathbf{R}^n to \mathbf{R}^n :

$$q = f(p),$$

where p is the price vector, and q is the demand vector. Linear models of the demand function are often used to analyze the effects of small price changes. Denoting the current price and current demand vectors by p^* and q^* , we have that $q^* = f(p^*)$, and the linear approximation is:

$$q^* + \delta q \approx f(p^*) + \left. \frac{df}{dp} \right|_{p^*} \delta p.$$

This is usually rewritten in term of the *elasticity matrix* E , with entries

$$e_{ij} = \left. \frac{df_i}{dp_j} \right|_{p_j^*} \frac{1/q_i^*}{1/p_j^*}$$

(i.e., relative change in demand per relative change in price.) Define the vector y of relative demand changes, and the vector x of relative price changes,

$$y_i = \frac{\delta q_i}{q_i^*}, \quad x_j = \frac{\delta p_j}{p_j^*},$$

and, finally, we have the linear model $y = Ex$. Here are the questions:

- (a) What is a reasonable assumption about the diagonal elements e_{ii} of the elasticity matrix?
- (b) Consider two goods. The off-diagonal elements of E describe how the demand for one good is affected by changes in the price of the other good. Assume $e_{11} = e_{22} = -1$ and $e_{12} = e_{21}$, that is,

$$E = \begin{bmatrix} -1 & e_{12} \\ e_{12} & -1 \end{bmatrix}.$$

Two goods are called *substitutes* if they provide a similar service or other satisfaction (for example: train tickets and bus tickets, cake and pie, etc.) Two goods are called *complements* if they tend to be used together (for example: automobiles and gasoline, left and right shoes, etc.) For each of these two generic situations, what can you say about e_{12} ?

- (c) Suppose the price elasticity of demand matrix is

$$E = \begin{bmatrix} -1 & -1 \\ -1 & -1 \end{bmatrix}.$$

Describe the *nullspace* of E , and give an interpretation (in one or two sentences.) What kind of goods could have such an elasticity matrix?

3.2 *Color perception.* Human color perception is based on the responses of three different types of color light receptors, called *cones*. The three types of cones have different spectral response characteristics and are called L, M, and, S because they respond mainly to long, medium, and short wavelengths, respectively. In this problem we will divide the visible spectrum into 20 bands, and model the cones' response as follows:

$$L_{\text{cone}} = \sum_{i=1}^{20} l_i p_i, \quad M_{\text{cone}} = \sum_{i=1}^{20} m_i p_i, \quad S_{\text{cone}} = \sum_{i=1}^{20} s_i p_i,$$

where p_i is the incident power in the i th wavelength band, and l_i , m_i and s_i are nonnegative constants that describe the spectral response of the different cones. The perceived color is a complex function of the three cone responses, *i.e.*, the vector $(L_{\text{cone}}, M_{\text{cone}}, S_{\text{cone}})$, with different cone response vectors perceived as different colors. (Actual color perception is a bit more complicated than this, but the basic idea is right.)

- (a) *Metamers.* When are two light spectra, p and \tilde{p} , visually indistinguishable? (Visually identical lights with different spectral power compositions are called *metamers*.)
- (b) *Visual color matching.* In a color matching problem, an observer is shown a test light and is asked to change the intensities of three primary lights until the sum of the primary lights looks like the test light. In other words, the observer is asked to find a spectrum of the form

$$p_{\text{match}} = a_1 u + a_2 v + a_3 w,$$

where u , v , w are the spectra of the primary lights, and a_i are the (nonnegative) intensities to be found, that is visually indistinguishable from a given test light spectrum p_{test} . Can this always be done? Discuss briefly.

- (c) *Visual matching with phosphors.* A computer monitor has three phosphors, R , G , and B . It is desired to adjust the phosphor intensities to create a color that looks like a reference test light. Find weights that achieve the match or explain why no such weights exist. The data for this problem is in the an m-file `color_perception.m`. Running `color_perception` will define and plot the vectors `wavelength`, `B_phosphor`, `G_phosphor`, `R_phosphor`, `L_coefficients`, `M_coefficients`, `S_coefficients`, and `test_light`.
- (d) *Effects of illumination.* An object's surface can be characterized by its reflectance (*i.e.*, the fraction of light it reflects) for each band of wavelengths. If the object is illuminated with a light spectrum characterized by I_i , and the reflectance of the object is r_i (which is between 0 and 1), then the reflected light spectrum is given by $I_i r_i$, where $i = 1, \dots, 20$ denotes the wavelength band. Now consider two objects illuminated (at different times) by two different light sources, say an incandescent bulb and sunlight. Sally argues that if the two objects look identical when illuminated by a tungsten bulb, they will look identical when illuminated by sunlight. Beth disagrees: she says that two objects can appear identical when illuminated by a tungsten bulb, but look different when lit by sunlight. Who is right? If Sally is right, explain why. If Beth is right give an example of two objects that appear identical under one light source and different under another. You can use the vectors `sunlight` and `tungsten` defined in `color_perception.m` as the light sources.

3.3 *Halfspace.* Suppose $a, b \in \mathbf{R}^n$ are two given points. Show that the set of points in \mathbf{R}^n that are closer to a than b is a halfspace, *i.e.*:

$$\{x \mid \|x - a\| \leq \|x - b\|\} = \{x \mid c^T x \leq d\}$$

for appropriate $c \in \mathbf{R}^n$ and $d \in \mathbf{R}$. Give c and d explicitly, and draw a picture showing a , b , c , and the halfspace.

3.4 *Some properties of the product of two matrices.* For each of the following statements, either show that it is true, or give a (specific) counterexample.

- If AB is full rank then A and B are full rank.
- If A and B are full rank then AB is full rank.
- If A and B have zero nullspace, then so does AB .
- If A and B are onto, then so is AB .

You can assume that $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{n \times p}$. Some of the false statements above become true under certain assumptions on the dimensions of A and B . As a trivial example, all of the statements above are true when A and B are scalars, *i.e.*, $n = m = p = 1$. For each of the statements above, find conditions on n , m , and p that make them true. Try to find the most general conditions you can. You can give your conditions as inequalities involving n , m , and p , or you can use more informal language such as “ A and B are both skinny.”

3.5 *Rank of a product.* Suppose that $A \in \mathbf{R}^{7 \times 5}$ has rank 4, and $B \in \mathbf{R}^{5 \times 7}$ has rank 3. What values can $\mathbf{Rank}(AB)$ possibly have? For each value r that is possible, give an example, *i.e.*, a specific A and B with the dimensions and ranks given above, for which $\mathbf{Rank}(AB) = r$. Please try to give simple examples, that make it easy for you to justify that the ranks of A , B , and AB are what you claim they are. We will deduct points for correct examples that are needlessly complicated. You can use Matlab to verify the ranks, but we don’t recommend it: numerical roundoff errors in Matlab’s calculations can sometimes give errors in computing the rank. (Matlab may still be useful; you just have to double check that the ranks it finds are correct.) Explain briefly why the rank of AB must be one of the values you give.

3.6 *Linearizing range measurements.* Consider a single (scalar) measurement y of the distance or range of $x \in \mathbf{R}^n$ to a fixed point or beacon at a , *i.e.*, $y = \|x - a\|$.

- (a) Show that the linearized model near x_0 can be expressed as $\delta y = k^T \delta x$, where k is the unit vector (*i.e.*, with length one) pointing from a to x_0 . Derive this analytically, and also draw a picture (for $n = 2$) to demonstrate it.
- (b) Consider the error e of the linearized approximation, *i.e.*,

$$e = \|x_0 + \delta x - a\| - \|x_0 - a\| - k^T \delta x.$$

The relative error of the approximation is given by $\eta = e/\|x_0 - a\|$. We know, of course, that the absolute value of the relative error is very small provided δx is small. In many specific applications, it is possible and useful to make a stronger statement, for example, to derive a bound on how large the error can be. You will do that here. In fact you will prove that

$$0 \leq \eta \leq \frac{\alpha^2}{2}$$

where $\alpha = \|\delta x\|/\|x_0 - a\|$ is the relative size of δx . For example, for a relative displacement of $\alpha = 1\%$, we have $\eta \leq 0.00005$, *i.e.*, the linearized model is accurate to about 0.005%. To prove this bound you can proceed as follows:

- Show that $\eta = -1 + \sqrt{1 + \alpha^2 + 2\beta} - \beta$ where $\beta = k^T \delta x / \|x_0 - a\|$.
- Verify that $|\beta| \leq \alpha$.
- Consider the function $g(\beta) = -1 + \sqrt{1 + \alpha^2 + 2\beta} - \beta$ with $|\beta| \leq \alpha$. By maximizing and minimizing g over the interval $-\alpha \leq \beta \leq \alpha$ show that

$$0 \leq \eta \leq \frac{\alpha^2}{2}.$$

3.7 *Orthogonal complement of a subspace.* If \mathcal{V} is a subspace of \mathbf{R}^n we define \mathcal{V}^\perp as the set of vectors orthogonal to every element in \mathcal{V} , *i.e.*,

$$\mathcal{V}^\perp = \{ x \mid \langle x, y \rangle = 0, \forall y \in \mathcal{V} \}.$$

- (a) Verify that \mathcal{V}^\perp is a subspace of \mathbf{R}^n .
- (b) Suppose \mathcal{V} is described as the span of some vectors v_1, v_2, \dots, v_r . Express \mathcal{V} and \mathcal{V}^\perp in terms of the matrix $V = \begin{bmatrix} v_1 & v_2 & \cdots & v_r \end{bmatrix} \in \mathbf{R}^{n \times r}$ using common terms (range, nullspace, transpose, etc.)

- (c) Show that every $x \in \mathbf{R}^n$ can be expressed uniquely as $x = v + v^\perp$ where $v \in \mathcal{V}$, $v^\perp \in \mathcal{V}^\perp$. *Hint:* let v be the projection of x on \mathcal{V} .
- (d) Show that $\dim \mathcal{V}^\perp + \dim \mathcal{V} = n$.
- (e) Show that $\mathcal{V} \subseteq \mathcal{U}$ implies $\mathcal{U}^\perp \subseteq \mathcal{V}^\perp$.
- 3.8 Consider the linearized navigation equations from the lecture notes. Find the conditions under which A has full rank. Describe the conditions geometrically (*i.e.*, in terms of the relative positions of the unknown coordinates and the beacons).
- 3.9 Suppose that $\angle(Ax, x) = 0$ for all $x \in \mathbf{R}^n$, *i.e.*, x and Ax always point in the same direction. What can you say about the matrix A ? Be very specific.
- 3.10 *Proof of Cauchy-Schwarz inequality.* You will prove the Cauchy-Schwarz inequality.
- (a) Suppose $a \geq 0$, $c \geq 0$, and for all $\lambda \in \mathbf{R}$, $a + 2b\lambda + c\lambda^2 \geq 0$. Show that $|b| \leq \sqrt{ac}$.
- (b) Given $v, w \in \mathbf{R}^n$ explain why $(v + \lambda w)^T(v + \lambda w) \geq 0$ for all $\lambda \in \mathbf{R}$.
- (c) Apply (a) to the quadratic resulting when the expression in (b) is expanded, to get the Cauchy-Schwarz inequality:
- $$|v^T w| \leq \sqrt{v^T v} \sqrt{w^T w}.$$
- (d) When does equality hold?
- 3.11 *Vector spaces over the Boolean field.* In this course the *scalar field*, *i.e.*, the components of vectors, will usually be the real numbers, and sometimes the complex numbers. It is also possible to consider vector spaces over other fields, for example \mathbf{Z}_2 , which consists of the two numbers 0 and 1, with Boolean addition and multiplication (*i.e.*, $1 + 1 = 0$). Unlike \mathbf{R} or \mathbf{C} , the field \mathbf{Z}_2 is finite, indeed, has only two elements. A vector in \mathbf{Z}_2^n is called a *Boolean vector*. Much of the linear algebra for \mathbf{R}^n and \mathbf{C}^n carries over to \mathbf{Z}_2^n . For example, we define a function $f : \mathbf{Z}_2^n \rightarrow \mathbf{Z}_2^m$ to be linear (over \mathbf{Z}_2) if $f(x + y) = f(x) + f(y)$ and $f(\alpha x) = \alpha f(x)$ for every $x, y \in \mathbf{Z}_2^n$ and $\alpha \in \mathbf{Z}_2$. It is easy to show that every linear function can be expressed as matrix multiplication, *i.e.*, $f(x) = Ax$, where $A \in \mathbf{Z}_2^{m \times n}$ is a Boolean matrix, and all the operations in the matrix multiplication are Boolean, *i.e.*, in \mathbf{Z}_2 . Concepts like nullspace, range, independence and rank are all defined in the obvious way for vector spaces over \mathbf{Z}_2 . Although we won't consider them in this course, there are many important applications of vector spaces and linear dynamical systems over \mathbf{Z}_2 . In this problem you will explore one simple example: block codes. *Linear block codes.* Suppose $x \in \mathbf{Z}_2^n$ is a Boolean vector we wish to transmit over an unreliable channel. In a linear block code, the vector $y = Gx$ is formed, where $G \in \mathbf{Z}_2^{m \times n}$ is the *coding matrix*, and $m > n$. Note that the vector y is 'redundant'; roughly speaking we have *coded* an n -bit vector as a (larger) m -bit vector. This is called an (n, m) code. The coded vector y is transmitted over the channel; the received signal \hat{y} is given by

$$\hat{y} = y + v,$$

where v is a noise vector (which usually is zero). This means that when $v_i = 0$, the i th bit is transmitted correctly; when $v_i = 1$, the i th bit is changed during transmission. In a *linear decoder*, the received signal is multiplied by another matrix: $\hat{x} = H\hat{y}$, where $H \in \mathbf{Z}_2^{n \times m}$. One reasonable requirement is that if the transmission is perfect, *i.e.*, $v = 0$, then the decoding is perfect, *i.e.*, $\hat{x} = x$. This holds if and only if H is a left inverse of G , *i.e.*, $HG = I_n$, which we assume to be the case.

- (a) What is the practical significance of $\mathcal{R}(G)$?
- (b) What is the practical significance of $\mathcal{N}(H)$?
- (c) A one-bit *error correcting code* has the property that for any noise v with one component equal to one, we still have $\hat{x} = x$. Consider $n = 3$. Either design a one-bit error correcting linear block code with the smallest possible m , or explain why it cannot be done. (By design we mean, give G and H explicitly and verify that they have the required properties.)

Remark: linear decoders are never used in practice; there are far better nonlinear ones.

- 3.12 *Quadratic extrapolation of a time series.* We are given a series z up to time n . Using a quadratic model, we want to extrapolate, or predict, $z(n+1)$ based on the three previous elements of the series: $z(n)$, $z(n-1)$, and $z(n-2)$. We'll denote the predicted value of $z(n+1)$ by y . Another way to describe this problem is: find a quadratic function $f(t) = a_2t^2 + a_1t + a_0$ which satisfies $f(n) = z(n)$, $f(n-1) = z(n-1)$, and $f(n-2) = z(n-2)$. The extrapolated value is then given by $y = f(n+1)$. Let the vector x denote the three previous elements of the series,

$$x = \begin{bmatrix} z(n) \\ z(n-1) \\ z(n-2) \end{bmatrix}.$$

You'll find a vector $c \in \mathbf{R}^3$ such that the extrapolated value is given by the linear transformation $y = c^T x$. You'll do this in two different ways.

- (a) *With matrix inversion.* Without loss of generality, make the problem independent of n by writing the extrapolating function as

$$f(n+k) = u_1k^2 + u_2k + u_3.$$

With the condition that f fits the series for $k = 0$, $k = -1$, and $k = -2$, find the matrix $A \in \mathbf{R}^{3 \times 3}$ that relates x to the vector u (of the coefficients of f),

$$x = Au.$$

Find the vector $b \in \mathbf{R}^3$ such that

$$y = b^T u.$$

Compute the inverse of A , and find c by eliminating u from the two previous equations.

- (b) *With basis functions.* Assume first that $x = e_1 = (1, 0, 0)$, and find y (*i.e.*, fit a quadratic to x and find its value at $n+1$.) Repeat this for $x = e_2 = (0, 1, 0)$, and for $x = e_3 = (0, 0, 1)$. Find the extrapolated value y for an arbitrary x as a linear combination of the three previous values.

- 3.13 *Right inverses.* This problem concerns the specific matrix

$$A = \begin{bmatrix} -1 & 0 & 0 & -1 & 1 \\ 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix}.$$

This matrix is full rank (*i.e.*, its rank is 3), so there exists at least one right inverse. In fact, there are many right inverses of A , which opens the possibility that we can seek right inverses that in addition have other properties. For each of the cases below, either find a specific matrix $B \in \mathbf{R}^{5 \times 3}$ that satisfies $AB = I$ and the given property, or explain why there is no such B . In cases where there is a right inverse B with the required property, you must briefly explain how you found your B . You must also attach a printout of some Matlab scripts that show the verification that $AB = I$. (We'll be very angry if we have to type in your 5×3 matrix into Matlab to check it.) When there is no right inverse with the given property, briefly explain why there is no such B .

- (a) The second row of B is zero.
- (b) The nullspace of B has dimension one.
- (c) The third column of B is zero.
- (d) The second and third rows of B are the same.
- (e) B is upper triangular, *i.e.*, $B_{ij} = 0$ for $i > j$.

(f) B is lower triangular, i.e., $B_{ij} = 0$ for $i < j$.

3.14 *Nonlinear unbiased estimators.* We consider the standard measurement setup:

$$y = Ax + v,$$

where $A \in \mathbf{R}^{m \times n}$, $x \in \mathbf{R}^n$ is the vector of parameters we wish to estimate, $y \in \mathbf{R}^m$ is the vector of measurements we take, and $v \in \mathbf{R}^m$ is the vector of measurement errors and noise. You may not assume anything about the dimensions of A , its rank, nullspace, etc. If the function $f : \mathbf{R}^m \rightarrow \mathbf{R}^n$ satisfies $f(Ax) = x$ for all $x \in \mathbf{R}^n$, then we say that f is an *unbiased estimator* (of x , given y). What this means is that if f is applied to our measurement vector, and $v = 0$, then f returns the true parameter value x . In EE263 we have studied linear unbiased estimators, which are unbiased estimators that are also linear functions. Here, though, we allow the possibility that f is nonlinear (which we take to mean, f is not linear). One of the following statements is true. Pick the statement that is true, and justify it completely. You can quote any result given in the lecture notes.

- A. *There is no such thing as a nonlinear unbiased estimator.* In other words, if f is any unbiased estimator, then f must be a linear function. (This statement is taken to be true if there are no unbiased estimators for a particular A .) If you believe this statement, explain why.
- B. *Nonlinear unbiased estimators do exist, but you don't need them.* In other words: it's possible to have a nonlinear unbiased estimator. But whenever there is a nonlinear unbiased estimator, there is also a linear unbiased estimator. If you believe this statement, then give a specific example of a matrix A , and an unbiased nonlinear estimator. Explain in the general case why a linear unbiased estimator exists whenever there is a nonlinear one.
- C. *There are cases for which nonlinear unbiased estimators exist, but no linear unbiased estimator exists.* If you believe this statement, give a specific example of a matrix A , and a nonlinear unbiased estimator, and also explain why no linear unbiased estimator exists.

3.15 *Channel equalizer with disturbance rejection.* A communication channel is described by $y = Ax + v$ where $x \in \mathbf{R}^n$ is the (unknown) transmitted signal, $y \in \mathbf{R}^m$ is the (known) received signal, $v \in \mathbf{R}^m$ is the (unknown) disturbance signal, and $A \in \mathbf{R}^{m \times n}$ describes the (known) channel. The disturbance v is known to be a linear combination of some (known) disturbance patterns,

$$d_1, \dots, d_k \in \mathbf{R}^m.$$

We consider linear equalizers for the channel, which have the form $\hat{x} = By$, where $B \in \mathbf{R}^{n \times m}$. (We'll refer to the matrix B as the equalizer; more precisely, you might say that B_{ij} are the equalizer coefficients.) We say the equalizer B *rejects* the disturbance pattern d_i if $\hat{x} = x$, no matter what x is, when $v = d_i$. If the equalizer rejects a set of disturbance patterns, for example, disturbances d_1 , d_3 , and d_7 (say), then it can reconstruct the transmitted signal exactly, when the disturbance v is any linear combination of d_1 , d_3 , and d_7 . Here is the problem. For the problem data given in `cedr_data.m`, find an equalizer B that rejects as many disturbance patterns as possible. (The disturbance patterns are given as an $m \times k$ matrix D , whose columns are the individual disturbance patterns.) Give the specific set of disturbance patterns that your equalizer rejects, as in 'My equalizer rejects three disturbance patterns: d_2 , d_3 , and d_6 .' (We only need *one* set of disturbances of the maximum size.) Explain how you know that there is no equalizer that rejects more disturbance patterns than yours does. Show the Matlab verification that your B does indeed reconstruct x , and rejects the disturbance patterns you claim it does. Show any other calculations needed to verify that your equalizer rejects the maximum number of patterns possible.

3.16 *Identifying a point on the unit sphere from spherical distances.* In this problem we consider the *unit sphere* in \mathbf{R}^n , which is defined as the set of vectors with norm one: $S^n = \{x \in \mathbf{R}^n \mid \|x\| = 1\}$. We define the *spherical distance* between two vectors on the unit sphere as the distance between them,

measured along the sphere, *i.e.*, as the angle between the vectors, measured in radians: If $x, y \in S^n$, the spherical distance between them is

$$\text{sphdist}(x, y) = \angle(x, y),$$

where we take the angle as lying between 0 and π . (Thus, the maximum distance between two points in S^n is π , which occurs only when the two points x, y are *antipodal*, which means $x = -y$.) Now suppose $p_1, \dots, p_k \in S^n$ are the (known) positions of some beacons on the unit sphere, and let $x \in S^n$ be an unknown point on the unit sphere. We have exact measurements of the (spherical) distances between each beacon and the unknown point x , *i.e.*, we are given the numbers

$$\rho_i = \text{sphdist}(x, p_i), \quad i = 1, \dots, k.$$

We would like to determine, without any ambiguity, the exact position of x , based on this information. Find the conditions on p_1, \dots, p_k under which we can unambiguously determine x , for any $x \in S^n$, given the distances ρ_i . You can give your solution algebraically, using any of the concepts used in class (*e.g.*, nullspace, range, rank), or you can give a geometric condition (involving the vectors p_i). You must justify your answer.

- 3.17 *Some true/false questions.* Determine if the following statements are true or false. No justification or discussion is needed for your answers. What we mean by “true” is that the statement is true for all values of the matrices and vectors given. You can’t assume anything about the dimensions of the matrices (unless it’s explicitly stated), but you can assume that the dimensions are such that all expressions make sense. For example, the statement “ $A + B = B + A$ ” is true, because no matter what the dimensions of A and B (which must, however, be the same), and no matter what values A and B have, the statement holds. As another example, the statement $A^2 = A$ is false, because there are (square) matrices for which this doesn’t hold. (There are also matrices for which it does hold, *e.g.*, an identity matrix. But that doesn’t make the statement true.)

a. If all coefficients (*i.e.*, entries) of the matrix A are positive, then A is full rank.

b. If A and B are onto, then $A + B$ must be onto.

c. If A and B are onto, then so is the matrix $\begin{bmatrix} A & C \\ 0 & B \end{bmatrix}$.

d. If A and B are onto, then so is the matrix $\begin{bmatrix} A \\ B \end{bmatrix}$.

e. If the matrix $\begin{bmatrix} A \\ B \end{bmatrix}$ is onto, then so are the matrices A and B .

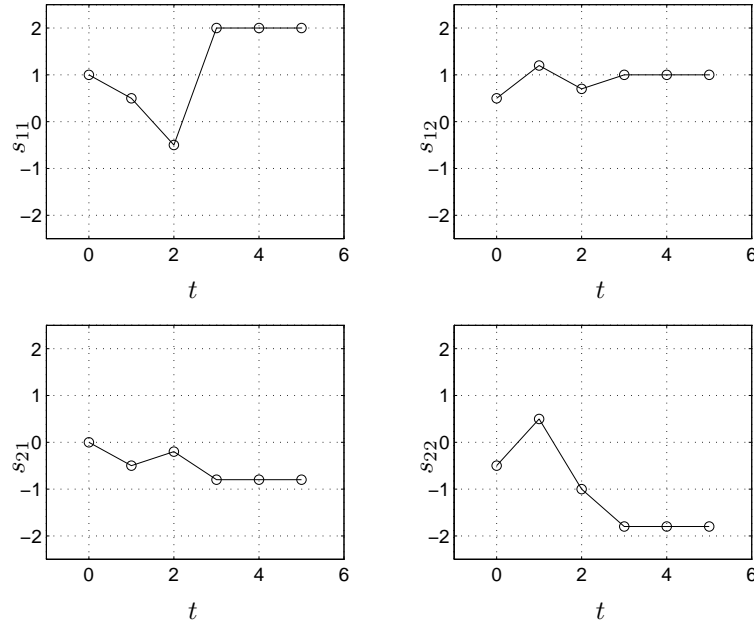
f. If A is full rank and skinny, then so is the matrix $\begin{bmatrix} A \\ B \end{bmatrix}$.

- 3.18 *Some true/false questions.* Determine if the following statements are true or false. What we mean by “true” is that the statement is true for all values of the matrices and vectors given. (You can assume the entries of the matrices and vectors are all real.) You can’t assume anything about the dimensions of the matrices (unless it’s explicitly stated), but you can assume that the dimensions are such that all expressions make sense. For example, the statement “ $A + B = B + A$ ” is true, because no matter what the dimensions of A and B (which must, however, be the same), and no matter what values A and B have, the statement holds. As another example, the statement $A^2 = A$ is false, because there are (square) matrices for which this doesn’t hold. (There are also matrices for which it does hold, *e.g.*, an identity matrix. But that doesn’t make the statement true.)

(a) If all coefficients (*i.e.*, entries) of the matrices A and B are nonnegative, and both A and B are onto, then $A + B$ is onto.

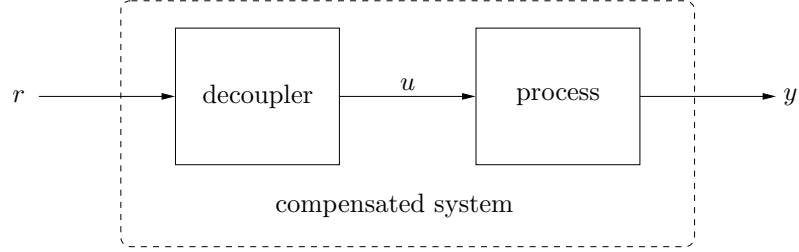
- (b) $\mathcal{N}\left(\begin{bmatrix} A \\ A+B \\ A+B+C \end{bmatrix}\right) = \mathcal{N}(A) \cap \mathcal{N}(B) \cap \mathcal{N}(C).$
- (c) $\mathcal{N}\left(\begin{bmatrix} A \\ AB \\ ABC \end{bmatrix}\right) = \mathcal{N}(A) \cap \mathcal{N}(B) \cap \mathcal{N}(C).$
- (d) $\mathcal{N}(B^T A^T A B + B^T B) = \mathcal{N}(B).$
- (e) If $\begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}$ is full rank, then so are the matrices A and B .
- (f) If $\begin{bmatrix} A & 0 \end{bmatrix}$ is onto, then A is full rank.
- (g) If A^2 is onto, then A is onto.
- (h) If $A^T A$ is onto, then A is onto.
- (i) Suppose $u_1, \dots, u_k \in \mathbf{R}^n$ are nonzero vectors such that $u_i^T u_j \geq 0$ for all i, j . Then the vectors are *nonnegative independent*, which means if $\alpha_1, \dots, \alpha_k \in \mathbf{R}$ are nonnegative scalars, and $\sum_{i=1}^k \alpha_i u_i = 0$, then $\alpha_i = 0$ for $i = 1, \dots, k$.
- (j) Suppose $A \in \mathbf{R}^{n \times k}$ and $B \in \mathbf{R}^{n \times m}$ are skinny, full rank matrices that satisfy $A^T B = 0$. Then $[A \ B]$ is skinny and full rank.

3.19 *Dynamic decoupling.* An industrial process is described by a 2-input 2-output discrete-time LDS with *finite impulse response* of length 4, which means that its impulse matrix h is nonzero only for $t = 0, 1, 2, 3$; $h(t) = 0$ for $t \geq 4$. This means that its step response matrix, defined as $s(t) = \sum_{\tau=0}^t h(\tau)$, converges to its final value by $t = 3$. If you want to think of this system in concrete terms, you can imagine it as a chemical process reactor, with u_1 a heater input, u_2 as a reactant flow rate, y_1 as the reactor temperature, and y_2 as the reactor pressure. The step response matrix of the system is shown below. The impulse response matrix of the system (for $t = 0, 1, 2, 3$) can be obtained from the class web page in `dynamic_dec.h.m`, where you will find the impulse response matrices `h0`, `h1`, `h2`, `h3` $\in \mathbf{R}^{2 \times 2}$ (which represent $h(0), h(1), h(2), h(3)$).



The plots show that u_1 has a substantial effect on y_2 , and that u_2 has a substantial effect on y_1 , neither of which we want. To eliminate them, you will explore the design of a *dynamic decoupler* for

this system, which is another 2-input, 2-output LDS with impulse matrix g . The decoupler is also FIR of length 4: $g(0), g(1), g(2), g(3) \in \mathbf{R}^{2 \times 2}$ can be nonzero, but $g(t) = 0$ for $t \geq 4$. The decoupler is used as a prefilter for the process: the input r (which is called the reference or command input) is applied as the input to the decoupler, and the output of the decoupler is u , the input to the industrial process. This is shown below.



We refer to this cascaded system as the *compensated system*. Let \tilde{s} denote the step response matrix of the compensated system, from input r to output y . The goal is to design the decoupler (*i.e.*, choose $g(0), g(1), g(2), g(3) \in \mathbf{R}^{2 \times 2}$) so that the compensated system satisfies the following specifications.

- $\lim_{t \rightarrow \infty} \tilde{s}(t) = I$. This means if the reference input is constant, the process output converges to the reference input.
- The off-diagonal entries of $\tilde{s}(t)$ are all zero (for all t). This means the compensated system is *decoupled*: r_1 has no effect on y_2 , and r_2 has no effect on y_1 .

Find such a decoupler, and plot the compensated system step response matrix. If there is no such decoupler (*i.e.*, the problem specifications are not feasible), say so, and explain why. If there are many decouplers that satisfy the given specifications, say so, and do something sensible with any extra degrees of freedom you may have.

- 3.20 *Temperatures in a multi-core processor.* We are concerned with the temperature of a processor at two critical locations. These temperatures, denoted $T = (T_1, T_2)$ (in degrees C), are affine functions of the power dissipated by three processor cores, denoted $P = (P_1, P_2, P_3)$ (in W). We make 4 measurements. In the first, all cores are idling, and dissipate 10W. In the next three measurements, one of the processors is set to full power, 100W, and the other two are idling. In each experiment we measure and note the temperatures at the two critical locations.

P_1	P_2	P_3	T_1	T_2
10W	10W	10W	27°	29°
100W	10W	10W	45°	37°
10W	100W	10W	41°	49°
10W	10W	100W	35°	55°

Suppose we operate all cores at the same power, p . How large can we make p , without T_1 or T_2 exceeding 70°?

You must fully explain your reasoning and method, in addition to providing the numerical solution.

- 3.21 *Relative deviation between vectors.* Suppose a and b are nonzero vectors of the same size. The relative deviation of b from a is defined as the distance between a and b , divided by the norm of a ,

$$\eta_{ab} = \frac{\|a - b\|}{\|a\|}.$$

This is often expressed as a percentage. The relative deviation is not a symmetric function of a and b ; in general, $\eta_{ab} \neq \eta_{ba}$.

Suppose $\eta_{ab} = 0.1$ (i.e., 10%). How big and how small can be η_{ba} be? How big and how small can $\angle(a, b)$ be? Explain your reasoning. For bounding $\angle(a, b)$, you can just draw some pictures; you don't have to give a formal argument.

- 3.22 *Single sensor failure detection and identification.* We have $y = Ax$, where $A \in \mathbf{R}^{m \times n}$ is known, and $x \in \mathbf{R}^n$ is to be found. Unfortunately, up to one sensor may have failed (but you don't know which one has failed, or even whether any has failed). You are given \tilde{y} and not y , where \tilde{y} is the same as y in all entries except, possibly, one (say, the k th entry). If all sensors are operating correctly, we have $y = \tilde{y}$. If the k th sensor fails, we have $\tilde{y}_i = y_i$ for all $i \neq k$.

The file `one_bad_sensor.m`, available on the course web site, defines A and \tilde{y} (as `A` and `ytilde`). Determine which sensor has failed (or if no sensors have failed). You must explain your method, and submit your code.

For this exercise, you can use the Matlab code `rank([F g])==rank(F)` to check if $g \in \mathcal{R}(F)$. (We will see later a much better way to check if $g \in \mathcal{R}(F)$.)

- 3.23 *Vector space multiple access (VSMA).* We consider a system of k transmitter-receiver pairs that share a common medium. The goal is for transmitter i to transmit a vector signal $x_i \in \mathbf{R}^{n_i}$ to the i th receiver, without interference from the other transmitters. All receivers have access to the same signal $y \in \mathbf{R}^m$, which includes the signals of all transmitters, according to

$$y = A_1 x_1 + \cdots + A_k x_k,$$

where $A_i \in \mathbf{R}^{m \times n_i}$. You can assume that the matrices A_i are skinny, i.e., $m \geq n_i$ for $i = 1, \dots, k$. (You can also assume that $n_i > 0$ and $A_i \neq 0$, for $i = 1, \dots, k$.) Since the k transmitters all share the same m -dimensional vector space, we call this *vector space multiple access*. Each receiver knows the received signal y , and the matrices A_1, \dots, A_k .

We say that the i th signal is *decodable* if the i th receiver can determine the value of x_i , no matter what values x_1, \dots, x_k have. Roughly speaking, this means that receiver i can process the received signal so as to perfectly recover the i th transmitted signal, while rejecting any interference from the other signals $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k$. Whether or not the i th signal is decodable depends, of course, on the matrices A_1, \dots, A_k .

Here are four statements about decodability:

- (a) Each of the signals x_1, \dots, x_k is decodable.
- (b) The signal x_1 is decodable.
- (c) The signals x_2, \dots, x_k are decodable, but x_1 isn't.
- (d) The signals x_2, \dots, x_k are decodable when x_1 is 0.

For each of these statements, you are to give the exact (i.e., necessary and sufficient) conditions under which the statement holds, in terms of A_1, \dots, A_k and n_1, \dots, n_k . Each answer, however, must have a very specific form: it must consist of a conjunction of one or more of the following properties:

- I. $\text{rank}(A_1) < n_1$.
- II. $\text{rank}([A_2 \cdots A_k]) = n_2 + \cdots + n_k$.
- III. $\text{rank}([A_1 \cdots A_k]) = n_1 + \text{rank}([A_2 \cdots A_k])$.
- IV. $\text{rank}([A_1 \cdots A_k]) = \text{rank}(A_1) + \text{rank}([A_2 \cdots A_k])$.

As examples, possible answers (for each statement) could be "I" or "I and II", or "I and II and IV". For some statements, there may be more than one correct answer; we will accept any correct one.

You can also give the response "My attorney has advised me not to respond to this question at this time." This response will receive partial credit.

For (just) this problem, we want only your answers. We do not want, and will not read, any further explanation or elaboration, or any other type of answers.

3.24 *Minimum distance and maximum correlation decoding.* We consider a simple communication system, in which a sender transmits one of N possible signals to a receiver, which receives a version of the signal sent that is corrupted by noise. Based on the corrupted received signal, the receiver has to estimate or guess which of the N signals was sent. We will represent the signals by vectors in \mathbf{R}^n . We will denote the possible signals as $a_1, \dots, a_N \in \mathbf{R}^n$. These signals, which collectively are called the *signal constellation*, are known to both the transmitter and receiver. When the signal a_k is sent, the received signal is $a_{\text{recd}} = a_k + v$, where $v \in \mathbf{R}^n$ is (channel or transmission) noise. In a communications course, the noise v is described by a statistical model, but here we'll just assume that it is 'small' (and in any case, it does not matter for the problem). The receiver must make a guess or estimate as to which of the signals was sent, based on the received signal a_{recd} . There are many ways to do this, but in this problem we explore two methods.

- *Minimum distance decoding.* Choose as the estimate of the decoded signal the one in the constellation that is closest to what is received, *i.e.*, choose a_k that minimizes $\|a_{\text{recd}} - a_i\|$. For example, if we have $N = 3$ and

$$\|a_{\text{recd}} - a_1\| = 2.2, \quad \|a_{\text{recd}} - a_2\| = 0.3, \quad \|a_{\text{recd}} - a_3\| = 1.1,$$

then the minimum distance decoder would guess that the signal a_2 was sent.

- *Maximum correlation decoding.* Choose as the estimate of the decoded signal the one in the constellation that has the largest inner product with the received signal, *i.e.*, choose a_k that maximizes $a_{\text{recd}}^T a_i$. For example, if we have $N = 3$ and

$$a_{\text{recd}}^T a_1 = -1.1, \quad a_{\text{recd}}^T a_2 = 0.2, \quad a_{\text{recd}}^T a_3 = 1.0,$$

then the maximum correlation decoder would guess that the signal a_3 was sent.

For both methods, let's not worry about breaking ties. You can just assume that ties never occur; one of the signals is always closest to, or has maximum inner product with, the received signal. Give some general conditions on the constellation (*i.e.*, the set of vectors a_1, \dots, a_N) under which these two decoding methods are the same. By 'same' we mean this: for any received signal a_{recd} , the decoded signal for the two methods is the same. Give the simplest condition you can. You must show how the decoding schemes always give the same answer, when your conditions hold. Also, give a specific counterexample, for which your conditions don't hold, and the methods differ. (We are *not* asking you to show that when your conditions don't hold, the two decoding schemes differ for some received signal.) You might want to check simple cases like $n = 1$ (scalar signals), $N = 2$ (only two messages in the constellation), or draw some pictures. But then again, you might not.

Lecture 4 – Orthonormal sets of vectors and QR factorization

4.1 *Bessel's inequality.* Suppose the columns of $U \in \mathbf{R}^{n \times k}$ are orthonormal. Show that $\|U^T x\| \leq \|x\|$. When do we have $\|U^T x\| = \|x\|$?

4.2 *Orthogonal matrices.*

- (a) Show that if U and V are orthogonal, then so is UV .
- (b) Show that if U is orthogonal, then so is U^{-1} .
- (c) Suppose that $U \in \mathbf{R}^{2 \times 2}$ is orthogonal. Show that U is either a rotation or a reflection. Make clear how you decide whether a given orthogonal U is a rotation or reflection.

4.3 *Projection matrices.* A matrix $P \in \mathbf{R}^{n \times n}$ is called a *projection matrix* if $P = P^T$ and $P^2 = P$.

- (a) Show that if P is a projection matrix then so is $I - P$.
- (b) Suppose that the columns of $U \in \mathbf{R}^{n \times k}$ are orthonormal. Show that UU^T is a projection matrix. (Later we will show that the converse is true: every projection matrix can be expressed as UU^T for some U with orthonormal columns.)
- (c) Suppose $A \in \mathbf{R}^{n \times k}$ is full rank, with $k \leq n$. Show that $A(A^T A)^{-1}A^T$ is a projection matrix.
- (d) If $S \subseteq \mathbf{R}^n$ and $x \in \mathbf{R}^n$, the point y in S closest to x is called the *projection of x on S* . Show that if P is a projection matrix, then $y = Px$ is the projection of x on $\mathcal{R}(P)$. (Which is why such matrices are called projection matrices ...)

4.4 *Reflection through a hyperplane.* Find the matrix $R \in \mathbf{R}^{n \times n}$ such that reflection of x through the hyperplane $\{z | a^T z = 0\}$ is given by Rx . Verify that the matrix R is orthogonal. (To reflect x through the hyperplane means the following: find the point z on the hyperplane closest to x . Starting from x , go in the direction $z - x$ through the hyperplane to a point on the opposite side, which has the same distance to z as x does.)

4.5 *Sensor integrity monitor.* A suite of m sensors yields measurement $y \in \mathbf{R}^m$ of some vector of parameters $x \in \mathbf{R}^n$. When the system is operating normally (which we hope is almost always the case) we have $y = Ax$, where $m > n$. If the system or sensors fail, or become faulty, then we no longer have the relation $y = Ax$. We can exploit the redundancy in our measurements to help us identify whether such a fault has occurred. We'll call a measurement y *consistent* if it has the form Ax for some $x \in \mathbf{R}^n$. If the system is operating normally then our measurement will, of course, be consistent. If the system becomes faulty, we hope that the resulting measurement y will become inconsistent, *i.e.*, not consistent. (If we are *really* unlucky, the system will fail in such a way that y is still consistent. Then we're out of luck.) A matrix $B \in \mathbf{R}^{k \times m}$ is called an *integrity monitor* if the following holds:

- $By = 0$ for any y which is consistent.
- $By \neq 0$ for any y which is inconsistent.

If we find such a matrix B , we can quickly check whether y is consistent; we can send an alarm if $By \neq 0$. Note that the first requirement says that every consistent y does not trip the alarm; the second requirement states that every inconsistent y does trip the alarm. Finally, the problem. Find an integrity monitor B for the matrix

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 1 & -1 & -2 \\ -2 & 1 & 3 \\ 1 & -1 & -2 \\ 1 & 1 & 0 \end{bmatrix}.$$

Your B should have the smallest k (i.e., number of rows) as possible. As usual, you have to explain what you're doing, as well as giving us your explicit matrix B . You must also verify that the matrix you choose satisfies the requirements. *Hints:*

- You might find one or more of the Matlab commands `orth`, `null`, or `qr` useful. Then again, you might not; there are many ways to find such a B .
- When checking that your B works, don't expect to have By exactly zero for a consistent y ; because of roundoff errors in computer arithmetic, it will be really, really small. That's OK.
- Be very careful typing in the matrix A . It's not just a random matrix.

4.6 *Householder reflections.* A *Householder matrix* is defined as

$$Q = I - 2uu^T,$$

where $u \in \mathbf{R}^n$ is normalized, that is, $u^T u = 1$.

- Show that Q is orthogonal.
- Show that $Qu = -u$. Show that $Qv = v$, for any v such that $u^T v = 0$. Thus, multiplication by Q gives reflection through the plane with normal vector u .
- Show that $\det Q = -1$.
- Given a vector $x \in \mathbf{R}^n$, find a unit-length vector u for which

$$Qx = (I - 2uu^T)x \in \text{span}\{e_1\},$$

where $e_1 = [1 \ 0 \ \dots \ 0]^T$. *Hint:* Try a u of the form $u = v/\|v\|$, with $v = x + \alpha e_1$ (find the appropriate α and show that such a u works ...) Compute such a u for $x = [3 \ 2 \ 4 \ 1 \ 5]^T$. Apply the corresponding Householder reflection to x : what is Qx ?

Note: Multiplication by an orthogonal matrix has very good numerical properties, in the sense that it does not accumulate much roundoff error. For this reason, Householder reflections are used as building blocks for fast, numerically sound algorithms. They are used, for example, in algorithms for the QR and SVD decompositions that progressively introduce zeros in the non-diagonal entries of a matrix.

4.7 *Finding a basis for the intersection of ranges.*

- Suppose you are given two matrices, $A \in \mathbf{R}^{n \times p}$ and $B \in \mathbf{R}^{n \times q}$. Explain how you can find a matrix $C \in \mathbf{R}^{n \times r}$, with independent columns, for which

$$\text{range}(C) = \text{range}(A) \cap \text{range}(B).$$

This means that the columns of C are a basis for $\text{range}(A) \cap \text{range}(B)$. Your solution can involve any standard methods or decompositions we have studied, e.g., compact and full SVD, QR factorization, pseudo-inverse, eigenvalue decomposition, and so on.

- Carry out the method described in part (a) for the particular matrices A and B defined in `intersect_range_data.m`. Be sure to give us your matrix C , as well as the Matlab (or other) code that generated it. Verify that $\text{range}(C) \subseteq \text{range}(A)$ and $\text{range}(C) \subseteq \text{range}(B)$, by showing that each column of C is in the range of A , and also in the range of B .

Please carefully separate your answers to part (a) (the general case) and part (b) (the specific case).

4.8 *Groups of equivalent statements.* In the list below there are 11 statements about two square matrices A and B in $\mathbf{R}^{n \times n}$.

- $\mathcal{R}(B) \subseteq \mathcal{R}(A)$.

- (b) there exists a matrix $Y \in \mathbf{R}^{n \times n}$ such that $B = YA$.
- (c) $AB = 0$.
- (d) $BA = 0$.
- (e) $\text{rank}(\begin{bmatrix} A & B \end{bmatrix}) = \text{rank}(A)$.
- (f) $\mathcal{R}(A) \perp \mathcal{N}(B^T)$.
- (g) $\text{rank}\left(\begin{bmatrix} A \\ B \end{bmatrix}\right) = \text{rank}(A)$.
- (h) $\mathcal{R}(A) \subseteq \mathcal{N}(B)$.
- (i) there exists a matrix $Z \in \mathbf{R}^{n \times n}$ such that $B = AZ$.
- (j) $\text{rank}(\begin{bmatrix} A & B \end{bmatrix}) = \text{rank}(B)$.
- (k) $\mathcal{N}(A) \subseteq \mathcal{N}(B)$.

Your job is to collect them into (the largest possible) groups of equivalent statements. Two statements are equivalent if each one implies the other. For example, the statement ‘ A is onto’ is equivalent to ‘ $\mathcal{N}(A) = \{0\}$ ’ (when A is square, which we assume here), because every square matrix that is onto has zero nullspace, and vice versa. Two statements are not equivalent if there exist (real) square matrices A and B for which one holds, but the other does not. A group of statements is equivalent if any pair of statements in the group is equivalent.

We want *just* your answer, which will consist of lists of mutually equivalent statements; we do not need any justification.

Put your answer in the following specific form. List each group of equivalent statements on a line, in (alphabetic) order. Each new line should start with the first letter not listed above. For example, you might give your answer as

a, c, d, h
b, i
e
f, g, j, k.

This means you believe that statements a, c, d, and h are equivalent; statements b and i are equivalent; and statements f, g, j, and k are equivalent. You also believe that the first group of statements is not equivalent to the second, or the third, and so on.

Lecture 5 – Least-squares

- 5.1 *Least-squares residuals.* Let the matrix A be skinny and full-rank. Let x_{ls} be the least-squares solution of $Ax = y$, and let $y_{\text{ls}} = Ax_{\text{ls}}$. Show that the residual vector $r = y - y_{\text{ls}}$ satisfies

$$\|r\|^2 = \|y\|^2 - \|y_{\text{ls}}\|^2.$$

Also, give a brief geometric interpretation of this equality (just a couple of sentences, and maybe a conceptual drawing).

- 5.2 *Complex linear algebra and least-squares.* Most of the linear algebra you have seen is unchanged when the scalars, matrices, and vectors are complex, *i.e.*, have complex entries. For example, we say a set of complex vectors $\{v_1, \dots, v_n\}$ is dependent if there exist complex scalars $\alpha_1, \dots, \alpha_n$, not all zero, such that $\alpha_1 v_1 + \dots + \alpha_n v_n = 0$. There are some slight differences when it comes to the inner product and other expressions that, in the real case, involve the transpose operator. For complex matrices (or vectors) we define the *Hermitian conjugate* as the complex conjugate of the transpose. We denote this as A^* , which is equal to $(\bar{A})^T$. Thus, the ij entry of the matrix A^* is given by (\bar{A}_{ji}) . The Hermitian conjugate of a matrix is sometimes called its *conjugate transpose* (which is a nice, explanatory name). Note that for a real matrix or vector, the Hermitian conjugate is the same as the transpose. We define the inner product of two complex vectors $u, v \in \mathbf{C}^n$ as

$$\langle u, v \rangle = u^* v,$$

which, in general, is a complex number. The norm of a complex vector is defined as

$$\|u\| = \sqrt{\langle u, u \rangle} = (|u_1|^2 + \dots + |u_n|^2)^{1/2}.$$

Note that these two expressions agree with the definitions you already know when the vectors are real. The complex least-squares problem is to find the $x \in \mathbf{C}^n$ that minimizes $\|Ax - y\|^2$, where $A \in \mathbf{C}^{m \times n}$ and $y \in \mathbf{C}^m$ are given. Assuming A is full rank and skinny, the solution is $x_{\text{ls}} = A^\dagger y$, where A^\dagger is the (complex) pseudo-inverse of A , given by

$$A^\dagger = (A^* A)^{-1} A^*.$$

(Which also reduces to the pseudo-inverse you've already seen when A is real). There are two general approaches to dealing with complex linear algebra problems. In the first, you simply generalize all the results to work for complex matrices, vectors, and scalars. Another approach is to represent complex matrices and vectors using real matrices and vectors of twice the dimensions, and then you apply what you already know about real linear algebra. We'll explore that idea in this problem. We associate with a complex vector $u \in \mathbf{C}^n$ a *real* vector $\tilde{u} \in \mathbf{R}^{2n}$, given by

$$\tilde{u} = \begin{bmatrix} \Re u \\ \Im u \end{bmatrix}.$$

We associate with a complex matrix $A \in \mathbf{C}^{m \times n}$ the *real* matrix $\tilde{A} \in \mathbf{R}^{2m \times 2n}$ given by

$$\tilde{A} = \begin{bmatrix} \Re A & -\Im A \\ \Im A & \Re A \end{bmatrix}.$$

- What is the relation between $\langle u, v \rangle$ and $\langle \tilde{u}, \tilde{v} \rangle$? Note that the first inner product involves complex vectors and the second involves real vectors.
- What is the relation between $\|u\|$ and $\|\tilde{u}\|$?
- What is the relation between Au (complex matrix-vector multiplication) and $\tilde{A}\tilde{u}$ (real matrix-vector multiplication)?
- What is the relation between \tilde{A}^T and A^* ?
- Using the results above, verify that $A^\dagger y$ solves the complex least-squares problem of minimizing $\|Ax - y\|$ (where A, x, y are complex). Express $A^\dagger y$ in terms of the real and imaginary parts of A and y . (You don't need to simplify your expression; you can leave block matrices in it.)

Lecture 6 – Least-squares applications

6.1 *Design of a robust heating element power vector.* Consider a heating problem in which x_i denotes the change in power of the i th heating element from some nominal setting (so $x_i < 0$ makes sense), and $y \in \mathbf{R}$ denotes the temperature rise at some critical point. We have n heating elements, so $x \in \mathbf{R}^n$. The process is linear, so $y = a^T x$ for some $a \in \mathbf{R}^n$. In fact, the system is operated in any of k modes, for which the linear function mapping x into y is not quite the same. We model this as follows: in mode i , we have $y = a_i^T x$, where $a_i \in \mathbf{R}^n$, $i = 1, \dots, k$. The vectors a_i are similar, but not quite the same.

- Suppose that $k > n$, *i.e.*, there are more operating modes than heating elements. Explain how to find a single heating element power vector x that results in $y \approx 1$ for each mode. Thus, this x yields about the same temperature rise for any of the operating modes. If some assumptions are necessary, state clearly what they are.
- Now suppose that $k < n$, *i.e.*, we have more heating elements than operating modes. Is it possible in this case to find an x that results in $y = 1$ for *every* operating mode? Suggest a reasonable choice for x , carefully stating any assumptions that are necessary.

6.2 *Optimal control of unit mass.* In this problem you will use Matlab to solve an optimal control problem for a force acting on a unit mass. Consider a unit mass at position $p(t)$ with velocity $\dot{p}(t)$, subjected to force $f(t)$, where $f(t) = x_i$ for $i - 1 < t \leq i$, for $i = 1, \dots, 10$.

- Assume the mass has zero initial position and velocity: $p(0) = \dot{p}(0) = 0$. Find x that minimizes

$$\int_{t=0}^{10} f(t)^2 dt$$

subject to the following specifications: $p(10) = 1$, $\dot{p}(10) = 0$, and $p(5) = 0$. Plot the optimal force f , and the resulting p and \dot{p} . Make sure the specifications are satisfied. Give a short intuitive explanation for what you see.

- Assume the mass has initial position $p(0) = 0$ and velocity $\dot{p}(0) = 1$. Our goal is to bring the mass near or to the origin at $t = 10$, at or near rest, *i.e.*, we want

$$J_1 = p(10)^2 + \dot{p}(10)^2,$$

small, while keeping

$$J_2 = \int_{t=0}^{10} f(t)^2 dt$$

small, or at least not too large. Plot the optimal trade-off curve between J_2 and J_1 . Check that the end points make sense to you. *Hint:* the parameter μ has to cover a very large range, so it usually works better in practice to give it a logarithmic spacing, using, *e.g.*, `logspace` in Matlab. You don't need more than 50 or so points on the trade-off curve.

Your solution to this problem should consist of a clear written narrative that explains what you are doing, and gives formulas symbolically; the Matlab source code you devise to find the numerical answers, along with comments explaining it all; and the final plots produced by Matlab.

6.3 *AR system identification.* In this problem you will use least-squares to develop and validate autoregressive (AR) models of a system from some input/output (I/O) records. You are given I/O records

$$u(1), \dots, u(N), \quad y(1), \dots, y(N),$$

which are the measured input and output of an unknown system. You will use least-squares to find approximate models of the form

$$y(t) = a_0 u(t) + b_1 y(t-1) + \dots + b_n y(t-n).$$

Specifically you will choose coefficients a_0, b_1, \dots, b_n that minimize

$$\sum_{t=n+1}^N (y(t) - a_0 u(t) - b_1 y(t-1) - \dots - b_n y(t-n))^2$$

where u, y are the given data record. The squareroot of this quantity is the residual norm (on the model data). Dividing by $\sqrt{\sum_{t=n+1}^N y(t)^2}$ yields the relative error. You'll plot this as a function of n for $n = 1, \dots, 35$. To validate or evaluate your models, you can try them on validation data records

$$\tilde{u}(1), \dots, \tilde{u}(N), \quad \tilde{y}(1), \dots, \tilde{y}(N).$$

To find the predictive ability of an AR model with coefficients a_0, b_1, \dots, b_n , you can form the signal

$$\hat{y}(t) = a_0 \tilde{u}(t) + b_1 \tilde{y}(t-1) + \dots + b_n \tilde{y}(t-n)$$

for $t = n+1, \dots, N$, and compare it to the actual output signal, \tilde{y} . You will plot the squareroot of the sum of squares of the difference, divided by the squareroot of the sum of squares of \tilde{y} , for $n = 1, \dots, 35$. Compare this to the plot above. Briefly discuss the results. To develop the models for different values of n , you can use inefficient code that just loops over n ; you do not have to try to use an efficient method based on one *QR* factorization. The file `I0data.m` contains the data for this problem and is available on the class web page. The `toeplitz()` command may be helpful.

6.4 *Least-squares model fitting.* In this problem you will use least-squares to fit several different types of models to a given set of input/output data. The data consists of a scalar input sequence u , and a scalar output sequence y , for $t = 1, \dots, N$. You will develop several different models that relate the signals u and y . The models range in complexity from a simple constant to a nonlinear dynamic model.

- *Memoryless models.* In a memoryless model, the output at time t , *i.e.*, $y(t)$, depends only the input at time t , *i.e.*, $u(t)$. Another common term for such a model is *static*.
 - constant model: $y(t) = a$
 - static linear: $y(t) = bu(t)$
 - static affine: $y(t) = a + bu(t)$
 - static quadratic: $y(t) = a + bu(t) + cu(t)^2$
- *Dynamic models.* In a dynamic model, $y(t)$ depends on $u(s)$ for some $s \neq t$. We consider simple dynamic models in which $y(t)$ depends on $u(t)$ and $u(t-1)$; in other words the current output depends on the current input and the previous input. Such models are said to have a *finite memory* of one sample. Another term is 2-tap system (the taps refer to taps on a delay line).
 - linear, 2-tap: $y(t) = bu(t) + du(t-1)$
 - affine, 2-tap: $y(t) = a + bu(t) + du(t-1)$
 - quadratic, 2-tap: $y(t) = a + bu(t) + cu(t)^2 + du(t-1) + eu(t-1)^2 + fu(t)u(t-1)$

Each of these models is specified by its parameters, *i.e.*, the scalars a, b, \dots, f . For each of these models, find the least-squares fit to the given data. In other words, find parameter values that minimize the sum-of-squares of the residuals. For example, for the affine 2-tap model, pick a, b , and d that minimize

$$\sum_{t=2}^N (y(t) - a - bu(t) - du(t-1))^2.$$

(Note that we start the sum at $t = 2$ so $u(t-1)$ is defined.) For each model, give the root-mean-square (RMS) residual, *i.e.*, the squareroot of the mean of the optimal residual squared. Plot the output \hat{y} predicted by your model, and plot the residual (which is $y - \hat{y}$). The data for this problem are available

in an m-file named `model_fitting_data.m`. Running `model_fitting_data` in Matlab will create the vectors u and y and the length N . Now you can plot u , y , etc. *Note:* the dataset u , y are *not* generated by any of the models above. It is generated by a nonlinear recursion, which has infinite (but rapidly fading) memory.

6.5 *Interpolation with rational functions.* In this problem we consider a function $f : \mathbf{R} \rightarrow \mathbf{R}$ of the form

$$f(x) = \frac{a_0 + a_1x + \cdots + a_mx^m}{1 + b_1x + \cdots + b_mx^m},$$

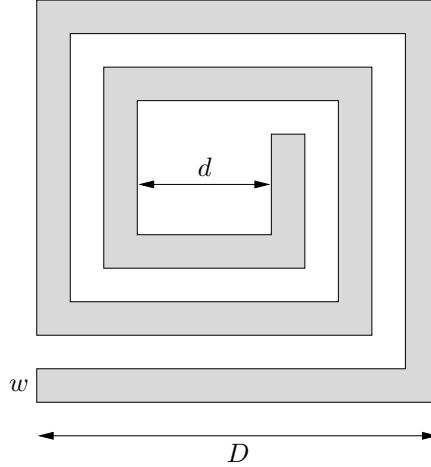
where a_0, \dots, a_m , and b_1, \dots, b_m are parameters, with either $a_m \neq 0$ or $b_m \neq 0$. Such a function is called a *rational function of degree m* . We are given data points $x_1, \dots, x_N \in \mathbf{R}$ and $y_1, \dots, y_N \in \mathbf{R}$, where $y_i = f(x_i)$. The problem is to find a rational function of smallest degree that is consistent with this data. In other words, you are to find m , which should be as small as possible, and $a_0, \dots, a_m, b_1, \dots, b_m$, which satisfy $f(x_i) = y_i$. Explain how you will solve this problem, and then carry out your method on the problem data given in `ri_data.m`. (This contains two vectors, `x` and `y`, that give the values x_1, \dots, x_N , and y_1, \dots, y_N , respectively.) Give the value of m you find, and the coefficients $a_0, \dots, a_m, b_1, \dots, b_m$. Please show us your verification that $y_i = f(x_i)$ holds (possibly with some small numerical errors).

6.6 *The middle inverse.* In this problem we consider the matrix equation

$$AXB = I,$$

where $A \in \mathbf{R}^{n \times p}$, $B \in \mathbf{R}^{q \times n}$, and $X \in \mathbf{R}^{p \times q}$. The matrices A and B are given, and the goal is to find a matrix X that satisfies the equation, or to determine that no such matrix exists. (When such an X exists, we call it a *middle inverse* of the pair A, B . It generalizes the notions of left-inverse and right-inverse: When $A = I$, X is a left-inverse of B , and when $B = I$, X is a right-inverse of A .) You will solve a specific instance of this problem, with data (*i.e.*, the matrices A and B) given in the mfile `axb_data.m`. If you think there is no X that satisfies $AXB = I$, explain why this is the case. Your explanation should be as concrete as possible. If you succeed in finding an X that satisfies $AXB = I$, please give it. You must explain how you found it, and you must submit the code that you used to generate your solution. You must also submit the Matlab code and output showing that you checked that $AXB = I$ holds (up to very small numerical errors). You can do this by typing `norm(A*X*B-eye(n))` in Matlab, and submitting a printout of what Matlab prints out. (We haven't yet studied the matrix norm, but it doesn't matter. Like the norm of a vector, it measures size, and here you are using it only to check that $AXB - I$ is small.) *The following interpretation is not needed to solve the problem.* We give it just to mention a concrete situation where a problem like this one might arise. One situation where this problem comes up is a nonstandard filtering or equalization problem. A vector signal $x \in \mathbf{R}^n$ is first processed by one channel, represented by B . At this point the signal is available for some filtering or processing by us, which is represented by the matrix X . After this processing, it is acted on by another channel, given by A . Our goal is to do the intermediate processing in such a way that it undoes the effect of the first and last channels.

6.7 *Approximate inductance formula.* The figure below shows a *planar spiral inductor*, implemented in CMOS, for use in RF circuits.



The inductor is characterized by four key parameters:

- n , the number of turns (which is a multiple of 1/4, but that needn't concern us)
- w , the width of the wire
- d , the inner diameter
- D , the outer diameter

The inductance L of such an inductor is a complicated function of the parameters n , w , d , and D . The inductance L can be found by solving Maxwell's equations, which takes considerable computer time, or by fabricating the inductor and measuring the inductance. In this problem you will develop a simple approximate inductance model of the form

$$\hat{L} = \alpha n^{\beta_1} w^{\beta_2} d^{\beta_3} D^{\beta_4},$$

where $\alpha, \beta_1, \beta_2, \beta_3, \beta_4 \in \mathbf{R}$ are constants that characterize the approximate model. (since L is positive, we have $\alpha > 0$, but the constants β_2, \dots, β_4 can be negative.) This simple approximate model, if accurate enough, can be used for design of planar spiral inductors. The file `inductor_data.m` on the course web site contains data for 50 inductors. (The data is real, not that it would affect how you solve the problem ...) For inductor i , we give parameters n_i , w_i , d_i , and D_i (all in μm), and also, the inductance L_i (in nH) obtained from measurements. (The data are organized as vectors of length 50. Thus, for example, w_{13} gives the wire width of inductor 13.) Your task, *i.e.*, the problem, is to find $\alpha, \beta_1, \dots, \beta_4$ so that

$$\hat{L}_i = \alpha n_i^{\beta_1} w_i^{\beta_2} d_i^{\beta_3} D_i^{\beta_4} \approx L_i \quad \text{for } i = 1, \dots, 50.$$

Your solution must include a clear description of how you found your parameters, as well as their actual numerical values. Note that we have not specified the criterion that you use to judge the approximate model (*i.e.*, the fit between \hat{L}_i and L_i); we leave that to your engineering judgment. But be sure to tell us what criterion you use. We define the *percentage error* between \hat{L}_i and L_i as

$$e_i = 100|\hat{L}_i - L_i|/L_i.$$

Find the average percentage error for your model, *i.e.*, $(e_1 + \dots + e_{50})/50$. (We are only asking you to find the average percentage error for your model; we do not require that your model minimize the average percentage error.) *Hint*: you might find it easier to work with $\log L$.

6.8 *Quadratic extrapolation of a time series, using least-squares fit.* This is an extension of problem 12, which concerns quadratic extrapolation of a time series. We are given a series z up to time n . Using a quadratic model, we want to extrapolate, or predict, $z(n+1)$ based on a least-squares fit to the

previous ten elements of the series: $z(n), z(n-1), \dots, z(n-9)$. We'll denote the predicted value of $z(n+1)$ by y . Another way to describe this problem is: find a quadratic function $f(t) = a_2 t^2 + a_1 t + a_0$ for which $f(n) \approx z(n)$, $f(n-1) \approx z(n-1)$, \dots , $f(n-9) \approx z(n-9)$. Specifically, the approximation should be such that

$$\sum_{k=-9}^0 (z(n+k) - f(n+k))^2$$

is minimized. The extrapolated value is then given by $y = f(n+1)$. Let the vector x denote the ten previous elements of the series,

$$x = \begin{bmatrix} z(n) \\ z(n-1) \\ \vdots \\ z(n-9) \end{bmatrix}.$$

Find a vector $c \in \mathbf{R}^{10}$ such that the extrapolated value is given by the linear transformation $y = c^T x$.

6.9 *Image reconstruction from line integrals.* In this problem we explore a simple version of a tomography problem. We consider a square region, which we divide into an $n \times n$ array of square pixels, as shown in Figure 1. The pixels are indexed column first, by a single index i ranging from 1 to n^2 , as shown

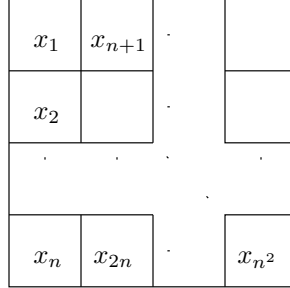


Figure 1: The square region and its division into pixels

in Figure 1. We are interested in some physical property such as density (say) which varies over the region. To simplify things, we'll assume that the density is constant inside each pixel, and we denote by x_i the density in pixel i , $i = 1, \dots, n^2$. Thus, $x \in \mathbf{R}^{n^2}$ is a vector that describes the density across the rectangular array of pixels. The problem is to estimate the vector of densities x , from a set of sensor measurements that we now describe. Each sensor measurement is a *line integral* of the density over a line L . In addition, each measurement is corrupted by a (small) noise term. In other words, the sensor measurement for line L is given by

$$\sum_{i=1}^{n^2} l_i x_i + v,$$

where l_i is the length of the intersection of line L with pixel i (or zero if they don't intersect), and v is a (small) measurement noise. Figure 2 gives an example, with a graphical explanation of l_i . Now suppose we have N line integral measurements, associated with lines L_1, \dots, L_N . From these measurements, we want to estimate the vector of densities x . The lines are characterized by the intersection lengths

$$l_{ij}, \quad i = 1, \dots, n^2, \quad j = 1, \dots, N,$$

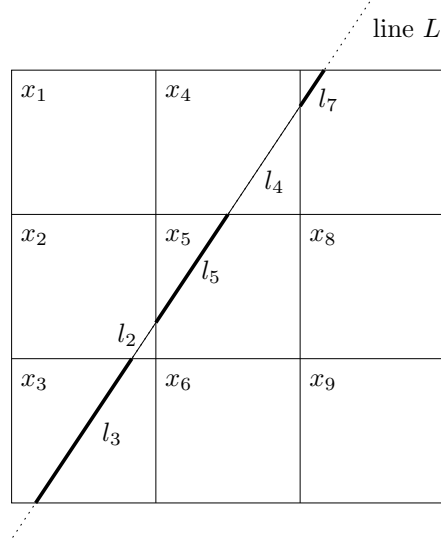


Figure 2: An example of a 3-by-3 pixel patch, with a line L and its intersections l_i with the pixels. Clearly, $l_1 = l_6 = l_8 = l_9 = 0$.

where l_{ij} gives the length of the intersection of line L_j with pixel i . Then, the whole set of measurements forms a vector $y \in \mathbf{R}^N$ whose elements are given by

$$y_j = \sum_{i=1}^{n^2} l_{ij} x_i + v_j, \quad j = 1, \dots, N.$$

And now the problem: you will reconstruct the pixel densities x from the line integral measurements y . The class webpage contains the M-file `tomodata.m`, which you should download and run in Matlab. It creates the following variables:

- `N`, the number of measurements (N),
- `n_pixels`, the side length in pixels of the square region (n),
- `y`, a vector with the line integrals y_j , $j = 1, \dots, N$,
- `lines_d`, a vector containing the displacement d_j , $j = 1, \dots, N$, (distance from the center of the region in pixels lengths) of each line, and
- `lines_theta`, a vector containing the angles θ_j , $j = 1, \dots, N$, of each line.

The file `tmeasure.m`, on the same webpage, shows how the measurements were computed, in case you're curious. You should take a look, but you don't need to understand it to solve the problem. We also provide the function `line_pixel_length.m` on the webpage, which you do need to use in order to solve the problem. This function computes the pixel intersection lengths for a given line. That is, given d_j and θ_j (and the side length n), `line_pixel_length.m` returns a $n \times n$ matrix, whose i, j th element corresponds to the intersection length for pixel i, j on the image. Use this information to find x , and display it as an image (of n by n pixels). You'll know you have it right when the image of x forms a familiar acronym... *Matlab hints:* Here are a few functions that you'll find useful to display an image:

- `A=reshape(v,n,m)`, converts the vector v (which must have $n*m$ elements) into an $n \times m$ matrix (the first column of A is the first n elements of v , etc.),

- `imagesc(A)`, displays the matrix **A** as an image, scaled so that its lowest value is black and its highest value is white,
- `colormap gray`, changes Matlab's image display mode to grayscale (you'll want to do this to view the pixel patch),
- `axis image`, redefines the axes of a plot to make the pixels square.

Note: While irrelevant to your solution, this is actually a simple version of *tomography*, best known for its application in medical imaging as the CAT scan. If an *x-ray* gets attenuated at rate x_i in pixel i (a little piece of a cross-section of your body), the j -th measurement is

$$z_j = \prod_{i=1}^{n^2} e^{-x_i l_{ij}},$$

with the l_{ij} as before. Now define $y_j = -\log z_j$, and we get

$$y_j = \sum_{i=1}^{n^2} x_i l_{ij}.$$

6.10 *Least-squares model fitting.* In this problem you will use least-squares to fit several different types of models to a given set of input/output data. The data consist of a scalar input sequence u , and a scalar output sequence y , for $t = 1, \dots, N$. You will develop several different models that relate the signals u and y .

- *Memoryless models.* In a memoryless model, the output at time t , *i.e.*, $y(t)$, depends only the input at time t , *i.e.*, $u(t)$. Another common term for such a model is *static*.

constant model:	$y(t) = c_0$
static linear:	$y(t) = c_1 u(t)$
static affine:	$y(t) = c_0 + c_1 u(t)$
static quadratic:	$y(t) = c_0 + c_1 u(t) + c_2 u(t)^2$

- *Dynamic models.* In a dynamic model, $y(t)$ depends on $u(s)$ for some $s \neq t$. We consider some simple time-series models (see problem 2 in the reader), which are linear dynamic models.

moving average (MA):	$y(t) = a_0 u(t) + a_1 u(t-1) + a_2 u(t-2)$
autoregressive (AR):	$y(t) = a_0 u(t) + b_1 y(t-1) + b_2 y(t-2)$
autoregressive moving average (ARMA):	$y(t) = a_0 u(t) + a_1 u(t-1) + b_1 y(t-1)$

Note that in the AR and ARMA models, $y(t)$ depends indirectly on all previous inputs, $u(s)$ for $s < t$, due to the recursive dependence on $y(t-1)$. For this reason, the AR and ARMA models are said to have *infinite memory*. The MA model, on the other hand, has a *finite memory*: $y(t)$ depends only on the current and two previous inputs. (Another term for this MA model is 3-tap system, where taps refer to taps on a delay line.)

Each of these models is specified by its parameters, *i.e.*, the scalars c_i, a_i, b_i . For each of these models, find the least-squares fit to the given data. In other words, find parameter values that minimize the sum-of-squares of the residuals. For example, for the ARMA model, pick a_0, a_1 , and b_1 that minimize

$$\sum_{t=2}^N (y(t) - a_0 u(t) - a_1 u(t-1) - b_1 y(t-1))^2.$$

(Note that we start the sum at $t = 2$ which ensures that $u(t-1)$ and $y(t-1)$ are defined.) For each model, give the root-mean-square (RMS) residual, *i.e.*, the squareroot of the mean of the optimal

residual squared. Plot the output \hat{y} predicted by your model, and plot the residual (which is $y - \hat{y}$). The data for this problem are available from the class web page, by following the link to Matlab datasets, in an m-file named `uy_data.m`. Copy this file to your working directory and type `uy_data` from within Matlab. This will create the vectors u and y and the scalar N (the length of the vectors). Now you can plot u , y , etc. *Note:* the dataset u , y is *not* generated by any of the models above. It is generated by a nonlinear recursion, which has infinite memory.

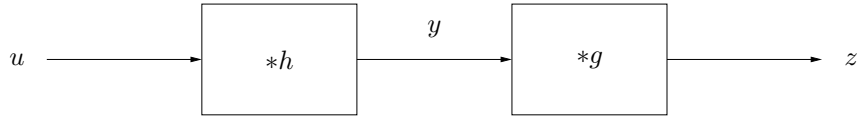
- 6.11 *Least-squares deconvolution.* A communications channel is modeled by a finite-impulse-response (FIR) filter:

$$y(t) = \sum_{\tau=0}^{n-1} u(t - \tau)h(\tau),$$

where $u : \mathbf{Z} \rightarrow \mathbf{R}$ is the channel input sequence, $y : \mathbf{Z} \rightarrow \mathbf{R}$ is the channel output, and $h(0), \dots, h(n-1)$ is the impulse response of the channel. In terms of discrete-time convolution we write this as $y = h * u$. You will design a *deconvolution filter* or *equalizer* which also has FIR form:

$$z(t) = \sum_{\tau=0}^{m-1} y(t - \tau)g(\tau),$$

where $z : \mathbf{Z} \rightarrow \mathbf{R}$ is the filter output, y is the channel output, and $g(0), \dots, g(m-1)$ is the impulse response of the filter, which we are to design. This is shown in the block diagram below.



The goal is to choose $g = (g(0), \dots, g(m-1))$ so that the filter output is approximately the channel input delayed by D samples, *i.e.*, $z(t) \approx u(t - D)$. Since $z = g * h * u$ (discrete-time convolution), this means that we'd like

$$(g * h)(t) \approx \begin{cases} 0 & t \neq D, \\ 1 & t = D \end{cases}$$

We will refer to $g * h$ as the *equalized impulse response*; the goal is to make it as close as possible to a D -sample delay. Specifically, we want the *least-squares* equalizer is g that minimizes the sum-of-squares error

$$\sum_{t \neq D} (g * h)(t)^2,$$

subject to the constraint

$$(g * h)(D) = 1.$$

To solve the problem below you'll need to get the file `deconv_data.m` from the class web page in the *Matlab files* section. It will define the channel impulse response h as a Matlab vector `h`. (Indices in Matlab run from 1 to n , while the argument of the channel impulse response runs from $t = 0$ to $t = n - 1$, so `h(3)` in Matlab corresponds to $h(2)$.)

- Find the least-squares equalizer g , of length $m = 20$, with delay $D = 12$. Plot the impulse responses of the channel (h) and the equalizer (g). Plot the equalized impulse response ($g * h$).
- The vector y (also defined in `deconv_data.m`) contains the channel output corresponding to a signal u passed through the channel (*i.e.*, $y = h * u$). The signal u is binary, *i.e.*, $u(t) \in \{-1, 1\}$, and starts at $t = 0$ (*i.e.*, $u(t) = 0$ for $t < 0$). Pass y through the least-squares equalizer found in part a, to form the signal z . Give a histogram plot of the amplitude distribution of both y and z . (You can remove the first and last D samples of z before making the histogram plot.) Comment on what you find.

Matlab hints: The command `conv` convolves two vectors; the command `hist` plots a histogram (of the amplitude distribution).

6.12 *Estimation with sensor offset and drift.* We consider the usual estimation setup:

$$y_i = a_i^T x + v_i, \quad i = 1, \dots, m,$$

where

- y_i is the i th (scalar) measurement
- $x \in \mathbf{R}^n$ is the vector of parameters we wish to estimate from the measurements
- v_i is the sensor or measurement error of the i th measurement

In this problem we assume the measurements y_i are taken at times evenly spaced, T seconds apart, starting at time $t = T$. Thus, y_i , the i th measurement, is taken at time $t = iT$. (This isn't really material; it just makes the interpretation simpler.) You can assume that $m \geq n$ and the measurement matrix

$$A = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_m^T \end{bmatrix}$$

is full rank (*i.e.*, has rank n). Usually we assume (often implicitly) that the measurement errors v_i are random, unpredictable, small, and centered around zero. (You don't need to worry about how to make this idea precise.) In such cases, least-squares estimation of x works well. In some cases, however, the measurement error includes some *predictable* terms. For example, each sensor measurement might include a (common) *offset* or *bias*, as well as a term that grows linearly with time (called a *drift*). We model this situation as

$$v_i = \alpha + \beta iT + w_i$$

where α is the sensor bias (which is unknown but the *same* for all sensor measurements), β is the drift term (again the same for all measurements), and w_i is part of the sensor error that is unpredictable, small, and centered around 0. If we knew the offset α and the drift term β we could just subtract the predictable part of the sensor signal, *i.e.*, $\alpha + \beta iT$ from the sensor signal. But we're interested in the case where we don't know the offset α or the drift coefficient β . Show how to use least-squares to *simultaneously* estimate the parameter vector $x \in \mathbf{R}^n$, the offset $\alpha \in \mathbf{R}$, and the drift coefficient $\beta \in \mathbf{R}$. Clearly explain your method. If your method always works, say so. Otherwise describe the conditions (on the matrix A) that must hold for your method to work, and give a simple example where the conditions don't hold.

6.13 *Estimating emissions from spot measurements.* There are n sources of a pollutant, at known locations $s_1, \dots, s_n \in \mathbf{R}^2$. Each source emits the pollutant at some emission rate; we let x_j denote the emission rate for source j . (These are positive, but to simplify things we won't concern ourselves with that.) The emission rates are to be determined, or estimated. We measure the total pollutant level at m spots, located at $t_1, \dots, t_m \in \mathbf{R}^2$, which are known. The total pollutant measured at spot i is the sum of the contributions from the n sources. The contribution from source j to measurement i is given by $\alpha x_j / \|s_j - t_i\|^2$, where α is a known (positive) constant. In other words, the pollutant concentration from a source follows an inverse square law, and is proportional to the emission rate. We assume that measurement spots do not coincide with the source locations, *i.e.*, we do not have $s_j = t_i$ for any i or j . We also assume that none of the spot locations is repeated (*i.e.*, we have $t_i \neq t_j$ for $i \neq j$) and that none of the source locations is repeated (*i.e.*, we have $s_i \neq s_j$ for $i \neq j$).

- (a) Give a *specific* example of source and spot measurement locations, with 4 sensors and 3 sources, for which it is *impossible* to find the emission rates given the spot measurements. In this part,

we ignore the issue of noise or sensor errors; we assume the spot measurements are exactly as described above. To show that your configuration is a valid example, give two *specific* different sets of emission rates that yield identical spot measurements. You are free to (briefly) explain your example using concepts such as range, nullspace, rank, and so on; but remember, we want a *specific numerical example*, such as $s_1 = [0 \ 1]^T, \dots, s_3 = [1 \ 2]^T, t_1 = [1 \ 1]^T, \dots, t_4 = [3 \ 2]^T$. (And similarly for the two emission rates that give the same spot measurements.)

- (b) Get the data from the file `emissions_data.m` that is available on the class web site. This file defines three source locations (given as a 2×3 matrix; the columns give the locations), and ten spot measurement locations (given as a 2×10 matrix). It also gives two sets of spot measurements: one for part (b), and one for part (c). Be careful to use the right set of measurements for each problem! The spot measurements are not perfect (as we assumed in part (a)); they contain small noise and errors. Estimate the pollutant emission rates. Explain your method, and give your estimate for the emissions rates of the three sources.
- (c) Now we suppose that *one* of the spot measurements is faulty, *i.e.*, its associated noise or error is far larger than the errors of the other spot measurements. Explain how you would identify or guess which one is malfunctioning, and then estimate the source emission rates. Carry out your method on the data given in the matlab file. Be sure to tell us which spot measurement you believe to be faulty, and what your guess of the emission rates is. (The emission rates are *not* the same as in part (b), but the source and spot measurement locations are.)

6.14 *Identifying a system from input/output data.* We consider the standard setup:

$$y = Ax + v,$$

where $A \in \mathbf{R}^{m \times n}$, $x \in \mathbf{R}^n$ is the input vector, $y \in \mathbf{R}^m$ is the output vector, and $v \in \mathbf{R}^m$ is the noise or disturbance. We consider here the problem of estimating the matrix A , given some input/output data. Specifically, we are given the following:

$$x^{(1)}, \dots, x^{(N)} \in \mathbf{R}^n, \quad y^{(1)}, \dots, y^{(N)} \in \mathbf{R}^m.$$

These represent N samples or observations of the input and output, respectively, possibly corrupted by noise. In other words, we have

$$y^{(k)} = Ax^{(k)} + v^{(k)}, \quad k = 1, \dots, N,$$

where $v^{(k)}$ are assumed to be small. The problem is to estimate the (coefficients of the) matrix A , based on the given input/output data. You will use a least-squares criterion to form an estimate \hat{A} of A . Specifically, you will choose as your estimate \hat{A} the matrix that minimizes the quantity

$$J = \sum_{k=1}^N \|Ax^{(k)} - y^{(k)}\|^2$$

over A .

- (a) Explain how to do this. If you need to make an assumption about the input/output data to make your method work, state it clearly. You may want to use the matrices $X \in \mathbf{R}^{n \times N}$ and $Y \in \mathbf{R}^{m \times N}$ given by

$$X = \begin{bmatrix} x^{(1)} & \dots & x^{(N)} \end{bmatrix}, \quad Y = \begin{bmatrix} y^{(1)} & \dots & y^{(N)} \end{bmatrix}$$

in your solution.

- (b) On the course web site you will find some input/output data for an instance of this problem in the mfile `sysid_data.m`. Executing this mfile will assign values to m , n , and N , and create two matrices that contain the input and output data, respectively. The $n \times N$ matrix variable X contains the input data $x^{(1)}, \dots, x^{(N)}$ (*i.e.*, the first column of X contains $x^{(1)}$, etc.). Similarly, the $m \times N$ matrix Y contains the output data $y^{(1)}, \dots, y^{(N)}$. You must give your final estimate \hat{A} , your source code, and also give an explanation of what you did.

- 6.15 *Robust least-squares estimation methods.* We consider a standard measurement setup, with $y = Ax + v$, where $x \in \mathbf{R}^n$ is a vector we'd like to estimate, $y \in \mathbf{R}^m$ is the vector of measurements, $v \in \mathbf{R}^m$ is the vector of measurement errors, and $A \in \mathbf{R}^{m \times n}$. We assume that $m > n$, *i.e.*, there are more measurements than parameters to be estimated. The measurement error v is not known, but is assumed to be small. The goal is to estimate x , given y . Here is the twist: we do not know the matrix A exactly. In fact we calibrated our sensor system on $k > 1$ different days, and found the values

$$A^{(1)}, \dots, A^{(k)}$$

for the matrix A , on the different days. These matrices are close to each other, but not exactly the same. There is no pattern to the (small) variations between the matrices; for example, there is no discernable drift; the variations appear to be small and random. You can assume that all of the matrices are full rank, *i.e.*, have rank n . Now suppose we have a measurement y taken on a day when we did not calibrate the sensor system. We want to form an estimate \hat{x} , based on this measurement. We don't know A exactly, but we can assume that it is close to the known values $A^{(1)}, \dots, A^{(k)}$ found during calibration. A method for guessing x in this situation is called a *robust estimation method*, since it attempts to take into account the uncertainty in the matrix A . Three very reasonable proposals for robust estimation are described below.

- *The average then estimate method.* First, we form the average of the calibration values,

$$A_{\text{avg}} = \frac{1}{k} \sum_{j=1}^k A^{(j)},$$

which is supposed to represent the most typical value of A . We then choose our estimate \hat{x} to minimize the least squares residual using A_{avg} , *i.e.*, to minimize $\|A_{\text{avg}}\hat{x} - y\|$. We refer to this value of \hat{x} as \hat{x}_{ae} , where the subscript stands for 'average (then) estimate'. (You can assume that A_{avg} is full rank.)

- *The estimate then average method.* First, we find the least-squares estimate of x for each of the calibration values, *i.e.*, we find $\hat{x}^{(j)}$ that minimizes $\|A^{(j)}\hat{x} - y\|$ over \hat{x} , for $j = 1, \dots, k$. Since the matrices $A^{(j)}$ are close but not equal, we find that the estimates $\hat{x}^{(j)}$ are also close but not equal. We find our final estimate of x as the average of these estimates:

$$\hat{x}_{\text{ea}} = \frac{1}{k} \sum_{j=1}^k \hat{x}^{(j)}.$$

(Here the subscript 'ea' stands for 'estimate (then) average'.)

- *Minimum RMS residuals method.* If we make the guess \hat{x} , then the residual, using the j th calibrated value of A , is given by $r^{(j)} = A^{(j)}\hat{x} - y$. The RMS value of the collection of residuals is given by

$$\left(\frac{1}{k} \sum_{j=1}^k \|r^{(j)}\|^2 \right)^{1/2}.$$

In the minimum RMS residual method, we choose \hat{x} to minimize this quantity. We denote this estimate of x as \hat{x}_{rms} .

Here is the problem:

- For each of these three methods, say whether the estimate \hat{x} is a linear function of y . If it is a linear function, give a formula for the matrix that gives \hat{x} in terms of y . For example, if you believe that \hat{x}_{ea} is a linear function of y , then you should give a formula for B_{ea} (in terms of $A^{(1)}, \dots, A^{(k)}$), where $\hat{x}_{\text{ea}} = B_{\text{ea}}y$.

- (b) Are the three methods described above different? If any two are the same (for all possible values of the data $A^{(1)}, \dots, A^{(k)}$ and y), explain why. If they are different, give a specific example in which the estimates differ.

6.16 *Estimating a signal with interference.* This problem concerns three proposed methods for estimating a signal, based on a measurement that is corrupted by a small noise and also by an interference, that need not be small. We have

$$y = Ax + Bv + w,$$

where $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{m \times p}$ are known. Here $y \in \mathbf{R}^m$ is the measurement (which is known), $x \in \mathbf{R}^n$ is the signal that we want to estimate, $v \in \mathbf{R}^p$ is the interference, and w is a noise. The noise is unknown, and can be assumed to be small. The interference is unknown, but cannot be assumed to be small. You can assume that the matrices A and B are skinny and full rank (*i.e.*, $m > n$, $m > p$), and that the ranges of A and B intersect only at 0. (If this last condition does not hold, then there is no hope of finding x , even when $w = 0$, since a nonzero interference can masquerade as a signal.) Each of the EE263 TAs proposes a method for estimating x . These methods, along with some informal justification from their proposers, are given below. Nikola proposes the **ignore and estimate method**. He describes it as follows:

We don't know the interference, so we might as well treat it as noise, and just ignore it during the estimation process. We can use the usual least-squares method, for the model $y = Ax + z$ (with z a noise) to estimate x . (Here we have $z = Bv + w$, but that doesn't matter.)

Almir proposes the **estimate and ignore method**. He describes it as follows:

We should simultaneously estimate both the signal x and the interference v , based on y , using a standard least-squares method to estimate $[x^T \ v^T]^T$ given y . Once we've estimated x and v , we simply ignore our estimate of v , and use our estimate of x .

Miki proposes the **estimate and cancel method**. He describes it as follows:

Almir's method makes sense to me, but I can improve it. We should simultaneously estimate both the signal x and the interference v , based on y , using a standard least-squares method, exactly as in Almir's method. In Almir's method, we then throw away \hat{v} , our estimate of the interference, but I think we should use it. We can form the "pseudo-measurement" $\tilde{y} = y - B\hat{v}$, which is our measurement, with the effect of the estimated interference subtracted off. Then, we use standard least-squares to estimate x from \tilde{y} , from the simple model $\tilde{y} = Ax + z$. (This is exactly as in Nikola's method, but here we have subtracted off or cancelled the effect of the estimated interference.)

These descriptions are a little vague; part of the problem is to translate their descriptions into more precise algorithms.

- Give an explicit formula for each of the three estimates. (That is, for each method give a formula for the estimate \hat{x} in terms of A , B , y , and the dimensions n, m, p .)
- Are the methods really different? Identify any pairs of the methods that coincide (*i.e.*, always give exactly the same results). If they are all three the same, or all three different, say so. Justify your answer. To show two methods are the same, show that the formulas given in part (a) are equal (even if they don't appear to be at first). To show two methods are different, give a specific numerical example in which the estimates differ.
- Which method or methods do you think work best? Give a very brief explanation. (If your answer to part (b) is "The methods are all the same" then you can simply repeat here, "The methods are all the same".)

- 6.17 *Vector time-series modeling.* This problem concerns a vector time-series, $y(1), \dots, y(T) \in \mathbf{R}^n$. The n components of $y(t)$ might represent measurements of different quantities, or prices of different assets, at time period t . Our goal is to develop a model that allows us to *predict* the next element in the time series, *i.e.*, to predict $y(t+1)$, given $y(1), \dots, y(t)$. A consultant proposes the following model for the time-series:

$$y(t) = Ay(t-1) + v(t), \quad t = 2, \dots, T,$$

where the matrix $A \in \mathbf{R}^{n \times n}$ is the parameter of the model, and $v(t) \in \mathbf{R}^n$ is a signal that is small and unpredictable. (We keep the definition of the terms ‘small’ and ‘unpredictable’ vague, since the exact meaning won’t matter.) This type of model has several names. It is called an VAR(1) model, which is short for *vector auto-regressive*, with one time lag. It is also called a Gauss-Markov model, which is a fancy name for a linear system driven by a noise. Once we have a model of this form, we can predict the next time-series sample using the formula

$$\hat{y}(t+1) = Ay(t), \quad t = 1, \dots, T.$$

The idea here is that $v(t)$ is unpredictable, so we simply replace it with zero when we estimate the next time-series sample. The prediction error is given by

$$e(t) = \hat{y}(t) - y(t), \quad t = 2, \dots, T.$$

The prediction error depends on the time-series data, and also A , the parameter in our model. There is one more twist. It is known that $y_1(t+1)$, the first component of the next time-series sample, does not depend on $y_2(t), \dots, y_n(t)$. The second component, $y_2(t+1)$, does not depend on $y_3(t), \dots, y_n(t)$. In general, the i th component, $y_i(t+1)$, does not depend on $y_{i+1}(t), \dots, y_n(t)$. Roughly speaking, this means that the current time-series component $y_i(t)$ only affects the next time-series components $y_1(t+1), \dots, y_i(t+1)$. This means that the matrix A is lower triangular, *i.e.*, $A_{ij} = 0$ for $i < j$. To find the parameter A that defines our model, you will use a least-squares criterion. You will pick A that minimizes the mean-square prediction error,

$$\frac{1}{T-1} \sum_{t=2}^T \|e(t)\|^2,$$

over all lower-triangular matrices. Carry out this method, using the data found in the `vts_data.m`, which contains an $n \times T$ matrix Y , whose columns are the vector time-series samples at $t = 1, \dots, T$. Explain your method, and submit the code that you use to solve the problem. Give your final estimated model parameter A , and the resulting mean-square error. Compare your mean-square prediction error to the mean-square value of y , *i.e.*,

$$\frac{1}{T} \sum_{t=1}^T \|y(t)\|^2.$$

Finally, predict what you think $y(T+1)$ is, based on the data given.

- 6.18 *Fitting a rational transfer function to frequency response data.* This problem concerns a rational function $H : \mathbf{C} \rightarrow \mathbf{C}$ of the form

$$H(s) = \frac{A(s)}{B(s)},$$

where A and B are the polynomials

$$A(s) = a_0 + a_1 s + \dots + a_m s^m, \quad B(s) = 1 + b_1 s + \dots + b_m s^m.$$

Here $a_0, \dots, a_m \in \mathbf{R}$ and $b_1, \dots, b_m \in \mathbf{R}$ are *real* parameters, and $s \in \mathbf{C}$ is the *complex* independent variable. We define $a = (a_0, \dots, a_m) \in \mathbf{R}^{m+1}$ and $b = (b_1, \dots, b_m) \in \mathbf{R}^m$, *i.e.*, a and b are vectors

containing the coefficients of A and B (not including the constant coefficient of B , which is fixed at one). We are given noisy measurements of H at some points on the imaginary axis, *i.e.*, some data

$$s_1 = j\omega_1, \dots, s_N = j\omega_N \in \mathbf{C}, \quad h_1, \dots, h_N \in \mathbf{C},$$

and hope to choose a and b so that we have $H(s_i) \approx h_i$. Here $\omega_1, \dots, \omega_N$ are real and nonnegative, and h_1, \dots, h_N are complex. To judge the quality of fit we use the mean-square error,

$$J = \frac{1}{N} \sum_{i=1}^N |H(s_i) - h_i|^2.$$

Interpretation. (Not needed to solve the problem.) You can think of H as a rational transfer function, with s the complex frequency variable. The data is a set of frequency response measurements, with some measurement errors. The goal is to find a rational transfer function that fits the measured frequency response. This problem explores a famous heuristic method, based on solving a sequence of (linear) least-squares problems, for finding coefficients a , b that approximately minimize J . We start by expressing J in the following (strange) way:

$$J = \frac{1}{N} \sum_{i=1}^N \left| \frac{A(s_i) - h_i B(s_i)}{z_i} \right|^2, \quad z_i = B(s_i), \quad i = 1, \dots, N.$$

The method works by choosing a and b that minimize the lefthand expression (with z_i fixed), then updating the numbers z_i using the righthand formula, and then repeating. More precisely, let k denote the iteration number, with $a^{(k)}$, $b^{(k)}$, and $z_i^{(k)}$ denoting the values of these parameters at iteration k , and $A^{(k)}$, $B^{(k)}$ denoting the associated polynomials. To update these parameters from iteration k to iteration $k+1$, we proceed as follows. First, we set $z_i^{(k+1)} = B^{(k)}(s_i)$, for $i = 1, \dots, N$. Then we choose $a^{(k+1)}$ and $b^{(k+1)}$ that minimize

$$\frac{1}{N} \sum_{i=1}^N \left| \frac{A^{(k+1)}(s_i) - h_i B^{(k+1)}(s_i)}{z_i^{(k+1)}} \right|^2.$$

(This can be done using ordinary linear least-squares.) We can start the iteration with $z_i^{(1)} = 1$, $i = 1, \dots, N$ (which is what would happen if we set $B^{(0)}(s) = 1$). The iteration is stopped when (or more accurately, if) successive iterates are very close, *i.e.*, we have $a^{(k+1)} \approx a^{(k)}$, and $b^{(k+1)} \approx b^{(k)}$. Several pathologies can arise in this algorithm. For example, we can end up with $z_i^{(k)} = 0$, or a certain matrix can be less than full rank, which complicates solving the least-squares problem to find $a^{(k)}$ and $b^{(k)}$. You can ignore these pathologies, however.

- (a) Explain how to obtain $a^{(k+1)}$ and $b^{(k+1)}$, given $z^{(k+1)}$. You must explain the math; you may not refer to any Matlab notation or operators (and especially, backslash) in your explanation. Please bear in mind that a_0, \dots, a_m and b_1, \dots, b_m are real, whereas many other variables and data in this problem are complex.
- (b) Implement the method, and apply it to the data given in `rat_data.m`. This file contains the data $\omega_1, \dots, \omega_N$, h_1, \dots, h_N , as well as m and N . Give the final coefficient vectors a , b , and the associated final value of J . Terminate the algorithm when

$$\left\| \begin{bmatrix} a^{(k+1)} - a^{(k)} \\ b^{(k+1)} - b^{(k)} \end{bmatrix} \right\| \leq 10^{-6}.$$

Plot J versus iteration k , with J on a logarithmic scale, and k on a linear scale, using the command `semilogy`. Plot $|H(j\omega)|$ on a logarithmic scale versus ω on a linear scale (using `semilogy`), for the first iteration, last iteration, and the problem data. To evaluate a polynomial in Matlab,

you can either write your own (short) code, or use the Matlab command `polyval`. This is a bit tricky, since `polyval` expects the polynomial coefficients to be listed in the reverse order than we use here. To evaluate $A(s)$ in Matlab you can use the command `polyval(a(m+1:-1:1),s)`. To evaluate $b(s)$ you can use `polyval([b(m:-1:1);1],s)`.

Note: no credit will be given for implementing any algorithm other than the one described in this problem.

- 6.19 *Quadratic placement.* We consider an integrated circuit (IC) that contains N cells or modules that are connected by K wires. We model a cell as a single point in \mathbf{R}^2 (which gives its location on the IC) and ignore the requirement that the cells must not overlap. The positions of the cells are

$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N),$$

i.e., x_i gives the x -coordinate of cell i , and y_i gives the y -coordinate of cell i . We have two types of cells: *fixed cells*, whose positions are fixed and given, and *free cells*, whose positions are to be determined. We will take the first n cells, at positions

$$(x_1, y_1), \dots, (x_n, y_n),$$

to be the free ones, and the remaining $N - n$ cells, at positions

$$(x_{n+1}, y_{n+1}), \dots, (x_N, y_N),$$

to be the fixed ones. The task of finding good positions for the free cells is called *placement*. (The fixed cells correspond to cells that are already placed, or external pins on the IC.) There are K wires that connect pairs of the cells. We will assign an orientation to each wire (even though wires are physically symmetric). Specifically, wire k goes *from* cell $I(k)$ *to* cell $J(k)$. Here I and J are functions that map wire number (*i.e.*, k) into the origination cell number (*i.e.*, $I(k)$), and the destination cell number (*i.e.*, $J(k)$), respectively. To describe the wire/cell topology and the functions I and J , we'll use the *node incidence matrix* A for the associated directed graph. The node incidence matrix $A \in \mathbf{R}^{K \times N}$ is defined as

$$A_{kj} = \begin{cases} 1 & \text{wire } k \text{ goes to cell } j, \text{ i.e., } j = J(k) \\ -1 & \text{wire } k \text{ goes from cell } j, \text{ i.e., } j = I(k) \\ 0 & \text{otherwise.} \end{cases}$$

Note that the k th row of A is associated with the k th wire, and the j th column of A is associated with the j th cell. The goal in placing the free cells is to use the smallest amount of interconnect wire, assuming that the wires are run as straight lines between the cells. (In fact, the wires in an IC are *not* run on straight lines directly between the cells, but that's another story. Pretending that the wires do run on straight lines seems to give good placements.) One common method, called *quadratic placement*, is to place the free cells in order to minimize the the total square wire length, given by

$$J = \sum_{k=1}^K ((x_{I(k)} - x_{J(k)})^2 + (y_{I(k)} - y_{J(k)})^2).$$

- (a) Explain how to find the positions of the free cells, *i.e.*,

$$(x_1, y_1), \dots, (x_n, y_n),$$

that minimize the total square wire length. You may make an assumption about the rank of one or more matrices that arise.

- (b) In this part you will determine the optimal quadratic placement for a specific set of cells and interconnect topology. The mfile `qplace_data.m` defines an instance of the quadratic placement

problem. Specifically, it defines the dimensions n , N , and K , and $N - n$ vectors `xfixed` and `yfixed`, which give the x - and y -coordinates of the fixed cells. The mfile also defines the node incidence matrix `A`, which is $K \times N$. Be sure to explain how you solve this problem, and to explain the matlab source code that solves it (which you must submit). Give the optimal locations of the free cells. Check your placement against various others, such as placing all free cells at the origin. You will also find an mfile that plots a proposed placement in a nice way:

`view_layout(xfree,yfree,xfixed,yfixed,A)`.

This mfile takes as argument the x - and y -coordinates of the free and fixed cells, as well as the node incidence matrix that describes the wires. It plots the proposed placement. Plot your optimal placement using `view_layout`.

- 6.20 *Least-squares state tracking.* Consider the system $x(t+1) = Ax(t) + Bu(t) \in \mathbf{R}^n$, with $x(0) = 0$. We *do not* assume it is controllable. Suppose $x_{\text{des}}(t) \in \mathbf{R}^n$ is given for $t = 1, \dots, N$ (and is meant to be the desired or target state trajectory). For a given input u , we define the mean-square tracking error as

$$E(u) = \frac{1}{N} \sum_{t=1}^N \|x(t) - x_{\text{des}}(t)\|^2.$$

- Explain how to find u_{opt} that minimizes E (in the general case). Your solution can involve a (general) pseudo-inverse.
- True or false:* If the system is controllable, there is a unique u_{opt} that minimizes $E(u)$. Briefly justify your answer.
- Find $E(u_{\text{opt}})$ for the specific system with

$$A = \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix},$$

$$x_{\text{des}}(t) = [t \ 0 \ 0]^T, \text{ and } N = 10.$$

- 6.21 *Time series prediction.* We consider an autonomous discrete-time linear system of the form

$$x(t+1) = Ax(t), \quad y(t) = Cx(t) + v(t),$$

where $x(t) \in \mathbf{R}^n$, $y(t) \in \mathbf{R}$ is the measured output signal, and $v(t) \in \mathbf{R}$ represents an output noise signal. In this problem, you *do not know* the matrices $A \in \mathbf{R}^{n \times n}$ or $C \in \mathbf{R}^{1 \times n}$, the state $x(t)$ (including the initial state $x(0)$), or even the system order n . You *do know* the measured output signal for $t = 1, \dots, p$:

$$y(1), \dots, y(p).$$

We give you two more pieces of information: the system order n is less than 20, and the RMS value of the noise, *i.e.*, $((1/p) \sum_{t=1}^p v(t)^2)^{1/2}$, is small (on the order of 0.001). The goal is to predict $y(t)$ for the next q time steps, *i.e.*, to predict what $y(p+1), \dots, y(p+q)$ will be. Here is the problem: get the time series data from the class web site in the file `timeseriesdata.m`, which gives $y(1), \dots, y(200)$. (We have $p = 200$ in this specific problem.) Then, predict what $y(201), \dots, y(220)$ are. Plot your estimates $\hat{y}(201), \dots, \hat{y}(220)$, and also, of course, give us the numbers. (You may also want to plot the whole set of data, from $t = 1$ to $t = 220$, just to make sure your prediction satisfies the ‘eyeball test’.) It is extremely important that you explain very clearly how you come up with your prediction. What is your method? If you make choices during your procedure, how do you make them?

- 6.22 *Extracting RC values from delay data.* We consider a CMOS digital gate that drives a load consisting of interconnect wires that connect to the inputs of other gates. To find the delay of the gate plus its load, we have to solve a complex, nonlinear ordinary differential equation that takes into account

circuit nonlinearities, parasitic capacitances, and so on. This can be done using a circuit simulator such as SPICE. A very simple approximation of the delay can be obtained by modeling the gate as a simple resistor with resistance R , and the load as a simple capacitor with capacitance C . In this simple model, the delay of the gate has the form ηRC , where η is a constant that depends on the precise definition of delay used. (One common choice is $\eta = 0.69$, which is based on the time it takes the voltage of a simple RC circuit to decay to $1/2$ its original value.) This simple RC delay model can be used for design, or approximate (but very fast) analysis. We address the question of determining a value of R for each of a set of gates, and a value of C for each of a set of loads, given accurate delay data (obtained by simulation) for each combination of a gate driving a load. We have n digital gates labeled $1, \dots, n$, and m loads labeled $1, \dots, m$. By simulation, we obtain the (accurate) delay D_{ij} for gate j driving load i . (D is given as an $m \times n$ matrix.) The goal is to find positive numbers R_1, \dots, R_n and C_1, \dots, C_m so that $D_{ij} \approx R_j C_i$. (For simplicity we'll take $\eta = 1$ in our delay model.) Finding good values of parameters for a simple model, given accurate data, is sometimes called *parameter extraction*. In this problem, we are extracting the gate resistances R_j and the load capacitances C_i , given accurate delay data D_{ij} (obtained by simulation). If we scale all resistances by a constant $\alpha > 0$, and scale all capacitances by $1/\alpha$, the approximate delays $R_j C_i$ don't change. To remove this ambiguity, we will fix $C_1 = 1$, *i.e.*, we will take the first load as our 'unit load'. Finally we get to the problem.

- (a) *Minimum mean-square absolute error.* Explain how to find $R_1^{\text{msa}}, \dots, R_n^{\text{msa}}$ and $C_1^{\text{msa}}, \dots, C_m^{\text{msa}}$ (positive, with $C_1^{\text{msa}} = 1$) that minimize the mean-square absolute error,

$$E_{\text{msa}} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n (D_{ij} - R_j C_i)^2.$$

- (b) *Minimum mean-square logarithmic error.* Explain how to find $R_1^{\text{msl}}, \dots, R_n^{\text{msl}}$ and $C_1^{\text{msl}}, \dots, C_m^{\text{msl}}$ (positive, with $C_1^{\text{msl}} = 1$) that minimize the mean-square logarithmic error,

$$E_{\text{msl}} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n (\log D_{ij} - \log(R_j C_i))^2.$$

(The logarithm here is base e , but it doesn't really matter.)

- (c) Find $R_1^{\text{msa}}, \dots, R_n^{\text{msa}}$ and $C_1^{\text{msa}}, \dots, C_m^{\text{msa}}$, as well as $R_1^{\text{msl}}, \dots, R_n^{\text{msl}}$ and $C_1^{\text{msl}}, \dots, C_m^{\text{msl}}$, for the particular delay data given in `rcextract_data.m` from the course web site. Submit the Matlab code used to calculate these values, as well as the values themselves. Also write down your minimum E_{msa} and E_{msl} values.

Please note the following:

- You do not need to know anything about digital circuits to solve this problem.
- The two criteria (absolute and logarithmic) are clearly close, so we expect the extracted Rs and Cs to be similar in the two cases. But they are not the same.
- In this problem we are more interested in your *approach* and *method* than the final numerical results. We will take points off if your method is not the best possible one, even if your answer is numerically close to the correct one.

6.23 *Reconstructing values from sums over subsets.* There are real numbers u_1, \dots, u_p that we do not know, but want to find. We do have information about sums of some subsets of the numbers. Specifically, we know v_1, \dots, v_q , where

$$v_i = \sum_{j \in S_i} u_j.$$

Here, S_i denotes the subset of $\{1, \dots, p\}$ that defines the partial sum used to form v_i . (We know both v_i and S_i , for $i = 1, \dots, q$.) We call the collection of subsets S_1, \dots, S_q *informative* if we can determine,

or reconstruct, u_1, \dots, u_p without ambiguity, from v_1, \dots, v_q . If the set of subsets is not informative, we say it is *uninformative*. As an example with $p = 3$ and $q = 4$,

$$v_1 = u_2 + u_3, \quad v_2 = u_1 + u_2 + u_3, \quad v_3 = u_1 + u_3, \quad v_4 = u_1 + u_2.$$

This corresponds to the subsets

$$S_1 = \{2, 3\}, \quad S_2 = \{1, 2, 3\}, \quad S_3 = \{1, 3\}, \quad S_4 = \{1, 2\}.$$

This collection of subsets is informative. To see this, we show how to reconstruct u_1, u_2, u_3 . First we note that $u_1 = v_2 - v_1$. Now that we know u_1 we can find u_2 from $u_2 = v_4 - u_1 = v_4 - v_2 + v_1$. In the same way we can get $u_3 = v_3 - u_1 = v_3 - v_2 + v_1$. **Note:** this is only an example to illustrate the notation.

- (a) This subproblem concerns the following specific case, with $p = 6$ numbers and $q = 11$ subsets. The subsets are

$$\begin{aligned} S_1 &= \{1, 2, 3\}, & S_2 &= \{1, 2, 4\}, & S_3 &= \{1, 2, 6\}, & S_4 &= \{1, 3, 5\}, & S_5 &= \{1, 4, 5\}, \\ S_6 &= \{2, 3, 6\}, & S_7 &= \{2, 4, 6\}, & S_8 &= \{3, 4, 5\}, & S_9 &= \{3, 5, 6\}, & S_{10} &= \{4, 5, 6\}, \\ S_{11} &= \{1, 2, 3, 4, 5, 6\}. \end{aligned}$$

The associated sums are

$$\begin{aligned} v_1 &= -2, & v_2 &= 14, & v_3 &= 6, & v_4 &= 4, & v_5 &= 20, & v_6 &= -5, \\ v_7 &= 11, & v_8 &= 9, & v_9 &= 1, & v_{10} &= 17, & v_{11} &= 15. \end{aligned}$$

Choose one of the following:

- *The collection of subsets S_1, \dots, S_{11} is informative.*
Justify why you believe this is the case, and reconstruct u_1, \dots, u_6 .
 - *The collection of subsets S_1, \dots, S_{11} is uninformative.*
To justify this, give two different sets of values u_1, \dots, u_6 , and $\tilde{u}_1, \dots, \tilde{u}_6$, whose given subset sums agree with the given v_1, \dots, v_{11} .
- (b) This subproblem concerns a general method for reconstructing $u = (u_1, \dots, u_p)$ given $v = (v_1, \dots, v_q)$ (and of course, the subsets S_1, \dots, S_q). We define the *subset count matrix* $Z \in \mathbf{R}^{q \times p}$ as follows: Z_{ij} is the number of subsets containing both i and j . (Thus, Z_{ii} is the number of subsets that contain i .) For each i , we define f_i as the sum of all v_j , over subsets that contain i :

$$f_i = \sum_{i \in S_j} v_j, \quad i = 1, \dots, p.$$

Then we reconstruct u as $u = Z^{-1}f$. (Of course, this requires that Z is invertible.) Choose one of the following:

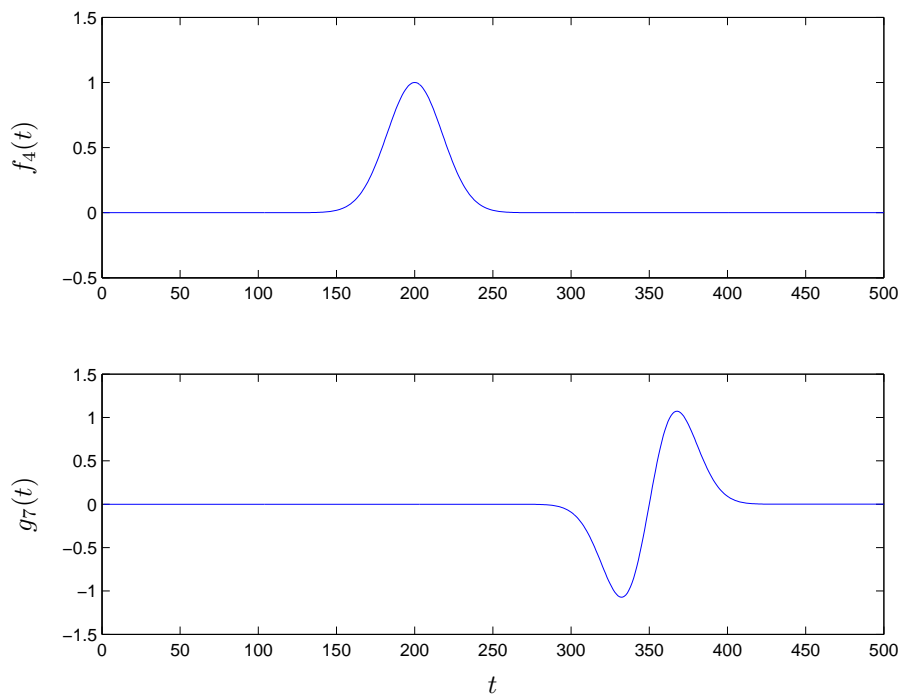
- *The method works, whenever the collection of subsets is informative.*
By ‘works’ we mean that Z is invertible, and that $Z^{-1}f$ is the unique u with subset sums v . If you believe this is the case, explain why.
- *The method can fail, even when the collection of subsets is informative.*
To convince us of this, give a *specific example*, where the collection of subsets is informative, but the method above fails, *i.e.*, either Z is singular, or $Z^{-1}f$ does not have the required subset sums. (Please give us the simplest example you can think of.)

6.24 *Signal estimation using least-squares.* This problem concerns discrete-time signals defined for $t = 1, \dots, 500$. We’ll represent these signals by vectors in \mathbf{R}^{500} , with the index corresponding to the time.

We are given a noisy measurement $y_{\text{meas}}(1), \dots, y_{\text{meas}}(500)$, of a signal $y(1), \dots, y(500)$ that is thought to be, at least approximately, a linear combination of the 22 signals

$$f_k(t) = e^{-(t-50k)^2/25^2}, \quad g_k(t) = \left(\frac{t-50k}{10}\right) e^{-(t-50k)^2/25^2},$$

where $t = 1, \dots, 500$ and $k = 0, \dots, 10$. Plots of f_4 and g_7 (as examples) are shown below.



As our estimate of the original signal, we will use the signal $\hat{y} = (\hat{y}(1), \dots, \hat{y}(500))$ in the span of $f_0, \dots, f_{10}, g_0, \dots, g_{10}$, that is closest to $y_{\text{meas}} = (y_{\text{meas}}(1), \dots, y_{\text{meas}}(500))$ in the RMS (root-mean-square) sense. Explain how to find \hat{y} , and carry out your method on the signal y_{meas} given in `sig_est_data.m` on the course web site. Plot y_{meas} and \hat{y} on the same graph. Plot the residual (the difference between these two signals) on a different graph, and give its RMS value.

6.25 *Point of closest convergence of a set of lines.* We have m lines in \mathbf{R}^n , described as

$$\mathcal{L}_i = \{p_i + tv_i \mid t \in \mathbf{R}\}, \quad i = 1, \dots, m,$$

where $p_i \in \mathbf{R}^n$, and $v_i \in \mathbf{R}^n$, with $\|v_i\| = 1$, for $i = 1, \dots, m$. We define the distance of a point $z \in \mathbf{R}^n$ to a line \mathcal{L} as

$$\mathbf{dist}(z, \mathcal{L}) = \min\{\|z - u\| \mid u \in \mathcal{L}\}.$$

(In other words, $\mathbf{dist}(z, \mathcal{L})$ gives the closest distance between the point z and the line \mathcal{L} .)

We seek a point $z^* \in \mathbf{R}^n$ that minimizes the sum of the squares of the distances to the lines,

$$\sum_{i=1}^m \mathbf{dist}(z, \mathcal{L}_i)^2.$$

The point z^* that minimizes this quantity is called the *point of closest convergence*.

- (a) Explain how to find the point of closest convergence, given the lines (*i.e.*, given p_1, \dots, p_m and v_1, \dots, v_m). If your method works provided some condition holds (such as some matrix being full rank), say so. If you can relate this condition to a simple one involving the lines, please do so.
- (b) Find the point z^* of closest convergence for the lines with data given in the Matlab file `line_conv_data.m`. This file contains $n \times m$ matrices P and V whose columns are the vectors p_1, \dots, p_m , and v_1, \dots, v_m , respectively. The file also contains commands to plot the lines and the point of closest convergence (once you have found it). Please include this plot with your solution.

6.26 *Estimating direction and amplitude of a light beam.* A light beam with (nonnegative) amplitude a comes from a direction $d \in \mathbf{R}^3$, where $\|d\| = 1$. (This means the beam travels in the direction $-d$.) The beam falls on $m \geq 3$ photodetectors, each of which generates a scalar signal that depends on the beam amplitude and direction, and the direction in which the photodetector is pointed. Specifically, photodetector i generates an output signal p_i , with

$$p_i = a\alpha \cos \theta_i + v_i,$$

where θ_i is the angle between the beam direction d and the outward normal vector q_i of the surface of the i th photodetector, and α is the photodetector sensitivity. You can interpret $q_i \in \mathbf{R}^3$, which we assume has norm one, as the direction the i th photodetector is pointed. We assume that $|\theta_i| < 90^\circ$, *i.e.*, the beam illuminates the top of the photodetectors. The numbers v_i are small measurement errors.

You are given the photodetector direction vectors $q_1, \dots, q_m \in \mathbf{R}^3$, the photodetector sensitivity α , and the noisy photodetector outputs, $p_1, \dots, p_m \in \mathbf{R}$. Your job is to estimate the beam direction $d \in \mathbf{R}^3$ (which is a unit vector), and a , the beam amplitude.

To describe unit vectors q_1, \dots, q_m and d in \mathbf{R}^3 we will use azimuth and elevation, defined as follows:

$$q = \begin{bmatrix} \cos \phi \cos \theta \\ \cos \phi \sin \theta \\ \sin \phi \end{bmatrix}.$$

Here ϕ is the elevation (which will be between 0° and 90° , since all unit vectors in this problem have positive 3rd component, *i.e.*, point upward). The azimuth angle θ , which varies from 0° to 360° , gives the direction in the plane spanned by the first and second coordinates. If $q = e_3$ (*i.e.*, the direction is directly up), the azimuth is undefined.

- (a) Explain how to do this, using a method or methods from this class. The simpler the method the better. If some matrix (or matrices) needs to be full rank for your method to work, say so.
- (b) Carry out your method on the data given in `beam_estim_data.m`. This mfile defines `p`, the vector of photodetector outputs, a vector `det_az`, which gives the azimuth angles of the photodetector directions, and a vector `det_el`, which gives the elevation angles of the photodetector directions. Note that both of these are given in *degrees*, not radians. Give your final estimate of the beam amplitude a and beam direction d (in azimuth and elevation, in degrees).

6.27 *Smooth interpolation on a 2D grid.* This problem concerns arrays of real numbers on an $m \times n$ grid. Such an array can represent an image, or a sampled description of a function defined on a rectangle. We can describe such an array by a matrix $U \in \mathbf{R}^{m \times n}$, where U_{ij} gives the real number at location i, j , for $i = 1, \dots, m$ and $j = 1, \dots, n$. We will think of the index i as associated with the y axis, and the index j as associated with the x axis.

It will also be convenient to describe such an array by a vector $u = \text{vec}(U) \in \mathbf{R}^{mn}$. Here `vec` is the function that stacks the columns of a matrix on top of each other:

$$\text{vec}(U) = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix},$$

where $U = [u_1 \cdots u_n]$. To go back to the array representation, from the vector, we have $U = \mathbf{vec}^{-1}(u)$. (This looks complicated, but isn't; \mathbf{vec}^{-1} just arranges the elements in a vector into an array.)

We will need two linear functions that operate on $m \times n$ arrays. These are simple approximations of partial differentiation with respect to the x and y axes, respectively. The first function takes as argument an $m \times n$ array U and returns an $m \times (n - 1)$ array V of forward (rightward) differences:

$$V_{ij} = U_{i,j+1} - U_{ij}, \quad i = 1, \dots, m, \quad j = 1, \dots, n - 1.$$

We can represent this linear mapping as multiplication by a matrix $D_x \in \mathbf{R}^{m(n-1) \times mn}$, which satisfies

$$\mathbf{vec}(V) = D_x \mathbf{vec}(U).$$

(This looks scarier than it is—each row of the matrix D_x has exactly one +1 and one -1 entry in it.)

The other linear function, which is a simple approximation of partial differentiation with respect to the y axis, maps an $m \times n$ array U into an $(m - 1) \times n$ array W , is defined as

$$W_{ij} = U_{i+1,j} - U_{ij}, \quad i = 1, \dots, m - 1, \quad j = 1, \dots, n.$$

We define the matrix $D_y \in \mathbf{R}^{(m-1)n \times mn}$, which satisfies $\mathbf{vec}(W) = D_y \mathbf{vec}(U)$.

We define the *roughness* of an array U as

$$R = \|D_x \mathbf{vec}(U)\|^2 + \|D_y \mathbf{vec}(U)\|^2.$$

The roughness measure R is the sum of the squares of the differences of each element in the array and its neighbors. Small R corresponds to smooth, or smoothly varying, U . The roughness measure R is zero precisely for constant arrays, *i.e.*, when U_{ij} are all equal.

Now we get to the problem, which is to interpolate some unknown values in an array in the smoothest possible way, given the known values in the array. To define this precisely, we partition the set of indices $\{1, \dots, mn\}$ into two sets: I_{known} and I_{unknown} . We let $k \geq 1$ denote the number of known values (*i.e.*, the number of elements in I_{known}), and $mn - k$ the number of unknown values (the number of elements in I_{unknown}). We are given the values u_i for $i \in I_{\text{known}}$; the goal is to guess (or estimate or assign) values for u_i for $i \in I_{\text{unknown}}$. We'll choose the values for u_i , with $i \in I_{\text{unknown}}$, so that the resulting U is as smooth as possible, *i.e.*, so it minimizes R . Thus, the goal is to fill in or interpolate missing data in a 2D array (an image, say), so the reconstructed array is as smooth as possible.

We give the k known values in a vector $w_{\text{known}} \in \mathbf{R}^k$, and the $mn - k$ unknown values in a vector $w_{\text{unknown}} \in \mathbf{R}^{mn-k}$. The complete array is obtained by putting the entries of w_{known} and w_{unknown} into the correct positions of the array. We describe these operations using two matrices $Z_{\text{known}} \in \mathbf{R}^{mn \times k}$ and $Z_{\text{unknown}} \in \mathbf{R}^{mn \times (mn-k)}$, that satisfy

$$\mathbf{vec}(U) = Z_{\text{known}} w_{\text{known}} + Z_{\text{unknown}} w_{\text{unknown}}.$$

(This looks complicated, but isn't: Each row of these matrices is a unit vector, so multiplication with either matrix just stuffs the entries of the w vectors into particular locations in $\mathbf{vec}(U)$. In fact, the matrix $[Z_{\text{known}} \ Z_{\text{unknown}}]$ is an $mn \times mn$ permutation matrix.)

In summary, you are given the problem data w_{known} (which gives the known array values), Z_{known} (which gives the locations of the known values), and Z_{unknown} (which gives the locations of the unknown array values, in some specific order). Your job is to find w_{unknown} that minimizes R .

- (a) Explain how to solve this problem. You are welcome to use any of the operations, matrices, and vectors defined above in your solution (*e.g.*, \mathbf{vec} , \mathbf{vec}^{-1} , D_x , D_y , Z_{known} , Z_{unknown} , w_{known} , \dots). If your solution is valid provided some matrix is (or some matrices are) full rank, say so.

- (b) Carry out your method using the data created by `smooth_interpolation.m`. The file gives m , n , w_{known} , Z_{known} and Z_{unknown} . This file also creates the matrices D_x and D_y , which you are welcome to use. (This was *very* nice of us, by the way.) You are welcome to look at the code that generates these matrices, but you do not need to understand it. For this problem instance, around 50% of the array elements are known, and around 50% are unknown.

The mfile also includes the original array `Uorig` from which we removed elements to create the problem. This is just so you can see how well your smooth reconstruction method does in reconstructing the original array. Of course, you cannot use `Uorig` to create your interpolated array `U`.

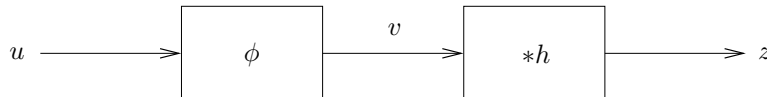
To visualize the arrays use the Matlab command `imagesc()`, with matrix argument. If you prefer a grayscale image, or don't have a color printer, you can issue the command `colormap gray`. The mfile that gives the problem data will plot the original image `Uorig`, as well as an image containing the known values, with zeros substituted for the unknown locations. This will allow you to see the pattern of known and unknown array values.

Compare `Uorig` (the original array) and `U` (the interpolated array found by your method), using `imagesc()`. Hand in complete source code, as well as the plots. Be sure to give the value of roughness R of U .

Hints:

- In Matlab, `vec(U)` can be computed as `U(:)`;
- `vec-1(u)` can be computed as `reshape(u,m,n)`.

- 6.28 *Designing a nonlinear equalizer from I/O data.* This problem concerns the discrete-time system shown below, which consists of a memoryless nonlinearity ϕ , followed by a convolution filter with finite impulse response h . The scalar signal u is the input, and the scalar signal z is the output.



What this means is

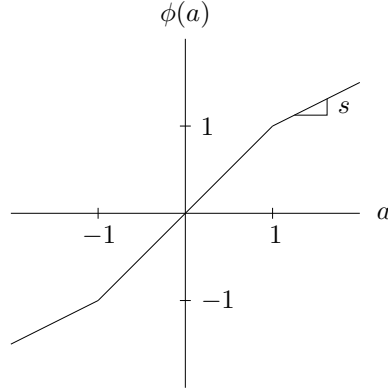
$$z(t) = \sum_{\tau=0}^{M-1} h(\tau)v(t-\tau), \quad v(t) = \phi(u(t)), \quad t \in \mathbf{Z}.$$

(Note that these signals are defined for all integer times, not just nonnegative times.)

Here $\phi : \mathbf{R} \rightarrow \mathbf{R}$, with the specific form

$$\phi(a) = \begin{cases} a & -1 \leq a \leq 1 \\ 1 - s + sa & a > 1 \\ -1 + s + sa & a < -1, \end{cases}$$

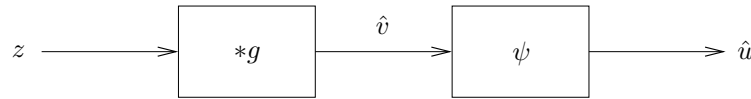
where $s > 0$ is a parameter. This function is shown below.



Here is an interpretation (that is not needed to solve the problem). The nonlinear function ϕ represents a power amplifier that is nonlinear for input signals larger than one in magnitude; s is called the *saturation gain* of the amplifier. The convolution system represents the transmission channel.

We are going to design an *equalizer* for the system, *i.e.*, another system that takes the signal z as input, and gives an output \hat{u} which is an approximation of the input signal u .

Our equalizer will have the form shown below.



This means

$$\hat{v}(t) = \sum_{\tau=0}^{M-1} g(\tau)z(t-\tau), \quad \hat{u}(t) = \psi(\hat{v}(t)), \quad t \in \mathbf{Z}.$$

This equalizer will work well provided $g * h \approx \delta$ (in which case $\hat{v}(t) \approx v(t)$), and $\psi = \phi^{-1}$ (*i.e.*, $\psi(\phi(a)) = a$ for all a).

To make sure our (standard) notation here is clear, recall that

$$(g * h)(t) = \sum_{\tau=\max\{0, t-M+1\}}^{\min\{M-1, t\}} g(\tau)h(t-\tau), \quad t = 0, \dots, 2M-1.$$

(Note: in matlab `conv(g,h)` gives the convolution of \mathbf{g} and \mathbf{h} , but these vectors are indexed from 1 to M , *i.e.*, $\mathbf{g}(1)$ corresponds to $g(0)$.) The term δ is the Kronecker delta, defined as $\delta(0) = 1$, $\delta(i) = 0$ for $i \neq 0$.

Now, finally, we come to the problem. You are given some input/output (I/O) data $u(1), \dots, u(N)$, $z(1), \dots, z(N)$, and M (the length of g , and also the length of h). You do *not* know the parameter s , or the channel impulse response $h(0), \dots, h(M-1)$. You also don't know $u(t)$, $z(t)$ for $t \leq 0$.

- (a) Explain how to find \hat{s} , an estimate of the saturation gain s , and $g(0), \dots, g(M-1)$, that minimize

$$J = \frac{1}{N-M+1} \sum_{i=M}^N \left(\hat{v}(i) - \phi(u(i)) \right)^2.$$

Here u refers to the given input data, and \hat{v} comes from the given output data z . Note that if $g * h = \delta$ and $s = \hat{s}$, we have $J = 0$.

We exclude $i = 1, \dots, M-1$ in the sum defining J because these terms depend (through \hat{v}) on $z(0), z(-1), \dots$, which are unknown.

- (b) Apply your method to the data given in the file `nleq_data.m`. Give the values of the parameters \hat{s} and $g(0), \dots, g(M-1)$ found, as well as J . Plot g using the matlab command `stem`.
- (c) Using the values of \hat{s} and $g(0), \dots, g(M-1)$ found in part (b), find the equalized signal $\hat{u}(t)$, for $t = 1, \dots, N$. For the purposes of finding $\hat{u}(t)$ you can assume that $z(t) = 0$ for $t \leq 0$. As a result, we can expect a large equalization error (*i.e.*, $\hat{u}(t) - u(t)$) for $t = 1, \dots, M-1$. Plot the input signal $u(t)$, the output signal $z(t)$, the equalized signal $\hat{u}(t)$, and the equalization error $\hat{u}(t) - u(t)$, for $t = 1, \dots, N$.
- 6.29 *Simple fitting.* You are given some data $x_1, \dots, x_N \in \mathbf{R}$ and $y_1, \dots, y_N \in \mathbf{R}$. These data are available in `simplefitdata.m` on the course web site.
- (a) Find the best affine fit, *i.e.*, $y_i \approx ax_i + b$, where ‘best’ means minimizing $\sum_{i=1}^N (y_i - (ax_i + b))^2$. (This is often called the ‘best linear fit’.) Set this up and solve it as a least-squares problem. Plot the data and the fit in the same figure. Give us a and b , and submit the code you used to find a and b .
- (b) Repeat for the best least-squares cubic fit, *i.e.*, $y_i \approx ax_i^3 + bx_i^2 + cx_i + d$.
- 6.30 *Estimating parameters from noisy measurements.* In this problem you will compare a least-squares estimate of a parameter vector (which uses all available measurements) with a method that uses just enough measurements. Carry out the following steps.
- (a) First we generate some problem data in Matlab. (You’re free to use any other software system instead.) Generate a 50×20 matrix A using `A=randn(50,20)`. (This chooses the entries from a normal distribution, but this doesn’t really matter for us.) Generate a noise vector v of length 50 using `v=0.1*randn(50,1)`. Generate a measurement vector x of length 20 using `x=randn(20,1)`. Finally, generate a measurement vector $y = Ax + v$.
- (b) Find the least-squares approximate solution of $y = Ax$, and call it x^{ls} . Find the relative error $\|x^{\text{ls}} - x\|/\|x\|$.
- (c) Now form a 20-long truncated measurement vector y^{trunc} which consists of the first 20 entries of y . Form an estimate of x from y^{trunc} . Call this estimate x^{jem} (‘Just Enough Measurements’). Find the relative error of x^{jem} .
- (d) Run your script (*i.e.*, (a)–(c)) several times. You’ll generate different different data each time, and you’ll get different numerical results in parts (b) and (c). Give a one sentence comment about what you observe.
- Note.* Since you are generating the data randomly, it is remotely possible that the second method will work better than the first, at least for one run. If this happens to you, quickly run your script again. Do not mention the incident to anyone.
- 6.31 *Signal reconstruction for a bandlimited signal.* In this problem we refer to *signals*, which are just vectors, with index interpreted as (discrete) time. It is common to write the index for a signal as an argument, rather than as a subscript; for example, if $y \in \mathbf{R}^N$ is a signal, we use $y(t)$ to denote y_t , with $t \in \{1, 2, \dots, N\}$. Another notation you’ll sometimes see in signal processing texts is $y[t]$ for y_t . The *discrete cosine transformation* (DCT) of the signal $y \in \mathbf{R}^N$ is another signal, typically denoted using the corresponding upper case symbol $Y \in \mathbf{R}^N$. It is defined as

$$Y(k) = \sum_{t=1}^N y(t)w(k) \cos \frac{\pi(2t-1)(k-1)}{2N}, \quad k = 1, \dots, N,$$

where $w(k)$ are weights, with

$$w(k) = \begin{cases} \sqrt{1/N}, & k = 1, \\ \sqrt{2/N}, & k = 2, \dots, N. \end{cases}$$

The number $Y(k)$ is referred to as the k th DCT coefficient of y . The DCT bandwidth of the signal y is the smallest K for which $Y(K) \neq 0$, and $Y(k) = 0$ for $k = K + 1, \dots, N$. When $K < N$, the signal is called *DCT bandlimited*. (The term is typically used to refer to the case when the DCT bandwidth, K , is significantly smaller than N .)

A signal y can be reconstructed from its DCT Y , via the *inverse DCT transform*, with

$$y(t) = \sum_{k=1}^N Y(k)w(k) \cos \frac{\pi(2t-1)(k-1)}{2N}, \quad t = 1, \dots, N,$$

where $w(k)$ are the same weights as those used above in the DCT.

Now for the problem. You are given noise-corrupted values of a DCT bandlimited signal y , at some (integer) times t_1, \dots, t_M , where $1 \leq t_1 < t_2 < \dots < t_M \leq N$:

$$y_i^{\text{samp}} = y(t_i) + v_i, \quad i = 1, \dots, M.$$

Here, v_i are small noises or errors. You don't know v , but you do know that its RMS value is approximately σ , a known constant. (In signal processing, y^{samp} would be called a non-uniformly sampled, noise corrupted version of y .)

Your job is to

- Determine the smallest DCT bandwidth (*i.e.*, the smallest K) that y could have.
- Find an estimate of y , \hat{y} , which has this bandwidth.

Your estimate \hat{y} must be consistent with the sample measurements y^{samp} . While it need not match exactly (you were told there was a small amount of noise in y^{samp}), you should ensure that the vector of differences,

$$(y_1^{\text{samp}} - \hat{y}(t_1^{\text{samp}}), \dots, y_M^{\text{samp}} - \hat{y}(t_M^{\text{samp}})),$$

has a small RMS value, on the order of σ (and certainly no more than 3σ).

- Clearly describe how to solve this problem. You can use any concepts we have used, to date, in EE263. *You cannot use (and do not need) any concepts from outside the class. This includes the Fourier transform and other signal processing methods you might know.*
- Carry out your method on the data in `bandlimit.m`. Running this script will define `N`, `ysamp`, `M`, `tsamp`, and `sigma`. It will also plot the sampled signal.

Give K , your estimate of the DCT bandwidth of y . Show \hat{y} on the same plot as the original sampled signal. (We have added the command to do this in `bandlimit.m`, but commented it out.)

Also, give us $\hat{y}(129)$, to four significant figures.

You might find the Matlab functions `dct` and `idct` useful; `dct(eye(N))` and `idct(eye(N))` will return matrices whose columns are the DCT, and inverse DCT transforms, respectively, of the unit vectors. Note, however, that you can solve the problem without using these functions.

- 6.32 *Fitting a model for hourly temperature.* You are given a set of temperature measurements (in degrees C), $y_t \in \mathbf{R}$, $t = 1, \dots, N$, taken hourly over one week (so $N = 168$). An expert says that over this week, an appropriate model for the hourly temperature is a trend (*i.e.*, a linear function of t) plus a diurnal component (*i.e.*, a 24-periodic component):

$$\hat{y}_t = at + p_t,$$

where $a \in \mathbf{R}$ and $p \in \mathbf{R}^N$ satisfies $p_{t+24} = p_t$, for $t = 1, \dots, N - 24$. We can interpret a (which has units of degrees C per hour) as the warming or cooling trend (for $a > 0$ or $a < 0$, respectively) over the week.

- (a) Explain how to find $a \in \mathbf{R}$ and $p \in \mathbf{R}^N$ (which is 24-periodic) that minimize the RMS value of $y - \hat{y}$.
- (b) Carry out the procedure described in part (a) on the data set found in `tempdata.m`. Give the value of the trend parameter a that you find. Plot the model \hat{y} and the measured temperatures y on the same plot. (The Matlab code to do this is in the file `tempdata.m`, but commented out.)
- (c) *Temperature prediction.* Use the model found in part (b) to predict the temperature for the next 24-hour period (*i.e.*, from $t = 169$ to $t = 192$). The file `tempdata.m` also contains a 24 long vector `ytom` with tomorrow's temperatures. Plot tomorrow's temperature and your prediction of it, based on the model found in part (b), on the same plot. What is the RMS value of your prediction error for tomorrow's temperatures?

6.33 *Empirical algorithm complexity.* The runtime T of an algorithm depends on its input data, which is characterized by three key parameters: k , m , and n . (These are typically integers that give the dimensions of the problem data.) A simple and standard model that shows how T scales with k , m , and n has the form

$$\hat{T} = \alpha k^\beta m^\gamma n^\delta,$$

where $\alpha, \beta, \gamma, \delta \in \mathbf{R}$ are constants that characterize the approximate runtime model. If, for example, $\delta \approx 3$, we say that the algorithm has (approximately) cubic complexity in n . (In general, the exponents β, γ , and δ need not be integers, or close to integers.)

Now suppose you are given measured runtimes for N executions of the algorithm, with different sets of input data. For each data record, you are given T_i (the runtime), and the parameters k_i, m_i , and n_i . It's possible (and often occurs) that two data records have identical values of k, m , and n , but different values of T . This means the algorithm was run on two different data sets that had the same dimensions; the corresponding runtimes can be (and often are) a little different.

We wish to find values of α, β, γ , and δ for which our model (approximately) fits our measurements. We define the fitting cost as

$$J = (1/N) \sum_{i=1}^N \left(\log(\hat{T}_i/T_i) \right)^2,$$

where $\hat{T}_i = \alpha k_i^\beta m_i^\gamma n_i^\delta$ is the runtime predicted by our model, using the given parameter values. This fitting cost can be (loosely) interpreted in terms of relative or percentage fit. If $(\log(\hat{T}_i/T_i))^2 \leq \epsilon$, then \hat{T}_i lies between $T_i/\exp \sqrt{\epsilon}$ and $T_i \exp \sqrt{\epsilon}$.

Your task is to find constants $\alpha, \beta, \gamma, \delta$ that minimize J .

- (a) Explain how to do this. If your method always finds the values that give the true global minimum value of J , say so. If your algorithm cannot guarantee finding the true global minimum, say so. If your method requires some matrix (or matrices) to be full rank, say so.
- (b) Carry out your method on the data found in `empac_data.m`. Give the values of α, β, γ , and δ you find, and the corresponding value of J .

Lecture 7 – Regularized least-squares and Gauss-Newton method

7.1 *Fitting a Gaussian function to data.* A Gaussian function has the form

$$f(t) = ae^{-(t-\mu)^2/\sigma^2}.$$

Here $t \in \mathbf{R}$ is the independent variable, and $a \in \mathbf{R}$, $\mu \in \mathbf{R}$, and $\sigma \in \mathbf{R}$ are parameters that affect its shape. The parameter a is called the *amplitude* of the Gaussian, μ is called its *center*, and σ is called the *spread* or *width*. We can always take $\sigma > 0$. For convenience we define $p \in \mathbf{R}^3$ as the vector of the parameters, *i.e.*, $p = [a \ \mu \ \sigma]^T$. We are given a set of data,

$$t_1, \dots, t_N, \quad y_1, \dots, y_N,$$

and our goal is to fit a Gaussian function to the data. We will measure the quality of the fit by the root-mean-square (RMS) fitting error, given by

$$E = \left(\frac{1}{N} \sum_{i=1}^N (f(t_i) - y_i)^2 \right)^{1/2}.$$

Note that E is a function of the parameters a , μ , σ , *i.e.*, p . Your job is to choose these parameters to minimize E . You'll use the Gauss-Newton method.

- (a) Work out the details of the Gauss-Newton method for this fitting problem. Explicitly describe the Gauss-Newton steps, including the matrices and vectors that come up. You can use the notation $\Delta p^{(k)} = [\Delta a^{(k)} \ \Delta \mu^{(k)} \ \Delta \sigma^{(k)}]^T$ to denote the update to the parameters, *i.e.*,

$$p^{(k+1)} = p^{(k)} + \Delta p^{(k)}.$$

(Here k denotes the k th iteration.)

- (b) Get the data t , y (and N) from the file `gauss_fit_data.m`, available on the class website. Implement the Gauss-Newton (as outlined in part (a) above). You'll need an initial guess for the parameters. You can visually estimate them (giving a short justification), or estimate them by any other method (but you must explain your method). Plot the RMS error E as a function of the iteration number. (You should plot enough iterations to convince yourself that the algorithm has nearly converged.) Plot the final Gaussian function obtained along with the data on the same plot. Repeat for another reasonable, but different initial guess for the parameters. Repeat for another set of parameters that is *not* reasonable, *i.e.*, not a good guess for the parameters. (It's possible, of course, that the Gauss-Newton algorithm doesn't converge, or fails at some step; if this occurs, say so.) Briefly comment on the results you obtain in the three cases.

7.2 *E-911.* The federal government has mandated that cellular network operators must have the ability to locate a cell phone from which an emergency call is made. This problem concerns a simplified version of an E-911 system that uses time of arrival information at a number of base stations to estimate the cell phone location. A cell phone at location $x \in \mathbf{R}^2$ (we assume that the elevation is zero for simplicity) transmits an emergency signal at time τ . This signal is received at n base stations, located at locations $s_1, \dots, s_n \in \mathbf{R}^2$. Each base station can measure the time of arrival of the emergency signal, within a few tens of nanoseconds. (This is possible because the base stations are synchronized using the Global Positioning System.) The measured times of arrival are

$$t_i = \frac{1}{c} \|s_i - x\| + \tau + v_i, \quad i = 1, \dots, n,$$

where c is the speed of light, and v_i is the noise or error in the measured time of arrival. You can assume that v_i is on the order of a few tens of nanoseconds. The problem is to estimate the cell phone

position $x \in \mathbf{R}^2$, as well as the time of transmission τ , based on the time of arrival measurements t_1, \dots, t_n . The mfile `e911_data.m`, available on the course web site, defines the data for this problem. Specifically, it defines a 2×9 matrix \mathbf{S} , whose columns give the positions of the 9 basestations, a 1×9 vector \mathbf{t} that contains the measured times of arrival, and the constant c , which is the speed of light. Distances are given in meters, times in nanoseconds, and the speed of light in meters/nanosecond. You can assume that the position x is somewhere in the box

$$|x_1| \leq 3000, \quad |x_2| \leq 3000,$$

and that $|\tau| \leq 5000$ (although all that really matters are the time differences). Your solution must contain the following:

- An explanation of your approach to solving this problem, including how you will check that your estimate is reasonable.
- The matlab source code you use to solve the problem, and check the results.
- The numerical results obtained by your method, including the results of any verification you do.

7.3 *Curve-smoothing*. We are given a function $F : [0, 1] \rightarrow \mathbf{R}$ (whose graph gives a curve in \mathbf{R}^2). Our goal is to find another function $G : [0, 1] \rightarrow \mathbf{R}$, which is a *smoothed* version of F . We'll judge the smoothed version G of F in two ways:

- *Mean-square deviation from F* , defined as

$$D = \int_0^1 (F(t) - G(t))^2 dt.$$

- *Mean-square curvature*, defined as

$$C = \int_0^1 G''(t)^2 dt.$$

We want *both* D and C to be small, so we have a problem with two objectives. In general there will be a trade-off between the two objectives. At one extreme, we can choose $G = F$, which makes $D = 0$; at the other extreme, we can choose G to be an affine function (*i.e.*, to have $G''(t) = 0$ for all $t \in [0, 1]$), in which case $C = 0$. The problem is to identify the optimal trade-off curve between C and D , and explain how to find smoothed functions G on the optimal trade-off curve. To reduce the problem to a finite-dimensional one, we will represent the functions F and G (approximately) by vectors $f, g \in \mathbf{R}^n$, where

$$f_i = F(i/n), \quad g_i = G(i/n).$$

You can assume that n is chosen large enough to represent the functions well. Using this representation we will use the following objectives, which approximate the ones defined for the functions above:

- *Mean-square deviation*, defined as

$$d = \frac{1}{n} \sum_{i=1}^n (f_i - g_i)^2.$$

- *Mean-square curvature*, defined as

$$c = \frac{1}{n-2} \sum_{i=2}^{n-1} \left(\frac{g_{i+1} - 2g_i + g_{i-1}}{1/n^2} \right)^2.$$

In our definition of c , note that

$$\frac{g_{i+1} - 2g_i + g_{i-1}}{1/n^2}$$

gives a simple approximation of $G''(i/n)$. You will only work with this approximate version of the problem, *i.e.*, the vectors f and g and the objectives c and d .

- (a) Explain how to find g that minimizes $d + \mu c$, where $\mu \geq 0$ is a parameter that gives the relative weighting of sum-square curvature compared to sum-square deviation. Does your method always work? If there are some assumptions you need to make (say, on rank of some matrix, independence of some vectors, etc.), state them clearly. Explain how to obtain the two extreme cases: $\mu = 0$, which corresponds to minimizing d without regard for c , and also the solution obtained as $\mu \rightarrow \infty$ (*i.e.*, as we put more and more weight on minimizing curvature).
- (b) Get the file `curve_smoothing.m` from the course web site. This file defines a specific vector f that you will use. Find and plot the optimal trade-off curve between d and c . Be sure to identify any critical points (such as, for example, any intersection of the curve with an axis). Plot the optimal g for the two extreme cases $\mu = 0$ and $\mu \rightarrow \infty$, and for three values of μ in between (chosen to show the trade-off nicely). On your plots of g , be sure to include also a plot of f , say with dotted line type, for reference. Submit your Matlab code.

7.4 *Optimal choice of initial temperature profile.* We consider a thermal system described by an n -element finite-element model. The elements are arranged in a line, with the temperature of element i at time t denoted $T_i(t)$. Temperature is measured in degrees Celsius above ambient; negative $T_i(t)$ corresponds to a temperature below ambient. The dynamics of the system are described by

$$c_1 \dot{T}_1 = -a_1 T_1 - b_1 (T_1 - T_2),$$

$$c_i \dot{T}_i = -a_i T_i - b_i (T_i - T_{i+1}) - b_{i-1} (T_i - T_{i-1}), \quad i = 2, \dots, n-1,$$

and

$$c_n \dot{T}_n = -a_n T_n - b_{n-1} (T_n - T_{n-1}).$$

where $c \in \mathbf{R}^n$, $a \in \mathbf{R}^n$, and $b \in \mathbf{R}^{n-1}$ are given and are all positive.

We can interpret this model as follows. The parameter c_i is the heat capacity of element i , so $c_i \dot{T}_i$ is the net heat flow into element i . The parameter a_i gives the thermal conductance between element i and the environment, so $a_i T_i$ is the heat flow from element i to the environment (*i.e.*, the direct heat loss from element i .) The parameter b_i gives the thermal conductance between element i and element $i+1$, so $b_i (T_i - T_{i+1})$ is the heat flow from element i to element $i+1$. Finally, $b_{i-1} (T_i - T_{i-1})$ is the heat flow from element i to element $i-1$.

The goal of this problem is to choose the initial temperature profile, $T(0) \in \mathbf{R}^n$, so that $T(t^{\text{des}}) \approx T^{\text{des}}$. Here, $t^{\text{des}} \in \mathbf{R}$ is a specific time when we want the temperature profile to closely match $T^{\text{des}} \in \mathbf{R}^n$. We also wish to satisfy a constraint that $\|T(0)\|$ should be not be too large.

To formalize these requirements, we use the objective $(1/\sqrt{n})\|T(t^{\text{des}}) - T^{\text{des}}\|$ and the constraint $(1/\sqrt{n})\|T(0)\| \leq T^{\text{max}}$. The first expression is the RMS temperature deviation, at $t = t^{\text{des}}$, from the desired value, and the second is the RMS temperature deviation from ambient at $t = 0$. T^{max} is the (given) maximum initial RMS temperature value.

- (a) Explain how to find $T(0)$ that minimizes the objective while satisfying the constraint.
- (b) Solve the problem instance with the values of n , c , a , b , t^{des} , T^{des} and T^{max} defined in the file `temp_prof_data.m`.
Plot, on one graph, your $T(0)$, $T(t^{\text{des}})$ and T^{des} . Give the RMS temperature error $(1/\sqrt{n})\|T(t^{\text{des}}) - T^{\text{des}}\|$, and the RMS value of initial temperature $(1/\sqrt{n})\|T(0)\|$.

Lecture 8 – Least-norm solutions of underdetermined equations

8.1 *Smallest input that drives a system to a desired steady-state output.* We start with the discrete-time model of the system used in pages 16-19 of lecture 1:

$$x(t+1) = A_d x(t) + B_d u(t), \quad y(t) = C_d x(t), \quad t = 1, 2, \dots,$$

where $A_d \in \mathbf{R}^{16 \times 16}$, $B_d \in \mathbf{R}^{16 \times 2}$, $C_d \in \mathbf{R}^{2 \times 16}$. The system starts from the zero state, *i.e.*, $x(1) = 0$. (We start from initial time $t = 1$ rather than the more conventional $t = 0$ since Matlab indexes vectors starting from 1, not 0.) The data for this problem can be found in `ss_small_input_data.m`.

The goal is to find an input u that results in $y(t) \rightarrow y_{\text{des}} = (1, -2)$ as $t \rightarrow \infty$ (*i.e.*, asymptotic convergence to a desired output) or, even better, an input u that results in $y(t) = y_{\text{des}}$ for $t = T+1, \dots$ (*i.e.*, exact convergence after T steps).

- Steady-state analysis for desired constant output.* Suppose that the system is in steady-state, *i.e.*, $x(t) = x_{\text{ss}}$, $u(t) = u_{\text{ss}}$ and $y(t) = y_{\text{des}}$ are constant (do not depend on t). Find u_{ss} and x_{ss} .
- Simple simulation.* Find $y(t)$, with initial state $x(1) = 0$, with $u(t) = u_{\text{ss}}$, for $t = 1, \dots, 20000$. Plot u and y versus t . If you've done everything right, you should observe that $y(t)$ appears to be converging to y_{des} .

You can use the following Matlab code to obtain plots that look like the ones in lecture 1.

```
figure;
subplot(411); plot(u(1,:));
subplot(412); plot(u(2,:));
subplot(413); plot(y(1,:));
subplot(414); plot(y(2,:));
```

Here we assume that u and y are 2×20000 matrices. There will be two differences between these plots and those in lecture 1: These plots start from $t = 1$, and the plots in lecture 1 scale t by a factor of 0.1.

- Smallest input.* Let $u^*(t)$, for $t = 1, \dots, T$, be the input with minimum RMS value

$$\left(\frac{1}{T} \sum_{t=1}^T \|u(t)\|^2 \right)^{1/2}$$

that yields $x(T+1) = x_{\text{ss}}$ (the value found in part (a)). Note that if $u(t) = u^*(t)$ for $t = 1, \dots, T$, and then $u(t) = u_{\text{ss}}$ for $t = T+1, T+2, \dots$, then $y(t) = y_{\text{des}}$ for $t \geq T+1$. In other words, we have exact convergence to the desired output in T steps.

For the three cases $T = 100$, $T = 200$, and $T = 500$, find u^* and its associated RMS value. For each of these three cases, plot u and y versus t .

- Plot the RMS value of u^* versus T for T between 100 and 1000 (for multiples of 10, if you like). The plot is probably better viewed on a log-log scale, which can be done using the command `loglog` instead of the command `plot`.

Minimum fuel and minimum peak input solutions. Suppose $A \in \mathbf{R}^{m \times n}$ is fat and full rank, so there are many x 's that satisfy $Ax = y$. In lecture we encountered the *least-norm* solution given by $x_{\text{ln}} = A^T(AA^T)^{-1}y$. This solution has the minimum (Euclidean) norm among all solutions of $Ax = y$. In many applications we want to minimize another norm of x (*i.e.*, measure of size of x) subject to $Ax = y$. Two common examples are the 1-norm and ∞ -norm, which are defined as

$$\|x\|_1 = \sum_{i=1}^n |x_i|, \quad \|x\|_\infty = \max_{i=1, \dots, n} |x_i|.$$

The 1-norm, for example, is often a good measure of fuel use; the ∞ -norm is the *peak* of the vector or signal x . There is no simple formula for the least 1-norm or ∞ -norm solution of $Ax = y$, like there is for the least (Euclidean) norm solution. They can be computed very easily, however. (That's one of the topics of EE364.) The analysis is a bit trickier as well, since we can't just differentiate to verify that we have the minimizer. For example, how would you *know* that a solution of $Ax = y$ has minimum 1-norm? In this problem you will explore this idea. First verify the following inequality, which is like the Cauchy-Schwarz inequality (but even easier to prove): for any $v, w \in \mathbf{R}^p$, the following inequality holds: $w^T v \leq \|v\|_\infty \|w\|_1$. From this inequality it follows that whenever $v \neq 0$,

$$\|w\|_1 \geq \frac{w^T v}{\|v\|_\infty}.$$

Now let z be any solution of $Az = y$, and let $\lambda \in \mathbf{R}^m$ be such that $A^T \lambda \neq 0$. Explain why we must have

$$\|z\|_1 \geq \frac{\lambda^T y}{\|A^T \lambda\|_\infty}.$$

Thus, any solution of $Az = y$ must have 1-norm at least as big as the righthand side expression. Therefore if you can find $x_{\text{mf}} \in \mathbf{R}^n$ (mf stands for minimum fuel) and $\lambda \in \mathbf{R}^m$ such that $Ax_{\text{mf}} = y$ and

$$\|x_{\text{mf}}\|_1 = \frac{\lambda^T y}{\|A^T \lambda\|_\infty},$$

then x_{mf} is a minimum fuel solution. (Explain why.) Methods for computing x_{mf} and the mysterious vector λ are described in EE364. In the rest of this problem, you'll use these ideas to verify a statement made during lecture. Now consider the problem from the lecture notes of a unit mass acted on by forces x_1, \dots, x_{10} for one second each. The mass starts at position $p(0) = 0$ with zero velocity and is required to satisfy $p(10) = 1, \dot{p}(10) = 0$. There are, of course, many force vectors that satisfy these requirements. In the lecture notes, you can see a plot of the least (Euclidean) norm force profile. In class I stated that the minimum fuel solution is given by $x_{\text{mf}} = (1/9, 0, \dots, 0, -1/9)$, *i.e.*, an accelerating force at the beginning, 8 seconds of coasting, and a (braking) force at the end to decelerate the mass to zero velocity at $t = 10$. Prove this. *Hint*: try $\lambda = (1, -5)$. Verify that the 1-norm of x_{mf} is less than the 1-norm of x_{ln} , the (Euclidean) least-norm solution. Feel free to use Matlab. There are several convenient ways to find the 1- and ∞ -norm of a vector z , *e.g.*, `norm(z,1)` and `norm(z,inf)` or `sum(abs(z))` and `max(abs(z))`. One last question, for fun: what do you think is the minimum peak force vector x_{mp} ? How would you verify that a vector x_{mp} (mp for minimum peak) is a minimum ∞ -norm solution of $Ax = y$? This input, by the way, is very widely used in practice. It is (basically) the input used in a disk drive to move the head from one track to another, while respecting a maximum possible current in the disk drive motor coil. *Hints*:

- The input is called *bang-bang*.
- Some people drive this way.

8.2 *Simultaneous left inverse of two matrices.* Consider a system where

$$y = Gx, \quad \tilde{y} = \tilde{G}x$$

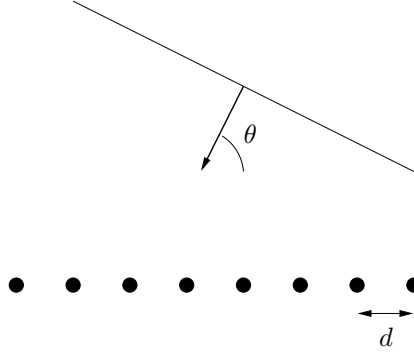
where $G \in \mathbf{R}^{m \times n}$, $\tilde{G} \in \mathbf{R}^{m \times n}$. Here x is some variable we wish to estimate or find, y gives the measurements with some set of (linear) sensors, and \tilde{y} gives the measurements with some *alternate* set of (linear) sensors. We want to find a *reconstruction matrix* $H \in \mathbf{R}^{n \times m}$ such that $HG = H\tilde{G} = I$. Such a reconstruction matrix has the nice property that it recovers x perfectly from *either* set of

measurements (y or \tilde{y}), *i.e.*, $x = Hy = H\tilde{y}$. Consider the specific case

$$G = \begin{bmatrix} 2 & 3 \\ 1 & 0 \\ 0 & 4 \\ 1 & 1 \\ -1 & 2 \end{bmatrix}, \quad \tilde{G} = \begin{bmatrix} -3 & -1 \\ -1 & 0 \\ 2 & -3 \\ -1 & -3 \\ 1 & 2 \end{bmatrix}.$$

Either find an explicit reconstruction matrix H , or explain why there is no such H .

8.3 *Phased-array antenna weight design.* We consider the phased-array antenna system shown below.



The array consists of n individual antennas (called *antenna elements*) spaced on a line, with spacing d between elements. A sinusoidal plane wave, with wavelength λ and angle of arrival θ , impinges on the array, which yields the output $e^{2\pi j(k-1)(d/\lambda)\cos\theta}$ (which is a complex number) from the k th element. (We've chosen the *phase center* as element 1, *i.e.*, the output of element 1 does not depend on the incidence angle θ .) A (complex) linear combination of these outputs is formed, and called the *combined array output*,

$$y(\theta) = \sum_{k=1}^n w_k e^{2\pi j(k-1)(d/\lambda)\cos\theta}.$$

The complex numbers w_1, \dots, w_n , which are the coefficients of the linear combination, are called the *antenna weights*. We can choose, *i.e.*, design, the weights. The combined array output depends on the angle of arrival of the wave. The function $|y(\theta)|$, for $0^\circ \leq \theta \leq 180^\circ$, is called the *antenna array gain pattern*. By choosing the weights w_1, \dots, w_n intelligently, we can shape the gain pattern to satisfy some specifications. As a simple example, if we choose the weights as $w_1 = 1, w_2 = \dots = w_n = 0$, then we get a uniform or omnidirectional gain pattern; $|y(\theta)| = 1$ for all θ . In this problem, we want a gain pattern that is one in a given ('target') direction θ_{target} , but small at other angles. Such a pattern would receive a signal coming from the direction θ_{target} , and attenuate signals (*e.g.*, 'jammers' or multipath reflections) coming from other directions. Here's the problem. You will design the weights for an array with $n = 20$ elements, and a spacing $d = 0.4\lambda$ (which is a typical value). We want $y(70^\circ) = 1$, and we want $|y(\theta)|$ small for $0^\circ \leq \theta \leq 60^\circ$ and $80^\circ \leq \theta \leq 180^\circ$. In other words, we want the antenna array to be relatively insensitive to plane waves arriving from angles more than 10° away from the target direction. (In the language of antenna arrays, we want a beamwidth of 20° around a target direction of 70° .) To solve this problem, you will first discretize the angles between 0° and 180° in 1° increments. Thus $y \in \mathbf{C}^{180}$ will be a (complex) vector, with y_k equal to $y(k^\circ)$, *i.e.*, $y(\pi k/180)$, for $k = 1, \dots, 180$. You are to choose $w \in \mathbf{C}^{20}$ that minimizes

$$\sum_{k=1}^{60} |y_k|^2 + \sum_{k=80}^{180} |y_k|^2$$

subject to the constraint $y_{70} = 1$. As usual, you must explain how you solve the problem. Give the weights you find, and also a plot of the antenna array response, *i.e.*, $|y_k|$, versus k (which, hopefully, will achieve the desired goal of being relatively insensitive to plane waves arriving at angles more than 10° from $\theta = 70^\circ$). *Hints:*

- You'll probably want to rewrite the problem as one involving real variables (*i.e.*, the real and imaginary parts of the antenna weights), and real matrices. You can then rewrite your solution in a more compact formula that uses complex matrices and vectors (if you like).
- **Very important:** in Matlab, the prime is actually the Hermitian conjugate operator. In other words, if A is a complex matrix or vector, A' gives the conjugate transpose, or Hermitian conjugate, of A .
- Although we don't require you to, you might find it fun to also plot your antenna gain pattern on a polar plot, which allows you to easily visualize the pattern. In Matlab, this is done using the `polar` command.

8.4 *Modifying measurements to satisfy known conservation laws.* A vector $y \in \mathbf{R}^n$ contains n measurements of some physical quantities $x \in \mathbf{R}^n$. The measurements are good, but not perfect, so we have $y \approx x$. From physical principles it is known that the quantities x must satisfy some linear equations, *i.e.*,

$$a_i^T x = b_i, \quad i = 1, \dots, m,$$

where $m < n$. As a simple example, if x_1 is the current in a circuit flowing into a node, and x_2 and x_3 are the currents flowing out of the node, then we must have $x_1 = x_2 + x_3$. More generally, the linear equations might come from various conservation laws, or balance equations (mass, heat, energy, charge ...). The vectors a_i and the constants b_i are known, and we assume that a_1, \dots, a_m are independent. Due to measurement errors, the measurement y won't satisfy the conservation laws (*i.e.*, linear equations above) exactly, although we would expect $a_i^T y \approx b_i$. An engineer proposes to adjust the measurements y by adding a correction term $c \in \mathbf{R}^n$, to get an adjusted estimate of x , given by

$$y_{\text{adj}} = y + c.$$

She proposes to find the smallest possible correction term (measured by $\|c\|$) such that the adjusted measurements y_{adj} satisfy the known conservation laws. Give an explicit formula for the correction term, in terms of y , a_i , b_i . If any matrix inverses appear in your formula, explain why the matrix to be inverted is nonsingular. Verify that the resulting adjusted measurement satisfies the conservation laws, *i.e.*, $a_i^T y_{\text{adj}} = b_i$.

8.5 *Estimator insensitive to certain measurement errors.* We consider the usual measurement setup: $y = Ax + v$, where

- $y \in \mathbf{R}^m$ is the vector of measurements
- $x \in \mathbf{R}^n$ is the vector of parameters we wish to estimate
- $v \in \mathbf{R}^m$ is the vector of measurement errors
- $A \in \mathbf{R}^{m \times n}$ is the coefficient matrix relating the parameters to the measurements

You can assume that $m > n$, and A is full rank. In this problem we assume that the measurement errors lie in the subspace

$$\mathcal{V} = \text{span}\{f_1, \dots, f_k\},$$

where $f_1, \dots, f_k \in \mathbf{R}^m$ are given, known vectors. Now consider a linear estimator of the form $\hat{x} = By$. Recall that the estimator is called *unbiased* if whenever $v = 0$, we have $\hat{x} = x$, for any $x \in \mathbf{R}^n$. In other words, an unbiased estimator predicts x perfectly when there is no measurement error. In this problem we consider the stronger condition that the estimator predicts x perfectly, for any measurement error

in \mathcal{V} . In other words, we have $\hat{x} = x$, for any $x \in \mathbf{R}^n$, and any $v \in \mathcal{V}$. If this condition holds, we say that the estimator is insensitive to measurement errors in \mathcal{V} . (Note that this condition is a stronger condition than the estimator being unbiased.)

- (a) Show that if $\mathcal{R}(A) \cap \mathcal{V} \neq \{0\}$, then there is no estimator insensitive to measurement errors in \mathcal{V} .
- (b) Now we consider a specific example, with

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & -1 \\ 2 & 1 \\ -1 & 2 \end{bmatrix}, \quad f_1 = \begin{bmatrix} 1 \\ 2 \\ -1 \\ 1 \\ 0 \end{bmatrix}, \quad f_2 = \begin{bmatrix} 3 \\ 3 \\ 2 \\ 2 \\ 1 \end{bmatrix}.$$

Either construct a specific $B \in \mathbf{R}^{2 \times 5}$ for which the linear estimator $\hat{x} = By$ is insensitive to measurement errors in \mathcal{V} , or explain in detail why none exists. If you find such a B , you must explain how you found it, and verify (say, in Matlab) that it satisfies the required properties. (We'll be really annoyed if you just give a matrix and leave the verification to us!)

8.6 Optimal flow on a data collection network. We consider a communications network with m nodes, plus a special destination node, and n communication links. Each communication link connects two (distinct) nodes and is bidirectional, *i.e.*, information can flow in either direction. We will assume that the network is connected, *i.e.*, there is a path, or sequence of links, from every node (including the special destination node) to every other node. With each communication link we associate a directed arc, which defines the direction of information flow that we will call positive. Using these reference directions, the flow or traffic on link j is denoted f_j . (The units are bits per second, but that won't matter to us.) The traffic on the network (*i.e.*, the flow in each communication link) is given by a vector $f \in \mathbf{R}^n$. A small example is shown in part 2 of this problem. In this example, nodes 1 and 3 are connected by communication link 4, and the associated arc points from node 1 to node 3. Thus $f_4 = 12$ means the flow on that link is 12 (bits per second), from node 1 to node 3. Similarly, $f_4 = -3$ means the flow on link 4 is 3 (bits per second), from node 3 to node 1. External information enters each of the m regular nodes and flows across links to the special destination node. In other words, the network is used to collect information from the nodes and route it through the links to the special destination node. (That explains why we call it a data collection network.) At node i , an external information flow s_i (which is nonnegative) enters. The vector $s \in \mathbf{R}^m$ of external flows is sometimes called the *source vector*. Information flow is conserved. This means that at each node (except the special destination node) the sum of all flows entering the node from communication links connected to that node, plus the external flow, equals the sum of the flows leaving that node on communication links. As an example, consider node 3 in the network of part 2. Links 4 and 5 enter this node, and link 6 leaves the node. Therefore, flow conservation at node 3 is given by

$$f_4 + f_5 + s_3 = f_6.$$

The first two terms on the left give the flow entering node 3 on links 4 and 5; the last term on the left gives the external flow entering node 3. The term on the righthand side gives the flow leaving over link 6. Note that this equation correctly expresses flow conservation regardless of the signs of f_4 , f_5 , and f_6 . Finally, here is the problem.

- (a) The vector of external flows, $s \in \mathbf{R}^m$, and the network topology, are given, and you must find the flow f that satisfies the conservation equations, and minimizes the mean-square traffic on the network, *i.e.*,

$$\frac{1}{n} \sum_{j=1}^n f_j^2.$$

Your answer should be in terms of the external flow s , and the *node incidence matrix* $A \in \mathbf{R}^{m \times n}$ that describes the network topology. The node incidence matrix is defined as

$$A_{ij} = \begin{cases} 1 & \text{arc } j \text{ enters (or points into) node } i \\ -1 & \text{arc } j \text{ leaves (or points out of) node } i \\ 0 & \text{otherwise.} \end{cases}$$

Note that each row of A is associated with a node on the network (not including the destination node), and each column is associated with an arc or link.

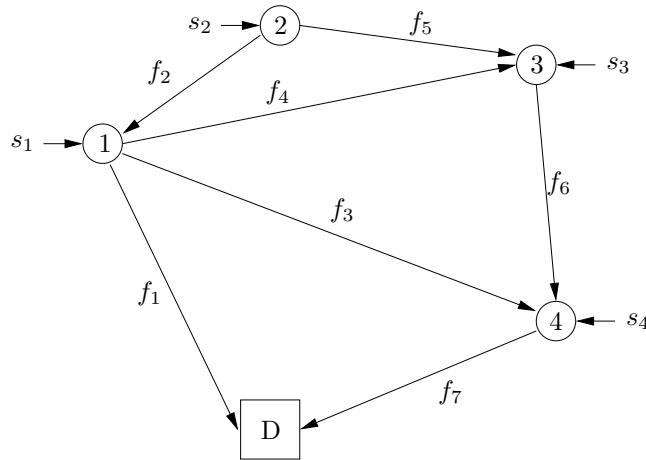
- (b) Now consider the specific (and very small) network shown below. The nodes are shown as circles, and the special destination node is shown as a square. The external flows are

$$s = \begin{bmatrix} 1 \\ 4 \\ 10 \\ 10 \end{bmatrix}.$$

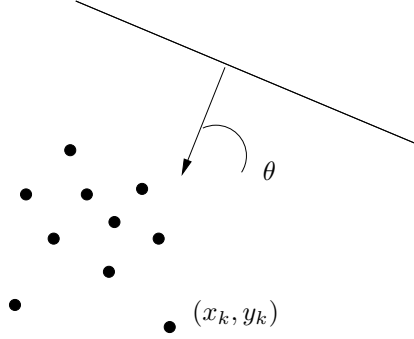
One simple feasible flow is obtained by routing all the external flow entering each node along a shortest path to the destination. For example, all the external flow entering node 2 goes to node 1, then to the destination node. For node 3, which has two shortest paths to the destination, we arbitrarily choose the path through node 4. This simple routing scheme results in the feasible flow

$$f_{\text{simple}} = \begin{bmatrix} 5 \\ 4 \\ 0 \\ 0 \\ 0 \\ 10 \\ 20 \end{bmatrix}.$$

Find the mean square optimal flow for this problem (as in part 1). Compare the mean square flow of the optimal flow with the mean square flow of f_{simple} .



8.7 *Random geometry antenna weight design.* We consider the phased-array antenna system shown below.



The array consists of n individual antennas (called *antenna elements*) randomly placed on $2d$ space, with the coordinates of the k th element being x_k and y_k . A sinusoidal plane wave, with wavelength λ and angle of arrival θ , impinges on the array, which yields the output $e^{j\frac{2\pi}{\lambda}(x_k \cos \theta + y_k \sin \theta)}$ (which is a complex number) from the k th element. A (complex) linear combination of these outputs is formed, and called the *combined array output*,

$$r(\theta) = \sum_{k=1}^n w_k e^{j\frac{2\pi}{\lambda}(x_k \cos \theta + y_k \sin \theta)}.$$

The complex numbers w_1, \dots, w_n , which are the coefficients of the linear combination, are called the *antenna weights*. We can choose, *i.e.*, design, the weights. The combined array output depends on the angle of arrival of the wave. The function $|r(\theta)|$, for $0^\circ \leq \theta < 360^\circ$, is called the antenna array *gain pattern*. By choosing the weights w_1, \dots, w_n intelligently, we can shape the gain pattern to satisfy some specifications. In this problem, we want a gain pattern that is one in a given ('target') direction θ_{target} , but small at other angles. Such a pattern would receive a signal coming from the direction θ_{target} , and attenuate signals (*e.g.*, 'jammers' or multipath reflections) coming from other directions. Here's the problem. Design the weights for the antenna array, whose elements have coordinates given in the file `antenna_geom.mat`. We want $r(70^\circ) = 1$, and we want $|r(\theta)|$ small for $0^\circ \leq \theta \leq 60^\circ$ and $80^\circ \leq \theta < 360^\circ$. In other words, we want the antenna array to be relatively insensitive to plane waves arriving from angles more than 10° away from the target direction. (In the language of antenna arrays, we want a beamwidth of 20° around a target direction of 70° .) You are told that $\lambda = 1$. To solve this problem, you will first discretize the angles between 1° and 360° in 1° increments. Thus $r \in \mathbb{C}^{360}$ will be a (complex) vector, with r_k equal to $r(k^\circ)$, *i.e.*, $r(\pi k/180)$, for $k = 1, \dots, 360$. You are to choose $w \in \mathbb{C}^n$ that minimizes

$$\sum_{k=1}^{60} |r_k|^2 + \sum_{k=80}^{360} |r_k|^2$$

subject to the constraint $r_{70} = 1$. As usual, you must explain how you solve the problem. Give the weights you find, and also a plot of the antenna array response, *i.e.*, $|r_k|$, versus k (which, hopefully, will achieve the desired goal of being relatively insensitive to plane waves arriving at angles more than 10° from $\theta = 70^\circ$). *Hints:*

- You'll probably want to rewrite the problem as one involving real variables (*i.e.*, the real and imaginary parts of the antenna weights), and real matrices. You can then rewrite your solution in a more compact formula that uses complex matrices and vectors (if you like).
- **Very important:** in Matlab, the prime is actually the Hermitian conjugate operator. In other words, if A is a complex matrix or vector, A' gives the conjugate transpose, or Hermitian conjugate, of A .
- Although we don't require you to, you might find it fun to also plot your antenna gain pattern on a polar plot, which allows you to easily visualize the pattern. In Matlab, this is done using the `polar` command.

8.8 *Estimation with known input norm.* We consider a standard estimation setup: $y = Ax + v$, where $A \in \mathbf{R}^{m \times n}$ is a full rank, skinny matrix, $x \in \mathbf{R}^n$ is the vector we wish to estimate, $v \in \mathbf{R}^m$ is an unknown noise vector, and $y \in \mathbf{R}^m$ is the measurement vector. As usual, we assume that smaller values of $\|v\|$ are more plausible than larger values. In this problem, we add one more piece of prior information: we know that $\|x\| = 1$. (In other words, the vector we are estimating is known ahead of time to have norm one.) This might occur in a communications system, where the transmitted signal power is known to be equal to one. (You may assume that the norm of the least-squares approximate solution exceeds one, *i.e.*, $\|(A^T A)^{-1} A^T y\| > 1$.)

- (a) Explain clearly how would you find the best estimate of x , taking into account the prior information $\|x\| = 1$. Explain how you would compute your estimate \hat{x} , given A and y . Is your estimate \hat{x} a linear function of y ?
- (b) On the EE263 webpage, you will find the file `mtprob4.m`, which gives the matrix A and the observed vector y . Carry out the estimation procedure you developed in part (a). Give your estimate \hat{x} , and verify that it satisfies $\|\hat{x}\| = 1$. Give the Matlab source you use to compute \hat{x} .

8.9 *Minimum energy rendezvous.* The dynamics of two vehicles, at sampling times $t = 0, 1, 2, \dots$, are given by

$$x(t+1) = Ax(t) + bu(t), \quad z(t+1) = Fz(t) + gv(t)$$

where

- $x(t) \in \mathbf{R}^n$ is the state of vehicle 1
- $z(t) \in \mathbf{R}^n$ is the state of vehicle 2
- $u(t) \in \mathbf{R}$ is the (scalar) input to vehicle 1
- $v(t) \in \mathbf{R}$ is the (scalar) input to vehicle 2

The initial states of the two vehicles are fixed and given:

$$x(0) = x_0, \quad z(0) = z_0.$$

We are interested in finding inputs for the two vehicles over the time interval $t = 0, 1, \dots, N-1$ so that they *rendezvous* at state $w \in \mathbf{R}^n$ at time $t = N$, *i.e.*, $x(N) = w$, $z(N) = w$. (The point $w \in \mathbf{R}^n$ is called the *rendezvous point*.) You can select the inputs to the two vehicles,

$$u(0), u(1), \dots, u(N-1), \quad v(0), v(1), \dots, v(N-1),$$

as well as the rendezvous point $w \in \mathbf{R}^n$. Among choices of u , v , and w that satisfy the rendezvous condition, we want the one that minimizes the total input energy defined as

$$E = \sum_{t=0}^{N-1} u(t)^2 + \sum_{t=0}^{N-1} v(t)^2.$$

Give explicit formulas for the optimal u , v , and w in terms of the problem data, *i.e.*, A , b , F , g , x_0 , and z_0 . If you need to assume that one or more matrices that arise in your solution are invertible, full rank, etc., that's fine, but be sure to make very clear what you're assuming.

8.10 *Least-norm solution of nonlinear equations.* Suppose $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ is a function, and $y \in \mathbf{R}^m$ is a vector, where $m < n$ (*i.e.*, x has larger dimension than y). We say that $x \in \mathbf{R}^n$ is a least-norm solution of $f(x) = y$ if for any $z \in \mathbf{R}^n$ that satisfies $f(z) = y$, we have $\|z\| \geq \|x\|$. When the function f is linear or affine (*i.e.*, linear plus a constant), the equations $f(x) = y$ are linear, and we know how to find the least-norm solution for such problems. In general, however, it is an extremely difficult problem to compute a least-norm solution to a set of nonlinear equations. There are, however, some good heuristic iterative methods that work well when the function f is not too far from affine, *i.e.*, its nonlinear terms are small compared to its linear and constant part. You may assume that you have a starting guess, which we'll call $x^{(0)}$. This guess doesn't necessarily satisfy the equations $f(x) = y$.

- (a) Suggest an iterative method for (approximately) solving the nonlinear least-norm problem, starting from the initial guess $x^{(0)}$. Use the notation $x^{(k)}$ to denote the k th iteration of your method. Explain clearly how you obtain $x^{(k+1)}$ from $x^{(k)}$. If you need to make any assumptions about rank of some matrix, do so. (You don't have to worry about what happens if the matrix is not full rank.) Your method should have the property that $f(x^{(k)})$ converges to y as k increases. (In particular, we don't need to have the iterates satisfy the nonlinear equations exactly.) Suggest a name for the method you invent. Your method should not be complicated or require a long explanation. You do not have to prove that the method converges, or that when it converges, it converges to a least-norm solution. All you have to do is suggest a sensible, simple method that ought to work well when f is not too nonlinear, and the starting guess $x^{(0)}$ is good.
- (b) Now we consider a specific example, with the function $f : \mathbf{R}^5 \rightarrow \mathbf{R}^2$ given by

$$\begin{aligned} f_1(x) &= 2x_1 - 3x_3 + x_5 + 0.1x_1x_2 - 0.5x_2x_5, \\ f_2(x) &= -x_2 + x_3 - x_4 + x_5 - 0.6x_1x_4 + 0.3x_3x_4. \end{aligned}$$

Note that each component of f consists of a linear part, and also a quadratic part. Use the method you invented in part a to find the least-norm solution of

$$f(x) = y = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

(We repeat that you do not have to prove that the solution you found is really the least-norm one.) As initial guess, you can use the least-norm solution of the linear equations resulting if you ignore the quadratic terms in f . Make sure to turn in your Matlab code as well as to identify the least-norm x you find, its norm, and the equation residual, *i.e.*, $f(x) - y$ (which should be very small).

- 8.11 *The smoothest input that takes the state to zero.* We consider the discrete-time linear dynamical system $x(t+1) = Ax(t) + Bu(t)$, with

$$A = \begin{bmatrix} 1.0 & 0.5 & 0.25 \\ 0.25 & 0 & 1.0 \\ 1.0 & -0.5 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1.0 \\ 0.1 \\ 0.5 \end{bmatrix}, \quad x(0) = \begin{bmatrix} 25 \\ 0 \\ -25 \end{bmatrix}.$$

The goal is to choose an input sequence $u(0), u(1), \dots, u(19)$ that yields $x(20) = 0$. Among the input sequences that yield $x(20) = 0$, we want the one that is *smoothest*, *i.e.*, that minimizes

$$J_{\text{smooth}} = \left(\frac{1}{20} \sum_{t=0}^{19} (u(t) - u(t-1))^2 \right)^{1/2},$$

where we take $u(-1) = 0$ in this formula. Explain how to solve this problem. Plot the smoothest input u_{smooth} , and give the associated value of J_{smooth} .

- 8.12 *Minimum energy input with way-point constraints.* We consider a vehicle that moves in \mathbf{R}^2 due to an applied force input. We will use a discrete-time model, with time index $k = 1, 2, \dots$; time index k corresponds to time $t = kh$, where $h > 0$ is the sample interval. The position at time index k is denoted by $p(k) \in \mathbf{R}^2$, and the velocity by $v(k) \in \mathbf{R}^2$, for $k = 1, \dots, K+1$. These are related by the equations

$$p(k+1) = p(k) + hv(k), \quad v(k+1) = (1 - \alpha)v(k) + (h/m)f(k), \quad k = 1, \dots, K,$$

where $f(k) \in \mathbf{R}^2$ is the force applied to the vehicle at time index k , $m > 0$ is the vehicle mass, and $\alpha \in (0, 1)$ models drag on the vehicle: In the absence of any other force, the vehicle velocity decreases

by the factor $1 - \alpha$ in each time index. (These formulas are approximations of more accurate formulas that we will see soon, but for the purposes of this problem, we consider them exact.) The vehicle starts at the origin, at rest, *i.e.*, we have $p(1) = 0$, $v(1) = 0$. (We take $k = 1$ as the initial time, to simplify indexing.)

The problem is to find forces $f(1), \dots, f(K) \in \mathbf{R}^2$ that minimize the cost function

$$J = \sum_{k=1}^K \|f(k)\|^2,$$

subject to *way-point constraints*

$$p(k_i) = w_i, \quad i = 1, \dots, M,$$

where k_i are integers between 1 and K . (These state that at the time $t_i = hk_i$, the vehicle must pass through the location $w_i \in \mathbf{R}^2$.) Note that there is no requirement on the vehicle velocity at the way-points.

- Explain how to solve this problem, given all the problem data (*i.e.*, h , α , m , K , the way-points w_1, \dots, w_M , and the way-point indices k_1, \dots, k_M).
- Carry out your method on the specific problem instance with data $h = 0.1$, $m = 1$, $\alpha = 0.1$, $K = 100$, and the $M = 4$ way-points

$$w_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \quad w_2 = \begin{bmatrix} -2 \\ 3 \end{bmatrix}, \quad w_3 = \begin{bmatrix} 4 \\ -3 \end{bmatrix}, \quad w_4 = \begin{bmatrix} -4 \\ -2 \end{bmatrix},$$

with way-point indices $k_1 = 10$, $k_2 = 30$, $k_3 = 40$, and $k_4 = 80$.

Give the optimal value of J .

Plot $f_1(k)$ and $f_2(k)$ versus k , using

```
subplot(211); plot(f(1,:));
subplot(212); plot(f(2,:));
```

We assume here that \mathbf{f} is a $2 \times K$ matrix, with columns $f(1), \dots, f(K)$.

Plot the vehicle trajectory, using `plot(p(1,:), p(2,:))`. Here \mathbf{p} is a $2 \times (K + 1)$ matrix with columns $p(1), \dots, p(K + 1)$.

- 8.13 In this problem you will show that, for *any* matrix A , and any positive number μ , the matrices $A^T A + \mu I$ and $AA^T + \mu I$ are both invertible, and

$$(A^T A + \mu I)^{-1} A^T = A^T (AA^T + \mu I)^{-1}.$$

- Let's first show that $A^T A + \mu I$ is invertible, assuming $\mu > 0$. (The same argument, with A^T substituted for A , will show that $AA^T + \mu I$ is invertible.) Suppose that $(A^T A + \mu I)z = 0$. Multiply on the left by z^T , and argue that $z = 0$. This is what we needed to show. (Your job is to fill all details of the argument.)
- Now let's establish the identity above. First, explain why

$$A^T (AA^T + \mu I) = (A^T A + \mu I) A^T$$

holds. Then, multiply on the left by $(A^T A + \mu I)^{-1}$, and on the right by $(AA^T + \mu I)^{-1}$. (These inverses exist, by part (a).)

- Now assume that A is fat and full rank. Show that as μ tends to zero from above (*i.e.*, μ is positive) we have

$$(A^T A + \mu I)^{-1} A^T \rightarrow A^T (AA^T)^{-1}.$$

(This is asserted, but not shown, in the lecture notes on page 8-12.)

8.14 *Singularity of the KKT matrix.* This problem concerns the general norm minimization with equality constraints problem (described in the lectures notes on pages 8-13),

$$\begin{array}{ll} \text{minimize} & \|Ax - b\| \\ \text{subject to} & Cx = d, \end{array}$$

where the variable is $x \in \mathbf{R}^n$, $A \in \mathbf{R}^{m \times n}$, and $C \in \mathbf{R}^{k \times n}$. We assume that C is fat ($k \leq n$), i.e., the number of equality constraints is no more than the number of variables.

Using Lagrange multipliers, we found that the solution can be obtained by solving the linear equations

$$\begin{bmatrix} A^T A & C^T \\ C & 0 \end{bmatrix} \begin{bmatrix} x \\ \lambda \end{bmatrix} = \begin{bmatrix} A^T b \\ d \end{bmatrix}$$

for x and λ . (The vector x gives the solution of the norm minimization problem above.) The matrix above, which we will call $K \in \mathbf{R}^{(n+k) \times (n+k)}$, is called the *KKT matrix* for the problem. (KKT are the initials of some of the people who came up with the optimality conditions for a more general type of problem.)

One question that arises is, when is the KKT matrix K nonsingular? The answer is: K is nonsingular if and only if C is full rank and $\mathcal{N}(A) \cap \mathcal{N}(C) = \{0\}$.

You will fill in all details of the argument below.

- (a) Suppose C is not full rank. Show that K is singular.
- (b) Suppose that there is a nonzero $u \in \mathcal{N}(A) \cap \mathcal{N}(C)$. Use this u to show that K is singular.
- (c) Suppose that K is singular, so there exists a nonzero vector $[u^T \ v^T]^T$ for which

$$\begin{bmatrix} A^T A & C^T \\ C & 0 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = 0.$$

Write this out as two block equations, $A^T A u + C^T v = 0$ and $C u = 0$. Conclude that $u \in \mathcal{N}(C)$. Multiply $A^T A u + C^T v = 0$ on the left by u^T , and use $C u = 0$ to conclude that $\|A u\| = 0$, which implies $u \in \mathcal{N}(A)$. Finish the argument that leads to the conclusion that either C is not full rank, or $\mathcal{N}(A) \cap \mathcal{N}(C) \neq \{0\}$.

8.15 *Minimum energy roundtrip.* We consider the linear dynamical system

$$x(t+1) = Ax(t) + Bu(t), \quad x(0) = 0,$$

with $u(t) \in \mathbf{R}$ and $x(t) \in \mathbf{R}^n$. We must choose $u(0), u(1), \dots, u(T-1)$ so that $x(T) = 0$ (i.e., after T steps we are back at the zero state), and $x(t_{\text{dest}}) = x_{\text{dest}}$ (i.e., at time t_{dest} the state is equal to x_{dest}). Here $x_{\text{dest}} \in \mathbf{R}^n$ is a given destination state. The time t_{dest} is *not* given; it can be any integer between 1 and $T-1$. The goal is to minimize the total input energy, defined as

$$E = \sum_{t=0}^{T-1} u(t)^2.$$

Note that you have to find t_{dest} , the time when the state hits the desired state, as well as the input trajectory $u(0), \dots, u(T-1)$.

- (a) Explain how to do this. For this problem, you may assume that $n \leq t_{\text{dest}} \leq T - n$. If you need some matrix or matrices that arise in your analysis to be full rank, you can just assume they are. But you must state this clearly.

(b) Carry out your method on the particular problem instance with data

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad T = 30, \quad x_{\text{dest}} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}.$$

Give the optimal value of t_{dest} and the associated value of E , and plot the optimal input trajectory u .

8.16 *Optimal dynamic purchasing.* You are to complete a large order to buy a certain number, B , of shares in some company. You are to do this over T time periods. (Depending on the circumstances, a single time period could be between tens of milliseconds and minutes.) We will let b_t denote the number of shares bought in time period t , for $t = 1, \dots, T$, so we have $b_1 + \dots + b_T = B$. (The quantities B, b_1, \dots, b_T can all be any real number; $b_t < 0$, for example, means we *sold* shares in the period t . We also don't require b_t to be integers.) We let p_t denote the price per share in period t , so the total cost of purchasing the B shares is $C = p_1 b_1 + \dots + p_T b_T$.

The amounts we purchase are large enough to have a noticeable effect on the price of the shares. The prices change according to the following equations:

$$p_1 = \bar{p} + \alpha b_1, \quad p_t = \theta p_{t-1} + (1 - \theta)\bar{p} + \alpha b_t, \quad t = 2, \dots, T.$$

Here \bar{p} is the base price of the shares and α and θ are parameters that determine how purchases affect the prices. The parameter α , which is positive, tells us how much the price goes up in the current period when we buy one share. The parameter θ , which lies between 0 and 1, measures the *memory*: If $\theta = 0$ the share price has no memory, and the purchase made in period t only affects the price in that period; if θ is 0.5 (say), the effect a purchase has on the price decays by a factor of two between periods. If $\theta = 1$, the price has perfect memory and the price change will persist for all future periods.

If purchases didn't increase the price, the cost of purchasing the shares would always be $\bar{p}B$. The difference between the total cost and this cost, $C - \bar{p}B$, is called the *transaction cost*.

Find the purchase quantities b_1, \dots, b_T that minimize the transaction cost $C - \bar{p}B$, for the particular problem instance with

$$B = 10000, \quad T = 10, \quad \bar{p} = 10, \quad \theta = 0.8, \quad \alpha = 0.00015.$$

Give the optimal transaction cost. Also give the transaction cost if all the shares were purchased in the first period, and the transaction cost if the purchases were evenly spread over the periods (*i.e.*, if 1000 shares were purchased in each period). Compare these three quantities.

You must explain your method clearly, using any concepts from this class, such as least-squares, pseudo-inverses, eigenvalues, singular values, etc. If your method requires that some rank or other conditions to hold, say so. You must also check, in your Matlab code, that these conditions are satisfied for the given problem instance.

Lecture 9 – Autonomous linear dynamical systems

9.1 *A simple population model.* We consider a certain population of fish (say) each (yearly) season. $x(t) \in \mathbf{R}^3$ will describe the population of fish at year $t \in \mathbf{Z}$, as follows:

- $x_1(t)$ denotes the number of fish less than one year old
- $x_2(t)$ denotes the number of fish between one and two years old
- $x_3(t)$ denotes the number of fish between two and three years

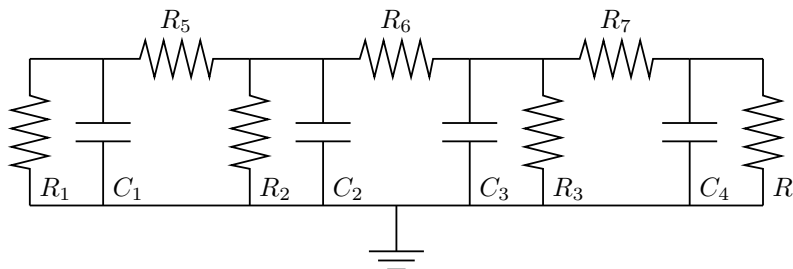
(We will ignore the fact that these numbers are integers.) The population evolves from year t to year $t + 1$ as follows.

- The number of fish less than one year old in the next year ($t + 1$) is equal to the total number of offspring born during the current year. Fish that are less than one year old in the current year (t) bear no offspring. Fish that are between one and two years old in the current year (t) bear an average of 2 offspring each. Fish that are between two and three years old in the current year (t) bear an average of 1 offspring each.
- 40% of the fish less than one year old in the current year (t) die; the remaining 60% live on to be between one and two years old in the next year ($t + 1$). fish, and 50% of the old fish.
- 30% of the one-to-two year old fish in the current year die, and 70% live on to be two-to-three year old fish in the next year.
- All of the two-to-three year old fish in the current year die.

Express the population dynamics as an autonomous linear system with state $x(t)$, *i.e.*, in the form $x(t + 1) = Ax(t)$. **Remark:** this example is silly, but more sophisticated population dynamics models are very useful and widely used.

9.2 *Tridiagonal systems.* A square matrix A is called tridiagonal if $A_{ij} = 0$ whenever $|i - j| > 1$. Tridiagonal matrices arise in many applications.

- Draw a pretty block diagram of $\dot{x} = Ax$, where $A \in \mathbf{R}^{4 \times 4}$ is tridiagonal.
- Consider a Markov chain with four states labeled 1,2,3,4. Let $z(k)$ denote the state at time k . The state transition probabilities are described as follows: when z is not 4, it increases by one with probability 0.3; when z is not 1, it decreases by one with probability 0.2. (If z neither increases nor decreases, it stays the same, *i.e.*, $z(k + 1) = z(k)$). Draw a graph of this Markov chain as in the lecture notes. Give the discrete time linear system equations that govern the evolution of the state distribution.
- Find the linear dynamical system description for the circuit shown below. Use state $x = [v_1 \ v_2 \ v_3 \ v_4]^T$, where v_i is the voltage across the capacitor C_i .



9.3 *A distributed congestion control scheme.* A data network is modeled as a set of l directed links that connect n nodes. There are p routes in the network, which is a path from a *source node*, along one or more links in the network, to the *destination node*. The routes are determined and known. Each route

has a *source rate* (in, say, bits per second). We denote the source rate for route j at time t as $x_j(t)$, $t = 0, 1, 2, \dots$ (We assume the system operates in discrete time.) The total *traffic* on a link is the sum of the source rates for the routes that pass through it. We use $T_i(t)$ to denote the total traffic on link i at time t , for $i = 1, \dots, l$. Each link has a *target traffic level*, which we denote T_i^{target} , $i = 1, \dots, l$. We define the *congestion* on link i as $T_i(t) - T_i^{\text{target}}$, $i = 1, \dots, l$. The congestion is positive if the traffic exceeds the target rate, and negative if it is below the target rate. The goal in *congestion control* is to adjust the source rates in such a way that the traffic levels converge to the target levels if possible, or close to the target levels otherwise. In this problem we consider a very simple congestion control protocol. Each route monitors the congestion for the links along its route. It then adjusts its source rate proportional to the sum of the congestion along its route. This can be expressed as:

$$x_j(t+1) = x_j(t) - \alpha (\text{sum of congestion along route } j), \quad j = 1, \dots, p,$$

where α is a positive scalar that determines how aggressively the source rates react to congestion. Note that this congestion control method is *distributed*; each source only needs to know the congestion along its own route, and does not directly coordinate its adjustments with the other routes. In real congestion control, the rates and traffic are nonnegative, and the traffic on each link must be below a maximum allowed level called the link capacity. In this problem, however, we ignore these effects; we do not take into account the link capacities, and allow the source rates and total traffic levels to become negative. Before we get to the questions, we define a matrix that may be useful. The *route-link matrix* $R \in \mathbf{R}^{l \times p}$, is defined as

$$R_{ij} = \begin{cases} 1 & \text{route } j \text{ utilizes link } i \\ 0 & \text{otherwise.} \end{cases}$$

- Show that $x(t)$, the vector of source rates, can be expressed as a linear dynamical system with constant input, *i.e.*, we have $x(t+1) = Ax(t) + b$. Be as explicit as you can about what A and b are. Try to use the simplest notation you can. *Hint*: use the matrix R .
- Simulate the congestion control algorithm for the network shown in figure 3, from two different initial source rates, using algorithm parameter $\alpha = 0.1$, and all target traffic levels equal to one. Plot the traffic level $T_i(t)$ for each link (on the same plot) versus t , for each of the initial source rates. (You are welcome to simulate the system from more than two initial source rates; we only ask you to hand in the plots for two, however.) Make a brief comment on the results.
- Now we come back to the general case (and *not* just the specific example from part (b)). Assume the congestion control update (*i.e.*, the linear dynamical system found in part (a)) has a unique equilibrium point \bar{x} , and that the rate $x(t)$ converges to it as $t \rightarrow \infty$. What can you say about \bar{x} ? Limit yourself to a few sentences. Does the rate \bar{x} always correspond to zero congestion on every link? Is it optimal in any way?

9.4 *Sequence transmission with constraints.* A communication system is based on 3 symbols: 1, 0, and -1. For this communication system, a valid sequence of symbols x_1, x_2, \dots, x_k , must satisfy several constraints:

- Transition constraint*: Consecutive symbols cannot be more than one apart from each other: we must have $|x_{i+1} - x_i| \leq 1$, for $i = 1, \dots, k-1$. Thus, for example, a 0 or a 1 can follow a 1, but a -1 cannot follow a 1.
- Power constraint*: The sum of the squares of any three consecutive symbols cannot exceed 2:

$$x_i^2 + x_{i+1}^2 + x_{i+2}^2 \leq 2,$$

for $i = 1, \dots, k-2$.

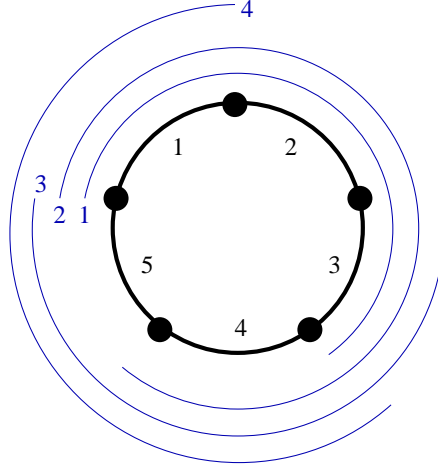


Figure 3: Data network for part (b), with links shown darker. Route 1 is (1, 2, 3), route 2 is (1, 2, 3, 4), route 3 is (3, 4, 5), and route 4 is (4, 5, 1), where routes are defined as sequences of links. All traffic and routes flow counterclockwise (although this doesn't matter).

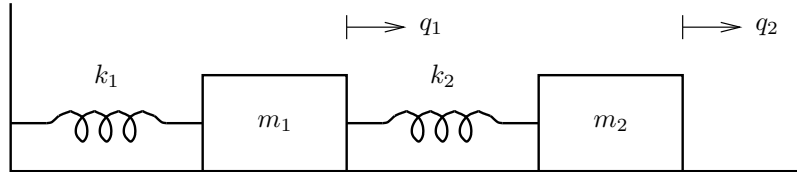
- *Average constraint:* The sum of any three consecutive symbols must not exceed one in absolute value:

$$|x_i + x_{i+1} + x_{i+2}| \leq 1,$$

for $i = 1, \dots, k - 2$. So, for example, a sequence that contains 1100 would not be valid, because the sum of the first three consecutive symbols is 2.

How many different (valid) sequences of length 20 are there?

9.5 Consider the mechanical system shown below:



Here q_i give the displacements of the masses, m_i are the values of the masses, and k_i are the spring stiffnesses, respectively. The dynamics of this system are

$$\dot{x} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -\frac{k_1+k_2}{m_1} & \frac{k_2}{m_1} & 0 & 0 \\ \frac{k_2}{m_2} & -\frac{k_2}{m_2} & 0 & 0 \end{bmatrix} x$$

where the state is given by

$$x = \begin{bmatrix} q_1 \\ q_2 \\ \dot{q}_1 \\ \dot{q}_2 \end{bmatrix}.$$

Immediately before $t = 0$, you are able to apply a strong impulsive force α_i to mass i , which results in initial condition

$$x(0) = \begin{bmatrix} 0 \\ 0 \\ \alpha_1/m_1 \\ \alpha_2/m_2 \end{bmatrix}.$$

(i.e., each mass starts with zero position and a velocity determined by the impulsive forces.) This problem concerns selection of the impulsive forces α_1 and α_2 . For parts a–c below, the parameter values are

$$m_1 = m_2 = 1, \quad k_1 = k_2 = 1.$$

Consider the following specifications:

- (a) $q_2(10) = 2$
- (b) $q_1(10) = 1, q_2(10) = 2$
- (c) $q_1(10) = 1, q_2(10) = 2, \dot{q}_1(10) = 0, \dot{q}_2(10) = 0$
- (d) $q_2(10) = 2$ when the parameters have the values used above (i.e., $m_1 = m_2 = 1, k_1 = k_2 = 1$), and also, $q_2(10) = 2$ when the parameters have the values $m_1 = 1, m_2 = 1.3, k_1 = k_2 = 1$.

Determine whether each of these specifications is feasible or not (i.e., whether there exist $\alpha_1, \alpha_2 \in \mathbf{R}$ that make the specification hold). If the specification is feasible, find the particular α_1, α_2 that satisfy the specification and minimize $\alpha_1^2 + \alpha_2^2$. If the specification is infeasible, find the particular α_1, α_2 that come closest, in a least-squares sense, to satisfying the specification. (For example, if you cannot find α_1, α_2 that satisfy $q_1(10) = 1, q_2(10) = 2$, then find α_i that minimize $(q_1(10) - 1)^2 + (q_2(10) - 2)^2$.) Be sure to be very clear about which alternative holds for each specification.

9.6 *Invariance of the unit square.* Consider the linear dynamical system $\dot{x} = Ax$ with $A \in \mathbf{R}^{2 \times 2}$. The unit square in \mathbf{R}^2 is defined by

$$S = \{ x \mid -1 \leq x_1 \leq 1, -1 \leq x_2 \leq 1 \}.$$

- (a) Find the exact conditions on A for which the unit square S is invariant under $\dot{x} = Ax$. Give the conditions as explicitly as possible.
- (b) Consider the following statement: if the eigenvalues of A are real and negative, then S is invariant under $\dot{x} = Ax$. Either show that this is true, or give an explicit counterexample.

9.7 *Iterative solution of linear equations.* In many applications we need to solve a set of linear equations $Ax = b$, where A is nonsingular (square) and x is very large (e.g., $x \in \mathbf{R}^{100000}$). We assume that Az can be computed at reasonable cost, for any z , but the standard methods for computing $x = A^{-1}b$ (e.g., LU decomposition) are not feasible. A common approach is to use an *iterative* method, which computes a sequence $x(1), x(2), \dots$ that *converges* to the solution $x = A^{-1}b$. These methods rely on another matrix \hat{A} , which is supposed to be ‘close’ to A . More importantly, \hat{A} has the property that $\hat{A}^{-1}z$ is easily or cheaply computed for any given z . As a simple example, the matrix \hat{A} might be the diagonal part of the matrix A (which, presumably, has relatively small off-diagonal elements). Obviously computing $\hat{A}^{-1}z$ is fast; it’s just scaling the entries of z . There are many, many other examples. A simple iterative method, sometimes called *relaxation*, is to set $\hat{x}(0)$ equal to some approximation of x (e.g., $\hat{x}(0) = \hat{A}^{-1}b$) and repeat, for $t = 0, 1, \dots$

$$r(t) = A\hat{x}(t) - b; \quad \hat{x}(t+1) = \hat{x}(t) - \hat{A}^{-1}r(t);$$

(The hat reminds us that $\hat{x}(t)$ is an approximation, after t iterations, of the true solution $x = A^{-1}b$.) This iteration uses only ‘cheap’ calculations: multiplication by A and \hat{A}^{-1} . Note that $r(t)$ is the residual after the t th iteration.

- (a) Let $\beta = \|\hat{A}^{-1}(A - \hat{A})\|$ (which is a measure of how close \hat{A} and A are). Show that if we choose $\hat{x}(0) = \hat{A}^{-1}b$, then $\|\hat{x}(t) - x\| \leq \beta^{t+1}\|x\|$. Thus if $\beta < 1$, the iterative method works, *i.e.*, for any b we have $\hat{x}(t) \rightarrow x$ as $t \rightarrow \infty$. (And if $\beta < 0.8$, say, then convergence is pretty fast.)
- (b) Find the exact conditions on A and \hat{A} such that the method works for any starting approximation $\hat{x}(0)$ and any b . Your condition can involve norms, singular values, condition number, and eigenvalues of A and \hat{A} , or some combination, etc. Your condition should be as explicit as possible; for example, it should not include any limits. Try to avoid the following two errors:
- Your condition guarantees convergence but is too restrictive. (For example: $\beta = \|\hat{A}^{-1}(A - \hat{A})\| < 0.8$)
 - Your condition doesn't guarantee convergence.

9.8 *Periodic solution of periodic linear dynamical system.* Consider the linear dynamical system $\dot{x} = A(t)x$ where

$$A(t) = \begin{cases} A_1 & 2k \leq t < 2k+1, \quad k = 0, 1, 2, \dots \\ A_2 & 2k+1 \leq t < 2k+2, \quad k = 0, 1, 2, \dots \end{cases}$$

In other words, $A(t)$ switches between the two values $A_1 \in \mathbf{R}^{n \times n}$ and $A_2 \in \mathbf{R}^{n \times n}$ every second. The matrix $A(t)$ is periodic with period 2, *i.e.*, $A(t+2) = A(t)$ for all $t \geq 0$.

- (a) *Existence of a periodic trajectory.* What are the conditions on A_1 and A_2 under which the system has a nonzero periodic trajectory, with period 2? By this we mean: there exists $x : \mathbf{R}_+ \rightarrow \mathbf{R}^n$, x not identically zero, with $x(t+2) = x(t)$ and $\dot{x} = A(t)x$.
- (b) *All trajectories are asymptotically periodic.* What are the conditions on A_1 and A_2 under which all trajectories of the system are asymptotically 2-periodic? By this we mean: for every $x : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ with $\dot{x} = A(t)x$, we have

$$\lim_{t \rightarrow \infty} \|x(t+2) - x(t)\| = 0.$$

(Note that this holds when x converges to zero ...)

Please note:

- Your conditions should be as explicit as possible. You can refer to the matrices A_1 and A_2 , or any matrices derived from them using standard matrix operations, their eigenvalues and eigenvectors or Jordan forms, singular values and singular vectors, etc.
- We do not want you to give us *a* condition under which the property described holds. We want you to give us the most general conditions under which the property holds.

9.9 *Analysis of a power control algorithm.* In this problem we consider again the power control method described in homework problem 1. Please refer to this problem for the setup and background. In that problem, you expressed the power control method as a discrete-time linear dynamical system, and simulated it for a specific set of parameters, with several values of initial power levels, and two target SINRs. You found that for the target SINR value $\gamma = 3$, the powers converged to values for which each SINR exceeded γ , no matter what the initial power was, whereas for the larger target SINR value $\gamma = 5$, the powers appeared to diverge, and the SINRs did not appear to converge. You are going to analyze this, now that you know a lot more about linear systems.

- (a) *Explain the simulations.* Explain your simulation results from the problem 1(b) for the given values of G , α , σ , and the two SINR threshold levels $\gamma = 3$ and $\gamma = 5$.
- (b) *Critical SINR threshold level.* Let us consider fixed values of G , α , and σ . It turns out that the power control algorithm works provided the SINR threshold γ is less than some critical value γ_{crit} (which might depend on G , α , σ), and doesn't work for $\gamma > \gamma_{\text{crit}}$. ('Works' means that no matter what the initial powers are, they converge to values for which each SINR exceeds γ .) Find an expression for γ_{crit} in terms of $G \in \mathbf{R}^{n \times n}$, α , and σ . Give the simplest expression you can. Of course you must explain how you came up with your expression.

9.10 Stability of a time-varying system. We consider a discrete-time linear dynamical system

$$x(t+1) = A(t)x(t),$$

where $A(t) \in \{A_1, A_2, A_3, A_4\}$. These 4 matrices, which are 4×4 , are given in `tv_data.m`.

Show that this system is stable, *i.e.*, for any trajectory x , we have $x(t) \rightarrow 0$ as $t \rightarrow \infty$. (This means that for any $x(0)$, and for any sequence $A(0), A(1), A(2), \dots$, we have $x(t) \rightarrow 0$ as $t \rightarrow \infty$.)

You may use any methods or concepts used in the class, *e.g.*, least-squares, eigenvalues, singular values, controllability, and so on. Your proof will consist of two parts:

- An explanation of how you are going to show that any trajectory converges to zero. Your argument of course will require certain conditions (that you will find) to hold for the given data A_1, \dots, A_4 .
- The numerical calculations that verify the conditions hold for the given data. You must provide the source code for these calculations, and show the results as well.

Lecture 10 – Solution via Laplace transform and matrix exponential

10.1 Suppose $\dot{x} = Ax$ and $\dot{z} = \sigma z + Az = (A + \sigma I)z$ where $\sigma \in \mathbf{R}$, and $x(0) = z(0)$. How are $z(t)$ and $x(t)$ related? Find the simplest possible expression for $z(t)$ in terms of $x(t)$. Justify your answer. When $\sigma < 0$, some people refer to the system $\dot{z} = \sigma z + Az$ as a *damped version* of $\dot{x} = Ax$. Another way to think of the damped system is in terms of *leaky integrators*. A leaky integrator satisfies $\dot{y} - \sigma y = u$; to get the damped system, you replace every integrator in the original system with a leaky integrator.

10.2 *Harmonic oscillator*. The system $\dot{x} = \begin{bmatrix} 0 & \omega \\ -\omega & 0 \end{bmatrix} x$ is called a *harmonic oscillator*.

- Find the eigenvalues, resolvent, and state transition matrix for the harmonic oscillator. Express $x(t)$ in terms of $x(0)$.
- Sketch the vector field of the harmonic oscillator.
- The state trajectories describe circular orbits, *i.e.*, $\|x(t)\|$ is constant. Verify this fact using the solution from part (a).
- You may remember that circular motion (in a plane) is characterized by the velocity vector being orthogonal to the position vector. Verify that this holds for any trajectory of the harmonic oscillator. Use only the differential equation; do not use the explicit solution you found in part (a).

10.3 *Properties of the matrix exponential*.

- Show that $e^{A+B} = e^A e^B$ if A and B commute, *i.e.*, $AB = BA$.
- Carefully show that $\frac{d}{dt} e^{At} = A e^{At} = e^{At} A$.

10.4 *Two-point boundary value problem*. Consider the system described by $\dot{x} = Ax$, where $A = \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}$.

- Find e^A .
- Suppose $x_1(0) = 1$ and $x_2(1) = 2$. Find $x(2)$. (This is called a *two-point boundary value problem*, since we are given conditions on the state at two time points instead of the usual single initial point.)

10.5 *Determinant of matrix exponential*.

- Suppose the eigenvalues of $A \in \mathbf{R}^{n \times n}$ are $\lambda_1, \dots, \lambda_n$. Show that the eigenvalues of e^A are $e^{\lambda_1}, \dots, e^{\lambda_n}$. You can assume that A is diagonalizable, although it is true in the general case.
- Show that $\det e^A = e^{\text{Tr } A}$. *Hint*: $\det X$ is the product of the eigenvalues of X , and $\text{Tr } Y$ is the sum of the eigenvalues of Y .

10.6 *Linear system with a quadrant detector*. In this problem we consider the specific system

$$\dot{x} = Ax = \begin{bmatrix} 0.5 & 1.4 \\ -0.7 & 0.5 \end{bmatrix} x.$$

We have a detector or sensor that gives us the sign of each component of the state $x = [x_1 \ x_2]^T$ each second:

$$y_1(t) = \text{sgn}(x_1(t)), \quad y_2(t) = \text{sgn}(x_2(t)), \quad t = 0, 1, 2, \dots$$

where the function $\text{sgn} : \mathbf{R} \rightarrow \mathbf{R}$ is defined by

$$\text{sgn}(a) = \begin{cases} 1 & a > 0 \\ 0 & a = 0 \\ -1 & a < 0 \end{cases}$$

There are several ways to think of these sensor measurements. You can think of $y(t) = [y_1(t) \ y_2(t)]^T$ as determining which quadrant the state is in at time t (thus the name *quadrant detector*). Or, you can think of $y(t)$ as a one-bit quantized measurement of the state at time t . Finally, the problem. You observe the sensor measurements

$$y(0) = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad y(1) = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

Based on these measurements, what values could $y(2)$ possibly take on? In terms of the quadrants, the problem can be stated as follows. $x(0)$ is in quadrant IV, and $x(1)$ is also in quadrant IV. The question is: which quadrant(s) can $x(2)$ possibly be in? You do not know the initial state $x(0)$. Of course, you must completely justify and explain your answer.

10.7 *Linear system with one-bit quantized output.* We consider the system

$$\dot{x} = Ax, \quad y(t) = \text{sign}(cx(t))$$

where

$$A = \begin{bmatrix} -0.1 & 1 \\ -1 & 0.1 \end{bmatrix}, \quad c = \begin{bmatrix} 1 & -1 \end{bmatrix},$$

and the sign function is defined as

$$\text{sign}(a) = \begin{cases} +1 & \text{if } a > 0 \\ -1 & \text{if } a < 0 \\ 0 & \text{if } a = 0 \end{cases}$$

Roughly speaking, the output of this autonomous linear system is quantized to one-bit precision. The following outputs are observed:

$$y(0.4) = +1, \quad y(1.2) = -1, \quad y(2.3) = -1, \quad y(3.8) = +1$$

What can you say (if anything) about the following:

$$y(0.7), \quad y(1.8), \quad \text{and } y(3.7)?$$

Your response might be, for example: “ $y(0.7)$ is definitely $+1$, and $y(1.8)$ is definitely -1 , but $y(3.7)$ can be anything (*i.e.*, -1 , 0 , or 1)”. Of course you *must* fully explain how you arrive at your conclusions. (What we mean by “ $y(0.7)$ is definitely $+1$ ” is: for any trajectory of the system for which $y(0.4) = +1$, $y(1.2) = -1$, $y(2.3) = -1$, and $y(3.8) = +1$, we also have $y(0.7) = +1$.)

10.8 *Some basic properties of eigenvalues.* Show that

- (a) the eigenvalues of A and A^T are the same
- (b) A is invertible if and only if A does not have a zero eigenvalue
- (c) if the eigenvalues of A are $\lambda_1, \dots, \lambda_n$ and A is invertible, then the eigenvalues of A^{-1} are $1/\lambda_1, \dots, 1/\lambda_n$,
- (d) the eigenvalues of A and $T^{-1}AT$ are the same.

Hint: you'll need to use the facts that $\det A = \det(A^T)$, $\det(AB) = \det A \det B$, and, if A is invertible, $\det A^{-1} = 1/\det A$.

10.9 *Characteristic polynomial.* Consider the characteristic polynomial $\mathcal{X}(s) = \det(sI - A)$ of the matrix $A \in \mathbf{R}^{n \times n}$.

- (a) Show that \mathcal{X} is *monic*, which means that its leading coefficient is one: $\mathcal{X}(s) = s^n + \dots$.

- (b) Show that the s^{n-1} coefficient of \mathcal{X} is given by $-\text{Tr } A$. ($\text{Tr } X$ is the *trace* of a matrix: $\text{Tr } X = \sum_{i=1}^n X_{ii}$.)
- (c) Show that the constant coefficient of \mathcal{X} is given by $\det(-A)$.
- (d) Let $\lambda_1, \dots, \lambda_n$ denote the eigenvalues of A , so that

$$\mathcal{X}(s) = s^n + a_{n-1}s^{n-1} + \dots + a_1s + a_0 = (s - \lambda_1)(s - \lambda_2) \cdots (s - \lambda_n).$$

By equating coefficients show that $a_{n-1} = -\sum_{i=1}^n \lambda_i$ and $a_0 = \prod_{i=1}^n (-\lambda_i)$.

10.10 *The adjoint system.* The adjoint system associated with the linear dynamical system $\dot{x} = Ax$ is $\dot{z} = A^T z$. Evidently the adjoint system and the system have the same eigenvalues.

- (a) How are the state-transition matrices of the system and the adjoint system related?
- (b) Show that $z(0)^T x(t) = z(t)^T x(0)$.

10.11 *Spectral resolution of the identity.* Suppose $A \in \mathbf{R}^{n \times n}$ has n linearly independent eigenvectors p_1, \dots, p_n , $p_i^T p_i = 1$, $i = 1, \dots, n$, with associated eigenvalues λ_i . Let $P = [p_1 \cdots p_n]$ and $Q = P^{-1}$. Let q_i^T be the i th row of Q .

- (a) Let $R_k = p_k q_k^T$. What is the range of R_k ? What is the rank of R_k ? Can you describe the null space of R_k ?
- (b) Show that $R_i R_j = 0$ for $i \neq j$. What is R_i^2 ?
- (c) Show that

$$(sI - A)^{-1} = \sum_{k=1}^n \frac{R_k}{s - \lambda_k}.$$

Note that this is a partial fraction expansion of $(sI - A)^{-1}$. For this reason the R_i 's are called the *residue* matrices of A .

- (d) Show that $R_1 + \dots + R_n = I$. For this reason the residue matrices are said to constitute a *resolution of the identity*.
- (e) Find the residue matrices for

$$A = \begin{bmatrix} 1 & 0 \\ 1 & -2 \end{bmatrix}$$

both ways described above (*i.e.*, find P and Q and then calculate the R 's, and then do a partial fraction expansion of $(sI - A)^{-1}$ to find the R 's).

10.12 *Using Matlab to find an invariant plane.* Consider the continuous-time system $\dot{x} = Ax$ with A given by

$$A = \begin{bmatrix} -0.1005 & 1.0939 & 2.0428 & 4.4599 \\ -1.0880 & -0.1444 & 5.9859 & -3.0481 \\ -2.0510 & -5.9709 & -0.1387 & 1.9229 \\ -4.4575 & 3.0753 & -1.8847 & -0.1164 \end{bmatrix}$$

You can verify that the eigenvalues of A are

$$\lambda_{1,2} = -0.10 \pm j5, \quad \lambda_{3,4} = -0.15 \pm j7.$$

- (a) Find an orthonormal basis (q_1, q_2) for the invariant plane associated with λ_1 and λ_2 .
- (b) Find $q_3, q_4 \in \mathbf{R}^4$ such that $Q = [q_1 \ q_2 \ q_3 \ q_4]$ is orthogonal. You might find the Matlab command `null` useful; it computes an orthonormal basis of the null space of a matrix.
- (c) Plot the individual states constituting the trajectory $x(t)$ of the system starting from an initial point in the invariant plane, say $x(0) = q_1$, for $0 \leq t \leq 40$.

- (d) If $x(t)$ is in the invariant plane what can you say about the components of the vector $Q^T x(t)$?
- (e) Using the result of part (12d) verify that the trajectory you found in part (12c) is in the invariant plane.

Note: The A matrix is available on the class web site in the file `inv_plane_matrix.m`.

10.13 *Positive quadrant invariance.* We consider a system $\dot{x} = Ax$ with $x(t) \in \mathbf{R}^2$ (although the results of this problem can be generalized to systems of higher dimension). We say the system is *positive quadrant invariant* (PQI) if whenever $x_1(T) \geq 0$ and $x_2(T) \geq 0$, we have $x_1(t) \geq 0$ and $x_2(t) \geq 0$ for all $t \geq T$. In other words, if the state starts inside (or enters) the positive (*i.e.*, first) quadrant, then the state remains indefinitely in the positive quadrant.

- (a) Find the precise conditions on A under which the system $\dot{x} = Ax$ is PQI. Try to express the conditions in the simplest form.
- (b) *True or False:* if $\dot{x} = Ax$ is PQI, then the eigenvalues of A are real.

10.14 *Some Matlab exercises.* Consider the continuous-time system $\dot{x} = Ax$ with

$$A = \begin{bmatrix} -0.1005 & 1.0939 & 2.0428 & 4.4599 \\ -1.0880 & -0.1444 & 5.9859 & -3.0481 \\ -2.0510 & -5.9709 & -0.1387 & 1.9229 \\ -4.4575 & 3.0753 & -1.8847 & -0.1164 \end{bmatrix}.$$

- (a) What are the eigenvalues of A ? Is the system stable? You can use the command `eig` in Matlab.
- (b) Plot a few trajectories of $x(t)$, *i.e.*, $x_1(t)$, $x_2(t)$, $x_3(t)$ and $x_4(t)$, for a few initial conditions. To do this you can use the matrix exponential command in Matlab `expm` (*not* `exp` which gives the element-by-element exponential of a matrix), or more directly, the Matlab command `initial` (use `help initial` for details.) Verify that the qualitative behavior of the system is consistent with the eigenvalues you found in part (14a).
- (c) Find the matrix Z such that $Zx(t)$ gives $x(t+15)$. Thus, Z is the ‘15 seconds forward predictor matrix’.
- (d) Find the matrix Y such that $Yx(t)$ gives $x(t-20)$. Thus Y reconstructs what the state was 20 seconds ago.
- (e) Briefly comment on the size of the elements of the matrices Y and Z .
- (f) Find $x(0)$ such that $x(10) = [1 \ 1 \ 1 \ 1]^T$.

10.15 *Volume preserving flows.* Suppose we have a set $S \subseteq \mathbf{R}^n$ and a linear dynamical system $\dot{x} = Ax$. We can propagate S along the ‘flow’ induced by the linear dynamical system by considering

$$S(t) = e^{At}S = \{ e^{At}s \mid s \in S \}.$$

Thus, $S(t)$ is the image of the set S under the linear transformation e^{tA} . What are the conditions on A so that the flow preserves volume, *i.e.*, $\text{vol}S(t) = \text{vol}S$ for all t ? Can the flow $\dot{x} = Ax$ be stable? *Hint:* if $F \in \mathbf{R}^{n \times n}$ then $\text{vol}(FS) = |\det F| \text{vol}S$, where $FS = \{ Fs \mid s \in S \}$.

10.16 *Stability of a periodic system.* Consider the linear dynamical system $\dot{x} = A(t)x$ where

$$A(t) = \begin{cases} A_1 & 2n \leq t < 2n+1, \quad n = 0, 1, 2, \dots \\ A_2 & 2n+1 \leq t < 2n+2, \quad n = 0, 1, 2, \dots \end{cases}$$

In other words, $A(t)$ switches between the two values A_1 and A_2 every second. We say that this (time-varying) linear dynamical system is stable if every trajectory converges to zero, *i.e.*, we have $x(t) \rightarrow 0$ as $t \rightarrow \infty$ for any $x(0)$. Find the conditions on A_1 and A_2 under which the periodic system is stable. Your conditions should be as explicit as possible.

10.17 *Computing trajectories of a continuous-time LDS.* We have seen in class that if $x(t)$ is the solution to the continuous-time, time-invariant, linear dynamical system

$$\dot{x} = Ax, \quad x(0) = x_0,$$

then the Laplace transform of $x(t)$ is given by

$$X(s) = (sI - A)^{-1} x_0.$$

Hence, we can obtain $x(t)$ from the inverse Laplace transform of the resolvent of A :

$$x(t) = \mathcal{L}^{-1} \left((sI - A)^{-1} \right) x_0.$$

- (a) Assuming that $A \in \mathbf{R}^{n \times n}$ has n independent eigenvectors, write $x(t)$ in terms of the residue matrices R_i and associated eigenvalues λ_i , $i = 1, \dots, n$. (The residue matrices are defined in the previous problem.)
- (b) Consider once again the matrix

$$A = \begin{bmatrix} 1 & 3 \\ 0 & -1 \end{bmatrix}.$$

Write the solution $x(t)$ for this dynamics matrix, with the initial condition $x_0 = [2 \ -1]^T$. Compute $x_1(2)$, *i.e.*, the value of the first entry of $x(t)$ at $t = 2$.

- (c) *Forward Euler approximation.* With this same A and x_0 , compute an approximation to the trajectory $x(t)$ by Euler approximation, with different step-sizes h . Run your simulation from $t = 0$ to $t = 2$, with N steps. For the number of steps N , use the values 10, 100, 1000, and 10000 (with the corresponding step-size $h = 2/N$). For each run, you'll obtain the sequence resulting from the discrete-time LDS

$$y(k+1) = (I + hA)y(k), \quad k = 0, \dots, N-1$$

with $y(0) = x_0$. On the same graph, plot the first entry, $y_1(k)$, of each of the four sequences you obtain (with hk on the horizontal axis).

- (d) *Error in Euler approximation.* For each of the four runs, compute the final error in x_1 , given by $\epsilon = y_1(N) - x_1(2)$. Plot ϵ as a function of N on a logarithmic scale (*hint*: use the Matlab function `loglog`). How many steps do you estimate you would need to achieve a precision of 10^{-6} ?
- (e) *Matrix exponential.* The matrix exponential is defined by the series

$$e^A = I + \sum_{k=1}^{+\infty} \frac{1}{k!} A^k.$$

With A as above and $h = 0.5$, compute an approximation of the matrix exponential of hA by adding the first ten term of the series:

$$B = I + \sum_{k=1}^{10} \frac{1}{k!} (hA)^k.$$

Compute 4 iterates of the discrete-time LDS

$$z(k+1) = Bz(k), \quad k = 0, \dots, 3,$$

with $z(0) = x_0$. Add $z_1(k)$ to the plot of the $y_1(k)$. What is the final error $\epsilon = z_1(4) - x_1(2)$? *Note*: The Matlab function `expm` uses a much more efficient algorithm to compute the matrix exponential. For this example, `expm` requires about the same computational effort as is needed to add the first ten terms of the series, but the result is much more accurate. (If you're curious, go ahead and compute the corresponding final error ϵ .)

- 10.18 Suppose $\dot{x} = Ax$ with $A \in \mathbf{R}^{n \times n}$. Two one-second experiments are performed. In the first, $x(0) = [1 \ 1]^T$ and $x(1) = [4 \ -2]^T$. In the second, $x(0) = [1 \ 2]^T$ and $x(1) = [5 \ -2]^T$.
- Find $x(1)$ and $x(2)$, given $x(0) = [3 \ -1]^T$.
 - Find A , by first computing the matrix exponential.
 - Either find $x(1.5)$ or explain why you cannot ($x(0) = [3 \ -1]^T$).
 - More generally, for $\dot{x} = Ax$ with $A \in \mathbf{R}^{n \times n}$, describe a procedure for finding A using experiments with different initial values. What conditions must be satisfied for your procedure to work?
- 10.19 *Output response envelope for linear system with uncertain initial condition.* We consider the autonomous linear dynamical system $\dot{x} = Ax$, $y(t) = Cx(t)$, where $x(t) \in \mathbf{R}^n$ and $y(t) \in \mathbf{R}$. We do not know the initial condition exactly; we only know that it lies in a ball of radius r centered at the point x_0 :

$$\|x(0) - x_0\| \leq r.$$

We call x_0 the *nominal* initial condition, and the resulting output, $y_{\text{nom}}(t) = Ce^{tA}x_0$, the *nominal output*. We define the *maximum output* or *upper output envelope* as

$$\bar{y}(t) = \max\{y(t) \mid \|x(0) - x_0\| \leq r\},$$

i.e., the maximum possible value of the output at time t , over all possible initial conditions. (Here you can choose a different initial condition for each t ; you are not required to find a single initial condition.) In a similar way, we define the *minimum output* or *lower output envelope* as

$$\underline{y}(t) = \min\{y(t) \mid \|x(0) - x_0\| \leq r\},$$

i.e., the minimum possible value of the output at time t , over all possible initial conditions.

- Explain how to find $\bar{y}(t)$ and $\underline{y}(t)$, given the problem data A , C , x_0 , and r .
 - Carry out your method on the problem data in `uie_data.m`. On the same axes, plot y_{nom} , \bar{y} , and \underline{y} , versus t , over the range $0 \leq t \leq 10$.
- 10.20 *Alignment of a fleet of vehicles.* We consider a fleet of vehicles, labeled $1, \dots, n$, which move along a line with (scalar) positions y_1, \dots, y_n . We let v_1, \dots, v_n denote the velocities of the vehicles, and u_1, \dots, u_n the net forces applied to the vehicles. The vehicle motions are governed by the equations

$$\dot{y}_i = v_i, \quad \dot{v}_i = u_i - v_i.$$

(Here we take each vehicle mass to be one, and include a damping term in the equations.) We assume that $y_1(0) < \dots < y_n(0)$, *i.e.*, the vehicles start out with vehicle 1 in the leftmost position, followed by vehicle 2 to its right, and so on, with vehicle n in the rightmost position. The goal is for the vehicles to converge to the configuration

$$y_i = i, \quad v_i = 0, \quad i = 1, \dots, n,$$

i.e., first vehicle at position 1, with unit spacing between adjacent vehicles, and all stationary. We call this configuration *aligned*, and the goal is to drive the vehicles to this configuration, *i.e.*, to align the vehicles. We define the spacing between vehicle i and $i+1$ as $s_i(t) = y_{i+1}(t) - y_i(t)$, for $i = 1, \dots, n-1$. (When the vehicles are aligned, these spacings are all one.) We will investigate three control schemes for aligning the fleet of vehicles.

- *Right looking control* is based on the spacing to the vehicle to the right. We use the control law

$$u_i(t) = s_i(t) - 1, \quad i = 1, \dots, n-1,$$

for vehicles $i = 1, \dots, n-1$. In other words, we apply a force on vehicle i proportional to its spacing error with respect to the vehicle to the right (*i.e.*, vehicle $i+1$). The rightmost vehicle uses the control law

$$u_n(t) = -(y_n(t) - n),$$

which applies a force proportional to its position error, in the opposite direction. This control law has the advantage that only the rightmost vehicle needs an absolute measurement sensor; the others only need a measurement of the distance to their righthand neighbor.

- *Left and right looking control* adjusts the input force based on the spacing errors to the vehicle to the left and the vehicle to the right:

$$u_i(t) = \frac{s_i(t) - 1}{2} - \frac{s_{i-1}(t) - 1}{2}, \quad i = 2, \dots, n-1,$$

The rightmost vehicle uses the same absolute error method as in right looking control, *i.e.*,

$$u_n(t) = -(y_n(t) - n),$$

and the first vehicle, which has no vehicle to its left, uses a right looking control scheme,

$$u_1(t) = s_1(t) - 1.$$

This scheme requires vehicle n to have an absolute position sensor, but the other vehicles only need to measure the distance to their neighbors.

- *Independent alignment* is based on each vehicle independently adjusting its position with respect to its required position:

$$u_i(t) = -(y_i(t) - i), \quad i = 1, \dots, n.$$

This scheme requires all vehicles to have absolute position sensors.

In the questions below, we consider the specific case with $n = 5$ vehicles.

- Which of the three schemes work? By ‘work’ we mean that the vehicles converge to the alignment configuration, no matter what the initial positions and velocities are. Among the schemes that do work, which one gives the fastest asymptotic convergence to alignment? (If there is a tie between two or three schemes, say so.) In this part of the problem you can ignore the issue of vehicle collisions, *i.e.*, spacings that pass through zero.
- Collisions.* In this problem we analyze vehicle collisions, which occur when any spacing between vehicles is equal to zero. (For example, $s_3(5.7) = 0$ means that vehicles 3 and 4 collide at $t = 5.7$.) We take the particular starting configuration

$$y = (0, 2, 3, 5, 7), \quad v = (0, 0, 0, 0, 0),$$

which corresponds to the vehicles with zero initial velocity, but not in the aligned positions. For each of the three schemes above (whether or not they work), determine if a collision occurs. If a collision does occur, find the earliest collision, giving the time and the vehicles involved. (For example, ‘Vehicles 3 and 4 collide at $t = 7.7$.’) If there is a tie, *i.e.*, two pairs of vehicles collide at the same time, say so. If the vehicles do not collide, find the point of closest approach, *i.e.*, the minimum spacing that occurs, between any pair of vehicles, for $t \geq 0$. (Give the time, the vehicles involved, and the minimum spacing.) If there is a tie, *i.e.*, two or more pairs of vehicles have the same distance of closest approach, say so. Be sure to give us times of collisions or closest approach with an absolute precision of at least 0.1.

- 10.21 *Scalar time-varying linear dynamical system.* Show that the solution of $\dot{x}(t) = a(t)x(t)$, where $x(t) \in \mathbf{R}$, is given by

$$x(t) = \exp\left(\int_0^t a(\tau) d\tau\right) x(0).$$

(You can just differentiate this expression, and show that it satisfies $\dot{x}(t) = a(t)x(t)$.) Find a specific example showing that the analogous formula does not hold when $x(t) \in \mathbf{R}^n$, with $n > 1$.

- 10.22 *Optimal initial conditions for a bioreactor.* The dynamics of a bioreactor are given by $\dot{x}(t) = Ax(t)$, where $x(t) \in \mathbf{R}^n$ is the state, with $x_i(t)$ representing the total mass of species or component i at time t . Component i has (positive) value (or cost) c_i , so the total value (or cost) of the components at time t is $c^T x(t)$. (We ignore any extra cost that would be incurred in separating the components.) Your job is to choose the initial state, under a budget constraint, that maximizes the total value at time T . More specifically, you are to choose $x(0)$, with all entries nonnegative, that satisfies $c^T x(0) \leq B$, where B is a given positive budget. The problem data (*i.e.*, things you know) are A , c , T , and B .

You can assume that A is such that, for any $x(0)$ with nonnegative components, $x(t)$ will also have all components nonnegative, for any $t \geq 0$. (This occurs, by the way, if and only if the off-diagonal entries of A are nonnegative.)

- (a) Explain how to solve this problem.
 (b) Carry out your method on the specific instance with data

$$A = \begin{bmatrix} 0.1 & 0.1 & 0.3 & 0 \\ 0 & 0.2 & 0.4 & 0.3 \\ 0.1 & 0.3 & 0.1 & 0 \\ 0 & 0 & 0.2 & 0.1 \end{bmatrix}, \quad c = \begin{bmatrix} 3.5 \\ 0.6 \\ 1.1 \\ 2.0 \end{bmatrix}, \quad T = 10, \quad B = 1.$$

Give the optimal $x(0)$, and the associated (optimal) terminal value $c^T x(T)$.

Give us the terminal value obtained when the initial state has equal mass in each component, *i.e.*, $x(0) = \alpha \mathbf{1}$, with α adjusted so that the total initial cost is B . Compare this with the optimal terminal value.

Also give us the terminal value obtained when the same amount, B/n , is spent on each initial state component (*i.e.*, $x(0)_i = B/(nc_i)$). Compare this with the optimal terminal value.

- 10.23 *Optimal espresso cup pre-heating.* At time $t = 0$ boiling water, at 100°C , is poured into an espresso cup; after P seconds (the ‘pre-heating time’), the water is poured out, and espresso, with initial temperature 95°C , is poured in. (You can assume this operation occurs instantaneously.) The espresso is then consumed exactly 15 seconds later (yes, instantaneously). The problem is to choose the pre-heating time P so as to maximize the temperature of the espresso when it is consumed.

We now give the thermal model used. We take the temperature of the liquid in the cup (water or espresso) as one state; for the cup we use an n -state finite element model. The vector $x(t) \in \mathbf{R}^{n+1}$ gives the temperature distribution at time t : $x_1(t)$ is the liquid (water or espresso) temperature at time t , and $x_2(t), \dots, x_{n+1}(t)$ are the temperatures of the elements in the cup. All of these are in degrees C, with t in seconds. The dynamics are

$$\frac{d}{dt}(x(t) - 20 \cdot \mathbf{1}) = A(x(t) - 20 \cdot \mathbf{1}),$$

where $A \in \mathbf{R}^{(n+1) \times (n+1)}$. (The vector $20 \cdot \mathbf{1}$, with all components 20, represents the ambient temperature.) The initial temperature distribution is

$$x(0) = \begin{bmatrix} 100 \\ 20 \\ \vdots \\ 20 \end{bmatrix}.$$

At $t = P$, the liquid temperature changes instantly from whatever value it has, to 95; the other states do not change. Note that the dynamics of the system are the same before and after pre-heating (because we assume that water and espresso behave in the same way, thermally speaking).

We have *very generously* derived the matrix A for you. You will find it in `espressodata.m`. In addition to A , the file also defines \mathbf{n} , and, respectively, the ambient, espresso and preheat water temperatures T_a (which is 20), T_e (95), and T_l (100).

Explain your method, submit your code, and give final answers, which must include the optimal value of P and the resulting optimal espresso temperature when it is consumed. Give both to an accuracy of one decimal place, as in

‘ $P = 23.5$ s, which gives an espresso temperature at consumption of 62.3°C .’

(This is not the correct answer, of course.)

Lecture 11 – Eigenvectors and diagonalization

11.1 *Left eigenvector properties.* Suppose w is a left eigenvector of $A \in \mathbf{R}^{n \times n}$ with real negative eigenvalue λ .

- Find a simple expression for $w^T e^{At}$.
- Let $\alpha < \beta$. The set $\{ z \mid \alpha \leq w^T z \leq \beta \}$ is referred to as a *slab*. Briefly explain this terminology. Draw a picture in \mathbf{R}^2 .
- Show that the slab $\{ z \mid 0 \leq w^T z \leq \beta \}$ is invariant for $\dot{x} = Ax$.

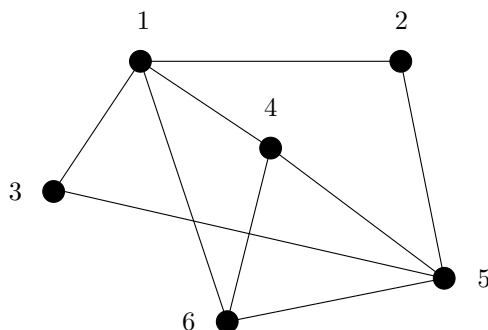
11.2 Consider the linear dynamical system $\dot{x} = Ax$ where $A \in \mathbf{R}^{n \times n}$ is diagonalizable with eigenvalues λ_i , eigenvectors v_i , and left eigenvectors w_i for $i = 1, \dots, n$. Assume that $\lambda_1 > 0$ and $\Re \lambda_i < 0$ for $i = 2, \dots, n$. Describe the trajectories qualitatively. Specifically, what happens to $x(t)$ as $t \rightarrow \infty$? Give the answer geometrically, in terms of $x(0)$.

11.3 *Another formula for the matrix exponential.* You might remember that for any complex number $a \in \mathbf{C}$, $e^a = \lim_{k \rightarrow \infty} (1 + a/k)^k$. You will establish the matrix analog: for any $A \in \mathbf{R}^{n \times n}$,

$$e^A = \lim_{k \rightarrow \infty} (I + A/k)^k.$$

To simplify things, you can assume A is diagonalizable. *Hint:* diagonalize.

11.4 *Synchronizing a communication network.* The graph below shows a communication network, with communication links shown as lines between the nodes, which are labeled 1, ..., 6. We refer to one node as a *neighbor* of another if they are connected by a link.



Each node has a clock. The clocks run at the same speed, but are not (initially) synchronized. The shift or offset of clock i , with respect to some absolute clock (*e.g.*, NIST's atomic clocks or the clock for the GPS system) will be denoted x_i . Thus $x_i > 0$ means the clock at node i is running in advance of the standard clock, while $x_i < 0$ means the i th clock is running behind the standard clock. The nodes *do not* know their own clock offsets (or the offsets of any of the other clocks); we introduce the numbers x_i only so we can analyze the system. At discrete intervals, which we denote $t = 0, 1, 2, \dots$, the nodes exchange communications messages. Through this exchange each node is able to find out the *relative time offset* of its own clock compared to the clocks of its neighboring nodes. For example, node 2 is able to find out the differences $x_1 - x_2$ and $x_5 - x_2$. (But remember, node 2 does not know any of the absolute clock offsets x_1 , x_2 , or x_5 .) While node i does not know its absolute offset x_i , it is able to *adjust it* by adding a delay or advance to it. The new offset takes effect at the next interval. Thus we have $x_i(t+1) = x_i(t) + a_i(t)$, where $a_i(t)$ is the adjustment made by the i th node to its clock in the t th interval. An engineer suggests the following scheme of adjusting the clock offsets. At each interval, each node determines its relative offset with each of its neighboring nodes. Then it computes

the average of these relative offsets. The node then adjusts its offset by this average. For example, for node 2 we would have the adjustment

$$a_2(t) = \frac{(x_1(t) - x_2(t)) + (x_5(t) - x_2(t))}{2}.$$

Finally, the question.

- (a) What happens?
- (b) Why?

We are interested in questions such as: do all the clocks become synchronized with the standard clock (*i.e.*, $x(t) \rightarrow 0$ as $t \rightarrow \infty$)? Do the clocks become synchronized with each other (*i.e.*, do all $x_i(t) - x_j(t)$ converge to zero as $t \rightarrow \infty$)? Does the system become synchronized no matter what the initial offsets are, or only for some initial offsets? You are welcome to use Matlab to do some relevant numerical computations, but you *must* explain what you are doing and why. We will *not* accept simulations of the network as an explanation. Another engineer suggests a modification of the scheme described above. She notes that if the scheme above were applied to a simple network consisting of two connected nodes, then the two nodes would just trade their offsets each time interval, so synchronization does not occur. To avoid this, she proposes to adjust each node's clock by only half the average offset with its neighbors. Thus, for node 2, this means:

$$a_2(t) = \frac{1}{2} \frac{(x_1(t) - x_2(t)) + (x_5(t) - x_2(t))}{2}.$$

- (c) Would you say this scheme is better or worse than the original one described above? If one is better than the other, how is it better? (For example, does it achieve synchronization from a bigger set of initial offsets, does it achieve synchronization faster, etc.)

11.5 Population dynamics. In this problem we will study how some population distribution (say, of people) evolves over time, using a discrete-time linear dynamical system model. Let $t = 0, 1, \dots$ denote time in years (since the beginning of the study). The vector $x(t) \in \mathbf{R}^n$ will give the population distribution at year t (on some fixed census date, *e.g.*, January 1). Specifically, $x_i(t)$ is the number of people at year t , of age $i - 1$. Thus $x_5(3)$ denotes the number of people of age 4, at year 3, and $x_1(t)$ (the number of 0 year-olds) denotes the number of people born since the last census. We assume n is large enough that no one lives to age n . We'll also ignore the fact that x_i are integers, and treat them as real numbers. (If $x_3(4) = 1.2$ bothers you, you can imagine the units as millions, say.) The total population at year t is given by $\mathbf{1}^T x(t)$, where $\mathbf{1} \in \mathbf{R}^n$ is the vector with all components 1.

- *Death rate.* The death rate depends only on age, and not on time t . The coefficient d_i is the fraction of people of age $i - 1$ who will die during the year. Thus we have, for $t = 0, 1, \dots$,

$$x_{k+1}(t+1) = (1 - d_k)x_k(t), \quad k = 1, \dots, n-1.$$

(As mentioned above, we assume that $d_n = 1$, *i.e.*, all people who make it to age $n - 1$ die during the year.) The death rate coefficients satisfy $0 < d_i < 1$, $i = 1, \dots, n - 1$. We define the *survival rate* coefficients as $s_k = 1 - d_k$, so $0 < s_k < 1$, $k = 1, \dots, n - 1$.

- *Birth rate.* The birth rate depends only on age, and not on time t . The coefficient b_i is the fraction of people of age $i - 1$ who will have a child during the year (taking into account multiple births). Thus the total births during a year is given by

$$x_1(t+1) = b_1x_1(t) + \dots + b_nx_n(t).$$

The birth rate coefficients satisfy $b_i \geq 0$, $i = 1, \dots, n$. We'll assume that at least one of the b_k 's is positive. (Of course you'd expect that b_i would be zero for non-fertile ages, *e.g.*, age below 11 and over 60, but we won't make that explicit assumption.)

The assumptions imply the following important property of our model: if $x_i(0) > 0$ for $i = 1, \dots, n$, then $x_i(t) > 0$ for $i = 1, \dots, n$. Therefore we don't have to worry about negative $x_i(t)$, so long as our initial population distribution $x(0)$ has all positive components. (To use fancy language we'd say the system is *positive orthant invariant*.)

- (a) Express the population dynamics model described above as a discrete-time linear dynamical system. That is, find a matrix A such that $x(t+1) = Ax(t)$.
- (b) Draw a block diagram of the system found in part (a).
- (c) Find the characteristic polynomial of the system explicitly in terms of the birth and death rate coefficients (or, if you prefer, the birth and survival rate coefficients).
- (d) *Survival normalized variables.* For each person born, s_1 make it to age 1, $s_1 s_2$ make it to age 2, and in general, $s_1 \cdots s_k$ make it to age k . We define

$$y_k(t) = \frac{x_k(t)}{s_1 \cdots s_{k-1}}$$

(with $y_1(t) = x_1(t)$) as new population variables that are normalized to the survival rate. Express the population dynamics as a linear dynamical system using the variable $y(t) \in \mathbf{R}^n$. That is, find a matrix \tilde{A} such $y(t+1) = \tilde{A}y(t)$.

Determine whether each of the next four statements is true or false. (Of course by 'true' we mean true for any values of the coefficients consistent with our assumptions, and by 'false' we mean false for some choice of coefficients consistent with our assumptions.)

- (e) Let x and z both satisfy our population dynamics model, *i.e.*, $x(t+1) = Ax(t)$ and $z(t+1) = Az(t)$, and assume that all components of $x(0)$ and $z(0)$ are positive. If $\mathbf{1}^T x(0) > \mathbf{1}^T z(0)$, then $\mathbf{1}^T x(t) > \mathbf{1}^T z(t)$ for $t = 1, 2, \dots$ (In words: we consider two populations that satisfy the same dynamics. Then the population that is initially larger will always be larger.)
- (f) All the eigenvalues of A are real.
- (g) If $d_k \geq b_k$ for $k = 1, \dots, n$, then $\mathbf{1}^T x(t) \rightarrow 0$ as $t \rightarrow \infty$, *i.e.*, the population goes extinct.
- (h) Suppose that $(b_1 + \cdots + b_n)/n \leq (d_1 + \cdots + d_n)/n$, *i.e.*, the 'average' birth rate is less than the 'average' death rate. Then $\mathbf{1}^T x(t) \rightarrow 0$ as $t \rightarrow \infty$.

11.6 *Rate of a Markov code.* Consider the Markov language described in exercise 13, with five symbols 1, 2, 3, 4, 5, and the following symbol transition rules:

- 1 must be followed by 2 or 3
- 2 must be followed by 2 or 5
- 3 must be followed by 1
- 4 must be followed by 4 or 2 or 5
- 5 must be followed by 1 or 3

- (a) *The rate of the code.* Let K_N denote the number of allowed sequences of length N . The number

$$R = \lim_{N \rightarrow \infty} \frac{\log_2 K_N}{N}$$

(if it exists) is called the *rate* of the code, in bits per symbol. Find the rate of this code. Compare it to the rate of the code which consists of all sequences from an alphabet of 5 symbols (*i.e.*, with no restrictions on which symbols can follow which symbols).

- (b) *Asymptotic fraction of sequences with a given starting or ending symbol.* Let $F_{N,i}$ denote the number of allowed sequences of length N that start with symbol i , and let $G_{N,i}$ denote the number of allowed sequences of length N that end with symbol i . Thus, we have

$$F_{N,1} + \cdots + F_{N,5} = G_{N,1} + \cdots + G_{N,5} = K_N.$$

Find the asymptotic fractions

$$f_i = \lim_{N \rightarrow \infty} F_{N,i}/K_N, \quad g_i = \lim_{N \rightarrow \infty} G_{N,i}/K_N.$$

Please don't find your answers by simple simulation or relatively mindless computation; we want to see (and understand) your method.

11.7 *Companion matrices.* A matrix A of the form

$$A = \begin{bmatrix} -a_1 & -a_2 & \cdots & -a_{n-1} & -a_n \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix}$$

is said to be a (top) *companion matrix*. There can be four forms of companion matrices depending on whether the a_i 's occur in the first or last row, or first or last column. These are referred to as top-, bottom-, left-, or right-companion matrices. Let $\dot{x} = Ax$ where A is top-companion.

- Draw a block diagram for the system $\dot{x} = Ax$.
 - Find the characteristic polynomial of the system using the block diagram and show that A is nonsingular if and only if $a_n \neq 0$.
 - Show that if A is nonsingular, then A^{-1} is a bottom-companion matrix with last row $-[1 \ a_1 \ \cdots \ a_{n-1}]/a_n$.
 - Find the eigenvector of A associated with the eigenvalue λ .
 - Suppose that A has distinct eigenvalues $\lambda_1, \dots, \lambda_n$. Find T such that $T^{-1}AT$ is diagonal.
- 11.8 *Squareroot and logarithm of a (diagonalizable) matrix.* Suppose that $A \in \mathbf{R}^{n \times n}$ is diagonalizable. Therefore, an invertible matrix $T \in \mathbf{C}^{n \times n}$ and diagonal matrix $\Lambda \in \mathbf{C}^{n \times n}$ exist such that $A = T\Lambda T^{-1}$. Let $\Lambda = \mathbf{diag}(\lambda_1, \dots, \lambda_n)$.
- We say $B \in \mathbf{R}^{n \times n}$ is a squareroot of A if $B^2 = A$. Let μ_i satisfy $\mu_i^2 = \lambda_i$. Show that $B = T \mathbf{diag}(\mu_1, \dots, \mu_n) T^{-1}$ is a squareroot of A . A squareroot is sometimes denoted $A^{1/2}$ (but note that there are in general many squareroots of a matrix). When λ_i are real and nonnegative, it is conventional to take $\mu_i = \sqrt{\lambda_i}$ (i.e., the nonnegative squareroot), so in this case $A^{1/2}$ is unambiguous.
 - We say B is a logarithm of A if $e^B = A$, and we write $B = \log A$. Following the idea of part a, find an expression for a logarithm of A (which you can assume is invertible). Is the logarithm unique? What if we insist on B being real?
- 11.9 *Separating hyperplane for a linear dynamical system.* A *hyperplane* (passing through 0) in \mathbf{R}^n is described by the equation $c^T x = 0$, where $c \in \mathbf{R}^n$ is nonzero. (Note that if $\beta \neq 0$, the vector $\tilde{c} = \beta c$ defines the same hyperplane.) Now consider the autonomous linear dynamic system $\dot{x} = Ax$, where $A \in \mathbf{R}^{n \times n}$ and $x(t) \in \mathbf{R}^n$. We say that the hyperplane defined by c is a *separating hyperplane* for this system if no trajectory of the system ever crosses the hyperplane. This means it is impossible to have $c^T x(t) > 0$ for some t , and $c^T x(\tilde{t}) < 0$ for some other \tilde{t} , for any trajectory x of the system. Explain how to find *all* separating hyperplanes for the system $\dot{x} = Ax$. In particular, give the conditions on A under which there is no separating hyperplane. (If you think there is always a separating hyperplane for a linear system, say so.) You can assume that A has distinct eigenvalues (and therefore is diagonalizable).

- 11.10 *Equi-angle sets.* Let $x_1, \dots, x_n \in \mathbf{R}^n$. We say that they form a (normalized) *equi-angle set*, with angle θ , if $\|x_i\| = 1$, $i = 1, \dots, n$, and

$$\angle(x_i, x_j) = \theta, \quad i, j = 1, \dots, n, \quad i \neq j.$$

In other words, each of the vectors has unit norm, and the angle between any pair of the vectors is θ . We'll take θ to be between 0 and π . An orthonormal set is a familiar example of an equi-angle set, with $\theta = \pi/2$. In \mathbf{R}^2 , there are equi-angle sets for every value of θ . It's easy to find such sets: just take

$$x_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad x_2 = \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}.$$

In \mathbf{R}^n , with $n > 2$, however, you can't have an equi-angle set with angle $\theta = \pi$. To see this, suppose that x_1, \dots, x_n is an equi-angle set in \mathbf{R}^n , with $n > 2$. Then we have $x_2 = -x_1$ (since $\angle(x_1, x_2) = \pi$), but also $x_3 = -x_1$ (since $\angle(x_1, x_3) = \pi$), so $\angle(x_2, x_3) = 0$. The question then arises, for what values of θ (between 0 and π) can you have an equi-angle set on \mathbf{R}^n ? The angle $\theta = 0$ always has an equi-angle set (just choose any unit vector u and set $x_1 = \dots = x_n = u$), and so does $\theta = \pi/2$ (just choose any orthonormal basis, *e.g.*, e_1, \dots, e_n). But what other angles are possible? For $n = 2$, we know the answer: any value of θ between 0 and π is possible, *i.e.*, for every value of θ there is an equi-angle set with angle θ .

- For general n , describe the values of θ for which there is an equi-angle set with angle θ . In particular, what is the maximum possible value θ can have?
 - Construct a specific equi-angle set in \mathbf{R}^4 for angle $\theta = 100^\circ = 5\pi/9$. Attach Matlab output to verify that your four vectors are unit vectors, and that the angle between any two of them is 100° . (Since $\angle(u, v) = \angle(v, u)$, you only have to check 6 angles. Also, you might find a clever way to find all the angles at once.)
- 11.11 *Optimal control for maximum asymptotic growth.* We consider the controllable linear system

$$x(t+1) = Ax(t) + Bu(t), \quad x(0) = 0,$$

where $x(t) \in \mathbf{R}^n$, $u(t) \in \mathbf{R}^m$. You can assume that A is diagonalizable, and that it has a single dominant eigenvalue (which here, means that there is one eigenvalue with largest magnitude). An input $u(0), \dots, u(T-1)$ is applied over time period $0, 1, \dots, T-1$; for $t \geq T$, we have $u(t) = 0$. The input is subject to a total energy constraint:

$$\|u(0)\|^2 + \dots + \|u(T-1)\|^2 \leq 1.$$

The goal is to choose the inputs $u(0), \dots, u(T-1)$ that maximize the norm of the state for large t . To be more precise, we're searching for $u(0), \dots, u(T-1)$, that satisfies the total energy constraint, and, for any other input sequence $\tilde{u}(0), \dots, \tilde{u}(T-1)$ that satisfies the total energy constraint, satisfies $\|x(t)\| \geq \|\tilde{x}(t)\|$ for t large enough. Explain how to do this. You can use any of the ideas from the class, *e.g.*, eigenvector decomposition, SVD, pseudo-inverse, etc. Be sure to summarize your final description of how to solve the problem. Unless you have to, you should *not* use limits in your solution. For example you cannot explain how to make $\|x(t)\|$ as large as possible (for a specific value of t), and then say, "Take the limit as $t \rightarrow \infty$ " or "Now take t to be really large".

- 11.12 *Estimating a matrix with known eigenvectors.* This problem is about estimating a matrix $A \in \mathbf{R}^{n \times n}$. The matrix A is not known, but we do have a noisy measurement of it, $A^{\text{meas}} = A + E$. Here the matrix E is measurement error, which is assumed to be small. While A is not known, we do know real, independent eigenvectors v_1, \dots, v_n of A . (Its eigenvalues $\lambda_1, \dots, \lambda_n$, however, are not known.)

We will combine our measurement of A with our prior knowledge to find an estimate \hat{A} of A . To do this, we choose \hat{A} as the matrix that minimizes

$$J = \frac{1}{n^2} \sum_{i,j=1}^n (A_{ij}^{\text{meas}} - \hat{A}_{ij})^2$$

among all matrices which have eigenvectors v_1, \dots, v_n . (Thus, \hat{A} is the matrix closest to our measurement, in the mean-square sense, that is consistent with the known eigenvectors.)

- (a) Explain how you would find \hat{A} . If your method is iterative, say whether you can guarantee convergence. Be sure to say whether your method finds the exact minimizer of J (except, of course, for numerical error due to roundoff), or an approximate solution. You can use any of the methods (least-squares, least-norm, Gauss-Newton, low rank approximation, *etc.*) or decompositions (QR, SVD, eigenvalue decomposition, *etc.*) from the course.
- (b) Carry out your method with the data

$$A^{\text{meas}} = \begin{bmatrix} 2.0 & 1.2 & -1.0 \\ 0.4 & 2.0 & -0.5 \\ -0.5 & 0.9 & 1.0 \end{bmatrix}, \quad v_1 = \begin{bmatrix} 0.7 \\ 0 \\ 0.7 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 0.3 \\ 0.6 \\ 0.7 \end{bmatrix}, \quad v_3 = \begin{bmatrix} 0.6 \\ 0.6 \\ 0.3 \end{bmatrix}.$$

Be sure to check that \hat{A} does indeed have v_1, v_2, v_3 as eigenvectors, by (numerically) finding its eigenvectors and eigenvalues. Also, give the value of J for \hat{A} . **Hint.** You might find the following useful (but then again, you might not.) In Matlab, if \mathbf{A} is a matrix, then $\mathbf{A}(:)$ is a (column) vector consisting of all the entries of \mathbf{A} , written out column by column. Therefore `norm(A(:))` gives the squareroot of the sum of the squares of entries of the matrix \mathbf{A} , *i.e.*, its Frobenius norm. The inverse operation, *i.e.*, writing a vector out as a matrix with some given dimensions, is done using the function `reshape`. For example, if \mathbf{a} is an mn vector, then `reshape(a,m,n)` is an $m \times n$ matrix, with elements taken from \mathbf{a} (column by column).

- 11.13 *Real modal form.* Generate a matrix A in $\mathbf{R}^{10 \times 10}$ using `A=randn(10)`. (The entries of A will be drawn from a unit normal distribution.) Find the eigenvalues of A . If by chance they are all real, please generate a new instance of A . Find the real modal form of A , *i.e.*, a matrix S such that $S^{-1}AS$ has the real modal form given in lecture 11. Your solution should include a clear explanation of how you will find S , the source code that you use to find S , and some code that checks the results (*i.e.*, computes $S^{-1}AS$ to verify it has the required form).

- 11.14 *Spectral mapping theorem.* Suppose $f : \mathbf{R} \rightarrow \mathbf{R}$ is analytic, *i.e.*, given by a power series expansion

$$f(u) = a_0 + a_1 u + a_2 u^2 + \dots$$

(where $a_i = f^{(i)}(0)/(i!)$). (You can assume that we only consider values of u for which this series converges.) For $A \in \mathbf{R}^{n \times n}$, we define $f(A)$ as

$$f(A) = a_0 I + a_1 A + a_2 A^2 + \dots$$

(again, we'll just assume that this converges).

Suppose that $Av = \lambda v$, where $v \neq 0$, and $\lambda \in \mathbf{C}$. Show that $f(A)v = f(\lambda)v$ (ignoring the issue of convergence of series). We conclude that if λ is an eigenvalue of A , then $f(\lambda)$ is an eigenvalue of $f(A)$. This is called the *spectral mapping theorem*.

To illustrate this with an example, generate a random 3×3 matrix, for example using `A=randn(3)`. Find the eigenvalues of $(I + A)(I - A)^{-1}$ by first computing this matrix, then finding its eigenvalues, and also by using the spectral mapping theorem. (You should get very close agreement; any difference is due to numerical round-off errors in the various computations.)

Lecture 12 – Jordan canonical form

12.1 *Some true/false questions.* Determine if the following statements are true or false. No justification or discussion is needed for your answers. What we mean by “true” is that the statement is true for all values of the matrices and vectors that appear in the statement. You can’t assume anything about the dimensions of the matrices (unless it’s explicitly stated), but you can assume that the dimensions are such that all expressions make sense. For example, the statement “ $A + B = B + A$ ” is true, because no matter what the dimensions of A and B are (they must, however, be the same), and no matter what values A and B have, the statement holds. As another example, the statement $A^2 = A$ is false, because there are (square) matrices for which this doesn’t hold. (There are also matrices for which it does hold, *e.g.*, an identity matrix. But that doesn’t make the statement true.) “False” means the statement isn’t true, in other words, it can fail to hold for some values of the matrices and vectors that appear in it.

- (a) If $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{n \times p}$ are both full rank, and $AB = 0$, then $n \geq m + p$.
- (b) If $A \in \mathbf{R}^{3 \times 3}$ satisfies $A + A^T = 0$, then A is singular.
- (c) If $A^k = 0$ for some integer $k \geq 1$, then $I - A$ is nonsingular.
- (d) If $A, B \in \mathbf{R}^{n \times n}$ are both diagonalizable, then AB is diagonalizable.
- (e) If $A, B \in \mathbf{R}^{n \times n}$, then every eigenvalue of AB is an eigenvalue of BA .
- (f) If $A, B \in \mathbf{R}^{n \times n}$, then every eigenvector of AB is an eigenvector of BA .
- (g) If A is nonsingular and A^2 is diagonalizable, then A is diagonalizable.

12.2 Consider the discrete-time system $x(t+1) = Ax(t)$, where $x(t) \in \mathbf{R}^n$.

- (a) Find $x(t)$ in terms of $x(0)$.
- (b) Suppose that $\det(zI - A) = z^n$. What are the eigenvalues of A ? What (if anything) can you say about $x(k)$ for $k < n$ and $k \geq n$, without knowing $x(0)$?

12.3 *Asymptotically periodic trajectories.* We say that $x : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ is *asymptotically T -periodic* if $\|x(t+T) - x(t)\|$ converges to 0 as $t \rightarrow \infty$. (We assume $T > 0$ is fixed.) Now consider the (time-invariant) linear dynamical system $\dot{x} = Ax$, where $A \in \mathbf{R}^{n \times n}$. Describe the precise conditions on A under which *all* trajectories of $\dot{x} = Ax$ are asymptotically T -periodic. Give your answer in terms of the Jordan form of A . (The period T can appear in your answer.) Make sure your answer works for ‘silly’ cases like $A = 0$ (for which all trajectories are constant, hence asymptotically T -periodic), or stable systems (for which all trajectories converge to 0, hence are asymptotically T -periodic). Mark your answer clearly, to isolate it from any (brief) discussion or explanation. You do not need to formally prove your answer; a brief explanation will suffice.

12.4 *Jordan form of a block matrix.* We consider the block 2×2 matrix

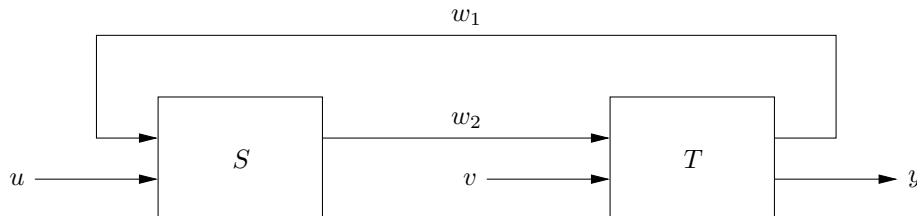
$$C = \begin{bmatrix} A & I \\ 0 & A \end{bmatrix}.$$

Here $A \in \mathbf{R}^{n \times n}$, and is diagonalizable, with real, distinct eigenvalues $\lambda_1, \dots, \lambda_n$. We’ll let v_1, \dots, v_n denote (independent) eigenvectors of A associated with $\lambda_1, \dots, \lambda_n$.

- (a) Find the Jordan form J of C . Be sure to explicitly describe its block sizes.
- (b) Find a matrix T such that $J = T^{-1}CT$.

Lecture 13 – Linear dynamical systems with inputs and outputs

13.1 *Interconnection of linear systems.* Often a linear system is described in terms of a block diagram showing the interconnections between components or subsystems, which are themselves linear systems. In this problem you consider the specific interconnection shown below:



Here, there are two subsystems S and T . Subsystem S is characterized by

$$\dot{x} = Ax + B_1u + B_2w_1, \quad w_2 = Cx + D_1u + D_2w_1,$$

and subsystem T is characterized by

$$\dot{z} = Fz + G_1v + G_2w_2, \quad w_1 = H_1z, \quad y = H_2z + Jw_2.$$

We don't specify the dimensions of the signals (which can be vectors) or matrices here. You can assume all the matrices are the correct (*i.e.*, compatible) dimensions. Note that the subscripts in the matrices above, as in B_1 and B_2 , refer to different matrices. Now the problem. Express the overall system as a single linear dynamical system with input, state, and output given by

$$\begin{bmatrix} u \\ v \end{bmatrix}, \quad \begin{bmatrix} x \\ z \end{bmatrix}, \quad y,$$

respectively. Be sure to explicitly give the input, dynamics, output, and feedthrough matrices of the overall system. If you need to make any assumptions about the rank or invertibility of any matrix you encounter in your derivations, go ahead. But be sure to let us know what assumptions you are making.

13.2 *Minimum energy control.* Consider the discrete-time linear dynamical system

$$x(t+1) = Ax(t) + Bu(t), \quad t = 0, 1, 2, \dots$$

where $x(t) \in \mathbf{R}^n$, and the input $u(t)$ is a scalar (hence, $A \in \mathbf{R}^{n \times n}$ and $B \in \mathbf{R}^{n \times 1}$). The initial state is $x(0) = 0$.

(a) Find the matrix \mathcal{C}_T such that

$$x(T) = \mathcal{C}_T \begin{bmatrix} u(T-1) \\ \vdots \\ u(1) \\ u(0) \end{bmatrix}.$$

(b) For the remainder of this problem, we consider a specific system with $n = 4$. The dynamics and input matrices are

$$A = \begin{bmatrix} 0.5 & 0.7 & -0.9 & -0.5 \\ 0.4 & -0.7 & 0.1 & 0.3 \\ 0.7 & 0.0 & -0.6 & 0.1 \\ 0.4 & -0.1 & 0.8 & -0.5 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}.$$

Suppose we want the state to be x_{des} at time T . Consider the desired state

$$x_{\text{des}} = \begin{bmatrix} 0.8 \\ 2.3 \\ -0.7 \\ -0.3 \end{bmatrix}.$$

What is the smallest T for which we can find inputs $u(0), \dots, u(T-1)$, such that $x(T) = x_{\text{des}}$? What are the corresponding inputs that achieve x_{des} at this minimum time? What is the smallest T for which we can find inputs $u(0), \dots, u(T-1)$, such that $x(T) = x_{\text{des}}$ for any $x_{\text{des}} \in \mathbf{R}^4$? We'll denote this T by T_{\min} .

- (c) Suppose the energy expended in applying inputs $u(0), \dots, u(T-1)$ is

$$E(T) = \sum_{t=0}^{T-1} (u(t))^2.$$

For a given T (greater than T_{\min}) and x_{des} , how can you compute the inputs which achieve $x(T) = x_{\text{des}}$ with the minimum expense of energy? Consider now the desired state

$$x_{\text{des}} = \begin{bmatrix} -1 \\ 1 \\ 0 \\ 1 \end{bmatrix}.$$

For each T ranging from T_{\min} to 30, find the minimum energy inputs that achieve $x(T) = x_{\text{des}}$. For each T , evaluate the corresponding input energy, which we denote by $E_{\min}(T)$. Plot $E_{\min}(T)$ as a function of T . (You should include in your solution a description of how you computed the minimum-energy inputs, and the plot of the minimum energy as a function of T . But you don't need to list the actual inputs you computed!)

- (d) You should observe that $E_{\min}(T)$ is non-increasing in T . Show that this is the case in general (*i.e.*, for any A , B , and x_{des}).

Note: There is a direct way of computing the asymptotic limit of the minimum energy as $T \rightarrow \infty$. We'll cover these ideas in more detail in ee363.

13.3 Output feedback for maximum damping. Consider the discrete-time linear dynamical system

$$\begin{aligned} x(t+1) &= Ax(t) + Bu(t), \\ y(t) &= Cx(t), \end{aligned}$$

with $A \in \mathbf{R}^{n \times n}$, $B \in \mathbf{R}^{n \times m}$, $C \in \mathbf{R}^{p \times n}$. In *output feedback control* we use an input which is a linear function of the output, that is,

$$u(t) = Ky(t),$$

where $K \in \mathbf{R}^{m \times p}$ is the *feedback gain matrix*. The resulting state trajectory is identical to that of an autonomous system,

$$x(t+1) = \bar{A}x(t).$$

- (a) Write \bar{A} explicitly in terms of A , B , C , and K .
 (b) Consider the single-input, single-output system with

$$A = \begin{bmatrix} 0.5 & 1.0 & 0.1 \\ -0.1 & 0.5 & -0.1 \\ 0.2 & 0.0 & 0.9 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad C = [0 \quad 1 \quad 0].$$

In this case, the feedback gain matrix K is a scalar (which we call simply the *feedback gain*.) The question is: find the feedback gain K_{opt} such that the feedback system is maximally damped. By maximally damped, we mean that the state goes to zero with the fastest asymptotic decay rate (measured for an initial state $x(0)$ with non-zero coefficient in the slowest mode.) *Hint:* You are only required to give your answer K_{opt} up to a precision of ± 0.01 , and you can assume that $K_{\text{opt}} \in [-2, 2]$.

- 13.4 *Affine dynamical systems.* A function $f : \mathbf{R}^n \rightarrow \mathbf{R}^m$ is called *affine* if it is a linear function plus a constant, i.e., of the form $f(x) = Ax + b$. Affine functions are more general than linear functions, which result when $b = 0$. We can generalize linear dynamical systems to *affine dynamical systems*, which have the form

$$\dot{x} = Ax + Bu + f, \quad y = Cx + Du + g.$$

Fortunately we don't need a whole new theory for (or course on) affine systems; a simple shift of coordinates converts it to a linear dynamical system. Assuming A is invertible, define $\tilde{x} = x + A^{-1}f$ and $\tilde{y} = y - g + CA^{-1}f$. Show that \tilde{x} , u , and \tilde{y} are the state, input, and output of a linear dynamical system.

- 13.5 Two separate experiments are performed for $t \geq 0$ on the single-input single-output (SISO) linear system

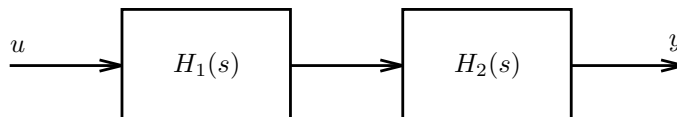
$$\dot{x} = Ax + Bu, \quad y = Cx + Du, \quad x(0) = [1 \ 2 \ -1]^T$$

(the initial condition is the same in each experiment). In the first experiment, $u(t) = e^{-t}$ and the resulting output is $y(t) = e^{-3t} + e^{-2t}$. In the second, $u(t) = e^{-3t}$ and the resulting output is $y(t) = 3e^{-3t} - e^{-2t}$.

- Can you determine the transfer function $C(sI - A)^{-1}B + D$ from this information? If it is possible, do so. If not, find two linear systems consistent with all the data given which have different transfer functions.
- Can you determine A , B , C , or D ?

- 13.6 *Cascade connection of systems.*

- Two linear systems (A_1, B_1, C_1, D_1) and (A_2, B_2, C_2, D_2) with states x_1 and x_2 (these are two *column vectors*, not two scalar components of one vector), have transfer functions $H_1(s)$ and $H_2(s)$, respectively. Find state equations for the cascade system:



Use the state $x = [x_1^T \ x_2^T]^T$.

- Use the state equations above to verify that the cascade system has transfer function $H_2(s)H_1(s)$. (To simplify, you can assume $D_1 = 0$, $D_2 = 0$.)
 - Find the dual of the LDS found in (a). Draw a block diagram of the dual system as a cascade connection of two systems. (To simplify, you can assume $D_1 = 0$, $D_2 = 0$.) *Remark:* quite generally, the block diagram corresponding to the dual system is the original block diagram, “turned around,” with all arrows reversed.
- 13.7 *Inverse of a linear system.* Suppose $H(s) = C(sI - A)^{-1}B + D$, where D is square and invertible. You will find a linear system with transfer function $H(s)^{-1}$.
- Start with $\dot{x} = Ax + Bu$, $y = Cx + Du$, and solve for \dot{x} and u in terms of x and y . Your answer will have the form: $\dot{x} = Ex + Fy$, $u = Gx + Hy$. Interpret the result as a linear system with state x , input y , and output u .

(b) Verify that

$$(G(sI - E)^{-1}F + H)(C(sI - A)^{-1}B + D) = I.$$

Hint: use the following “resolvent identity:”

$$(sI - X)^{-1} - (sI - Y)^{-1} = (sI - X)^{-1}(X - Y)(sI - Y)^{-1}$$

which can be verified by multiplying by $sI - X$ on the left and $sI - Y$ on the right.

13.8 *Offset or skewed discretization.* In the lecture notes we considered sampling a continuous-time system in which the input update and output sampling occur at the same time, *i.e.*, are synchronized. In this problem we consider what happens when there is a constant time offset or skew between them (which often happens in practice). Consider the continuous-time LDS $\dot{x} = Ax + Bu$, $y = Cx + Du$. We define the sequences x_d and y_d as

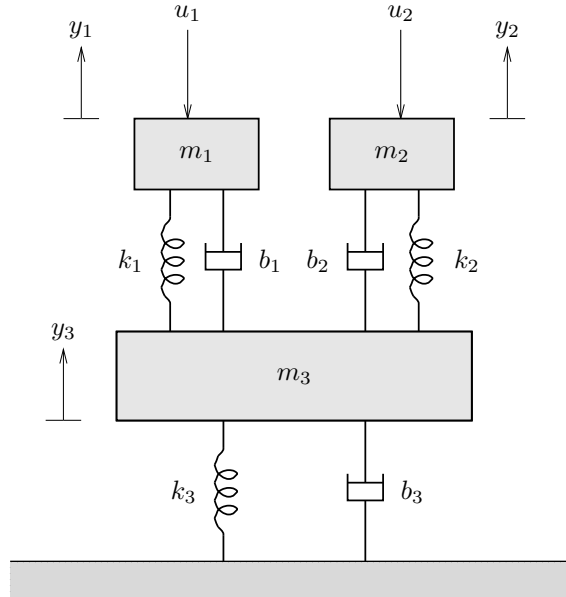
$$x_d(k) = x(kh), \quad y_d(k) = y(kh), \quad k = 0, 1, \dots$$

where $h > 0$ (*i.e.*, the state and output are sampled every h seconds). The input u is given by

$$u(t) = u_d(k) \quad \text{for} \quad kh + \delta \leq t < (k+1)h + \delta, \quad k = 0, 1, \dots$$

where δ is a delay or offset in the input update, with $0 \leq \delta < h$. Find a discrete-time LDS with u_d as input and y_d as output. Give the matrices that describe this LDS.

13.9 *Static decoupling.* Consider the mechanical system shown below.



Two masses with values $m_1 = 1$ and $m_2 = 2$ are attached via spring/damper suspensions with stiffnesses $k_1 = 1$, $k_2 = 2$ and damping $b_1 = 1$, $b_2 = 2$ to a platform, which is another mass of value $m_3 = 3$. The platform is attached to the ground by a spring/damper suspension with stiffness $k_3 = 3$ and damping $b_3 = 3$. The displacements of the masses (with respect to ground) are denoted y_1 , y_2 , and y_3 . Forces u_1 and u_2 are applied to the first two masses.

(a) Find matrices $A \in \mathbf{R}^{6 \times 6}$ and $B \in \mathbf{R}^{6 \times 2}$ such that the dynamics of the mechanical system is given by $\dot{x} = Ax + Bu$ where

$$x = [y_1 \ y_2 \ y_3 \ \dot{y}_1 \ \dot{y}_2 \ \dot{y}_3]^T, \quad u = [u_1 \ u_2]^T.$$

Ignore the effect of gravity (or you can assume the effect of gravity has already been taken into account in the definition of y_1 , y_2 and y_3).

- (b) Plot the step responses matrix, *i.e.*, the step responses from inputs u_1 and u_2 to outputs y_1 , y_2 and y_3 . Briefly interpret and explain your plots.
- (c) Find the DC gain matrix $H(0)$ from inputs u_1 and u_2 to outputs y_1 and y_2 .
- (d) *Design of an asymptotic decoupler.* In order to make the steady-state deflections of masses 1 and 2 independent of each other, we let $u = H(0)^{-1}y_{\text{cmd}}$, where $y_{\text{cmd}} : \mathbf{R}_+ \rightarrow \mathbf{R}^2$. Plot the step responses from y_{cmd} to y_1 and y_2 , and compare with the original ones found in part b.

13.10 *A method for rapidly driving the state to zero.* We consider the discrete-time linear dynamical system

$$x(t+1) = Ax(t) + Bu(t),$$

where $A \in \mathbf{R}^{n \times n}$ and $B \in \mathbf{R}^{n \times k}$, $k < n$, is full rank. The goal is to choose an input u that causes $x(t)$ to converge to zero as $t \rightarrow \infty$. An engineer proposes the following simple method: at time t , choose $u(t)$ that minimizes $\|x(t+1)\|$. The engineer argues that this scheme will work well, since the norm of the state is made as small as possible at every step. In this problem you will analyze this scheme.

- (a) Find an explicit expression for the proposed input $u(t)$ in terms of $x(t)$, A , and B .
- (b) Now consider the linear dynamical system $x(t+1) = Ax(t) + Bu(t)$ with $u(t)$ given by the proposed scheme (*i.e.*, as found in (10a)). Show that x satisfies an autonomous linear dynamical system equation $x(t+1) = Fx(t)$. Express the matrix F explicitly in terms of A and B .
- (c) Now consider a specific case:

$$A = \begin{bmatrix} 0 & 3 \\ 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

Compare the behavior of $x(t+1) = Ax(t)$ (*i.e.*, the original system with $u(t) = 0$) and $x(t+1) = Fx(t)$ (*i.e.*, the original system with $u(t)$ chosen by the scheme described above) for a few initial conditions. Determine whether each of these systems is stable.

13.11 *Analysis of investment allocation strategies.* Each year or period (denoted $t = 0, 1, \dots$) an investor buys certain amounts of one-, two-, and three-year certificates of deposit (CDs) with interest rates 5%, 6%, and 7%, respectively. (We ignore minimum purchase requirements, and assume they can be bought in any amount.)

- $B_1(t)$ denotes the amount of one-year CDs bought at period t .
- $B_2(t)$ denotes the amount of two-year CDs bought at period t .
- $B_3(t)$ denotes the amount of three-year CDs bought at period t .

We assume that $B_1(0) + B_2(0) + B_3(0) = 1$, *i.e.*, a total of 1 is to be invested at $t = 0$. (You can take $B_j(t)$ to be zero for $t < 0$.) The total payout to the investor, $p(t)$, at period t is a sum of six terms:

- $1.05B_1(t-1)$, *i.e.*, principle plus 5% interest on the amount of one-year CDs bought one year ago.
- $1.06B_2(t-2)$, *i.e.*, principle plus 6% interest on the amount of two-year CDs bought two years ago.
- $1.07B_3(t-3)$, *i.e.*, principle plus 7% interest on the amount of three-year CDs bought three years ago.
- $0.06B_2(t-1)$, *i.e.*, 6% interest on the amount of two-year CDs bought one year ago.
- $0.07B_3(t-1)$, *i.e.*, 7% interest on the amount of three-year CDs bought one year ago.
- $0.07B_3(t-2)$, *i.e.*, 7% interest on the amount of three-year CDs bought two years ago.

The total wealth held by the investor at period t is given by

$$w(t) = B_1(t) + B_2(t) + B_2(t-1) + B_3(t) + B_3(t-1) + B_3(t-2).$$

Two re-investment allocation strategies are suggested.

- *The 35-35-30 strategy.* The total payout is re-invested 35% in one-year CDs, 35% in two-year CDs, and 30% in three-year CDs. The initial investment allocation is the same: $B_1(0) = 0.35$, $B_2(0) = 0.35$, and $B_3(0) = 0.30$.
 - *The 60-20-20 strategy.* The total payout is re-invested 60% in one-year CDs, 20% in two-year CDs, and 20% in three-year CDs. The initial investment allocation is $B_1(0) = 0.60$, $B_2(0) = 0.20$, and $B_3(0) = 0.20$.
- (a) Describe the investments over time as a linear dynamical system $x(t+1) = Ax(t)$, $y(t) = Cx(t)$ with $y(t)$ equal to the total wealth at time t . Be *very clear* about what the state $x(t)$ is, and what the matrices A and C are. You will have two such linear systems: one for the 35-35-30 strategy and one for the 60-20-20 strategy.
- (b) *Asymptotic wealth growth rate.* For each of the two strategies described above, determine the asymptotic growth rate, defined as $\lim_{t \rightarrow \infty} w(t+1)/w(t)$. (If this limit doesn't exist, say so.) *Note:* simple numerical simulation of the strategies (*e.g.*, plotting $w(t+1)/w(t)$ versus t to guess its limit) is *not* acceptable. (You can, of course, simulate the strategies to *check* your answer.)
- (c) *Asymptotic liquidity.* The total wealth at time t can be divided into three parts:
- $B_1(t) + B_2(t-1) + B_3(t-2)$ is the amount that matures in one year (*i.e.*, the amount of principle we could get back next year)
 - $B_2(t) + B_3(t-1)$ is the amount that matures in two years
 - $B_3(t)$ is the amount that matures in three years (*i.e.*, is least liquid)

We define liquidity ratios as the ratio of these amounts to the total wealth:

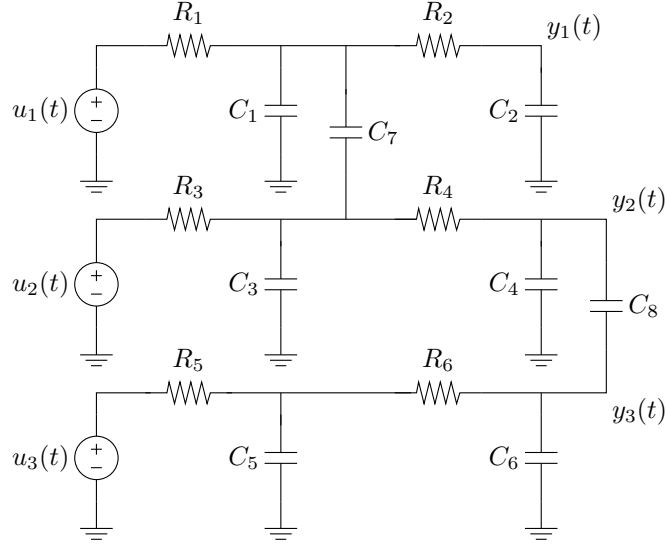
$$\begin{aligned} L_1(t) &= (B_1(t) + B_2(t-1) + B_3(t-2))/w(t), \\ L_2(t) &= (B_2(t) + B_3(t-1))/w(t), \\ L_3(t) &= B_3(t)/w(t). \end{aligned}$$

For the two strategies above, do the liquidity ratios converge as $t \rightarrow \infty$? If so, to what values?

Note: as above, simple numerical simulation alone is *not* acceptable.

- (d) Suppose you could change the *initial* investment allocation for the 35-35-30 strategy, *i.e.*, choose some other nonnegative values for $B_1(0)$, $B_2(0)$, and $B_3(0)$ that satisfy $B_1(0) + B_2(0) + B_3(0) = 1$. What allocation would you pick, and how would it be better than the (0.35, 0.35, 0.30) initial allocation? (For example, would the asymptotic growth rate be larger?) How much better is your choice of initial investment allocations? *Hint for part d:* think *very carefully* about this one. *Hint for whole problem:* watch out for nondiagonalizable, or nearly nondiagonalizable, matrices. Don't just blindly type in Matlab commands; check to make sure you're computing what you think you're computing.

13.12 *Analysis of cross-coupling in interconnect wiring.* In integrated circuits, wires which connect the output of one gate to the inputs of one (or more) other gates are called *nets*. As feature sizes shrink to well below a micron (*i.e.*, 'deep submicron') the capacitance of a wire to the substrate (which in a simple analysis can be approximated as ground), as well as to neighboring wires, must be taken into account. A simple lumped model of three nets is shown below. The inputs are the voltage sources u_1 , u_2 , u_3 , and the outputs are the three voltages labeled y_1 , y_2 , y_3 . The resistances R_1, \dots, R_6 represent the resistance of the wire segments. The capacitances C_1, \dots, C_6 are capacitances from the interconnect wires to the substrate; the capacitances C_7 and C_8 are capacitances between wires 1 and 2, and wires 2 and 3, respectively. (The different locations of these cross-coupling capacitances models the wire 1 crossing over wire 2 near the driving gate, and wire 2 crossing over wire 3 near the end of the wire, but you don't need to know this to do the problem ...) In static conditions, the circuit reduces to three wires (with resistance $R_1 + R_2$, $R_3 + R_4$, and $R_5 + R_6$, respectively) connecting the inputs to the outputs.



To simplify the problem we'll assume that all resistors have value 1 and all capacitors have value 1. We recognize that some of you don't know how to write the equations that govern this circuit, so we've done it for you. (If you're an EE student in this category, then shame on you.) The equations are

$$C\dot{v} + Gv = Fu, \quad y = Kv,$$

where

$$C = \begin{bmatrix} 2 & 0 & -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & -1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 & 0 & 2 \end{bmatrix}, \quad G = \begin{bmatrix} 2 & -1 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & -1 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & -1 \\ 0 & 0 & 0 & 0 & -1 & 1 \end{bmatrix},$$

$$F = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}, \quad K = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

and $v \in \mathbf{R}^6$ is the vector of voltages at capacitors C_1, \dots, C_6 , respectively. To save you the trouble of typing these in, we've put an mfile `interconn.m` on the course web page, which defines these matrices. The inputs (which represent the gates that drive the three nets) are Boolean valued, *i.e.*, $u_i(t) \in \{0, 1\}$ for all t . In this problem we will only consider inputs that switch (change value from 0 to 1 or 1 to 0) at most once.

- (a) *50%-threshold delay.* For $t < 0$, the system is in static condition, and the inputs have values $u(t) = f$ for $t < 0$, where $f_i \in \{0, 1\}$. At $t = 0$, the input switches to the Boolean vector g , *i.e.*, for $t \geq 0$, $u(t) = g$, where $g_i \in \{0, 1\}$. Since the DC gain matrix of this system is I , and the system is stable, the output converges to the input value: $y(t) \rightarrow g$ as $t \rightarrow \infty$. We define the 50%-threshold delay of the transition as smallest T such that $|y_i(t) - g_i| \leq 0.5$ for $t \geq T$, and for $i = 1, 2, 3$. (If the following gate thresholds were set at 0.5, then this would be first time after which the outputs would be guaranteed correct.) Among the 64 possible transitions, find the largest (*i.e.*, worst) 50%-threshold delay. Give the largest delay, and also describe which transition gives the largest delay (*e.g.*, the transition with $f = (0, 0, 1)$ to $g = (1, 0, 0)$).

- (b) *Maximum bounce due to cross-coupling.* Now suppose that input 2 remains zero, but inputs 1 and 3 undergo transitions at times $t = T_1$ and $t = T_3$, respectively. (In part 1, in contrast, all transitions occurred at $t = 0$.) To be more precise (and also so nobody can say we weren't clear),

$$u_1(t) = \begin{cases} f_1 & \text{for } t < T_1 \\ g_1 & \text{for } t \geq T_1 \end{cases}, \quad u_3(t) = \begin{cases} f_3 & \text{for } t < T_3 \\ g_3 & \text{for } t \geq T_3 \end{cases}, \quad u_2(t) = 0 \text{ for all } t,$$

where $f_1, f_3, g_1, g_3 \in \{0, 1\}$. The transitions in inputs 1 and 3 induce a nonzero response in output 2. (But y_2 does converge back to zero, since $u_2 = 0$.) This phenomenon of y_2 deviating from zero (which is what it would be if there were no cross-coupling capacitance) is called *bounce* (induced by the cross-coupling between the nets). If for any t , $y_2(t)$ is large enough to trigger the following gate, things can get very, very ugly. What is the maximum possible bounce? In other words, what is the maximum possible value of $y_2(t)$, over all possible $t, T_1, T_3, f_1, f_3, g_1, g_3$? Be sure to give not only the maximum value, but also the times t, T_1 , and T_3 , and the transitions f_1, f_3, g_1, g_3 , which maximize $y(t)$.

Note: in this problem we don't consider multiple transitions, but it's not hard to do so.

- 13.13 *Periodic solution with intermittent input.* We consider the *stable* linear dynamical system $\dot{x} = Ax + Bu$, where $x(t) \in \mathbf{R}^n$, and $u(t) \in \mathbf{R}$. The input has the specific form

$$u(t) = \begin{cases} 1 & kT \leq t < (k + \theta)T, \quad k = 0, 1, 2, \dots \\ 0 & (k + \theta)T \leq t < (k + 1)T, \quad k = 0, 1, 2, \dots \end{cases}$$

Here $T > 0$ is the *period*, and $\theta \in [0, 1]$ is called the *duty cycle* of the input. You can think of u as a constant input value one, that is applied over a fraction θ of each cycle, which lasts T seconds. Note that when $\theta = 0$, the input is $u(t) = 0$ for all t , and when $\theta = 1$, the input is $u(t) = 1$ for all t .

- Explain how to find an initial state $x(0)$ for which the resulting state trajectory is T -periodic, *i.e.*, $x(t + T) = x(t)$ for all $t \geq 0$. Give a formula for $x(0)$ in terms of the problem data, *i.e.*, A, B, T , and θ . Try to give the simplest possible formula.
- Explain why there is always *exactly one* value of $x(0)$ that results in $x(t)$ being T -periodic. In addition, explain why the formula you found in part (a) always makes sense and is valid. (For example, if your formula involves a matrix inverse, explain why the matrix to be inverted is nonsingular.)
- We now consider the specific system with

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -1 & -2 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 8 \\ 2 \\ -14 \end{bmatrix}, \quad T = 5.$$

Plot J , the mean-square norm of the state,

$$J = \frac{1}{T} \int_0^T \|x(t)\|^2 dt,$$

versus θ , for $0 \leq \theta \leq 1$, where $x(0)$ is the periodic initial condition that you found in part (a). You may approximate J as

$$J \approx \frac{1}{N} \sum_{i=0}^{N-1} \|x(iT/N)\|^2,$$

for N large enough (say 1000). Estimate the value of θ that maximizes J .

- 13.14 *System identification of a linear dynamical system.* In *system identification*, we are given some time series values for a discrete-time input vector signal,

$$u(1), u(2), \dots, u(N) \in \mathbf{R}^m,$$

and also a discrete-time state vector signal,

$$x(1), x(2), \dots, x(N) \in \mathbf{R}^n,$$

and we are asked to find matrices $A \in \mathbf{R}^{n \times n}$ and $B \in \mathbf{R}^{n \times m}$ such that we have

$$x(t+1) \approx Ax(t) + Bu(t), \quad t = 1, \dots, N-1. \quad (2)$$

We use the symbol \approx since there may be small measurement errors in the given signal data, so we don't expect to find matrices A and B for which the linear dynamical system equations hold exactly. Let's give a quantitative measure of how well the linear dynamical system model (2) holds, for a particular choice of matrices A and B . We define the RMS (root-mean-square) value of the residuals associated with our signal data and a candidate pair of matrices A, B as

$$R = \left(\frac{1}{N-1} \sum_{t=1}^{N-1} \|x(t+1) - Ax(t) - Bu(t)\|^2 \right)^{1/2}.$$

We define the RMS value of x , over the same period, as

$$S = \left(\frac{1}{N-1} \sum_{t=1}^{N-1} \|x(t+1)\|^2 \right)^{1/2}.$$

We define the *normalized residual*, denoted ρ , as $\rho = R/S$. If we have $\rho = 0.05$, for example, it means that the state equation (2) holds, roughly speaking, to within 5%. Given the signal data, we will choose the matrices A and B to minimize the RMS residual R (or, equivalently, the normalized residual ρ).

- (a) Explain how to do this. Does the method always work? If some conditions have to hold, specify them.
- (b) Carry out this procedure on the data in `lds_sysid.m` on the course web site. Give the matrices A and B , and give the associated value of the normalized residual. Of course you must show your Matlab source code and the output it produces.

- 13.15 *System identification with selection of inputs & states.* This problem continues, or rather extends, the previous one on system identification, problem 14. Here too we need to fit a linear dynamical system model to some given signal data. To complicate things, though, we are not told which of the scalar signals are input components and which are state components. That's part of what we have to decide. We are given the time series data, *i.e.*, a vector signal,

$$z(1), z(2), \dots, z(N) \in \mathbf{R}^p.$$

We will assign each component of z as either an input component, or a state component. For example, if z has four components we might assign its first and third to be the state, and its second and fourth to be the input, *i.e.*,

$$x(t) = \begin{bmatrix} z_1(t) \\ z_3(t) \end{bmatrix}, \quad u(t) = \begin{bmatrix} z_2(t) \\ z_4(t) \end{bmatrix}.$$

You can assume that we always assign at least one component to the state, so the dimension of the state is always at least one. Once we assign components of z to either x or u , we then proceed as in problem (14): we find matrices A and B that minimize the RMS residuals as defined in problem (14).

One measure of the complexity of the model is the number of components assigned to the input u ; the larger the dimension of u , the more complex the model. If the dimension of u is small, then we have a compact model, in the sense that the data are explained by a linear dynamical system driven by only a few inputs. As an extreme case, if all components of z are assigned to x , then we have an autonomous linear dynamical system model for the data, *i.e.*, one with no inputs at all. Finally, here is the problem. Get the data given in `lds_sysid2.m` on the class web server, which contains a vector $z(t) \in \mathbf{R}^8$ for $t = 1, \dots, 100$. Assign the components of z to either state or input, and develop a linear dynamical system model (*i.e.*, find matrices A and B) for your choice of x and u . We seek the simplest model, *i.e.*, the one with the smallest dimension of u , for which the normalized RMS residuals is smaller than around 5%. Your solution should consist of the following:

- *Your approach.* Explain how you solved the problem.
- *Your assignments to state and input.* Give a clear description of what x and u are. Please order the components in x and u in the same order as in z .
- Your matrices A and B .
- The relative RMS residuals obtained by your matrices.
- The Matlab code used to solve the problem, and its output.

13.16 *A greedy control scheme.* Our goal is to choose an input $u : \mathbf{R}_+ \rightarrow \mathbf{R}^m$, that is not too big, and drives the state $x : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ of the system $\dot{x} = Ax + Bu$ to zero quickly. To do this, we will choose $u(t)$, for each t , to minimize the quantity

$$\frac{d}{dt} \|x(t)\|^2 + \rho \|u(t)\|^2,$$

where $\rho > 0$ is a given parameter. The first term gives the rate of decrease (if it is negative) of the norm-squared of the state vector; the second term is a penalty for using a large input.

This scheme is greedy because at each instant t , $u(t)$ is chosen to minimize the compositive objective above, without regard for the effects such an input might have in the future.

- Show that $u(t)$ can be expressed as $u(t) = Kx(t)$, where $K \in \mathbf{R}^{m \times n}$. Give an explicit formula for K . (In other words, the control scheme has the form of a constant linear state feedback.)
- What are the conditions on A , B , and ρ under which we have $(d/dt)\|x(t)\|^2 < 0$ whenever $x(t) \neq 0$, using the scheme described above? (In other words, when does this control scheme result in the norm squared of the state always decreasing?)
- Find an example of a system (*i.e.*, A and B), for which the open-loop system $\dot{x} = Ax$ is stable, but the closed-loop system $\dot{x} = Ax + Bu$ (with u as above) is unstable, when $\rho = 1$. Try to find the simplest example you can, and be sure to show us verification that the open-loop system is stable and that the closed-loop system is not. (We will *not* check this for you. You must explain how to check this, and attach code and associated output.)

13.17 *FIR filter with small feedback.* Consider a cascade of 100 one-sample delays:

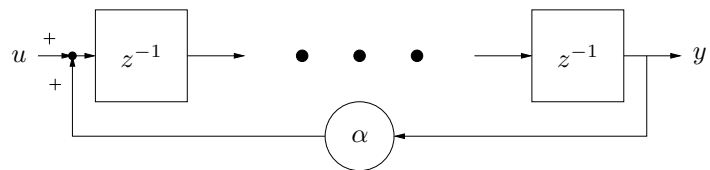


- Express this as a linear dynamical system

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t)$$

- What are the eigenvalues of A ?

- (c) Now we add simple feedback, with gain $\alpha = 10^{-5}$, to the system:



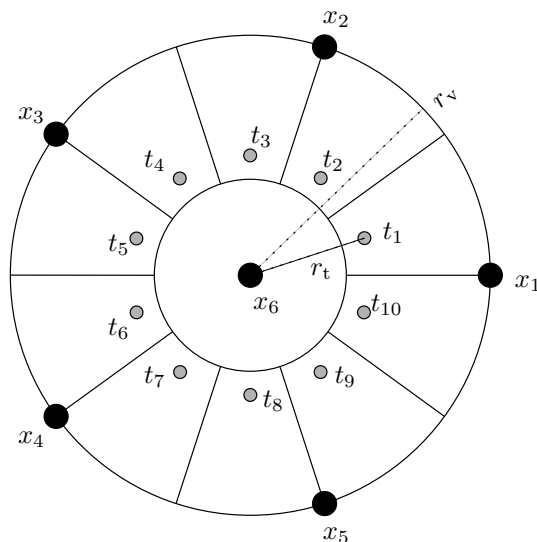
Express this as a linear dynamical system

$$x(t+1) = A_f x(t) + B_f u(t), \quad y(t) = C_f x(t) + D_f u(t)$$

- (d) What are the eigenvalues of A_f ?
- (e) How different is the impulse response of the system with feedback ($\alpha = 10^{-5}$) and without feedback ($\alpha = 0$)?

Lecture 15 – Symmetric matrices, quadratic forms, matrix norm, and SVD

- 15.1 *Simplified temperature control.* A circular room consists of 10 identical cubicles around a circular shaft. There are 6 temperature-control systems in the room. Of those, 5 have vents evenly distributed around the perimeter of the room, and one is located in the centre. Each vent j blows a stream of air at temperature x_j , measured relative to the surrounding air (ambient air temperature.) The temperatures may be hotter ($x_j > 0$) or colder ($x_j < 0$) than the ambient air temperature. The temperature in each



cubicle (measured at its center as shown in the figure) is t_i , and the effect of vent j on temperature t_i is given by

$$A_{ij} = \frac{1}{r_{ij}^2}$$

where r_{ij} is the distance between vent j and measuring point i . So the system can be described by $t = Ax$ (where A is tall.) The temperature preferences differ among the inhabitants of the 10 cubicles. More precisely, the inhabitant of cubicle i wants the temperature to be y_i hotter than the surrounding air (which is colder if $y_i < 0$!) The objective is then to choose the x_j to best match these preferences (*i.e.*, obtain *exactly* the least possible sum of squares error in temperature), with minimal cost. Here, “cost” means the total power spent on the temperature-control systems, which is the sum of the power consumed by each heater/cooler, which in turn is proportional to x_j^2 .

- How would you choose the x_j to best match the given preferences y_i , with minimal power consumption?
- Temperature measurement points are at distance r_t from the center, and vents are at distance r_v . Vent 1 lies exactly on the horizontal. The file `temp_control.m` on the course webpage defines `r_t`, `r_v`, and a preferences vector `y`. It also provides code for computing the distances from each vent to each desired temperature location. Using these data, find the optimal vent temperatures x , and the corresponding RMS error in temperature, as well as the power usage.

Comment: In this problem we ignore the fact that, in general, cooling requires more power (per unit of temperature difference) than heating ... But this was not meant to be an entirely realistic problem to start with!

- 15.2 *Norm expressions for quadratic forms.* Let $f(x) = x^T Ax$ (with $A = A^T \in \mathbf{R}^{n \times n}$) be a quadratic form.

- (a) Show that f is positive semidefinite (i.e., $A \geq 0$) if and only if it can be expressed as $f(x) = \|Fx\|^2$ for some matrix $F \in \mathbf{R}^{k \times n}$. Explain how to find such an F (when $A \geq 0$). What is the size of the smallest such F (i.e., how small can k be)?
- (b) Show that f can be expressed as a difference of squared norms, in the form $f(x) = \|Fx\|^2 - \|Gx\|^2$, for some appropriate matrices F and G . How small can the sizes of F and G be?

15.3 *Congruences and quadratic forms.* Suppose $A = A^T \in \mathbf{R}^{n \times n}$.

- (a) Let $Z \in \mathbf{R}^{n \times p}$ be any matrix. Show that $Z^T A Z \geq 0$ if $A \geq 0$.
- (b) Suppose that $T \in \mathbf{R}^{n \times n}$ is invertible. Show that $T^T A T \geq 0$ if and only if $A \geq 0$. When T is invertible, $T A T^T$ is called a *congruence* of A and $T A T^T$ and A are said to be *congruent*. This problem shows that congruences preserve positive semidefiniteness.

15.4 *Positive semidefinite (PSD) matrices.*

- (a) Show that if A and B are PSD and $\alpha \in \mathbf{R}$, $\alpha \geq 0$, then so are αA and $A + B$.
- (b) Show that any (symmetric) submatrix of a PSD matrix is PSD. (To form a symmetric submatrix, choose any subset of $\{1, \dots, n\}$ and then throw away all other columns and rows.)
- (c) Show that if $A \geq 0$, $A_{ii} \geq 0$.
- (d) Show that if $A \geq 0$, $|A_{ij}| \leq \sqrt{A_{ii} A_{jj}}$. In particular, if $A_{ii} = 0$, then the entire i th row and column of A are zero.

15.5 *A Pythagorean inequality for the matrix norm.* Suppose that $A \in \mathbf{R}^{m \times n}$ and $B \in \mathbf{R}^{p \times n}$. Show that

$$\left\| \begin{bmatrix} A \\ B \end{bmatrix} \right\| \leq \sqrt{\|A\|^2 + \|B\|^2}.$$

Under what conditions do we have equality?

15.6 *Gram matrices.* Given functions $f_i : [a, b] \rightarrow \mathbf{R}$, $i = 1, \dots, n$, the *Gram matrix* $G \in \mathbf{R}^{n \times n}$ associated with them is defined by

$$G_{ij} = \int_a^b f_i(t) f_j(t) dt.$$

- (a) Show that $G = G^T \geq 0$.
- (b) Show that G is singular if and only if the functions f_1, \dots, f_n are linearly dependent.

15.7 *Properties of symmetric matrices.* In this problem P and Q are symmetric matrices. For each statement below, either give a proof or a specific counterexample.

- (a) If $P \geq 0$ then $P + Q \geq Q$.
- (b) If $P \geq Q$ then $-P \leq -Q$.
- (c) If $P > 0$ then $P^{-1} > 0$.
- (d) If $P \geq Q > 0$ then $P^{-1} \leq Q^{-1}$.
- (e) If $P \geq Q$ then $P^2 \geq Q^2$.

Hint: you might find it useful for part (d) to prove $Z \geq I$ implies $Z^{-1} \leq I$.

15.8 Express $\sum_{i=1}^{n-1} (x_{i+1} - x_i)^2$ in the form $x^T P x$ with $P = P^T$. Is $P \geq 0$? $P > 0$?

15.9 Suppose A and B are symmetric matrices that yield the same quadratic form, i.e., $x^T A x = x^T B x$ for all x . Show that $A = B$. *Hint:* first try $x = e_i$ (the i th unit vector) to conclude that the entries of A and B on the diagonal are the same; then try $x = e_i + e_j$.

- 15.10 *A power method for computing $\|A\|$.* The following method can be used to compute the largest singular value (σ_1), and also the corresponding left and right singular vectors (u_1 and v_1) of $A \in \mathbf{R}^{m \times n}$. You can assume (to simplify) that the largest singular value of A is isolated, i.e., $\sigma_1 > \sigma_2$. Let $z(0) = a \in \mathbf{R}^n$ be nonzero, and then repeat the iteration

$$w(t) = Az(t); \quad z(t+1) = A^T w(t);$$

for $t = 1, 2, \dots$. For large t , $w(t)/\|w(t)\| \approx u_1$ and $z(t)/\|z(t)\| \approx v_1$. Analyze this algorithm. Show that it ‘usually’ works. Be very explicit about when it fails. (In practice it always works.)

- 15.11 *An invertibility criterion.* Suppose that $A \in \mathbf{R}^{n \times n}$. Show that $\|A\| < 1$ implies $I - A$ is invertible. *Interpretation:* every matrix whose distance to the identity is less than one is invertible.

- 15.12 *A bound on maximum eigenvalue for a matrix with entries smaller than one.* Suppose $A = A^T \in \mathbf{R}^{n \times n}$, with $|A_{ij}| \leq 1$, $i, j = 1, \dots, n$. How large can $\lambda_{\max}(A)$ be?

- 15.13 *Some problems involving matrix inequalities.* In the following problems you can assume that $A = A^T \in \mathbf{R}^{n \times n}$ and $B = B^T \in \mathbf{R}^{n \times n}$. We *do not*, however, assume that A or B is positive semidefinite. For $X = X^T \in \mathbf{R}^{n \times n}$, $\lambda_i(X)$ will denote its i th eigenvalue, sorted so $\lambda_1(X) \geq \lambda_2(X) \geq \dots \geq \lambda_n(X)$. As usual, the symbol \leq between symmetric matrices denotes matrix inequality (e.g., $A \leq B$ means $B - A$ is positive semidefinite). Decide whether each of the following statements is true or false. (‘True’ means the statement holds for all A and B ; ‘false’ means there is at least one pair A, B for which the statement does not hold.)

- (a) $A \geq B$ if $\lambda_i(A) \geq \lambda_i(B)$ for $i = 1, \dots, n$.
- (b) If $\{x|x^T A x \leq 1\} \subseteq \{x|x^T B x \leq 1\}$, then $A \geq B$.
- (c) If $A \leq B$, then $\{x|x^T A x \leq 1\} \subseteq \{x|x^T B x \leq 1\}$.
- (d) If the eigenvalues of A and B are the same, i.e., $\lambda_i(A) = \lambda_i(B)$ for $i = 1, \dots, n$, then there is an orthogonal matrix Q such that $A = Q^T B Q$.
- (e) If there is an orthogonal matrix Q such that $A = Q^T B Q$, then the eigenvalues of A and B are the same, i.e., $\lambda_i(A) = \lambda_i(B)$ for $i = 1, \dots, n$.
- (f) If $A \geq B$ then for all $t \geq 0$, $e^{At} \geq e^{Bt}$.
- (g) If $A \geq B$ then $A_{ij} \geq B_{ij}$ for $i, j = 1, \dots, n$.
- (h) If $A_{ij} \geq B_{ij}$ for $i, j = 1, \dots, n$, then $A \geq B$.

- 15.14 *Eigenvalues and singular values of a symmetric matrix.* Let $\lambda_1, \dots, \lambda_n$ be the eigenvalues, and let $\sigma_1, \dots, \sigma_n$ be the singular values of a matrix $A \in \mathbf{R}^{n \times n}$, which satisfies $A = A^T$. (The singular values are based on the full SVD: If $\text{Rank}(A) < n$, then some of the singular values are zero.) You can assume the eigenvalues (and of course singular values) are sorted, i.e., $\lambda_1 \geq \dots \geq \lambda_n$ and $\sigma_1 \geq \dots \geq \sigma_n$. How are the eigenvalues and singular values related?

- 15.15 *More facts about singular values of matrices.* For each of the following statements, prove it if it is true; otherwise give a *specific* counterexample. Here $X, Y, Z \in \mathbf{R}^{n \times n}$.

- (a) $\sigma_{\max}(X) \geq \max_{1 \leq i \leq n} \sqrt{\sum_{1 \leq j \leq n} |X_{ij}|^2}$.
- (b) $\sigma_{\min}(X) \geq \min_{1 \leq i \leq n} \sqrt{\sum_{1 \leq j \leq n} |X_{ij}|^2}$.
- (c) $\sigma_{\max}(XY) \leq \sigma_{\max}(X)\sigma_{\max}(Y)$.
- (d) $\sigma_{\min}(XY) \geq \sigma_{\min}(X)\sigma_{\min}(Y)$.
- (e) $\sigma_{\min}(X + Y) \geq \sigma_{\min}(X) - \sigma_{\max}(Y)$.

- 15.16 A matrix can have all entries large and yet have small gain in some directions, that is, it can have a small σ_{\min} . For example,

$$A = \begin{bmatrix} 10^6 & 10^6 \\ 10^6 & 10^6 \end{bmatrix}$$

has “large” entries while $\|A[1 \ -1]^T\| = 0$. Can a matrix have all entries small and yet have a large gain in some direction, that is, a large σ_{\max} ? Suppose, for example, that $|A_{ij}| \leq \epsilon$ for $1 \leq i, j \leq n$. What can you say about $\sigma_{\max}(A)$?

- 15.17 *Frobenius norm of a matrix.* The Frobenius norm of a matrix $A \in \mathbf{R}^{n \times n}$ is defined as $\|A\|_F = \sqrt{\text{Tr } A^T A}$. (Recall Tr is the trace of a matrix, *i.e.*, the sum of the diagonal entries.)

(a) Show that

$$\|A\|_F = \left(\sum_{i,j} |A_{ij}|^2 \right)^{1/2}.$$

Thus the Frobenius norm is simply the Euclidean norm of the matrix when it is considered as an element of \mathbf{R}^{n^2} . Note also that it is much easier to compute the Frobenius norm of a matrix than the (spectral) norm (*i.e.*, maximum singular value).

(b) Show that if U and V are orthogonal, then $\|UA\|_F = \|AV\|_F = \|A\|_F$. Thus the Frobenius norm is not changed by a pre- or post- orthogonal transformation.

(c) Show that $\|A\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_r^2}$, where $\sigma_1, \dots, \sigma_r$ are the singular values of A . Then show that $\sigma_{\max}(A) \leq \|A\|_F \leq \sqrt{r} \sigma_{\max}(A)$. In particular, $\|Ax\| \leq \|A\|_F \|x\|$ for all x .

- 15.18 *Drawing a graph.* We consider the problem of drawing (in two dimensions) a graph with n vertices (or nodes) and m undirected edges (or links). This just means assigning an x - and a y - coordinate to each node. We let $x \in \mathbf{R}^n$ be the vector of x - coordinates of the nodes, and $y \in \mathbf{R}^n$ be the vector of y - coordinates of the nodes. When we draw the graph, we draw node i at the location $(x_i, y_i) \in \mathbf{R}^2$. The problem, of course, is to make the drawn graph look good. One goal is that neighboring nodes on the graph (*i.e.*, ones connected by an edge) should not be too far apart as drawn. To take this into account, we will choose the x - and y -coordinates so as to minimize the objective

$$J = \sum_{i < j, i \sim j} ((x_i - x_j)^2 + (y_i - y_j)^2),$$

where $i \sim j$ means that nodes i and j are connected by an edge. The objective J is precisely the sum of the squares of the lengths (in \mathbf{R}^2) of the drawn edges of the graph. We have to introduce some other constraints into our problem to get a sensible solution. First of all, the objective J is not affected if we shift all the coordinates by some fixed amount (since J only depends on differences of the x - and y -coordinates). So we can assume that

$$\sum_{i=1}^n x_i = 0, \quad \sum_{i=1}^n y_i = 0,$$

i.e., the sum (or mean value) of the x - and y -coordinates is zero. These two equations ‘center’ our drawn graph. Another problem is that we can minimize J by putting all the nodes at $x_i = 0, y_i = 0$, which results in $J = 0$. To force the nodes to spread out, we impose the constraints

$$\sum_{i=1}^n x_i^2 = 1, \quad \sum_{i=1}^n y_i^2 = 1, \quad \sum_{i=1}^n x_i y_i = 0.$$

The first two say that the variance of the x - and y - coordinates is one; the last says that the x - and y -coordinates are uncorrelated. (You don’t have to know what variance or uncorrelated mean; these are

just names for the equations given above.) The three equations above enforce ‘spreading’ of the drawn graph. Even with these constraints, the coordinates that minimize J are not unique. For example, if x and y are any set of coordinates, and $Q \in \mathbf{R}^{2 \times 2}$ is any orthogonal matrix, then the coordinates given by

$$\begin{bmatrix} \tilde{x}_i \\ \tilde{y}_i \end{bmatrix} = Q \begin{bmatrix} x_i \\ y_i \end{bmatrix}, \quad i = 1, \dots, n$$

satisfy the centering and spreading constraints, and have the same value of J . This means that if you have a proposed set of coordinates for the nodes, then by rotating or reflecting them, you get another set of coordinates that is just as good, according to our objective. We’ll just live with this ambiguity. Here’s the question:

- (a) Explain how to solve this problem, *i.e.*, how to find x and y that minimize J subject to the centering and spreading constraints, given the graph topology. You can use any method or ideas we’ve encountered in the course. Be clear as to whether your approach solves the problem exactly (*i.e.*, finds a set of coordinates with J as small as it can possibly be), or whether it’s just a good heuristic (*i.e.*, a choice of coordinates that achieves a reasonably small value of J , but perhaps not the absolute best). In describing your method, you may not refer to any Matlab commands or operators; your description must be entirely in mathematical terms.
- (b) Implement your method, and carry it out for the graph given in `dg_data.m`. This mfile contains the *node adjacency matrix* of the graph, denoted A , and defined as $A_{ij} = 1$ if nodes i and j are connected by an edge, and $A_{ij} = 0$ otherwise. (The graph is undirected, so A is symmetric. Also, we do not have self-loops, so $A_{ii} = 0$, for $i = 1, \dots, n$.) Draw your final result using the commands `gplot(A, [x y], 'o-')`, which plots the graph and `axis('square')`, which sets the x - y -scales to be equal. Give the value of J achieved by your choice of x and y . Verify that your x and y satisfy the centering and spreading conditions, at least approximately. If your method is iterative, plot the value of J versus iteration. The mfile `dg_data.m` also contains the vectors `x_circ` and `y_circ`. These coordinates are obtained using a standard technique for drawing a graph, by placing the nodes in order on a circle. The radius of the circle is chosen so that `x_circ` and `y_circ` satisfy the centering and spread constraints. You can draw the given graph this way using `gplot(A, [x_circ y_circ], 'o-'); axis('square');`.

Hint. You are welcome to use the results described below, without proving them. Let $A \in \mathbf{R}^{n \times n}$ be symmetric, with eigenvalue decomposition $A = \sum_{i=1}^n \lambda_i q_i q_i^T$, with $\lambda_1 \geq \dots \geq \lambda_n$, and $\{q_1, \dots, q_n\}$ orthonormal. You know that a solution of the problem

$$\begin{array}{ll} \text{minimize} & x^T A x \\ \text{subject to} & x^T x = 1, \end{array}$$

where the variable is $x \in \mathbf{R}^n$, is $x = q_n$. The related maximization problem is

$$\begin{array}{ll} \text{maximize} & x^T A x \\ \text{subject to} & x^T x = 1, \end{array}$$

with variable $x \in \mathbf{R}^n$. A solution to this problem is $x = q_1$. Now consider the following generalization of the first problem:

$$\begin{array}{ll} \text{minimize} & \text{Tr}(X^T A X) \\ \text{subject to} & X^T X = I_k, \end{array}$$

where the variable is $X \in \mathbf{R}^{n \times k}$, and I_k denotes the $k \times k$ identity matrix, and we assume $k \leq n$. The constraint means that the columns of X , say, x_1, \dots, x_k , are orthonormal; the objective can be written in terms of the columns of X as $\text{Tr}(X^T A X) = \sum_{i=1}^k x_i^T A x_i$. A solution of this problem is $X = [q_{n-k+1} \dots q_n]$. Note that when $k = 1$, this reduces to the first problem above. The related maximization problem is

$$\begin{array}{ll} \text{maximize} & \text{Tr}(X^T A X) \\ \text{subject to} & X^T X = I_k, \end{array}$$

with variable $X \in \mathbf{R}^{n \times k}$. A solution of this problem is $X = [q_1 \cdots q_k]$.

- 15.19 *Approximate left inverse with norm constraints.* Suppose $A \in \mathbf{R}^{m \times n}$ is full rank with $m \geq n$. We seek a matrix $F \in \mathbf{R}^{n \times m}$ that minimizes $\|I - FA\|$ subject to the constraint $\|F\| \leq \alpha$, where $\alpha > 0$ is given. Note that $\|I - FA\|$ gives a measure of how much F fails to be a left inverse of A . Give an explicit description of an optimal F . Your description can involve standard matrix operations and decompositions (eigenvector/eigenvalue, QR, SVD, ...).
- 15.20 *Finding worst-case inputs.* The single-input, single output system $x(t+1) = Ax(t) + Bu(t)$, $y(t) = Cx(t)$, $x(0) = 0$, where

$$A = \begin{bmatrix} 0.9 & 0.5 \\ -0.5 & 0.7 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 2 \end{bmatrix},$$

is a very simple (discretized and lumped) dynamical model of a building. The input u is ground displacement (during an earthquake), and y gives the displacement of the top of the building. The

input u is known to satisfy $\sum_{t=0}^{49} u(t)^2 \leq 1$ and $u(t) = 0$ for $t \geq 50$, i.e., the earthquake has energy less than one, and only lasts 50 samples.

- (a) How large can $\sum_{t=0}^{99} y(t)^2$ be? Plot an input u that maximizes $\sum_{t=0}^{99} y(t)^2$, along with the resulting output y .
- (b) How large can $|y(100)|$ be? Plot an input u that maximizes $|y(100)|$, along with the resulting output y .

As usual, you must explain how you solve the problem, as well as give explicit numerical answers and plots.

- 15.21 *Worst and best direction of excitation for a suspension system.* A suspension system is connected at one end to a base (that can move or vibrate) and at the other to the load (that it is supposed to isolate from vibration of the base). Suitably discretized, the system is described by

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t), \quad x(0) = 0,$$

where $u(t) \in \mathbf{R}^3$ represents the (x-, y-, and z- coordinates of the) displacement of base, and $y(t) \in \mathbf{R}^3$ represents the (x-, y-, and z- coordinates of the) displacement of the load. The input u has the form $u(t) = qv(t)$, where $q \in \mathbf{R}^3$ is a (constant) vector with $\|q\| = 1$, and $v(t) \in \mathbf{R}$ gives the displacement amplitude versus time. In other words, the driving displacement u is always in the direction q , with amplitude given by the (scalar) signal v . The response of the system is judged by the RMS deviation of the load over a 100 sample interval, i.e.,

$$D = \left(\frac{1}{100} \sum_{t=1}^{100} \|y(t)\|^2 \right)^{1/2}.$$

The data $A, B, C, v(0), \dots, v(99)$ are known (and available in the mfile `worst_susp_data.m` on the course web site). The problem is to find the direction $q_{\max} \in \mathbf{R}^3$ that maximizes D , and the direction $q_{\min} \in \mathbf{R}^3$ that minimizes D . Give the directions and the associated values of D (D_{\max} and D_{\min} , respectively).

- 15.22 *Two representations of an ellipsoid.* In the lectures, we saw two different ways of representing an ellipsoid, centered at 0, with non-zero volume. The first uses a quadratic form:

$$\mathcal{E}_1 = \{x \mid x^T S x \leq 1\},$$

with $S^T = S > 0$. The second is as the image of a unit ball under a linear mapping:

$$\mathcal{E}_2 = \{ y \mid y = Ax, \|x\| \leq 1 \},$$

with $\det(A) \neq 0$.

- (a) Given S , explain how to find an A so that $\mathcal{E}_1 = \mathcal{E}_2$.
- (b) Given A , explain how to find an S so that $\mathcal{E}_1 = \mathcal{E}_2$.
- (c) What about uniqueness? Given S , explain how to find *all* A that yield $\mathcal{E}_1 = \mathcal{E}_2$. Given A , explain how to find *all* S that yield $\mathcal{E}_1 = \mathcal{E}_2$.

15.23 *Determining initial bacteria populations.* We consider a population that consists of three strains of a bacterium, called strain 1, strain 2, and strain 3. The vector $x(t) \in \mathbf{R}^3$ will denote the amount, or *biomass* (in grams) of the strains present in the population at time t , measured in hours. For example, $x_2(3.4)$ denotes the amount of strain 2 (in grams) in the sample at time $t = 3.4$ hours. Over time, the biomass of each strain changes through several mechanisms including cell division, cell death, and mutation. (But you don't need to know any biology to answer this question!) The population dynamics is given by $\dot{x} = Ax$, where

$$A = \begin{bmatrix} -0.1 & 0.3 & 0 \\ 0 & -0.2 & 0.1 \\ 0.1 & 0 & -0.1 \end{bmatrix}.$$

You can assume that we always have $x_i(t) > 0$, *i.e.*, the biomass of each strain is always positive. The total biomass at time t is given by $\mathbf{1}^T x(t) = x_1(t) + x_2(t) + x_3(t)$, where $\mathbf{1} \in \mathbf{R}^3$ denotes the vector with all components one.

- (a) Give a very brief interpretation of the entries of the matrix A . For example, what is the significance of $a_{13} = 0$? What is the significance of the sign of a_{11} ? Limit yourself to 100 words. You may use phrases such as ‘the presence of strain i enhances (or inhibits) growth of strain j ’.
- (b) As $t \rightarrow \infty$, does the total biomass converge to ∞ (*i.e.*, grow without bound), converge to zero, or not converge at all (for example, oscillate)? Explain how you arrive at your conclusion and show any calculations (by hand or Matlab) that you need to do. You can assume that $x_i(0) > 0$ for $i = 1, 2, 3$. *Posterior intuitive explanation.* In 100 words or less, give a plausible story that explains, intuitively, the result you found.
- (c) *Selection of optimal assay time.* A biologist wishes to estimate the original biomass of each of the three strains, *i.e.*, the vector $x(0) \in \mathbf{R}^3$, based on measurements of the total biomass taken at $t = 0$, $t = 10$, and $t = T$, where T satisfies $0 < T < 10$. The three measurements of total biomass (which are called *assays*) will include a small additive error, denoted v_1 (for the assay at $t = 0$), v_2 (for the assay at $t = T$ and v_3 (for the assay at $t = 10$). You can assume that $v_1^2 + v_2^2 + v_3^2 \leq 0.01^2$, *i.e.*, the sum of the squares of the measurement errors is not more than 0.01^2 . You can also assume that a good method for computing the estimate of $x(0)$, given the measurements, will be used. (The estimation method won't make any use of the information that $x_i(0) > 0$.) The problem here is to choose the time T of the intermediate assay in such a way that the estimate of $x(0)$, denoted $\hat{x}(0)$, is as accurate as possible. We'll judge accuracy by the maximum value that $\|\hat{x}(0) - x(0)\|$ can have, over all measurement errors that satisfy $v_1^2 + v_2^2 + v_3^2 \leq 0.01^2$. Find the optimal value for T (of course, between 0 and 10), *i.e.*, the value of T that minimizes the maximum value $\|\hat{x}(0) - x(0)\|$ can have. We are looking for an answer that is accurate to within ± 0.1 . Of course you must explain exactly what you are doing, and submit your Matlab code as well the output it produces. Be sure to say what the optimal T is, and what the optimal accuracy is (*i.e.*, what the maximum value $\|\hat{x}(0) - x(0)\|$ is, for the T you choose).

- 15.24 *A measure of connectedness in a graph.* We consider an undirected graph with n nodes, described by its adjacency matrix $A \in \mathbf{R}^{n \times n}$, defined by

$$A_{ij} = \begin{cases} 1 & \text{if there is a link connecting nodes } i \text{ and } j \\ 0 & \text{otherwise.} \end{cases}$$

We assume the graph has no self-loops, *i.e.*, $A_{ii} = 0$. Note that $A = A^T$. We assume that the graph has at least one link, so $A \neq 0$. A *path* from node i to node j , of length $m > 0$, is an $m + 1$ -long sequence of nodes, connected by links, that start at i and end at j . More precisely it is a sequence $i = k_1, k_2, \dots, k_{m+1} = j$, with the property that $A_{k_1, k_2} = \dots = A_{k_m, k_{m+1}} = 1$. Note that a path can include loops; there is no requirement that each node be visited only once. For example, if node 3 and node 5 are connected by a link (*i.e.*, $A_{35} = 1$), then the sequence 3, 5, 3, 5 is a path between node 3 and node 5 of length 3. We say that each node is connected to itself by a path of length zero. Let $P_m(i, j)$ denote the total number of paths of length m from node i to node j . We define

$$C_{ij} = \lim_{m \rightarrow \infty} \frac{P_m(i, j)}{\sum_{i, j=1}^n P_m(i, j)},$$

when the limits exist. When the limits don't, we say that C_{ij} isn't defined. In the fraction in this equation, the numerator is the number of paths of length m between nodes i and j , and the denominator is the total number of paths of length m , so the ratio gives the fraction of all paths of length m that go between nodes i and j . When C_{ij} exists, it gives the asymptotic fraction of all (long) paths that go from node i to node j . The number C_{ij} gives a good measure of how "connected" nodes i and j are in the graph. You can make *one* of the following assumptions:

- (a) A is full rank.
- (b) A has distinct eigenvalues.
- (c) A has distinct singular values.
- (d) A is diagonalizable.
- (e) A has a dominant eigenvalue, *i.e.*, $|\lambda_1| > |\lambda_i|$ for $i = 2, \dots, n$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A .

(Be very clear about which one you choose.) Using your assumption, explain why C_{ij} exists, and derive an expression for C_{ij} . You can use any of the concepts from the class, such as singular values and singular vectors, eigenvalues and eigenvectors, pseudo-inverse, etc., but you cannot leave a limit in your expression. You must explain why your expression is, in fact, equal to C_{ij} .

- 15.25 *Recovering an ellipsoid from boundary points.* You are given a set of vectors $x^{(1)}, \dots, x^{(N)} \in \mathbf{R}^n$ that are thought to lie on or near the surface of an ellipsoid centered at the origin, which we represent as

$$\mathcal{E} = \{x \in \mathbf{R}^n \mid x^T A x = 1\},$$

where $A = A^T \in \mathbf{R}^{n \times n} \geq 0$. Your job is to recover, at least approximately, the matrix A , given the observed data $x^{(1)}, \dots, x^{(N)}$. Explain your approach to this problem, and then carry it out on the data given in the mfile `ellip_bdry_data.m`. Be sure to explain how you check that the ellipsoid you find is reasonably consistent with the given data, and also that the matrix A you find does, in fact, correspond to an ellipsoid. To simplify the explanation, you can give it for the case $n = 4$ (which is the dimension of the given data). But it should be clear from your discussion how it works in general.

- 15.26 *Predicting zero crossings.* We consider a linear system of the form

$$\dot{x} = Ax, \quad y = Cx,$$

where $x(t) \in \mathbf{R}^n$ and $y(t) \in \mathbf{R}$. We know A and C , but we do not know $x(0)$. We cannot directly observe the output, but we do know the times at which the output is zero, *i.e.*, we are given the *zero-crossing times* t_1, \dots, t_p at which $y(t_i) = 0$. You can assume these times are given in increasing order, *i.e.*, $0 \leq t_1 < \dots < t_p$, and that $y(t) \neq 0$ for $0 \leq t < t_p$ and $t \neq t_1, \dots, t \neq t_p$. (Note that this definition of zero-crossing times doesn't require the output signal to *cross* the value zero; it is enough to just have the value zero.) We are interested in the following question: given A , C , and the zero-crossing times t_1, \dots, t_p , can we *predict* the next zero-crossing time t_{p+1} ? (This means, of course, that $y(t) \neq 0$ for $t_p < t < t_{p+1}$, and $y(t_{p+1}) = 0$.) You will answer this question for the specific system

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -18 & -11 & -12 & -2 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix},$$

and zero-crossing times

$$t_1 = 0.000, \quad t_2 = 1.000, \quad t_3 = 2.000, \quad t_4 = 3.143.$$

(So here we have $p = 4$.) Note that the zero-crossing times are given to three significant digits. Specifically, you must do one of the following:

- If you think that you can determine the next zero-crossing time t_5 , explain in detail how to do it, and find the next time t_5 (to at least two significant figures).
- If you think that you cannot determine the next zero-crossing time t_5 , explain in detail why, and find two trajectories of the system which have t_1, \dots, t_4 as the first 4 zero-crossings, but have different 5th zero-crossings. (The zero-crossings should differ substantially, and not just in the last significant digit.)

Be sure to make it clear which one of these options you choose. *Hint:* Be careful estimating rank or determining singularity, if that's part of your procedure; remember that the zero-crossing times are only given to three significant figures.

15.27 *Optimal time compression equalizer.* We are given the (finite) impulse response of a communications channel, *i.e.*, the real numbers

$$c_1, c_2, \dots, c_n.$$

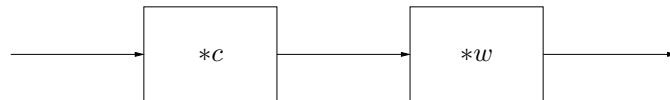
Our goal is to design the (finite) impulse response of an equalizer, *i.e.*, the real numbers

$$w_1, w_2, \dots, w_n.$$

(To make things simple, the equalizer has the same length as the channel.) The equalized channel response h is given by the convolution of w and c , *i.e.*,

$$h_i = \sum_{j=1}^{i-1} w_j c_{i-j}, \quad i = 2, \dots, 2n,$$

where we take w_i and c_i to be zero for $i \leq 0$ or $i > n$. This is shown below.



The goal is to choose w so that most of the energy of the equalized impulse response h is *concentrated* within k samples of $t = n + 1$, where $k < n - 1$ is given. To define this formally, we first define the total energy of the equalized response as

$$E_{\text{tot}} = \sum_{i=2}^{2n} h_i^2,$$

and the energy in the desired time interval as

$$E_{\text{des}} = \sum_{i=n+1-k}^{n+1+k} h_i^2.$$

For any w for which $E_{\text{tot}} > 0$, we define the *desired to total energy ratio*, or DTE, as $\text{DTE} = E_{\text{des}}/E_{\text{tot}}$. Thus number is clearly between 0 and 1; it tells us what fraction of the energy in h is contained in the time interval $t = n + 1 - k, \dots, t = n + 1 + k$. You can assume that h is such that for any $w \neq 0$, we have $E_{\text{tot}} > 0$.

- (a) How do you find a $w \neq 0$ that maximizes DTE? You must give a very clear description of your method, and explain why it works. Your description and justification must be *very clear*. You can appeal to any concepts used in the class, *e.g.*, least-squares, least-norm, eigenvalues and eigenvectors, singular values and singular vectors, matrix exponential, and so on.
- (b) Carry out your method for time compression length $k = 1$ on the data found in `time_comp_data.m`. Plot your solution w , the equalized response h , and give the DTE for your w .

Please note: You do not need to know anything about equalizers, communications channels, or even convolution; everything you need to solve this problem is clearly defined in the problem statement.

- 15.28 *Minimum energy required to leave safe operating region.* We consider the stable controllable system $\dot{x} = Ax + Bu$, $x(0) = 0$, where $x(t) \in \mathbf{R}^n$ and $u(t) \in \mathbf{R}^m$. The input u is beyond our control, but we have some idea of how large its total energy

$$\int_0^\infty \|u(\tau)\|^2 d\tau$$

is likely to be. The *safe operating region* for the state is the unit ball

$$\mathcal{B} = \{ x \mid \|x\| \leq 1 \}.$$

The hope is that input u will not drive the state outside the safe operating region. One measure of system security that is used is the minimum energy E_{min} that is required to drive the state outside the safe operating region:

$$E_{\text{min}} = \min \left\{ \int_0^t \|u(\tau)\|^2 d\tau \mid x(t) \notin \mathcal{B} \right\}.$$

(Note that we do not specify t , the time at which the state is outside the safe operating region.) If E_{min} is much larger than the energy of the u 's we can expect, we can be fairly confident that the state will not leave the safe operating region. (E_{min} can also be justified as a system security measure on statistical grounds, but we won't go into that here.)

- (a) Find E_{min} explicitly. Your solution should be in terms of the matrices A , B , or other matrices derived from them such as the controllability matrix \mathcal{C} , the controllability Gramian W_c , and its inverse $P = W_c^{-1}$. Make sure you give the simplest possible expression for E_{min} .
- (b) Suppose the safe operating region is the unit cube $C = \{ x \mid |x_i| \leq 1, i = 1, \dots, n \}$ instead of the unit ball \mathcal{B} . Let $E_{\text{min}}^{\text{cube}}$ denote the minimum energy required to drive the state outside the unit cube C . Repeat part (a) for $E_{\text{min}}^{\text{cube}}$.

- 15.29 *Energy storage efficiency in a linear dynamical system.* We consider the discrete-time linear dynamic system

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t),$$

where $x(t) \in \mathbf{R}^n$, and $u(t), y(t) \in \mathbf{R}$. The initial state is zero, *i.e.*, $x(0) = 0$. We apply an input sequence $u(0), \dots, u(N-1)$, and are interested in the output over the next N samples, *i.e.*, $y(N), \dots, y(2N-1)$. (We take $u(t) = 0$ for $t \geq N$.) We define the *input energy* as

$$\mathcal{E}_{\text{in}} = \sum_{t=0}^{N-1} u(t)^2,$$

and similarly, the output energy is defined as

$$\mathcal{E}_{\text{out}} = \sum_{t=N}^{2N-1} y(t)^2.$$

How would you choose the (nonzero) input sequence $u(0), \dots, u(N-1)$ to maximize the ratio of output energy to input energy, *i.e.*, to maximize $\mathcal{E}_{\text{out}}/\mathcal{E}_{\text{in}}$? What is the maximum value the ratio $\mathcal{E}_{\text{out}}/\mathcal{E}_{\text{in}}$ can have?

- 15.30 *Energy-optimal evasion.* A vehicle is governed by the following discrete-time linear dynamical system equations:

$$x(t+1) = Ax(t) + Bu(t), \quad y(t) = Cx(t), \quad x(0) = 0.$$

Here $x(t) \in \mathbf{R}^n$ is the vehicle state, $y(t) \in \mathbf{R}^3$ is the vehicle position, and $u(t) \in \mathbf{R}^m$ is the input signal. (The vehicle dynamics are really continuous; the equation above is the result of a suitable sampling.) The system is controllable.

- (a) *Minimum energy to reach a position.* Find the input $u(0), \dots, u(T-1)$ that reaches position $f \in \mathbf{R}^3$ at time T (where $T \geq n$), *i.e.*, $y(T) = f$, and minimizes the input ‘energy’

$$\|u(0)\|^2 + \dots + \|u(T-1)\|^2.$$

The input u is the (energy) optimal input for the vehicle to arrive at the position f at time T . Give an expression for E , the energy of the minimum energy input. (Of course E will depend on the data A, B, C , and f .)

- (b) *Energy-optimal evasion.* Now consider a second vehicle governed by

$$z(t+1) = Fz(t) + Gv(t), \quad w(t) = Hz(t), \quad z(0) = 0$$

where $z(t) \in \mathbf{R}^n$ is the state of the vehicle, $w(t) \in \mathbf{R}^3$ is the vehicle position, and $v(t) \in \mathbf{R}^m$ is the input signal. This vehicle is to be overtaken (intercepted) by the first vehicle at time T , where $T \geq n$. This means that $w(T) = y(T)$. How would you find $v(0), \dots, v(T-1)$ that maximizes the minimum energy the first vehicle must expend to intercept the second vehicle at time T , subject to a limit on input energy,

$$\|v(0)\|^2 + \dots + \|v(T-1)\|^2 \leq 1?$$

The input v is maximally evasive, in the sense that it requires the first vehicle to expend the largest amount of input energy to overtake it, given the limit on input energy the second vehicle is allowed to use. Express your answer in terms of the data A, B, C, F, G, H , and standard matrix functions (inverse, transpose, SVD, ...). *Remark:* This problem is obviously not a very realistic model of a real pursuit-evasion situation, for several reasons: both vehicles start from the zero initial state, the time of intercept (T) is known to the second vehicle, and the place of intercept ($w(T)$) is known ahead of time to the first vehicle. Still, it’s possible to extend the results of this problem to handle a realistic model of a pursuit/evasion.

15.31 *Worst-case analysis of impact.* We consider a (time-invariant) linear dynamical system

$$\dot{x} = Ax + Bu, \quad x(0) = x_{\text{init}},$$

with state $x(t) \in \mathbf{R}^n$, and input $u(t) \in \mathbf{R}^m$. We are interested in the state trajectory over the time interval $[0, T]$. In this problem the input u represents an *impact* on the system, so it has the form

$$u(t) = g\delta(t - T_{\text{imp}}),$$

where $g \in \mathbf{R}^m$ is a vector that gives the direction and magnitude of the impact, and T_{imp} is the time of the impact. We assume that $0 \leq T_{\text{imp}} \leq T_-$. ($T_{\text{imp}} = T_-$ means that the impact occurs right at the end of the period of interest, and does affect $x(T)$.) We let $x_{\text{nom}}(T)$ denote the state, at time $t = T$, of the linear system $\dot{x}_{\text{nom}} = Ax_{\text{nom}}$, $x_{\text{nom}}(0) = x_{\text{init}}$. The vector $x_{\text{nom}}(T)$ is what the final state $x(T)$ of the system above would have been at time $t = T$, had the impact not occurred (*i.e.*, with $u = 0$). We are interested in the deviation D between $x(T)$ and $x_{\text{nom}}(T)$, as measured by the norm:

$$D = \|x(T) - x_{\text{nom}}(T)\|.$$

D measures how far the impact has shifted the state at time T . We would like to know how large D can be, over all possible impact directions and magnitudes no more than one (*i.e.*, $\|g\| \leq 1$), and over all possible impact times between 0 and T_- . In other words, we would like to know the maximum possible state deviation, at time T , due to an impact of magnitude no more than one. We'll call the choices of T_{imp} and g that maximize D the *worst-case impact time* and *worst-case impact vector*, respectively.

- (a) Explain how to find the worst-case impact time, and the worst-case impact vector, given the problem data A , B , x_{init} , and T . Your explanation should be as short and clear as possible. You can use any of the concepts we have encountered in the class. Your approach can include a simple numerical search (such as plotting a function of one variable to find its maximum), if needed. If either the worst-case impact time or the worst-case impact vector do not depend on some of the problem data (*i.e.*, A , B , x_{init} , and T) say so.
- (b) Get the data from `wc_impact_data.m`, which defines A , B , x_{init} , and T , and carry out the procedure described in part (a). Be sure to give us the worst-case impact time (with absolute precision of 0.01), the worst-case impact vector, and the corresponding value of D .

15.32 *Worst time for control system failure.* In this problem we consider a system that under normal circumstances is governed by the equations

$$\dot{x}(t) = Ax(t) + Bu(t), \quad u(t) = Kx(t). \quad (3)$$

(This is called *state feedback*, and is very commonly used in automatic control systems.) Here the application is a regulator, which means that input u is meant to drive the state to zero as $t \rightarrow \infty$, no matter what $x(0)$ is. At time $t = T_f$, however, a fault occurs, and the input signal becomes zero. The fault is cleared (*i.e.*, corrected) T_c seconds after it occurs. Thus, for $T_f \leq t \leq T_f + T_c$, we have $\dot{x}(t) = Ax(t)$; for $t < T_f$ and $t > T_f + T_c$, the system is governed by the equations (3). You'll find the specific matrices A , B , and K , in the mfile `fault_ctrl_sys.m` on the class web site. Here's the problem: suppose the system fails for one second, some time in the time interval $[0, 9]$. In other words, we have $0 \leq T_f \leq 9$, and $T_c = 1$. We don't know what $x(0)$ is, but you can assume that $\|x(0)\| \leq 1$. We also don't know the time of failure T_f . The problem is to find the time of failure T_f (in $[0, 9]$) and the initial condition $x(0)$ (with $\|x(0)\| \leq 1$) that maximizes $\|x(10)\|$. In essence, you are carrying out a worst-case analysis of the effects of a one second control system failure. As usual, you must explain your approach clearly and completely. You must also give your source code, and the results, *i.e.*, the worse possible $x(0)$, the worst failure time T_f , and the resulting value of $\|x(10)\|$. An accuracy of 0.01 for T_f is fine.

15.33 *Some proof or counterexample questions.* Determine if the following statements are true or false. If the statement is true, prove it; if you think it is false, provide a *specific* (numerical) counterexample. You get five points for the correct solution (*i.e.*, the right answer and a valid proof or counterexample), two points for the right answer (*i.e.*, true or false), and zero points otherwise. What we mean by “true” is that the statement is true for all values of the matrices and vectors given. (You can assume the entries of the matrices and vectors are all real.) You can’t assume anything about the dimensions of the matrices (unless it’s explicitly stated), but you can assume that the dimensions are such that all expressions make sense. For example, the statement “ $A + B = B + A$ ” is true, because no matter what the dimensions of A and B (which must, however, be the same), and no matter what values A and B have, the statement holds. As another example, the statement $A^2 = A$ is false, because there are (square) matrices for which this doesn’t hold. In such a case, provide a specific counterexample, for example, $A = 2$ (which is a matrix in $\mathbf{R}^{1 \times 1}$).

- (a) Suppose $A = A^T \in \mathbf{R}^{n \times n}$ satisfies $A \geq 0$, and $A_{kk} = 0$ for some k (between 1 and n). Then A is singular.
- (b) Suppose $A, B \in \mathbf{R}^{n \times n}$, with $\|A\| > \|B\|$. Then, for all $k \geq 1$, $\|A^k\| \geq \|B^k\|$.
- (c) Suppose \tilde{A} is a submatrix of a matrix $A \in \mathbf{R}^{m \times n}$. (This means \tilde{A} is obtained from A by removing some rows and columns from A ; as an extreme case, any element of A is a (1×1) submatrix of A .) Then $\|\tilde{A}\| \leq \|A\|$.
- (d) For any A, B, C, D with compatible dimensions (see below),

$$\left\| \begin{bmatrix} A & B \\ C & D \end{bmatrix} \right\| \leq \left\| \begin{bmatrix} \|A\| & \|B\| \\ \|C\| & \|D\| \end{bmatrix} \right\|.$$

Compatible dimensions means: A and B have the same number of rows, C and D have the same number of rows, A and C have the same number of columns, and B and D have the same number of columns.

- (e) For any A and B with the same number of columns, we have

$$\max\{\|A\|, \|B\|\} \leq \left\| \begin{bmatrix} A \\ B \end{bmatrix} \right\| \leq \sqrt{\|A\|^2 + \|B\|^2}.$$

- (f) Suppose the fat (including, possibly, square) and full rank matrices A and B have the same number of rows. Then we have $\kappa(A) \leq \kappa(\begin{bmatrix} A & B \end{bmatrix})$, where $\kappa(Z)$ denotes, as usual, the condition number of the matrix Z , *i.e.*, the ratio of the largest to the smallest singular value.

15.34 *Uncovering a hidden linear explanation.* Consider a set of vectors $y_1, \dots, y_N \in \mathbf{R}^n$, which might represent a collection of measurements or other data. Suppose we have

$$y_i \approx Ax_i + b, \quad i = 1, \dots, N,$$

where $A \in \mathbf{R}^{n \times m}$, $x_i \in \mathbf{R}^m$, and $b \in \mathbf{R}^n$, with $m < n$. (Our main interest is in the case when N is much larger than n , and m is smaller than n .) Then we say that $y = Ax + b$ is a *linear explanation* of the data y . We refer to x as the vector of *factors* or *underlying causes* of the data y . For example, suppose $N = 500$, $n = 30$, and $m = 5$. In this case we have 500 vectors; each vector y_i consists of 30 scalar measurements or data points. But these 30-dimensional data points can be ‘explained’ by a much smaller set of 5 ‘factors’ (the components of x_i). This problem is about uncovering, or discovering, a linear explanation of a set of data, given only the data. In other words, we are given y_1, \dots, y_N , and the goal is to find m , A , b , and x_1, \dots, x_N so that $y_i \approx Ax_i + b$. To judge the accuracy of a proposed explanation, we’ll use the RMS explanation error, *i.e.*,

$$J = \left(\frac{1}{N} \sum_{i=1}^N \|y_i - Ax_i - b\|^2 \right)^{1/2}.$$

One rather simple linear explanation of the data is obtained with $x_i = y_i$, $A = I$, and $b = 0$. In other words, the data explains itself! In this case, of course, we have $y_i = Ax_i + b$, so the RMS explanation error is zero. But this is not what we're after. To be a useful explanation, we want to have m substantially smaller than n , *i.e.*, substantially fewer factors than the dimension of the original data (and for this smaller dimension, we'll accept a nonzero, but hopefully small, value of J .) Generally, we want m , the number of factors in the explanation, to be as small as possible, subject to the constraint that J is not too large. Even if we fix the number of factors as m , a linear explanation of a set of data is not unique. Suppose A , b , and x_1, \dots, x_N is a linear explanation of our data, with $x_i \in \mathbf{R}^m$. Then we can multiply the matrix A by two (say), and divide each vector x_i by two, and we have another linear explanation of the original data. More generally, let $F \in \mathbf{R}^{m \times m}$ be invertible, and $g \in \mathbf{R}^m$. Then we have

$$y_i \approx Ax_i + b = (AF^{-1})(Fx_i + g) + (b - AF^{-1}g).$$

Thus,

$$\tilde{A} = AF^{-1}, \quad \tilde{b} = b - AF^{-1}g, \quad \tilde{x}_1 = Fx_1 + g, \quad \dots, \quad \tilde{x}_N = Fx_N + g$$

is another equally good linear explanation of the data. In other words, we can apply any affine (*i.e.*, linear plus constant) mapping to the underlying factors x_i , and generate another equally good explanation of the original data by appropriately adjusting A and b . To standardize or normalize the linear explanation, it is usually assumed that

$$\frac{1}{N} \sum_{i=1}^N x_i = 0, \quad \frac{1}{N} \sum_{i=1}^N x_i x_i^T = I.$$

In other words, the underlying factors have an average value zero, and unit sample covariance. (You don't need to know what covariance is — it's just a definition here.) Finally, the problem.

- (a) Explain clearly how you would find a hidden linear explanation for a set of data y_1, \dots, y_N . Be sure to say how you find m , the dimension of the underlying causes, the matrix A , the vector b , and the vectors x_1, \dots, x_N . Explain clearly why the vectors x_1, \dots, x_N have the required properties.
 - (b) Carry out your method on the data in the file `linexp_data.m` available on the course web site. The file gives the matrix $Y = [y_1 \cdots y_N]$. Give your final A , b , and x_1, \dots, x_N , and verify that $y_i \approx Ax_i + b$ by calculating the norm of the error vector, $\|y_i - Ax_i - b\|$, for $i = 1, \dots, N$. Sort these norms in descending order and plot them. (This gives a good picture of the distribution of explanation errors.) By explicit computation verify that the vectors x_1, \dots, x_N obtained, have the required properties.
- 15.35 *Some bounds on singular values.* Suppose $A \in \mathbf{R}^{6 \times 3}$, with singular values 7, 5, 3, and $B \in \mathbf{R}^{6 \times 3}$, with singular values 2, 2, 1. Let $C = [A \ B] \in \mathbf{R}^{6 \times 6}$, with full SVD $C = U\Sigma V^T$, with $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_6)$. (We allow the possibility that some of these singular values are zero.)
- (a) How large can σ_1 be?
 - (b) How small can σ_1 be?
 - (c) How large can σ_6 be?
 - (d) How small can σ_6 be?

What we mean is, how large (or small) can the specified quantity be, for any A and B with the given sizes and given singular values.

Give your answers with 3 digits after the decimal place, as in

$$(a) \ 12.420, \quad (b) \ 10.000, \quad (c) \ 0.552, \quad (d) \ 0.000.$$

(This is just an example.) Briefly justify your answers. find A and B that achieve the values you give.

15.36 *Some simple matrix inequality counter-examples.*

- (a) Find a (square) matrix A , which has all eigenvalues real and positive, but there is a vector x for which $x^T A x < 0$. (Give A and x , and the eigenvalues of A .)

Moral: You cannot use positivity of the eigenvalues of A as a test for whether $x^T A x \geq 0$ holds for all x .

What is the correct way to check whether $x^T A x \geq 0$ holds for all x ? (You are allowed to find eigenvalues in this process.)

- (b) Find symmetric matrices A and B for which neither $A \geq B$ nor $B \geq A$ holds.

Of course, we'd like the simplest examples in each case.

15.37 *Some true-false questions.* In the following statements, $A \in \mathbf{R}^{n \times n}$, σ_{\min} refers to σ_n (the n th largest singular value), and κ refers to the condition number. Tell us whether each statement is true or false. 'True' means that the statement holds for any matrix $A \in \mathbf{R}^{n \times n}$, for any n . 'False' means that the statement is not true. The only answers we will read are 'True', 'False', and 'My attorney has advised me to not answer this question at this time'. (This last choice will receive partial credit.) If you write anything else, you will receive no credit for that statement. In particular, do not write justification for any answer, or provide any counter-examples.

- (a) $\|e^A\| \leq e^{\|A\|}$.
 (b) $\sigma_{\min}(e^A) \geq e^{\sigma_{\min}(A)}$.
 (c) $\kappa(e^A) \leq e^{\kappa(A)}$.
 (d) $\kappa(e^A) \leq e^{2\|A\|}$.
 (e) $\mathbf{Rank}(e^A) \geq \mathbf{Rank}(A)$.
 (f) $\mathbf{Rank}(e^A - I) \leq \mathbf{Rank}(A)$.

Lecture 16 – SVD applications

- 16.1 *Smallest matrix with given row and column sums.* Explain how to find the matrix $A \in \mathbf{R}^{m \times n}$ that minimizes

$$J = \sum_{i=1}^m \sum_{j=1}^n A_{ij}^2,$$

subject to the constraints

$$\sum_{j=1}^n A_{ij} = r_i, \quad i = 1, \dots, m, \quad \sum_{i=1}^m A_{ij} = c_j, \quad j = 1, \dots, n.$$

Here, r_i (which give the rows sums) are given, as are c_j (which give the column sums). You can assume that $\sum_{i=1}^m r_i = \sum_{j=1}^n c_j$; if this doesn't hold, there is no A that can satisfy the constraints. Using matrix notation, the objective can be written as $J = \mathbf{Tr}(A^T A)$, and the constraints are

$$A\mathbf{1} = r, \quad A^T\mathbf{1} = c,$$

where $\mathbf{1}$ denotes a vector (of appropriate size in each case) with all components one. The data $r \in \mathbf{R}^m$ and $c \in \mathbf{R}^n$ must satisfy $\mathbf{1}^T r = \mathbf{1}^T c$. Explain your method in the general case. If you can give a nice formula for the optimal A , do so. In addition, carry out your method for the specific data

$$r = \begin{bmatrix} 30 \\ 18 \\ 26 \\ 22 \\ 14 \\ 34 \end{bmatrix}, \quad c = \begin{bmatrix} 24 \\ 20 \\ 16 \\ 8 \\ 28 \\ 32 \\ 4 \\ 12 \end{bmatrix},$$

with $m = 6$ and $n = 8$. (Entries in A do not have to be integers.)

- 16.2 *Condition number.* Show that $\kappa(A) = 1$ if and only if A is a multiple of an orthogonal matrix. Thus the best conditioned matrices are precisely (scaled) orthogonal matrices.
- 16.3 *Tightness of the condition number sensitivity bound.* Suppose A is invertible, $Ax = y$, and $A(x + \delta x) = y + \delta y$. In the lecture notes we showed that $\|\delta x\|/\|x\| \leq \kappa(A)\|\delta y\|/\|y\|$. Show that this bound is not conservative, *i.e.*, there are x , y , δx , and δy such that equality holds. *Conclusion:* the bound on relative error can be taken on, if the data x is in a particularly unlucky direction and the data error δx is in (another) unlucky direction.
- 16.4 *Sensitivity to errors in A .* This problem concerns the relative error incurred in solving a set of linear equations when there are errors in the *matrix* A (as opposed to error in the data vector b). Suppose A is invertible, $Ax = b$, and $(A + \delta A)(x + \delta x) = b$. Show that $\|\delta x\|/\|x + \delta x\| \leq \kappa(A)\|\delta A\|/\|A\|$.
- 16.5 *Minimum and maximum RMS gain of an FIR filter.* Consider an FIR filter with impulse response

$$h_1 = 1, \quad h_2 = 0.6, \quad h_3 = 0.2, \quad h_4 = -0.2, \quad h_5 = -0.1.$$

We'll consider inputs of length 50 (*i.e.*, input signal that are zero for $t > 50$), so the output will have length (up to) 54, since the FIR filter has length 5. Let $u \in \mathbf{R}^{50}$ denote the input signal, and $y \in \mathbf{R}^{54}$ denote the resulting output signal. The RMS gain of the filter for a signal u is defined as

$$g = \frac{\frac{1}{\sqrt{54}}\|y\|}{\frac{1}{\sqrt{50}}\|u\|},$$

which is the ratio of the RMS value of the output to the RMS value of the input. Obviously, the gain g depends on the input signal.

- (a) Find the maximum RMS gain g_{\max} of the FIR filter, *i.e.*, the largest value g can have. Plot the corresponding input signal whose RMS gain is g_{\max} .
- (b) Find the minimum RMS gain g_{\min} of the FIR filter, *i.e.*, the smallest value g can have. Plot the corresponding input signal whose RMS gain is g_{\min} .
- (c) Plot the magnitude of the transfer function of the FIR filter, *i.e.*, $|H(e^{j\Omega})|$, where

$$H(e^{j\Omega}) = \sum_{k=1}^5 h_k e^{-jk\Omega}.$$

Find the maximum and minimum absolute values of the transfer function, and the frequencies at which they are attained. Compare to the results from parts a and b. *Hint:* To plot the magnitude of the transfer function, you may want to use the `freqz` Matlab command. Make sure you understand what `freqz` does (using `help freqz`, for example).

- (d) (This part is for fun.) Make a conjecture about the maximum and minimum singular values of a Toeplitz matrix, and the associated left and right singular vectors.

16.6 *Detecting linear relations.* Suppose we have N measurements y_1, \dots, y_N of a vector signal $x_1, \dots, x_N \in \mathbf{R}^n$:

$$y_i = x_i + d_i, \quad i = 1, \dots, N.$$

Here d_i is some small measurement or sensor noise. We hypothesize that there is a linear relation among the components of the vector signal x , *i.e.*, there is a nonzero vector q such that $q^T x_i = 0$, $i = 1, \dots, N$. The geometric interpretation is that all of the vectors x_i lie in the hyperplane $q^T x = 0$. We will assume that $\|q\| = 1$, which does not affect the linear relation. Even if the x_i 's do lie in a hyperplane $q^T x = 0$, our measurements y_i will not; we will have $q^T y_i = q^T d_i$. These numbers are small, assuming the measurement noise is small. So the problem of determining whether or not there is a linear relation among the components of the vectors x_i comes down to finding out whether or not there is a unit-norm vector q such that $q^T y_i$, $i = 1, \dots, N$, are all small. We can view this problem geometrically as well. Assuming that the x_i 's all lie in the hyperplane $q^T x = 0$, and the d_i 's are small, the y_i 's will all lie close to the hyperplane. Thus a scatter plot of the y_i 's will reveal a sort of flat cloud, concentrated near the hyperplane $q^T x = 0$. Indeed, for any z and $\|q\| = 1$, $|q^T z|$ is the distance from the vector z to the hyperplane $q^T x = 0$. So we seek a vector q , $\|q\| = 1$, such that all the measurements y_1, \dots, y_N lie close to the hyperplane $q^T x = 0$ (that is, $q^T y_i$ are all small). How can we determine if there is such a vector, and what is its value? We define the following normalized measure:

$$\rho = \sqrt{\frac{1}{N} \sum_{i=1}^N (q^T y_i)^2} \bigg/ \sqrt{\frac{1}{N} \sum_{i=1}^N \|y_i\|^2}.$$

This measure is simply the ratio between the *root mean square distance* of the vectors to the hyperplane $q^T x = 0$ and the *root mean square length* of the vectors. If ρ is small, it means that the measurements lie close to the hyperplane $q^T x = 0$. Obviously, ρ depends on q . Here is the problem: explain how to find the minimum value of ρ over all unit-norm vectors q , and the unit-norm vector q that achieves this minimum, given the data set y_1, \dots, y_N .

16.7 *Stability conditions for the distributed congestion control scheme.* We consider the congestion control scheme in problem 3, and will use the notation from that problem. In this problem, we study the dynamics and convergence properties of the rate adjustment scheme. To simplify things, we'll assume that the route matrix R is skinny and full rank. You can also assume that $\alpha > 0$. Let

$\bar{x}_{\text{ls}} = (R^T R)^{-1} R^T T^{\text{target}}$ denote the least-squares approximate solution of the (over-determined) equations $Rx \approx T^{\text{target}}$. (The source rates given by \bar{x}_{ls} minimize the sum of the squares of the congestion measures on all paths.)

- Find the conditions on the update parameter α under which the rate adjustment scheme converges to \bar{x}_{ls} , no matter what the initial source rate is.
- Find the value of α that results in the fastest possible asymptotic convergence of the rate adjustment algorithm. Find the associated asymptotic convergence rate. We define the convergence rate as the smallest number c for which $\|x(t) - \bar{x}_{\text{ls}}\| \leq ac^t$ holds for all trajectories and all t (the constant a can depend on the trajectory).

You can give your solutions in terms of any of the concepts we have studied, *e.g.*, matrix exponential, eigenvalues, singular values, condition number, and so on. Your answers can, of course, depend on R , T^{target} , and \bar{x}_{ls} . If your answer doesn't depend on some of these (or even all of them) be sure to point this out. We'll take points off if you give a solution that is correct, but not as simple as it can be.

16.8 Consider the system $\dot{x} = Ax$ with

$$A = \begin{bmatrix} 0.3132 & 0.3566 & 0.2545 & 0.2579 & 0.2063 \\ -0.0897 & 0.2913 & 0.1888 & 0.4392 & 0.1470 \\ 0.0845 & 0.2433 & -0.5888 & -0.0407 & 0.1744 \\ 0.2478 & -0.1875 & 0.2233 & 0.3126 & -0.6711 \\ 0.1744 & 0.2315 & -0.1004 & -0.2111 & 0.0428 \end{bmatrix}.$$

- Find the initial state $x(0) \in \mathbf{R}^5$ satisfying $\|x(0)\| = 1$ such that $\|x(3)\|$ is maximum. In other words, find an initial condition of unit norm that produces the *largest* state at $t = 3$.
- Find the initial state $x(0) \in \mathbf{R}^5$ satisfying $\|x(0)\| = 1$ such that $\|x(3)\|$ is minimum.

To save you the trouble of typing in the matrix A , you can find it on the web page in the file `max_min_init_state.m`.

16.9 *Regularization and SVD*. Let $A \in \mathbf{R}^{n \times n}$ be full rank, with SVD

$$A = \sum_{i=1}^n \sigma_i u_i v_i^T.$$

(We consider the square, full rank case just for simplicity; it's not too hard to consider the general nonsquare, non-full rank case.) Recall that the regularized approximate solution of $Ax = y$ is defined as the vector $x_{\text{reg}} \in \mathbf{R}^n$ that minimizes the function

$$\|Ax - y\|^2 + \mu \|x\|^2,$$

where $\mu > 0$ is the regularization parameter. The regularized solution is a linear function of y , so it can be expressed as $x_{\text{reg}} = By$ where $B \in \mathbf{R}^{n \times n}$.

- Express the SVD of B in terms of the SVD of A . To be more specific, let

$$B = \sum_{i=1}^n \tilde{\sigma}_i \tilde{u}_i \tilde{v}_i^T$$

denote the SVD of B . Express $\tilde{\sigma}_i$, \tilde{u}_i , \tilde{v}_i , for $i = 1, \dots, n$, in terms of σ_i , u_i , v_i , $i = 1, \dots, n$ (and, possibly, μ). Recall the convention that $\tilde{\sigma}_1 \geq \dots \geq \tilde{\sigma}_n$.

- Find the norm of B . Give your answer in terms of the SVD of A (and μ).

(c) Find the worst-case relative inversion error, defined as

$$\max_{y \neq 0} \frac{\|AB y - y\|}{\|y\|}.$$

Give your answer in terms of the SVD of A (and μ).

16.10 *Optimal binary signalling.* We consider a communication system given by

$$y(t) = Au(t) + v(t), \quad t = 0, 1, \dots$$

Here

- $u(t) \in \mathbf{R}^n$ is the transmitted (vector) signal at time t
- $y(t) \in \mathbf{R}^m$ is the received (vector) signal at time t
- $v(t) \in \mathbf{R}^m$ is noise at time t
- $t = 0, 1, \dots$ is the (discrete) time

Note that the system has no memory: $y(t)$ depends only on $u(t)$. For the noise, we assume that $\|v(t)\| < V_{\max}$. Other than this maximum value for the norm, we know nothing about the noise (for example, we do not assume it is random). We consider binary signalling, which means that at each time t , the transmitter sends one of two signals, *i.e.*, we have either $u(t) = s_1 \in \mathbf{R}^n$ or $u(t) = s_2 \in \mathbf{R}^n$. The receiver then guesses which of the two signals was sent, based on $y(t)$. The process of guessing which signal was sent, based on the received signal $y(t)$, is called *decoding*. In this problem we are only interested in the case when the communication is completely reliable, which means that the receiver's estimate of which signal was sent is always correct, no matter what $v(t)$ is (provided $\|v(t)\| < V_{\max}$, of course). Intuition suggests that this is possible when V_{\max} is small enough.

- (a) Your job is to design the signal constellation, *i.e.*, the vectors $s_1 \in \mathbf{R}^n$ and $s_2 \in \mathbf{R}^n$, and the associated (reliable) decoding algorithm used by the receiver. Your signal constellation should minimize the maximum transmitter power, *i.e.*,

$$P_{\max} = \max\{\|s_1\|, \|s_2\|\}.$$

You must describe:

- your analysis of this problem,
 - how you come up with s_1 and s_2 ,
 - the exact decoding algorithm used,
 - how you know that the decoding algorithm is reliable, *i.e.*, the receiver's guess of which signal was sent is always correct.
- (b) The file `opt_bin_data.m` contains the matrix A and the scalar V_{\max} . Using your findings from part 1, determine the optimal signal constellation.

16.11 *Some input optimization problems.* In this problem we consider the system $x(t+1) = Ax(t) + Bu(t)$, with

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad x(0) = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 1 \end{bmatrix}.$$

- (a) *Least-norm input to steer state to zero in minimum time.* Find the minimum T , T_{\min} , such that $x(T) = 0$ is possible. Among all $(u(0), u(1), \dots, u(T_{\min} - 1))$ that steer $x(0)$ to $x(T_{\min}) = 0$, find the one of minimum norm, *i.e.*, the one that minimizes

$$J_{T_{\min}} = \|u(0)\|^2 + \dots + \|u(T_{\min} - 1)\|^2.$$

Give the minimum value of $J_{T_{\min}}$ achieved.

- (b) *Least-norm input to achieve $\|x(10)\| \leq 0.1$.* In lecture we worked out the least-norm input that drives the state exactly to zero at $t = 10$. Suppose instead we only require the state to be *small* at $t = 10$, for example, $\|x(10)\| \leq 0.1$. Find $u(0), u(1), \dots, u(9)$ that minimize

$$J_9 = \|u(0)\|^2 + \dots + \|u(9)\|^2$$

subject to the condition $\|x(10)\| \leq 0.1$. Give the value of J_9 achieved by your input.

- 16.12 *Determining the number of signal sources.* The signal transmitted by n sources is measured at m receivers. The signal transmitted by each of the sources at sampling period k , for $k = 1, \dots, p$, is denoted by an n -vector $x(k) \in \mathbf{R}^n$. The gain from the j -th source to the i -th receiver is denoted by $a_{ij} \in \mathbf{R}$. The signal measured at the receivers is then

$$y(k) = Ax(k) + v(k), \quad k = 1, \dots, p,$$

where $v(k) \in \mathbf{R}^m$ is a vector of sensor noises, and $A \in \mathbf{R}^{m \times n}$ is the matrix of source to receiver gains. However, we do not know the gains a_{ij} , nor the transmitted signal $x(k)$, nor even the number of sources present n . We only have the following additional *a priori* information:

- We expect the number of sources to be less than the number of receivers (*i.e.*, $n < m$, so that A is skinny);
- A is full-rank and well-conditioned;
- All sources have roughly the same average power, the signal $x(k)$ is unpredictable, and the source signals are unrelated to each other; Hence, given enough samples (*i.e.*, p large) the vectors $x(k)$ will ‘point in all directions’;
- The sensor noise $v(k)$ is small relative to the received signal $Ax(k)$.

Here’s the question:

- (a) You are given a large number of vectors of sensor measurements $y(k) \in \mathbf{R}^m$, $k = 1, \dots, p$. How would you estimate the number of sources, n ? Be sure to clearly describe your proposed method for determining n , and to explain when and why it works.
- (b) Try your method on the signals given in the file `nsources.m`. Running this script will define the variables:
- `m`, the number of receivers;
 - `p`, the number of signal samples;
 - `Y`, the receiver sensor measurements, an array of size `m` by `p` (the k -th column of `Y` is $y(k)$.)

What can you say about the number of signal sources present?

Note: Our problem description and assumptions are not precise. An important part of this problem is to explain your method, and clarify the assumptions.

- 16.13 *The EE263 search engine.* In this problem we examine how linear algebra and low-rank approximations can be used to find matches to a search query in a set of documents. Let’s assume we have four documents: **A**, **B**, **C**, and **D**. We want to search these documents for three terms: *piano*, *violin*, and *drum*. We know that:

in **A**, the word *piano* appears 4 times, *violin* 3 times, and *drum* 1 time;

in **B**, the word *piano* appears 6 times, *violin* 1 time, and *drum* 0 times;

in **C**, the word *piano* appears 7 time, *violin* 4 times, and *drum* 39 times; and

in **D**, the word *piano* appears 0 times, *violin* 0 times, and *drum* 5 times.

We can tabulate this as follows:

	A	B	C	D
<i>piano</i>	4	6	7	0
<i>violin</i>	3	1	4	0
<i>drum</i>	1	0	39	5

This information is used to form a *term-by-document* matrix A , where A_{ij} specifies the frequency of the i th term in the j th document, *i.e.*,

$$A = \begin{bmatrix} 4 & 6 & 7 & 0 \\ 3 & 1 & 4 & 0 \\ 1 & 0 & 39 & 5 \end{bmatrix}.$$

Now let q be a *query vector*, with a non-zero entry for each term. The query vector expresses a criterion by which to select a document. Typically, q will have 1 in the entries corresponding to the words we want to search for, and 0 in all other entries (but other weighting schemes are possible.) A simple measure of how relevant document j is to the query is given by the inner product of the j th column of A with q :

$$a_j^T q.$$

However, this criterion is biased towards large documents. For instance, a query for *piano* ($q = [1\ 0\ 0]^T$) by this criterion would return document **C** as most relevant, even though document **B** (and even **A**) is probably much more relevant. For this reason, we use the inner product normalized by the norm of the vectors,

$$\frac{a_j^T q}{\|a_j\| \|q\|}.$$

Note that our criterion for measuring how well a document matches the query is now the cosine of the angle between the document and query vectors. Since all entries are non-negative, the cosine is in $[0, 1]$ (and the angle is in $[-\pi/2, \pi/2]$.) Define \tilde{A} and \tilde{q} as normalized versions of A and q (A is normalized column-wise, *i.e.*, each column is divided by its norm.) Then,

$$c = \tilde{A}^T \tilde{q}$$

is a column vector that gives a measure of the relevance of each document to the query. And now, the question. In the file `term_by_doc.m` you are given m search terms, n documents, and the corresponding term-by-document matrix $A \in \mathbf{R}^{m \times n}$. (They were obtained randomly from Stanford's *Portfolio* collection of internal documents.) The variables `term` and `document` are lists of strings. The string `term{i}` contains the i th word. Each document is specified by its URL, *i.e.*, if you point your web browser to the URL specified by the string `document{j}` you can inspect the contents of the j th document. The matrix entry `A(i,j)` specifies how many times term i appears in document j .

- Compute \tilde{A} , the normalized term-by-document matrix. Compute and plot the singular values of \tilde{A} .
- Perform a query for the word *students* ($i = 53$) on \tilde{A} . What are the 5 top results?
- We will now consider low-rank approximations of \tilde{A} , that is

$$\hat{A}_r = \min_{\hat{A}, \text{rank}(\hat{A}) \leq r} \|\tilde{A} - \hat{A}\|.$$

Compute \hat{A}_{32} , \hat{A}_{16} , \hat{A}_8 , and \hat{A}_4 . Perform a query for the word *students* on these matrices. Comment on the results.

- Are there advantages of using low-rank approximations over using the full-rank matrix? (You can assume that a very large number of searches will be performed before the term-by-document matrix is updated.)

Note: Variations and extensions of this idea are actually used in commercial search engines (although the details are closely guarded secrets ...) Issues in real search engines include the fact that m and n are enormous and change with time. These methods are very interesting because they can recover documents that don't include the term searched for. For example, a search for *automobile* could retrieve

a document with no mention of *automobile*, but many references to cars (can you give a justification for this?) For this reason, this approach is sometimes called *latent semantic indexing*. *Matlab hints:* You may find the command `sort` useful. It sorts a vector in *ascending* order, and can also return a vector with the original indexes of the sorted elements. Here's a sample code that sorts the vector c in *descending* order, and prints the URL of the top document and its corresponding c_j .

```
[c,j]=sort(-c);
c=-c;
disp(document{j(1)})
disp(c(1))
```

- 16.14 *Condition number and angles between columns.* Suppose $A \in \mathbf{R}^{n \times n}$ has columns $a_1, \dots, a_n \in \mathbf{R}^n$, each of which has unit norm:

$$A = [a_1 \ a_2 \ \cdots \ a_n], \quad \|a_i\| = 1, \quad i = 1, \dots, n.$$

Suppose that two of the columns have an angle less than 10° between them, *i.e.*, $a_k^T a_l \geq \cos 10^\circ$. Show that $\kappa(A) \geq 10$, where κ denotes the condition number. (If A is singular, we take $\kappa(A) = \infty$, and so $\kappa(A) \geq 10$ holds.) *Interpretation:* If the columns were orthogonal, *i.e.*, $\angle(a_i, a_j) = 90^\circ$ for $i \neq j$, $i, j = 1, \dots, n$, then A would be an orthogonal matrix, and therefore its condition number would be one (which is the smallest a condition number can be). At the other extreme, if two columns are the same (*i.e.*, have zero angle between them), the matrix is singular, and the condition number is infinite. Intuition suggests that if some pair of columns has a *small* angle, such as 10° , then the condition number must be large. (Although in many applications, a condition number of 10 is not considered large.)

- 16.15 *Analysis and optimization of a communication network.* A communication network is modeled as a set of m directed links connecting nodes. There are n routes in the network. A route is a path, along one or more links in the network, from a *source node* to a *destination node*. In this problem, the routes are fixed, and are described by an $m \times n$ route-link matrix A , defined as

$$A_{ij} = \begin{cases} 1 & \text{route } j \text{ passes through link } i \\ 0 & \text{otherwise.} \end{cases}$$

Over each route we have a nonnegative *flow*, measured in (say) bits per second. We denote the flow along route j as f_j , and we call $f \in \mathbf{R}^n$ the *flow vector*. The *traffic* on a link i , denoted t_i , is the sum of the flows on all routes passing through link i . The vector $t \in \mathbf{R}^m$ is called the *traffic vector*. handle. We're

Each link has an associated nonnegative *delay*, measured in (say) seconds. We denote the delay for link i as d_i , and refer to $d \in \mathbf{R}^m$ as the *link delay vector*. The *latency* on a route j , denoted l_j , is the sum of the delays along each link constituting the route, *i.e.*, the time it takes for bits entering the source to emerge at the destination. The vector $l \in \mathbf{R}^n$ is the *route latency vector*.

The total number of bits in the network at an instant in time is given by $B = f^T l = t^T d$.

- (a) *Worst-case flows and delays.* Suppose the flows and link delays satisfy

$$(1/n) \sum_{j=1}^n f_j^2 \leq F^2, \quad (1/m) \sum_{i=1}^m d_i^2 \leq D^2,$$

where F and D are given. What is the maximum possible number of bits in the network? What values of f and d achieve this maximum value? (For this problem you can ignore the constraint that the flows and delays must be nonnegative. It turns out, however, that the worst-case flows and delays can always be chosen to be nonnegative.)

- (b) *Utility maximization.* For a flow f_j , the network operator derives income at a rate $p_j f_j$, where p_j is the price per unit flow on route j . The network operator's total rate of income is thus $\sum_{j=1}^n p_j f_j$. (The route prices are known and positive.)

The network operator is charged at a rate $c_i t_i$ for having traffic t_i on link i , where c_i is the cost per unit of traffic on link i . The total charge rate for link traffic is $\sum_{i=1}^m t_i c_i$. (The link costs are known and positive.) The net income rate (or utility) to the network operator is therefore

$$U^{\text{net}} = \sum_{j=1}^n p_j f_j - \sum_{i=1}^m c_i t_i.$$

Find the flow vector f that maximizes the operator's net income rate, subject to the constraint that each f_j is between 0 and F^{max} , where F^{max} is a given positive maximum flow value.

- 16.16 *A heuristic for MAXCUT.* Consider a graph with n nodes and m edges, with the nodes labeled $1, \dots, n$ and the edges labeled $1, \dots, m$. We partition the nodes into two groups, B and C , i.e., $B \cap C = \emptyset$, $B \cup C = \{1, \dots, n\}$. We define the number of *cuts* associated with this partition as the number of edges between pairs of nodes when one of the nodes is in B and the other is in C . A famous problem, called the MAXCUT problem, involves choosing a partition (i.e., B and C) that maximizes the number of cuts for a given graph. For any partition, the number of cuts can be no more than m . If the number of cuts is m , nodes in group B connect only to nodes in group C and the graph is bipartite.

The MAXCUT problem has many applications. We describe one here, although you do not need it to solve this problem. Suppose we have a communication system that operates with a two-phase clock. During periods $t = 0, 2, 4, \dots$, each node in group B transmits data to nodes in group C that it is connected to; during periods $t = 1, 3, 5, \dots$, each node in group C transmits to the nodes in group B that it is connected to. The number of cuts, then, is exactly the number of successful transmissions that can occur in a two-period cycle. The MAXCUT problem is to assign nodes to the two groups so as to maximize the overall efficiency of communication.

It turns out that the MAXCUT problem is hard to solve exactly, at least if we don't want to resort to an exhaustive search over all, or most of, the 2^{n-1} possible partitions. In this problem we explore a sophisticated heuristic method for finding a good (if not the best) partition in a way that scales to large graphs.

We will encode the partition as a vector $x \in \mathbf{R}^n$, with $x_i \in \{-1, 1\}$. The associated partition has $x_i = 1$ for $i \in B$ and $x_i = -1$ for $i \in C$. We describe the graph by its node adjacency matrix $A \in \mathbf{R}^{n \times n}$ (C), with

$$A_{ij} = \begin{cases} 1 & \text{there is an edge between node } i \text{ and node } j \\ 0 & \text{otherwise} \end{cases}$$

Note that A is symmetric and $A_{ii} = 0$ (since we do not have self-loops in our graph).

- (a) Find a symmetric matrix P , with $P_{ii} = 0$ for $i = 1, \dots, n$, and a constant d , for which $x^T P x + d$ is the number of cuts encoded by any partitioning vector x . Explain how to calculate P and d from A . Of course, P and d cannot depend on x .

The MAXCUT problem can now be stated as the optimization problem

$$\begin{aligned} & \text{maximize} && x^T P x + d \\ & \text{subject to} && x_i^2 = 1, \quad i = 1, \dots, n, \end{aligned}$$

with variable $x \in \mathbf{R}^n$.

- (b) A famous heuristic for approximately solving MAXCUT is to replace the n constraints $x_i^2 = 1$, $i = 1, \dots, n$, with a single constraint $\sum_{i=1}^n x_i^2 = n$, creating the so-called *relaxed* problem

$$\begin{aligned} & \text{maximize} && x^T P x + d \\ & \text{subject to} && \sum_{i=1}^n x_i^2 = n. \end{aligned}$$

Explain how to solve this relaxed problem (even if you could not solve part (a)).

Let x^* be a solution to the relaxed problem. We generate our candidate partition with $x_i = \text{sign}(x_i^*)$. (This means that $x_i = 1$ if $x_i^* \geq 0$, and $x_i = -1$ if $x_i^* < 0$.) really

Remark: We can give a geometric interpretation of the relaxed problem, which will also explain why it's called relaxed. The constraints in the problem in part (a), that $x_i^2 = 1$, require x to lie on the vertices of the unit hypercube. In the relaxed problem, the constraint set is the unit ball of unit radius. Because this constraint set is larger than the original constraint set (*i.e.*, it includes it), we say the constraints have been relaxed.

- (c) Run the MAXCUT heuristic described in part (b) on the data given in `mc_data.m`. How many cuts does your partition yield?

A simple alternative to MAXCUT is to generate a large number of random partitions, using the random partition that maximizes the number of cuts as an approximate solution. Carry out this method with 1000 random partitions generated by `x = sign(rand(n,1)-0.5)`. What is the largest number of cuts obtained by these random partitions?

Note: There are many other heuristics for approximately solving the MAXCUT problem. However, we are not interested in them. In particular, please do not submit any other method for approximately solving MAXCUT.

- 16.17 *Simultaneously estimating student ability and exercise difficulty.* Each of n students takes an exam that contains m questions. Student j receives (nonnegative) grade G_{ij} on question i . One simple model for predicting the grades is to estimate $G_{ij} \approx \hat{G}_{ij} = a_j/d_i$, where a_j is a (nonnegative) number that gives the *ability* of student j , and d_i is a (positive) number that gives the *difficulty* of exam question i . Given a particular model, we could simultaneously scale the student abilities and the exam difficulties by any positive number, without affecting \hat{G}_{ij} . Thus, to ensure a unique model, we will *normalize* the exam question difficulties d_i , so that the mean exam question difficulty across the m questions is 1.

In this problem, you are given a complete set of grades (*i.e.*, the matrix $G \in \mathbf{R}^{m \times n}$). Your task is to find a set of nonnegative student abilities, and a set of positive, normalized question difficulties, so that $G_{ij} \approx \hat{G}_{ij}$. In particular, choose your model to minimize the RMS error, J ,

$$J = \left(\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n (G_{ij} - \hat{G}_{ij})^2 \right)^{1/2}.$$

This can be compared to the RMS value of the grades,

$$\left(\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n G_{ij}^2 \right)^{1/2}.$$

- (a) Explain how to solve this problem, using any concepts from EE263. If your method is approximate, or not guaranteed to find the global minimum value of J , say so. If carrying out your method requires some rank or other conditions to hold, say so.

Note: You do not have to concern yourself with the requirement that a_j are nonnegative and d_i are positive. You can just assume this works out, or is easily corrected.

- (b) Carry out your method on the data found in `grademodeldata.m`. Give the optimal value of J , and also express it as a fraction of the RMS value of the grades. Give the difficulties of the 7 problems on the exam.

- 16.18 *Angle between two subspaces.* The angle between two nonzero vectors v and w in \mathbf{R}^n is defined as

$$\angle(v, w) = \cos^{-1} \left(\frac{v^T w}{\|v\| \|w\|} \right),$$

where we take $\cos^{-1}(a)$ as being between 0 and π . We define the angle between a nonzero vector $v \in \mathbf{R}^n$ and a (nonzero) subspace $\mathcal{W} \subseteq \mathbf{R}^n$ as

$$\angle(v, \mathcal{W}) = \min_{w \in \mathcal{W}, w \neq 0} \angle(v, w).$$

Thus, $\angle(v, \mathcal{W}) = 10^\circ$ means that the smallest angle between v and any vector in \mathcal{W} is 10° . If $v \in \mathcal{W}$, we have $\angle(v, \mathcal{W}) = 0$.

Finally, we define the angle between two nonzero subspaces \mathcal{V} and \mathcal{W} as

$$\angle(\mathcal{V}, \mathcal{W}) = \max \left\{ \max_{v \in \mathcal{V}, v \neq 0} \angle(v, \mathcal{W}), \max_{w \in \mathcal{W}, w \neq 0} \angle(w, \mathcal{V}) \right\}.$$

This angle is zero if and only if the two subspaces are equal. If $\angle(\mathcal{V}, \mathcal{W}) = 10^\circ$, say, it means that either there is a vector in \mathcal{V} whose minimum angle to any vector of \mathcal{W} is 10° , or there is a vector in \mathcal{W} whose minimum angle to any vector of \mathcal{V} is 10° .

- (a) Suppose you are given two matrices $A \in \mathbf{R}^{n \times r}$, $B \in \mathbf{R}^{n \times r}$, each of rank r . Let $\mathcal{V} = \text{range}(A)$ and $\mathcal{W} = \text{range}(B)$. Explain how you could find or compute $\angle(\mathcal{V}, \mathcal{W})$. You can use any of the concepts in the class, *e.g.*, least-squares, QR factorization, pseudo-inverse, norm, SVD, Jordan form, etc.
 - (b) Carry out your method for the matrices found in `angsubdata.m`. Give the numerical value for $\angle(\text{range}(A), \text{range}(B))$.
- 16.19 *Extracting the faintest signal.* An n -vector valued signal, $x(t) \in \mathbf{R}^n$, is defined for $t = 1, \dots, T$. We'll refer to its i th component, $x_i(t)$, for $t = 1, \dots, T$, as the i th scalar signal. The scalar signals x_1, \dots, x_{n-1} have an RMS value substantially larger than x_n . In other words, x_n is the faintest scalar signal. It is also the signal of interest for this problem. We will assume that the scalar signals x_1, \dots, x_n are unrelated to each other, and so are nearly uncorrelated (*i.e.*, nearly orthogonal).

We aren't given the vector signal $x(t)$, but we are given a linear transformation of it,

$$y(t) = Ax(t), \quad t = 1, \dots, T,$$

where $A \in \mathbf{R}^{n \times n}$ is invertible. If we knew A , we could easily recover the original signal (and therefore also the faintest scalar signal $x_n(t)$), using $x(t) = A^{-1}y(t)$, $t = 1, \dots, T$. But, sadly, we don't know A . Here is a heuristic method for guessing $x_n(t)$. We will form our estimate as

$$\hat{x}_n(t) = w^T y(t), \quad t = 1, \dots, T,$$

where $w \in \mathbf{R}^n$ is a vector of weights. Note that if w were chosen so that $w^T A = \alpha e_n^T$, with $\alpha \neq 0$ a constant, then we would have $\hat{x}_n(t) = \alpha x_n(t)$, *i.e.*, a perfect reconstruction except for the scale factor α .

Now, the important part of our heuristic: we choose w to minimize the RMS value of \hat{x}_n , subject to $\|w\| = 1$. *Very roughly*, one idea behind the heuristic is that, in general, $w^T y$ is a linear combination of the scalar signals x_1, \dots, x_n . If the linear combination has a small norm, that's because the linear combination is 'rich in x_n ', and has only a small amount of energy contributed by x_1, \dots, x_{n-1} . That, in fact, is exactly what we want. In any case, you don't need to worry about why the heuristic works (or doesn't work)—it's the method you are going to use in this problem.

- (a) Explain how to find a w that minimizes the RMS value of \hat{x}_n , using concepts from the class (*e.g.*, range, rank, least-squares, QR factorization, eigenvalues, singular values, and so on).

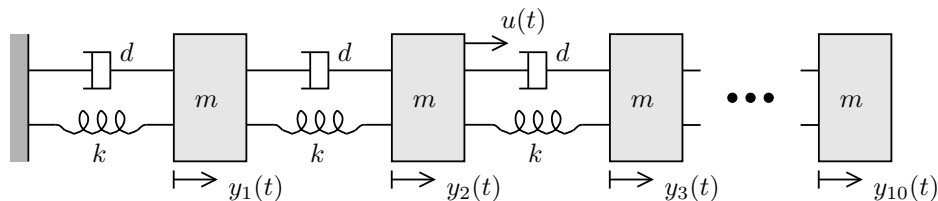
- (b) Carry out your method on the problem instance with $n = 4$, $T = 26000$, described in the Matlab file `faintestdata.m`. This file will define an $n \times T$ matrix Y , where the t th column of Y is the vector $y(t)$. The file will also define n and T . Submit your code, and give us the optimal weight vector $w \in \mathbf{R}^4$ you find, along with the associated RMS value of \hat{x}_n .

The following is not needed to solve the problem. The signals are actually audio tracks, each 3.25 seconds long and sampled at 8 kHz. The Matlab file `faintestaudio.m` contains commands to generate wave files of the linear combinations y_1, \dots, y_4 , and a wave file of your estimate \hat{x}_n . You are welcome to generate and listen to these files.

Lecture 18 – Controllability and state transfer

18.1 This problem has two parts that are mostly independent.

(a) *Actuator placement (10 masses).*



Ten masses are connected in series by springs and light dampers, as shown in the figure above. The mass positions (deviation from rest) are denoted by y_1, \dots, y_{10} . The masses, spring constants, and damping constants are all identical and given by

$$m = 1, \quad k = 1, \quad d = 0.01.$$

An actuator is used to apply a force $u(t)$ to one of the masses. In the figure, the actuator is shown located on the second mass from the left, but it could also have been placed in any of the other nine masses. Use state $x = [y^T \dot{y}^T]^T$.

- For which of the ten possible actuator placements is the system controllable?
- You are given the following design specification: any state should be reachable without the need for very large actuation forces. Where would you place the actuator? (Since the design specification is somewhat vague, you should clearly explain and justify your decision.)

Note: To avoid error propagation in solutions, use the Matlab script `spring_series.m`, available at the course web site, which constructs the dynamics and input matrices A and B .

(b) *Optimal control (4 masses).* Consider now a system with the same characteristics, but with only four masses. Four unit masses are connected in series by springs and light dampers (with $k = 1$, and $d = 0.01$.) A force actuator is placed on the *third* mass from the left. As before, use state $x = [y^T \dot{y}^T]^T$.

- Is the system controllable?
- You are given the initial state of the system, $x(0) = e_8 = [0 \ \cdots \ 0 \ 1]^T$, and asked to drive the state to as close to zero as possible at time $t_f = 20$ (*i.e.*, a velocity disturbance in the fourth mass is to be attenuated as much as possible in 20 seconds.) In other words, you are to choose an input $u(t)$, $t \in [0, t_f]$, that minimizes $\|x(t_f)\|^2$. Furthermore, from among all inputs that achieve the minimum $\|x(t_f)\|^2$, we want the smallest one, *i.e.*, the one for which the energy

$$\mathcal{E}_u = \int_0^{t_f} u(t)^2 dt$$

is minimized. Your answer should include (i) a plot of the minimizing input $u_{\text{opt}}(t)$; (ii) the corresponding energy $\mathcal{E}_{u,\text{min}}$; and (iii) the resulting $\|x(t_f)\|^2$. You must explain and justify how you obtained your solution. *Notes:*

- We will be happy with an approximate solution (by using, say, an input that is piece-wise constant in small intervals.) You may want to discretize the system, in which case we suggest you use 100 discretization intervals (or more.)
- You may (or may not) want to use the result

$$A \int_0^h e^{At} B dt = (e^{Ah} - I) B.$$

18.2 *Horizon selection.* Consider the (scalar input) system

$$x(t+1) = \begin{bmatrix} 0 & 0 & 0.8 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} x(t) + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} u(t), \quad x(0) = 0.$$

For $N \geq 3$ let $E_N(z)$ denote the minimum input energy, *i.e.*, the minimum value of

$$u(0)^2 + \cdots + u(N-1)^2,$$

required to reach $x(N) = z$. Let $E_\infty(z)$ denote the minimum energy required to reach the state $x(N) = z$, without fixing the final time N , *i.e.*, $E_\infty(z) = \lim_{N \rightarrow \infty} E_N(z)$. Find the minimum value of N such that $E_N(z) \leq 1.1E_\infty(z)$ for all z . (This is the shortest horizon that requires no more than 10% more control energy than infinite horizon control, for any final state). *Hint:* the Matlab command `P=dlyap(A,W)` computes the solution of the Lyapunov equation $APA^T + W = P$.

18.3 *Minimum energy required to steer the state to zero.* Consider a controllable discrete-time system $x(t+1) = Ax(t) + Bu(t)$, $x(0) = x_0$. Let $E(x_0)$ denote the minimum energy required to drive the state to zero, *i.e.*

$$E(x_0) = \min \left\{ \sum_{\tau=0}^{t-1} \|u(\tau)\|^2 \mid x(t) = 0 \right\}.$$

An engineer argues as follows:

This problem is like the minimum energy reachability problem, but ‘turned backwards in time’ since here we steer the state from a given state to zero, and in the reachability problem we steer the state from zero to a given state. The system $z(t+1) = A^{-1}z(t) - A^{-1}Bu(t)$ is the same as the given one, except time is running backwards. Therefore $E(x_0)$ is the same as the minimum energy required for z to reach x_0 (a formula for which can be found in the lecture notes).

Either justify or refute the engineer’s statement. You can assume that A is invertible.

18.4 *Minimum energy inputs with coasting.* We consider the controllable system $\dot{x} = Ax + Bu$, $x(0) = 0$, where $A \in \mathbf{R}^{n \times n}$ and $B \in \mathbf{R}^{n \times m}$. You are to determine an input u that results in $x(t_f) = x_{\text{des}}$, where t_f and x_{des} are given. You are also given t_a , where $0 < t_a \leq t_f$, and have the constraint that $u(t) = 0$ for $t > t_a$. Roughly speaking, you are allowed to apply a (nonzero) input u during the ‘controlled portion’ of the trajectory, *i.e.*, from $t = 0$ until $t = t_a$; from $t = t_a$ until the final time t_f , the system ‘coasts’ or ‘drifts’ with $u(t) = 0$. Among all u that satisfy these specifications, u_{ln} will denote the one that minimizes the ‘energy’

$$\int_0^{t_a} \|u(t)\|^2 dt.$$

- Give an explicit formula for $u_{\text{ln}}(t)$.
- Now suppose that t_a is increased (but still less than t_f). An engineer asserts that the minimum energy required will decrease. Another engineer disagrees, pointing out that the final time has not changed. Who is right? Justify your answer. (It is possible that neither is right.)
- Matlab exercise.* Consider the mechanical system on page 11-9 of the notes. Let $x_{\text{des}} = [1 \ 0 \ -1 \ 0 \ 0 \ 0]^T$ and $t_f = 6$. Plot the minimum energy required as a function of t_a , for $0 < t_a < t_f$. You can use a simple method to numerically approximate any integrals you encounter. You must explain what you are doing; just submitting some code and a plot is not enough.

18.5 *Some True/False questions.* By ‘True’, of course, we mean that the statement holds for all values of the matrices, vectors, dimensions, etc., mentioned in the statement. ‘False’ means that the statement fails to hold in at least one case.

- (a) Suppose $A \in \mathbf{R}^{n \times n}$ and $p(s) = s^n + a_1 s^{n-1} + \cdots + a_n$ is polynomial of degree n , with leading coefficient one, that satisfies $p(A) = 0$. Then p is the characteristic polynomial of A .
- (b) Suppose $x : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ is a trajectory of the linear dynamical system $\dot{x} = Ax$, which is stable. Then for any $t \geq 0$, we have $\|x(t)\| \leq \|x(0)\|$.
- (c) Let $A \in \mathbf{R}^{p \times q}$ and let $a_i \in \mathbf{R}^p$ denote the i th column of A . Then we have

$$\|A\| \geq \max_{i=1, \dots, q} \|a_i\|.$$

- (d) Suppose the two linear dynamical systems $\dot{x} = Fx$ and $\dot{z} = Gz$, where $F, G \in \mathbf{R}^{n \times n}$, are both stable. Then the linear dynamical system $\dot{w} = (F + G)w$ is stable.
- (e) Suppose P and Q are symmetric $n \times n$ matrices, and let $\{v_1, v_2, \dots, v_n\}$ be a basis for \mathbf{R}^n . Then if we have $v_i^T P v_i \geq v_i^T Q v_i$ for $i = 1, \dots, n$, we must have $P \geq Q$.
- (f) Let $A \in \mathbf{R}^{n \times n}$, and suppose $v \in \mathbf{R}^n$, $v \neq 0$, satisfies $v^T A = \lambda v^T$, where $\lambda \in \mathbf{R}$. Let $x : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ be any trajectory of the linear dynamical system $\dot{x} = Ax$. Then at least one of the following statements hold:
- $v^T x(t) \geq v^T x(0)$ for all $t \geq 0$
 - $v^T x(t) \leq v^T x(0)$ for all $t \geq 0$
- (g) Suppose $A \in \mathbf{R}^{p \times q}$ is fat (*i.e.*, $p \leq q$) and full rank, and $B \in \mathbf{R}^{q \times r}$ is skinny (*i.e.*, $q \geq r$) and full rank. Then AB is full rank.
- (h) Suppose $A \in \mathbf{R}^{n \times n}$ has all eigenvalues equal to zero, and the nullspace of A is the same as the nullspace of A^2 . Then $A = 0$.
- (i) Consider the discrete-time linear dynamical system $x(t+1) = Ax(t) + Bu(t)$, where $A \in \mathbf{R}^{n \times n}$. Suppose there is an input that steers the state from a particular initial state x_{init} at time $t = 0$ to a particular final state x_{final} at time $t = T$, where $T > n$. Then there is an input that steers the state from x_{init} at time $t = 0$ to x_{final} at time $t = n$.

Lecture 19 – Observability and state estimation

19.1 *Sensor selection and observer design.* Consider the system $\dot{x} = Ax$, $y = Cx$, with

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

(This problem concerns observer design so we've simplified things by not even including an input.) (The matrix A is the same as in problem 11, just to save you typing; there is no other connection between the problems.) We consider observers that (exactly and instantaneously) reconstruct the state from the output and its derivatives. Such observers have the form

$$x(t) = F_0 y(t) + F_1 \frac{dy}{dt}(t) + \cdots + F_k \frac{d^k y}{dt^k}(t),$$

where F_0, \dots, F_k are matrices that specify the observer. (Of course we require this formula to hold for any trajectory of the system and any t , *i.e.*, the observer has to work!) Consider an observer defined by F_0, \dots, F_k . We say the *degree* of the observer is the largest j such that $F_j \neq 0$. The degree gives the highest derivative of y used to reconstruct the state. If the i th columns of F_0, \dots, F_k are all zero, then the observer doesn't use the i th sensor signal $y_i(t)$ to reconstruct the state. We say the observer *uses* or *requires* the sensor i if at least one of the i th columns of F_0, \dots, F_k is nonzero.

- (a) What is the minimum number of sensors required for such an observer? List all combinations (*i.e.*, sets) of sensors, of this minimum number, for which there is an observer using only these sensors.
- (b) What is the minimum degree observer? List all combinations of sensors for which an observer of this minimum degree can be found.