

# Bayesian Scientific Computing Day 1

Daniela Calvetti, Erkki Somersalo

Case Western Reserve University  
Department of Mathematics, Applied Mathematics and Statistics

Lyngby, December 2019

# Introduction

Mathematical model building has often a *natural* direction (causality, locality):

**Forward** problem:

Cause  $\longrightarrow$  Effect

Conceptually tractable, mathematically often the easier (but not necessarily easy) direction.

**Inverse** problem:

Effect  $\longrightarrow$  Cause

Conceptually and technically challenging (may not even have a solution).

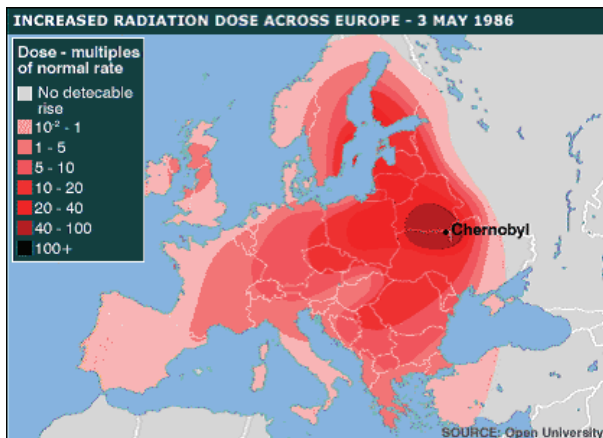
# Introduction



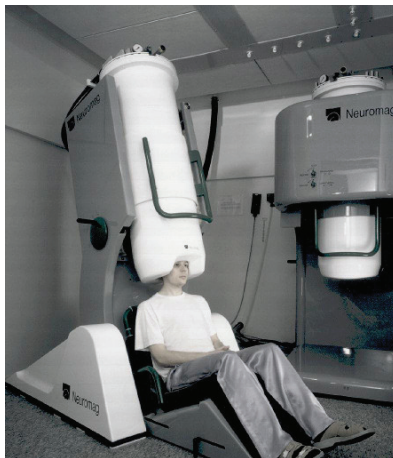
**Forward problem:** Predict the nuclear fallout

**Inverse problem:** Find the source of the leak

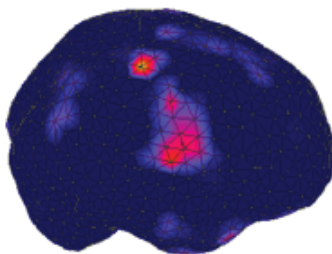
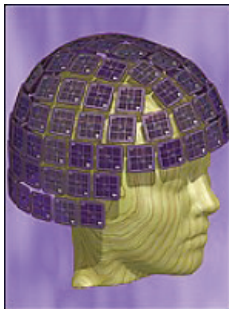
# Example: Non-cooperative Target



# Example: Magnetoencephalography (MEG)



# Example: Magnetoencephalography (MEG)



Biot-Savart law for a current dipole at  $\vec{r} = \vec{r}_0$ :

$$b(\vec{r}) = \vec{n}(\vec{r}) \cdot \vec{B}(\vec{r}) = \frac{\mu_0}{4\pi} \frac{\vec{n}(\vec{r}) \cdot \vec{q} \times (\vec{r} - \vec{r}_0)}{|\vec{r} - \vec{r}_0|^3}.$$

# Example: Magnetoencephalography (MEG)

Many dipoles, many measurements:

$$b_j = \sum_{n=1}^N \sum_{k=1}^3 a_{jn}^k q_k^n + \text{noise},$$

or in matrix form,

$$b = Ax + e.$$

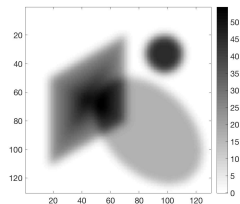
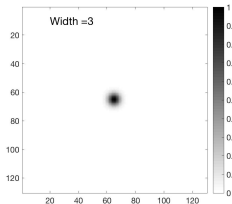
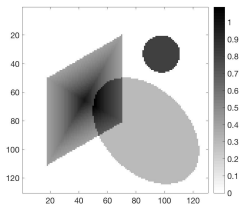
where

$$e = \text{noise}.$$

# Inverse problems and ill-posedness

Blurring by convolution,

$$g(r) = \int K(r - r')f(r')dr' + \varepsilon \xrightarrow{\text{discretize}} b = Ax + \varepsilon.$$





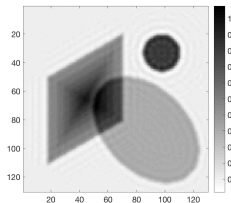
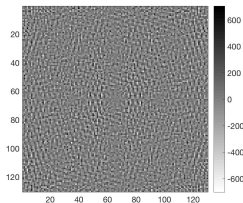
# Solving with or without regularization

Naive inversion,

$$x_0 = A^\dagger b,$$

versus Tikhonov regularized inversion,

$$x_\alpha = \operatorname{argmin}\{\|b - Ax\|^2 + \alpha^2\|x\|^2\}.$$



even with data with 5 digits of accuracy.

# Causality and Loss of Information

- Second Law of Thermodynamics: “The entropy of a closed system can never decrease:  $\Delta S \geq 0$ ”.
- Interpretation: The arrow of time points in larger entropy  $\Rightarrow$  In causal processes, information is lost (Entropy = lack of information)
- Solving inverse problems is a fight against the Second Law of Thermodynamics.
- To be successful, lost information needs to be replaced by extra information.

Bayesian methods are a systematic way  
to merge information from different sources.

# Example from Thermodynamics

Consider a thin rod of unit length, identified with the interval  $[0, 1]$ .

$$u(x, t) = \text{temperature at } x \in [0, 1] \text{ at time } t \geq 0.$$

End points held at fixed temperature,

$$u(0, t) = u(1, t) = a = \text{known constant}.$$

The temperature satisfies the heat equation,

$$\frac{\partial u}{\partial t} = D \frac{\partial^2 u}{\partial x^2},$$

where  $D$  is the heat conductivity of the material.

# Inverse Problem: Against Causality

**Inverse Problem:** *Given the final temperature distribution  $u(x, T)$ , is it possible to recover  $u(x, 0)$ ?*

$$u(x, 0) \xrightarrow{\text{forward}} u(x, T) \xrightarrow{\text{inverse}} u(x, 0).$$

**Question:** *What if we simply reverse time:  $t \rightarrow -t$ ? Define  $v(x, t) = u(x, T - t)$ ,*

$$\frac{\partial u}{\partial t} = D \frac{\partial^2 u}{\partial x^2}, \quad \longrightarrow \quad -\frac{\partial v}{\partial t} = D \frac{\partial^2 v}{\partial x^2},$$

with initial value

$$v(x, 0) = u(x, T).$$

Obviously  $v(x, T) = u(x, 0)$ . Can we solve it numerically? Let's try (Exercise)!

# Inverse Problems and Bayesian Inference

- In **inverse problems**, the goal is to estimate a quantity that is not directly observable, by using indirect observations
- In **statistical inference**, the goal is to estimate a probability distribution based on random draws from it

What is the connection?

- Noisy measurements are random draws from a distribution
- The distribution depends on the unknown quantity that we are interested in

**Bayesian approach:** *Develop a formalism to express what one can **believe** about the values of the unknown, given the measured data.*

# Probability

Two classical ways of *understanding* probability:

- *Frequentist's definition*: Probability can be understood in terms of frequencies in repeated experiments
- *Bayesian definition*: Probability is a subject's expression of degree of belief.

Bayesian probability = **subjective probability**

# Bayesian Subjective Probability

**Example:** Predict the continuation of a sequence:

26535	89793	23846	26433	83279	50288	41971
69399	37510	58209	74944	59230	78164	...

Subject 1:

- The sequence looks random
- Can be modeled as a realization of a random process
- The distribution is estimated by counting the frequencies of the numbers  $0 \dots 9$

# Bayesian Subjective Probability

**Example:** Predict the continuation of a sequence:

3, 14159	26535	89793	23846	26433	83279	50288	41971
	69399	37510	58209	74944	59230	78164	...

Subject 2:

- This looks like the decimal approximation of  $\pi$
- Prediction from a reference source
- Minimal or no uncertainty (although Subject 2 may be wrong!).



# Bayesian Subjective Probability

Bayesian probability:

- Expresses a subject's level of uncertainty (lack of information).
- Asserts that randomness is not the object's but the *subject's* property.
- May be subjective, but needs to be defensible.

**Note:** Subjective is *not* the same as arbitrary.

# Deterministic Quantities as Random Variables

*“Why should we interpret a fully deterministic quantity as a random variable?”*

Consider the following two “experiments”:

- ① Predict the outcome of a coin toss I’m going to make.
  - ② Guess the outcome of a coin toss I already made without showing you the outcome.
- In the first experiment, the outcome is “naturally” a random variable (result of a random process).
  - In the second one, the outcome is fully deterministic; however, you don’t know the outcome.

The same concepts of probability apply to both experiments.

# Bayes' Formula

Reverend Thomas Bayes (1701–1761), philosopher, theologian.

Bayes' formula appeared in his paper

- “An Essay towards solving a Problem in the Doctrine of Chances”, read posthumously in the Royal Society in 1763.
- Bayes's motivation: Studied the proofs for existence of God.

# Bayes' Formula

Pierre-Simon Laplace (1749 - 1827), French mathematician and natural scientist, one of the developers of modern probability and statistics.

- Gambling
- Astronomy
- “Mémoires sur la probabilité des causes par les évènements” (1774): *Inverse probability*.
- “Mémor on comets” (1813) Contains Bayes' formula in its present form, and a scientific application.

# Inverse Probability

Inverse probability (through an example)

- Given an urn with a known number of black and white balls, we can compute the probability of drawing (say) a white ball (“forward probability”).
- The inverse probability is to figure out the ratio between the numbers of white and black balls, given observed draws from the urn.

# Subjective Probability

BRUNO DE FINETTI 1906–1985



“Probability does not exist!”  
(B. de Finetti: *Theory of Probability*)

# Subjective Probability

Probability = uncertainty, or subject's lack of information.

*"The only relevant thing is uncertainty - the extent of our knowledge and ignorance. The actual fact of whether or not the events considered are in some sense determined, or known by other people, and so on, is of no consequence."*

– Bruno de Finetti

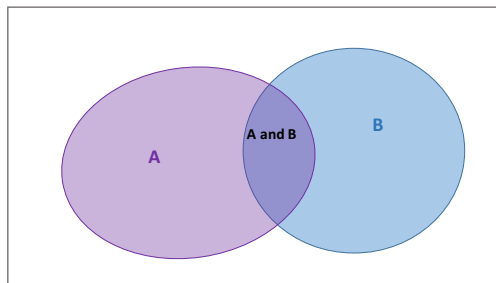
# Elementary Version of Bayes' Formula

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

- $A, B$  are events. The event  $B$  is observed.
- $P(A)$  = the probability of the event  $A$ .
- $P(B)$  = the probability of the event  $B$ .
- $P(B | A)$  = probability of  $B$  assuming that  $A$  happens.
- $P(A | B)$  = probability of  $A$  assuming that  $B$  happens.



# Visual Justification of Bayes' Formula



$$\begin{aligned}P(A, B) &= P(A | B)P(B) \\ &= P(B | A)P(A),\end{aligned}$$

$$P(A) = P(A | B)P(B) + P(A | \neg B)P(\neg B).$$

# Example

*A woman detects a lump in her breast and gets a mammography. Later, she is recalled for further diagnostics due to a positive mammography result. What is her probability of having breast cancer?*

- Event  $B$  = positive finding that makes the doctor suspect cancer.
- Event  $A$  = the patient has breast cancer.

# Example

Assume the following statistical data, based on long-term patient records:

	Malignant tumor	Benign tumor
Positive (recall)	0.8	0.6
Negative (no recall)	0.2	0.4

Interpretation in terms of conditional probabilities:

$$\begin{aligned}
 P(B \mid A) &= 0.8, & P(B \mid \neg A) &= 0.6 \\
 P(\neg B \mid A) &= 0.2, & P(\neg B \mid \neg A) &= 0.4
 \end{aligned}$$

Here,  $\neg A$  = “not A.”

# Example

More information: According medical records, for one out of four patients reporting a lump in the breast, the tumor is a malignant lesion.

Interpretation:

$$P(A) = 0.25, \quad P(\neg A) = 0.75.$$

We still need  $P(B)$ . This is obtained as follows.

$$P(B) = P(B \mid A)P(A) + P(B \mid \neg A)P(\neg A),$$

that is, “ $B$  may happen with  $A$  or with  $\neg A$ .” In numbers:

$$P(B) = 0.8 \times 0.25 + 0.6 \times 0.75 = 0.65.$$

# Example

Use Bayes' formula:

$$\begin{aligned}P(A | B) &= \frac{P(B | A)P(A)}{P(B)} \\&= \frac{0.8 \times 0.25}{0.65} \approx 0.31.\end{aligned}$$

We see that the spontaneous reaction “*The patient has cancer with probability 0.8*” is way too pessimistic!

# Prosecutor's Fallacy

*A newborn dies without a visible reason (SIDS). The probability of such event is very low; should one suspect that the parents have killed the baby?*

Sally Clark case, UK 1998: Two newborns died of SIDS, and the mother was accused of murder.

Argument: "1/73 million chance that SIDS happens independently twice in the same family."

A faulty reasoning leads to those suspicions, and a probabilistic one cleared the mother.

- Event  $B$  = Newborn dies with no visible reason.
- Event  $A$  = The parents have killed the baby.

Then

$$P(B) = P(B \mid A)P(A) + P(B \mid \neg A)P(\neg A),$$

How big must  $P(A)$  to support the murder conclusion?

# Random Variables, Continuum State Space

Probability distribution of a real valued random variable  $X$ ,

$$\mu_X((a, b)) = P\{a < X < b\}.$$

If a function  $\pi_X(x) \geq 0$  exists such that

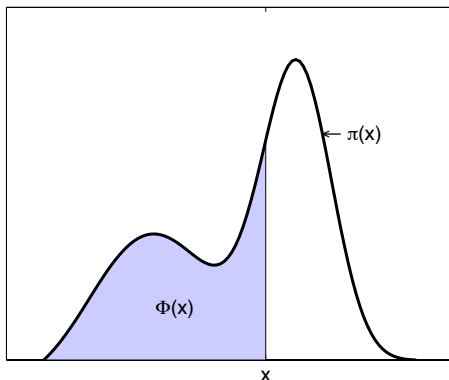
$$\mu_X(B) = \int_B \pi_X(x) dx, \quad B \subset \mathbb{R},$$

we refer to  $\pi_X$  as the *probability density of  $X$* .

Normalization:

$$\int_{\mathbb{R}} \pi_X(x) dx = P\{X \in \mathbb{R}\} = 1.$$

# Random Variables, Continuum State Space



Cumulative distribution

$$\Phi_X(x) = P\{X < x\} = \int_{-\infty}^x \pi_X(x') dx'.$$



# Joint Probability Density

For two random variables  $X$  and  $Y$ , define

$$\mu_{XY}(A \times B) = P\{X \in A, Y \in B\},$$

where  $\mu_{XY}$  is the joint probability distribution.

Assume the existence of the joint probability density  $\pi_{XY}$ ,

$$\mu_{XY}(A \times B) = \int_A \int_B \pi_{XY}(x, y) dy dx.$$

Again,

$$\int_{\mathbb{R}} \int_{\mathbb{R}} \pi_{XY}(x, y) dx dy = 1.$$

# Multivariate Random Variables

Real valued random variables  $X_1, X_2, \dots, X_n$ .

Define

$$X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} \in \mathbb{R}^n.$$

Joint probability density,

$$\pi_X(x_1, x_2, \dots, x_n) \geq 0,$$

defines the probability density in  $\mathbb{R}^n$ : For  $B \in \mathbb{R}^n$ ,

$$P\{X \in B\} = \int_B \pi_X(x) dx = \int \cdots \int_B \pi_X(x_1, \dots, x_n) dx_1 \dots dx_n.$$

# Marginal Densities

Define the marginal densities,

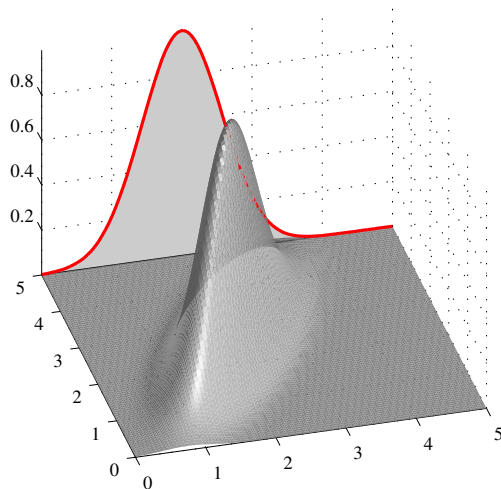
$$\pi_X(x) = \int_{\mathbb{R}} \pi_{XY}(x, y) dy, \quad \pi_Y(y) = \int_{\mathbb{R}} \pi_{XY}(x, y) dx,$$

which satisfy

$$P\{X \in A\} = \int_A \pi_X(x) dx = \int_A \int_{\mathbb{R}} \pi_{XY}(x, y) dy dx = P\{X \in A, Y \in \mathbb{R}\},$$

$$P\{Y \in B\} = \int_B \pi_Y(y) dy = \int_B \int_{\mathbb{R}} \pi_{XY}(x, y) dx dy = P\{Y \in B, X \in \mathbb{R}\}.$$

# Marginal Density Visualized



# Conditional Densities

Since

$$\pi_Y(y) = \int_{\mathbb{R}} \pi_{XY}(x, y) dx,$$

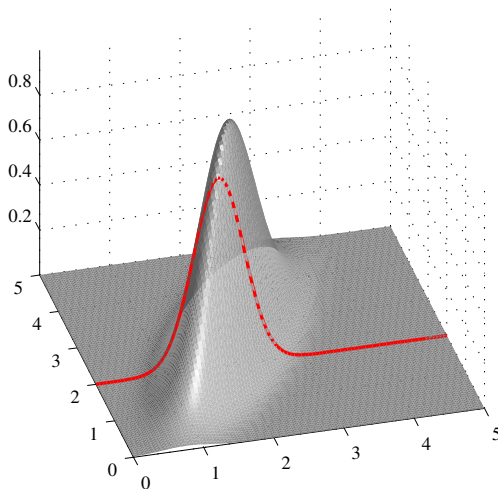
so if  $\pi_Y(y) \neq 0$ , we have

$$\frac{\pi_{XY}(x, y)}{\pi_Y(y)} \geq 0, \quad \int_{\mathbb{R}} \frac{\pi_{XY}(x, y)}{\pi_Y(y)} dx = 1.$$

Define a probability density  $\pi_{X|Y}(x | y)$  through the formula

$$\pi_{X|Y}(x | y) = \frac{\pi_{XY}(x, y)}{\pi_Y(y)}.$$

# Conditional Density Visualized



# Independency

Two random variables are **independent if and only if** the joint probability density can be factored as

$$\pi_{XY}(x, y) = \pi_X(x)\pi_Y(y).$$

**Observe:** For independent variables

$$\pi_{X|Y}(x | y) = \frac{\pi_{XY}(x, y)}{\pi_Y(y)} = \frac{\pi_X(x)\pi_Y(y)}{\pi_Y(y)} = \pi_X(x),$$

- By observing  $Y$  we learn nothing about  $X$ .

**Converse:** If

$$\pi_{X|Y}(x | y) = \pi_X(x),$$

the joint probability density factors and independency follows.

# Bayes' Formula

Interpretation: The conditional density is the probability density of  $X$ , assuming that the variable  $Y$  takes on value  $Y = y$ .

Symmetrically,

$$\pi_{Y|X}(y | x) = \frac{\pi_{XY}(x, y)}{\pi_X(x)},$$

assuming that  $\pi_X(x) \neq 0$ .

Solving for the joint density leads to

$$\pi_{XY}(x, y) = \pi_{X|Y}(x | y)\pi_Y(y) = \pi_{Y|X}(y | x)\pi_X(x),$$



# Bayes' Formula

The equation

$$\pi_{XY}(x, y) = \pi_{X|Y}(x | y)\pi_Y(y) = \pi_{Y|X}(y | x)\pi_X(x),$$

implies

**Bayes' formula for probability densities:**

$$\pi_{X|Y}(x | y) = \frac{\pi_{Y|X}(y | x)\pi_X(x)}{\pi_Y(y)}.$$

This is the key formula in Bayesian scientific computing

$\pi_X$  = prior density,

$\pi_{Y|X}$  = likelihood density,

$\pi_{X|Y}$  = posterior density.

posterior  $\propto$  prior  $\times$  likelihood.

# Expectation, Variance

Given a  $\mathbb{R}$ -valued random variable  $X$ , the *expectation* is the center of mass of the probability distribution,

$$E\{X\} = \int_{\mathbb{R}} x\pi_X(x)dx = \bar{x}.$$

The *variance* is the expectation of the squared deviation from the expectation,

$$\text{var}(X) = E\{(X - \bar{x})^2\} = \int_{\mathbb{R}} (x - \bar{x})^2\pi_X(x)dx,$$

assuming that the integrals converge.

# Example: Gaussian Distributions

A random variable  $X \in \mathbb{R}$  is normally distributed, or Gaussian,

$$X \sim \mathcal{N}(\mu, \sigma^2),$$

if

$$\mathbb{P}\{X \leq t\} = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^t \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) dx.$$

Probability density:

$$\pi_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right).$$

# Expectation, Variance

**Example:** Gaussian distribution:

$$\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} x \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) dx = \mu,$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{\infty} (x - \mu)^2 \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) dx = \sigma^2,$$

# Expectation, Variance

**Example:** Cauchy distribution:

$$\pi_X(x) = \frac{\alpha}{\pi} \frac{1}{1 + \alpha^2 x^2},$$

$$\int_{-\infty}^{\infty} |x| \pi_X(x) dx = \infty,$$

and therefore, the integral defining the expectation is non-convergent.

# Multivariate vs. Univariate

Define

$$X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

with each component  $X_i$  being an  $\mathbb{R}$ -valued variable.

Probability density of  $X$  = joint probability density  $\pi = \pi_X : \mathbb{R}^n \rightarrow \mathbb{R}_+$  of its components.

Expectation is

$$\bar{x} = \int_{\mathbb{R}^n} x \pi_X(x) dx \in \mathbb{R}^n,$$

or, componentwise,

$$\bar{x}_i = \int_{\mathbb{R}^n} x_i \pi_X(x) dx \in \mathbb{R}, \quad 1 \leq i \leq n.$$

# Multivariate vs. Univariate

Observe:

$$\begin{aligned} \int_{\mathbb{R}^n} x_1 \pi_X(x) dx &= \int_{\mathbb{R}} x_1 \underbrace{\int_{\mathbb{R}^{n-1}} \pi_X(x_1, x_2, \dots, x_n) dx_2 \dots dx_n}_{\text{marginal density of } x_1} \\ &= \int_{\mathbb{R}} x_1 \pi_{X_1}(x_1) dx_1. \end{aligned}$$

Similarly the other components.

# Multivariate vs. Univariate

The *covariance matrix* is defined as

$$\text{cov}(X) = \int_{\mathbb{R}^n} (x - \bar{x})(x - \bar{x})^T \pi_X(x) dx \in \mathbb{R}^{n \times n},$$

or, componentwise,

$$\text{cov}(X)_{ij} = \int_{\mathbb{R}^n} (x_i - \bar{x}_i)(x_j - \bar{x}_j) \pi_X(x) dx \in \mathbb{R}^{n \times n}, \quad 1 \leq i, j \leq n.$$



# Multivariate vs. Univariate

Recall: A matrix  $C \in \mathbb{R}^{m \times m}$  is *symmetric, positive definite* (SPD), if

1

$$C^T = C, \quad (\text{symmetry})$$

2

$$v^T C v > 0 \text{ for all } v \neq 0. \quad (\text{positive definiteness})$$

Equivalent conditions :

(2a.)  $C$  admits a Cholesky factorization,

$$C = R^T R, \quad R \text{ is upper triangular, } R_{jj} > 0.$$

(2b.)  $C$  admits an eigenvalue decomposition

$$C = V D V^T, \quad V \text{ orthogonal, } D \text{ diagonal, } D_{jj} > 0.$$

# Multivariate vs. Univariate

Covariance matrix is *symmetric* and *positive semi-definite*: For any  $v \in \mathbb{R}^n$ ,  $v \neq 0$ ,

$$\begin{aligned} v^T \text{cov}(X) v &= \int_{\mathbb{R}^n} [v^T (x - \bar{x})] [(x - \bar{x})^T v] \pi_X(x) dx \\ &= \int_{\mathbb{R}^n} (v^T (x - \bar{x}))^2 \pi_X(x) dx \geq 0. \end{aligned} \quad (1)$$

$v^T \text{cov}(X) v =$  variance of  $X$  into the direction  $v$ .

Usually, it is assumed that  $X$  has non-vanishing variance in all directions, and the covariance is SPD.

# Multivariate vs. Univariate

Denote by  $x'_i \in \mathbb{R}^{n-1}$  the vector  $x$  with the  $i$ th component deleted:

$$\begin{aligned}
 \text{cov}(X)_{ii} &= \int_{\mathbb{R}^n} (x_i - \bar{x}_i)^2 \pi_X(x) dx \\
 &= \int_{\mathbb{R}} (x_i - \bar{x}_i)^2 \underbrace{\left( \int_{\mathbb{R}^{n-1}} \pi_X(x_i, x'_i) dx'_i \right)}_{=\pi_{X_i}(x_i)} dx_i \\
 &= \int_{\mathbb{R}} (x_i - \bar{x}_i)^2 \pi_{X_i}(x_i) dx_i = \text{var}(X_i).
 \end{aligned}$$

**Conclusion:** Diagonal of the covariance matrix gives the variances of the individual components.

# Multivariate vs. Univariate

Consider two random variables,  $X_1$  and  $X_2$ , and assume that they are independent and Gaussian:

$$X_j \sim \mathcal{N}(0, \sigma_j^2), \quad j = 1, 2.$$

$$\begin{aligned} P\{a_1 < x_1 < b_1, a_2 < x_2 < b_2\} &= P\{a_1 < x_1 < b_1\}P\{a_2 < x_2 < b_2\} \\ &= \frac{1}{\sqrt{2\pi\sigma_1^2}} \int_{a_1}^{b_1} \exp\left(-\frac{1}{2\sigma_1^2}x_1^2\right) dx_1 \times \frac{1}{\sqrt{2\pi\sigma_2^2}} \int_{a_2}^{b_2} \exp\left(-\frac{1}{2\sigma_2^2}x_2^2\right) dx_2 \\ &= \int \int_Q \pi_{X_1 X_2}(x_1, x_2) dx_1, dx_2, \end{aligned}$$

# Multivariate vs. Univariate

where

$$Q = [a_1, b_1] \times [a_2, b_2] \subset \mathbb{R}^2,$$

and  $\pi_{X_1 X_2}(x_1, x_2) = \pi_X(x)$  is a two-dimensional Gaussian density,

$$\pi_X(x) = \frac{1}{\sqrt{(2\pi)^2 \sigma_1^2 \sigma_2^2}} \exp\left(-\frac{1}{2} \left(\frac{x_1^2}{\sigma_1^2} + \frac{x_2^2}{\sigma_2^2}\right)\right)$$

# Multivariate vs. Univariate

Defining

$$C = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix},$$

the two-dimensional Gaussian density can be written as

$$\begin{aligned} \pi_X(x) &= \frac{1}{\sqrt{(2\pi)^2 \sigma_1^2 \sigma_2^2}} \exp \left( -\frac{1}{2} \left( \frac{x_1^2}{\sigma_1^2} + \frac{x_2^2}{\sigma_2^2} \right) \right) \\ &= \frac{1}{\sqrt{(2\pi)^2 |C|}} \exp \left( -\frac{1}{2} x^T C^{-1} x \right), \end{aligned}$$

where

$$|C| = |\det(C)| = \sigma_1^2 \sigma_2^2.$$

# Multivariate vs. Univariate

This leads to the definition of a general *multivariate Gaussian probability density with independent components*:

Let

$$C = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}.$$

Define

$$\pi_X(x) = \frac{1}{\sqrt{(2\pi)^m |C|}} \exp\left(-\frac{1}{2} x^T C^{-1} x\right),$$

where  $|C| = |\det(C)| = \sigma_1^2 \sigma_2^2 \cdots \sigma_n^2$ .

# Multivariate vs. Univariate

Geometric intuition in two dimensions: The curve

$$\frac{x_1^2}{\sigma_1^2} + \frac{x_2^2}{\sigma_2^2} = \text{constant}$$

represents an ellipse with principal axes along the coordinate axes.

Add a rotation: Define

$$U = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix},$$

and rotated random variable  $X'$ ,

$$X' = U^T X.$$

Assume that the components  $X'_1$  and  $X'_2$  are independent Gaussian,

$$X'_j \sim \mathcal{N}(0, \sigma_j^2), \quad j = 1, 2.$$



# Multivariate vs. Univariate

If

$$D = \begin{bmatrix} \sigma_1^2 & \\ & \sigma_2^2 \end{bmatrix},$$

then the probability density of  $X'$  is

$$\pi_{X'}(x') = \frac{1}{\sqrt{(2\pi)^m |D|}} \exp\left(-\frac{1}{2}(x')^T D^{-1} x'\right),$$

Probability density of  $X = Ux'$ : Jacobian determinant of a rotation = 1, so

$$x = Ux' \Rightarrow dx = dx'.$$

## Multivariate vs. Univariate

$$\begin{aligned}
 \pi_X(x) &= \frac{1}{\sqrt{(2\pi)^m |D|}} \exp \left( -\frac{1}{2} (U^T x)^T D^{-1} U^T x \right) \\
 &= \frac{1}{\sqrt{(2\pi)^m |D|}} \exp \left( -\frac{1}{2} x^T (U D^{-1} U^T) x \right) \\
 &= \frac{1}{\sqrt{(2\pi)^m |C|}} \exp \left( -\frac{1}{2} x^T C^{-1} x \right),
 \end{aligned}$$

where

$$C = U D U^T, \quad |C| = |D|.$$

# Multivariate vs. Univariate

**Geometric interpretation:** The curve

$$x^T C^{-1} x = \text{constant}, \quad C \in \mathbb{R}^{2 \times 2}$$

is a rotated ellipse.

More generally, in  $\mathbb{R}^n$ ,

$$x^T C^{-1} x = \text{constant}, \quad C \in \mathbb{R}^{n \times n}$$

is an ellipsoidal hypersurface.

# Multivariate Gaussian

Multivariate extension of Gaussian densities in  $\mathbb{R}$ :  $X \in \mathbb{R}^n$  is Gaussian, if its probability density is

$$\pi_X(x) = \left( \frac{1}{(2\pi)^n |C|} \right)^{1/2} \exp \left( -\frac{1}{2} (x - \mu)^T C^{-1} (x - \mu) \right),$$

where  $\mu \in \mathbb{R}^n$ ,  $C \in \mathbb{R}^{n \times n}$  is symmetric positive definite, that is,

$$C^T = C, \quad v^T C v > 0 \quad \text{for all } v \neq 0.$$

$$P\{X \in B\} = \int_B \pi_X(x) dx.$$

Notation,

$$X \sim \mathcal{N}(\mu, C).$$

# Multivariate Gaussian

A straightforward computation shows: If  $X \sim \mathcal{N}(\mu, C)$ , then

$$\mathbb{E}\{X\} = \mu,$$

$$\text{cov}(X) = C.$$

# Inverse problems are ill-posed

A problem is called *well-posed* if

- It has a solution (existence)
- The solution is unique (uniqueness)
- Small perturbations in the problem setting lead to small perturbation in the solution (stability)

Problems that are not well-posed are called *ill-posed*

*Inverse problems are practically always ill-posed; they fail to satisfy at least one, but often several of the conditions defining well-posedness.*

# Linear algebra: Basic concepts

Inner product: If  $x, y \in \mathbb{R}^n$ ,

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

the inner product is defined as

$$x^T y = \sum_{j=1}^n x_j y_j = y^T x.$$

The vectors are said to be *orthogonal* if and only if

$$x^T y = 0.$$

# Linear algebra: Basic concepts

Matrix-vector product: If  $A \in \mathbb{R}^{m \times n}$ ,  $x \in \mathbb{R}^n$ , then

$$b = Ax \in \mathbb{R}^m, \quad b_k = \sum_{j=1}^n A_{kj}x_j, \quad 1 \leq k \leq m.$$

Denoting

$$A = \begin{bmatrix} | & | & \cdots & | \\ a_1 & a_2 & \cdots & a_n \\ | & | & \cdots & | \end{bmatrix}, \quad a_j \in \mathbb{R}^m,$$

we have

$$Ax = \sum_{j=1}^n x_j a_j = \text{linear combination of } a_1, \dots, a_n.$$



# Linear algebra: Basic concepts

Matrix-matrix product: Let  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times k}$ . Then

$$C = AB \in \mathbb{R}^{m \times k}, \quad C_{ij} = \sum_{\ell=1}^n A_{i\ell} B_{\ell j}, \quad 1 \leq i \leq m, \quad 1 \leq j \leq k.$$

Denoting

$$B = \begin{bmatrix} b_1 & b_2 & \cdots & b_k \end{bmatrix},$$

we have

$$AB = \begin{bmatrix} Ab_1 & Ab_2 & \cdots & Ab_k \end{bmatrix},$$

Recall:

$$(AB)^T = B^T A^T.$$

# Linear algebra: Basic concepts

Given vectors  $a_1, a_2, \dots, a_k \in \mathbb{R}^n$ , we say that the vectors are *linearly independent* if

$$c_1 a_1 + c_2 a_2 + \dots + c_k a_k = 0 \Leftrightarrow c_1 = c_2 = \dots = c_k = 0.$$

In matrix notation: Denoting

$$A = \begin{bmatrix} a_1 & a_2 & \dots & a_k \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad c = \begin{bmatrix} c_1 \\ \vdots \\ c_k \end{bmatrix},$$

the condition reads

$$Ac = 0_n \Leftrightarrow c = 0_k.$$

Here  $0_n$  ( $0_k$ ) is the null vector in  $\mathbb{R}^n$  ( $\mathbb{R}^k$ ).

# Linear algebra: Basic concepts

Given  $k$  vectors  $a_1, \dots, a_k$ , the subspace spanned by these vectors is defined as

$$H = \text{span}\{a_1, \dots, a_k\} = \{x \in \mathbb{R}^n \mid x = c_1 a_1 + \dots c_k a_k \\ \text{for some } c_1, \dots, c_k \in \mathbb{R}\}.$$

If the vectors  $a_1, \dots, a_k$  are linearly independent, they form a *basis* of  $H$ .  
The *dimension* of a subspace  $H \subset \mathbb{R}^n$  is the maximum number of independent vectors that span the subspace  $H$ .

# Linear algebra: Basic concepts

Given a matrix  $A \in \mathbb{R}^{m \times n}$ ,

- The *null space* of  $A$ , denoted as  $\mathcal{N}(A)$  is defined as

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

- The *range* of  $A$ , denoted as  $\mathcal{R}(A)$  is defined as

$$\mathcal{R}(A) = \{y \in \mathbb{R}^m \mid y = Ax \text{ for some } x \in \mathbb{R}^n\}.$$

# Linear algebra: Basic concepts

Another characterization for  $\mathcal{R}(A)$ : Denote the column vectors of  $A$  as  $a_1, \dots, a_n$ ,

$$A = \left[ \begin{array}{c|c|ccc|c} & & & & & \\ & a_1 & a_2 & \cdots & a_n & \\ & & & & & \end{array} \right], \quad a_j \in \mathbb{R}^m.$$

We have

$$Ax = \sum_{j=1}^n x_j a_j = \text{linear combination of } a_1, \dots, a_n.$$

Therefore,

$$\mathcal{R}(A) = \text{span}\{a_1, \dots, a_n\} = \begin{array}{l} \text{the linear space of all} \\ \text{linear combinations of columns of } A. \end{array}$$

# Linear algebra: Basic concepts

Consider the linear equation



$$Ax = b, \quad A \in \mathbb{R}^{m \times n}. \quad (2)$$

- If  $\mathcal{N}(A) \neq \{0\}$ , the solution of (??), if it exists, is non-unique: If  $x^* \in \mathbb{R}^n$  is a solution and  $0 \neq x_0 \in \mathcal{N}(A)$ , then

$$A(x^* + x_0) = Ax^* + Ax_0 = Ax^* = b,$$

and so  $x^* + x_0$  is also a solution.

- If  $\mathcal{R}(A) \neq \mathbb{R}^m$ , the problem (??) does not always have a solution: If  $b \notin \mathcal{R}(A)$ , then, by definition, there is no  $x \in \mathbb{R}^n$  such that (??) holds.

# Linear algebra: Basic concepts

The *rank* of a matrix  $A \in \mathbb{R}^{m \times n}$ , denoted as  $\text{rank}(A)$ , has two *equivalent* definitions:

- ①  $\text{rank}(A) = \max.$  number of linearly independent columns of  $A$ ,
- ②  $\text{rank}(A) = \max.$  number of linearly independent rows of  $A$ .

# Linear algebra: Basic concepts

Observations:

①  $\text{rank}(A) \leq \min\{m, n\}.$

② We have

$$\dim(\mathcal{R}(A)) = \dim(\text{span}\{a_1, \dots, a_n\}) = \text{rank}(A).$$

Therefore, if  $\text{rank}(A) < m$ , the problem  $Ax = b$  may not have a solution.



# Linear algebra: Basic concepts

Closer look at the null space: Let  $x \in \mathcal{N}(A)$ , that is  $Ax = 0$ . Then, for any  $y \in \mathbb{R}^m$ , we have

$$0 = y^T Ax = (y^T A)x = (A^T y)^T x \text{ for all } y \in \mathbb{R}^m.$$

The vectors  $A^T y$  span the subspace  $\mathcal{R}(A^T) \subset \mathbb{R}^n$ . Conclusion:

$$\mathcal{N}(A) \perp \mathcal{R}(A^T).$$

In fact, we can conclude that

$$\begin{aligned} \mathcal{N}(A) &= \mathcal{R}(A^T)^\perp = \{x \in \mathbb{R}^n \mid x \perp z \text{ for every } z \in \mathcal{R}(A^T)\} \\ &= \text{orthocomplement of } \mathcal{R}(A^T). \end{aligned}$$

# Linear algebra: Basic concepts

Denote

$$A = \begin{bmatrix} \alpha_1^\top \\ \alpha_2^\top \\ \vdots \\ \alpha_m^\top \end{bmatrix}, \quad \alpha_j \in \mathbb{R}^n,$$

or, equivalently,

$$A^\top = \begin{bmatrix} \alpha_1 & \alpha_2 & \cdots & \alpha_m \end{bmatrix} \in \mathbb{R}^{n \times m}.$$

Since

$$\mathcal{R}(A^\top) = \text{span}\{\alpha_1, \dots, \alpha_m\},$$

we conclude that

$$\dim(\mathcal{R}(A^\top)) = \text{rank}(A).$$

# Linear algebra: Basic concepts

Conclusion: If  $\text{rank}(A) < n$ , then  $\mathcal{R}(A^T) \neq \mathbb{R}^n$ , and its orthocomplement  $\mathcal{N}(A)$  is *at least* one-dimensional. In fact,

$$\dim(\mathcal{N}(A)) = n - \dim(\mathcal{R}(A^T)) = n - \text{rank}(A).$$

We add a third observation to the previous ones:

# Linear algebra: Basic concepts

Observations:

①  $\text{rank}(A) \leq \min\{m, n\}.$

② We have

$$\dim(\mathcal{R}(A)) = \dim(\text{span}\{a_1, \dots, a_n\}) = \text{rank}(A).$$

Therefore, if  $\text{rank}(A) < m$ , the problem  $Ax = b$  may not have a solution.

③ We have

$$\dim(\mathcal{N}(A)) = n - \dim(\mathcal{R}(A^T)) = n - \text{rank}(A),$$

Therefore, if  $\text{rank}(A) < n$ , the problem  $Ax = b$  cannot have a unique solution.

# Linear algebra: Conclusions

- ① To guarantee the existence of the solution  $Ax = b$ ,  $A \in \mathbb{R}^{m \times n}$ , we need to have

$$\text{rank}(A) = m.$$

- ② To guarantee uniqueness of the solution, we need to have

$$\text{rank}(A) = n.$$

- ③ The problem cannot be well-posed unless  $m = n$  and the matrix  $A$  is full rank,

$$\det(A) \neq 0,$$

that is, the matrix is square and invertible.

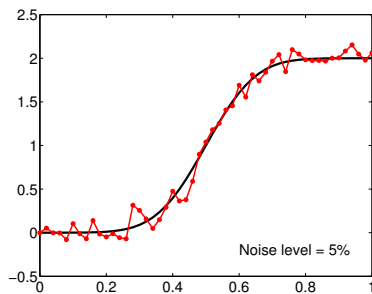
Therefore, *almost every linear problem you encounter is ill-posed!*

# Example: Numerical Differentiation

Let  $f : [0, 1] \rightarrow \mathbb{R}$  be a differentiable function,  $f(0) = 0$ .

$$\text{Data} = f(t_j) + \text{noise}, \quad t_j = \frac{j}{n}, \quad j = 1, 2, \dots, n.$$

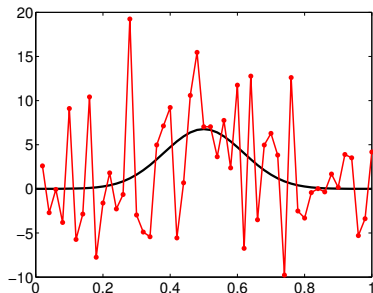
**Problem:** Estimate  $f'(t_j)$ .



# Naive Solution: Finite Difference Approximation

Approximate

$$f'(t_j) \approx \frac{f(t_j) - f(t_{j-1}))}{h}, \quad h = \frac{1}{n}.$$



# Naive Solution: Finite Difference Approximation

Where is the problem?

$$b_j = f(t_j) + e_j,$$

$$x_j = \frac{b_j - b_{j-1}}{h} = \underbrace{n(f(t_j) - f(t_{j-1}))}_{\approx f'(t_j)} + n(e_j - e_{j-1}).$$

The noise is amplified by a factor of the order  $\sim n = 50$ .



# Formulation as an Inverse Problem

Denote  $g(t) = f'(t)$ . Then,

$$f(t) = \int_0^t g(\tau) d\tau.$$

Linear model:

$$\text{Data} = b_j = f(t_j) + e_j = \int_0^{t_j} g(\tau) d\tau + e_j,$$

where  $e_j$  is the noise.

# Discretization

Write

$$\int_0^{t_j} g(\tau) d\tau \approx \frac{1}{n} \sum_{k=1}^j g(t_k).$$

By denoting  $g(t_k) = x_k$ ,

$$b = Ax + e,$$

where

$$A = \frac{1}{n} \begin{bmatrix} 1 & & & & \\ 1 & 1 & & & \\ \vdots & & \ddots & & \\ 1 & 1 & \dots & 1 & \end{bmatrix}.$$

# Properties of the Matrix

The matrix  $A$  is invertible: For a triangular matrix,

$$\det(A) = a_{11}a_{22} \cdots a_{nn} = \left(\frac{1}{n}\right)^n \neq 0.$$

Therefore,

$$\text{rank}(A) = m = n,$$

and the problem has a unique solution. One can check that

$$A^{-1} = n \begin{bmatrix} 1 & & & & \\ -1 & 1 & & & \\ & \ddots & \ddots & & \\ & & -1 & 1 & \\ & & & & \end{bmatrix},$$

corresponding to the finite difference method.

# Properties of the Matrix

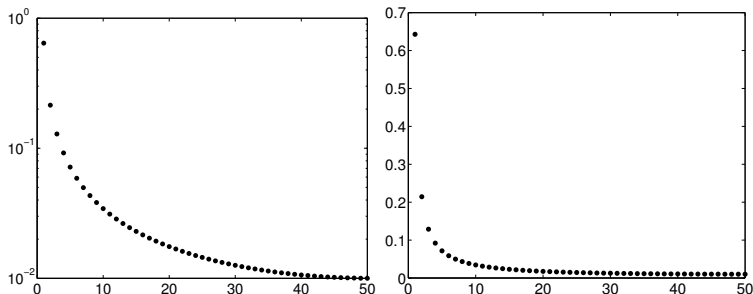
Singular values (revisited soon):

$$A = UDV^T,$$

where  $U, V \in \mathbb{R}^{n \times n}$  are orthogonal matrices,

$$D = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{bmatrix}, \quad d_1 \geq d_2 \geq \dots d_n > 0.$$

# Singular Value Analysis



Ratio of smallest and largest singular values (= condition number of A):

$$r(A) = \frac{d_1}{d_n} = 64.27$$

In inverse problems,  $r(A)$  is often *much* larger than this ( $10^6$ – $10^8$ )!

# SVD Revisited

*Fundamental result:* Any rectangular real matrix  $A \in \mathbb{R}^{m \times n}$  can be decomposed as

$$A = UDV^T,$$

where  $U \in \mathbb{R}^{m \times m}$ ,  $V \in \mathbb{R}^{n \times n}$  are orthogonal matrices, that is,

$$UU^T = U^TU = I_m, \quad VV^T = V^TV = I_n,$$

and  $D \in \mathbb{R}^{m \times n}$  is a diagonal matrix. The diagonal entries of  $D$  are non-negative, and are usually ordered in a non-decreasing order,

$$d_1 \geq d_2 \geq \dots \geq d_{\min(m,n)} \geq 0.$$

# SVD Revisited

*Remark 1:* Denote

$$V = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix}, \quad v_j \in \mathbb{R}^n.$$

Orthogonality of  $V$  is *equivalent* to saying that the column vectors  $v_j$  form an orthonormal basis of  $\mathbb{R}^n$ .

In particular, for every  $x \in \mathbb{R}^n$ ,

$$x = \sum_{j=1}^n (v_j^T x) v_j,$$

where

$v_j^T x$  = orthogonal projection of  $x$   
onto the subspace spanned by  $v_j$ .

# SVD Revisited

*Remark 2:* Orthogonality of a matrix  $V$  means that the action  $x \mapsto Vx$  is

- a rotation, or
- a reflection, or
- a permutation.

that is, a conformal (= angle-preserving) linear transformation.

$$(Vx)^T Vy = x^T \underbrace{V^T V}_{=I_n} y = x^T y,$$

that is,

$$\cos \angle(Vx, Vy) = \cos \angle(x, y), \quad \|Vx\| = \|u\|.$$



# SVD Revisited

*Remark 3:* A rectangular (=non-square) diagonal matrix looks like

$$D = \begin{bmatrix} d_1 & & 0 & \cdots & 0 \\ & \ddots & \vdots & & \vdots \\ & & d_m & 0 & \cdots & 0 \end{bmatrix}, \quad m < n,$$

and

$$D = \begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_n \\ 0 & \cdots & 0 \\ \vdots & & \vdots \\ 0 & \cdots & 0 \end{bmatrix}, \quad m > n.$$

# SVD in Action

Given the SVD of a matrix  $A \in \mathbb{R}^{m \times n}$ ,

$$Ax = \sum_{j=1}^{\max(m,n)} u_j d_j (v_j^T x) = \sum_{j=1}^r u_j d_j (v_j^T x),$$

where  $r$  is the number of non-zero singular values,

$$d_1 \geq d_2 \geq \dots \geq d_r > d_{r+1} = \dots d_{\min(m,n)} = 0.$$

$$r = \text{rank}(A).$$

# SVD in Action

Null space of a matrix:



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



$$Ax = \sum_{j=1}^r u_j d_j (v_j^T x) = 0,$$

if and only if

$$x \perp v_1, v_2, \dots, v_r,$$

that is,

$$x = \sum_{j=r+1}^n x_j v_j \in \text{span}\{v_{r+1}, \dots, v_n\}.$$

# SVD in Action

Range of a matrix

$$\mathcal{R}(A) = \{y \in \mathbb{R}^m \mid y = Ax \text{ for some } x \in \mathbb{R}^n\} = ?$$

$$b = Ax = \sum_{j=1}^r u_j d_j (v_j^T x)$$

for some  $x$  if and only if

$$b = \sum_{j=1}^r (u_j^T b) u_j \in \text{span}\{u_1, u_2, \dots, u_r\}.$$

# SVD in Action

Solving a linear equations in terms of the SVD:

$$b = Ax,$$



or in terms of SVD,

$$\sum_{j=1}^m (u_j^T b) u_j = \sum_{j=1}^r d_j (v_j^T x) u_j.$$

# SVD in Action

- Solution exists if and only if  $b \in \mathcal{R}(A)$ , that is,

$$b \perp u_{r+1}, \dots, u_m.$$

- For the solution to exist, we must have

$$d_j(v_j^T x) = (u_j^T b), \quad 1 \leq j \leq r.$$

- For any  $x_0 \in \mathcal{N}(A)$ , the vector

$$x = \sum_{j=1}^r \frac{(u_j^T b)}{d_j} v_j + x_0 = \sum_{j=1}^r \frac{(u_j^T b)}{d_j} v_j + \sum_{j=r+1}^n x_j v_j$$

is a solution.

# Four Fundamental Subspaces

Given a matrix  $A = UDV^T \in \mathbb{R}^{m \times n}$ ,

$\text{Rank}(A) = r = \# \text{ of non-zero singular values} \leq \min(m, n)$ ,

$$U = \left[ \begin{array}{c|c} \leftarrow r \rightarrow & \leftarrow (m-r) \rightarrow \\ \mathcal{R}(A) & \mathcal{N}(A^T) \end{array} \right] \in \mathbb{R}^{m \times m},$$

$$V = \left[ \begin{array}{c|c} \leftarrow r \rightarrow & \leftarrow (n-r) \rightarrow \\ \mathcal{R}(A^T) & \mathcal{N}(A) \end{array} \right] \in \mathbb{R}^{n \times n},$$

# Backwards Heat Equation and Linear Algebra

Heat equation:

$$\frac{\partial u}{\partial t} = D \frac{\partial^2 u}{\partial x^2}, \quad u(0, t) = u(1, t) = a, \quad u(x, 0) = u_0(x).$$

Look for the solution of the form

$$u(x, t) = a + \sum_{j=1}^{\infty} u_j(t) \sin \pi j x,$$

and initial value as

$$u_0(x) = a + \sum_{j=1}^{\infty} \alpha_j \sin \pi j x,$$



# Backwards Heat Equation and Linear Algebra

Substitute into the heat equation,

$$\sum_{j=1}^{\infty} u_j'(t) \sin \pi j x = - \sum_{j=1}^{\infty} D(\pi j)^2 u_j(t) \sin \pi j x,$$

which is satisfied if

$$u_j'(t) = -D(\pi j)^2 u_j(t), \quad u_j(0) = \alpha_j.$$

Solution:

$$u_j(t) = \alpha_j e^{-D(\pi j)^2 t}.$$

# Backwards Heat Equation and Linear Algebra

Forward map:

$$A : u(x, 0) \mapsto u(x, T),$$

$$a + \sum_{j=1}^{\infty} \alpha_j \sin \pi j x \mapsto a + \sum_{j=1}^{\infty} \alpha_j e^{-D(\pi j)^2 T} \sin \pi j x.$$

Inverse map (formally)

$$A^{-1} : u(x, T) \mapsto u(x, 0),$$

$$a + \sum_{j=1}^{\infty} \beta_j \sin \pi j x \mapsto a + \sum_{j=1}^{\infty} \beta_j e^{D(\pi j)^2 T} \sin \pi j x.$$

# In Matrix Language

Finite approximation (set  $a = 0$ ): Discretize,

$$x_k = \frac{k}{n}, \quad 1 \leq k \leq n-1, \quad (\text{interior points})$$

then truncate the trigonometric series,

$$u(x_k, 0) = \sqrt{\frac{2}{n}} \sum_{j=1}^{n-1} \alpha_j \sin \pi \frac{jk}{n},$$

or

$$U_0 = \begin{bmatrix} u(x_1, 0) \\ \vdots \\ u(x_{n-1}, 0) \end{bmatrix} = S \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_{n-1} \end{bmatrix},$$

where  $S \in \mathbb{R}^{(n-1) \times (n-1)}$  is the matrix with entries

$$S_{jk} = \sqrt{\frac{2}{n}} \sin \pi \frac{jk}{n}, \quad 1 \leq j, k \leq n-1.$$

Similarly, we write

$$u(x_k, T) = \sqrt{\frac{2}{n}} \sum_{j=1}^{n-1} \beta_j \sin \pi \frac{jk}{n},$$

or, in matrix form,

$$U_T = S\beta.$$

From the spectral analysis, we know that

$$\beta_j = e^{-D(\pi j)^2 T} \alpha_j = \lambda_j \alpha_j,$$

or, in matrix form,

$$\beta = \Lambda \alpha,$$

where

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_{n-1} \end{bmatrix} \in \mathbb{R}^{(n-1) \times (n-1)}.$$

You can verify (using trigonometry or numerically) that

$$S^{-1} = S^T = S.$$

Therefore:

$$U_T = S\beta = S\Lambda\alpha = S\Lambda S^T U_0.$$

We have found a propagation matrix

$$A = S\Lambda S^T : U_0 \mapsto U_T.$$

It is invertible (in theory),

$$A^{-1} = S\Lambda^{-1}S^T,$$

but this formula is useless in practice.

# Summary

- The ill-posedness of a linear problem depends on the dimensions of the matrix and the rank of it
- SVD gives the means to determine the rank of the matrix, the basis of the null space and the basis of the range
- Sensitivity to noise depends on the distribution of the singular values (condition number)
- SVD is a useful tool for analyzing linear problems, but may be impractical to compute