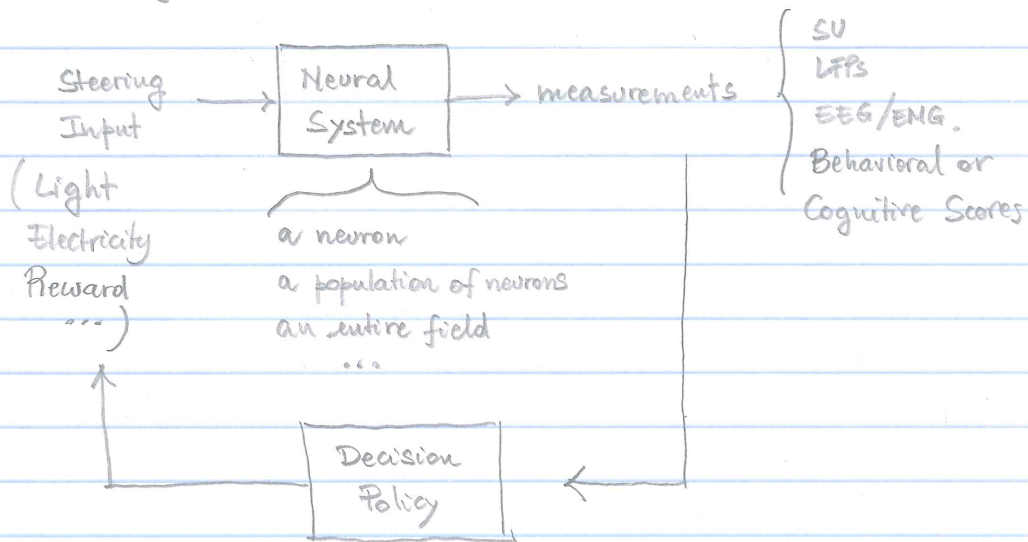


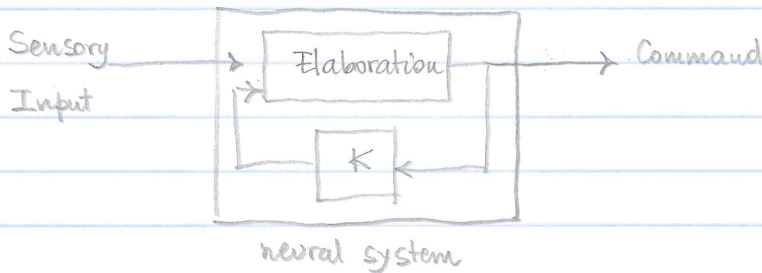
LECTURE 1

* Neural Control : what do we mean ?

1) Control = Exogenous Intervention

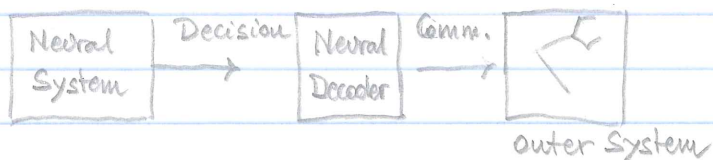


2) Control = Internal Regulation



Ex: Motor Commands
Learning, etc.

3) Control = Actuation of Outer System



②

In all three cases we need to provide:

- A representation of what the neural system does in response to an input
- A characterization of the uncertainty that affects such response



We need to introduce models (to represent) and estimation tools (to characterize) - We also need objectives to constrain the models



The topics of the course are:

- Modeling (Input-State-Output)
- Estimation (State, Parameters)
- Applications to:

- Regulation of Neural Activity
- Decoding of Movements
- Execution of Movements

□

* Tools that we need to carry on the estimation

1) Least-Squares

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \begin{array}{l} \text{- measurements} \\ \text{(e.g.; firing rates} \\ \text{in numerous repetitions)} \end{array}$$

$$X \triangleq \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \begin{array}{l} \text{- explicative} \\ \text{variables} \\ \text{(} m < n \text{)} \end{array}$$

For instance, $m=1$ and x_1 is the actual firing rate (which is

supposed constant) - We have:

$$\begin{aligned} \text{error}_1 &= y_1 - x_1 \\ \text{error}_2 &= y_2 - x_1 \\ &\vdots \\ \text{error}_n &= y_n - x_1 \end{aligned}$$

$\underbrace{\hspace{10em}}_E$

$$\Leftrightarrow E = Y - \begin{matrix} \text{w times} \\ \left[\begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right] x_1 = Y - A x_1 \\ \underline{\underline{A}} \end{matrix}$$

In general, if $m > 1$, then A is $n \times m$. Moreover coefficients of A can be arbitrary:

$$E = Y - AX$$

Problem: What is the estimation \hat{X} that minimizes the norm of E ?

Solution:

$$\begin{aligned} \|E\|^2 &\triangleq \sum_{i=1}^n \text{error}_i^2 = (Y - AX)^T (Y - AX) = \\ &= Y^T Y - X^T A^T Y - Y^T A X + X^T A^T A X = f(X) \quad (*) \end{aligned}$$

The min of $f(X)$ (which is a scalar function) has derivatives equal to zero \Rightarrow We impose:

$$\left. \begin{aligned} \frac{\partial f}{\partial x_1} &= 0 \\ \frac{\partial f}{\partial x_2} &= 0 \\ &\vdots \\ \frac{\partial f}{\partial x_m} &= 0 \end{aligned} \right\} \text{If } X \text{ were scalar, we would derive from } (*):$$

$$-A^T Y - Y^T A + 2A^T A X = 0$$

\Updownarrow Since it's scalar

$$2A^T A X - 2A^T Y = 0$$

\Updownarrow

$$\hat{X} = (A^T A)^{-1} A^T Y$$

4

The same result holds if X is a $m \times 1$ vector:

$$\hat{X} = (A^T A)^{-1} A^T Y \quad (**)$$

Note: $(**)$ is a projection from the n -dim space of Y to the m -dim space of X (\Rightarrow We are linearly combining all measurements in Y to extract an information about X)



What if not all the measurements in Y are equally good?

We introduce a matrix of weights W , e.g., a diagonal matrix:

$$W \triangleq \begin{bmatrix} \sigma_1 & & 0 \\ & \sigma_2 & \\ 0 & & \ddots \\ & & & \sigma_n \end{bmatrix} \quad \text{and we assign each entry to one specific error:}$$

$$\|W\epsilon\|_2^2 = (WY - WAX)^T (WY - WAX)$$

By repeating the argument above, we have:

$$2A^T W^T W Y - 2A^T W^T W A X = 0$$



$$\hat{X} = (A^T W^T W A)^{-1} A^T W^T W Y$$

Note: $W^T W = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \sigma_2^2 & \\ 0 & & \ddots \\ & & & \sigma_n^2 \end{bmatrix} \Rightarrow$ We are weighting each measurement a quadratic value

Note: ϵ is a $n \times 1$ vector $\Rightarrow R = \text{cov}(\epsilon)$ is a $n \times n$ matrix \Rightarrow It can be

shown that, by choosing $W^T W = R^{-1}$, the solution \hat{X} is the best unbiased estimator of X :

$$\hat{X} = (A^T R^{-1} A)^{-1} A^T R^{-1} Y$$

This choice of $W^T W$ is convenient also for another reason:

\hat{X} is an estimation \Rightarrow There is an error between \hat{X} and X of the actual X

\Rightarrow The covariance of this error is given by:

$$P \triangleq \underset{\substack{\uparrow \\ \text{expected} \\ \text{value}}}{E} [(X - \hat{X})(X - \hat{X})^T] = (A^T R^{-1} A)^{-1}$$

Example: $R \triangleq \begin{bmatrix} \sigma^2 & & & \\ & \sigma^2 & & \\ & & \ddots & \\ & & & \sigma^2 \end{bmatrix} \begin{matrix} \\ \\ \\ \underbrace{\hspace{2cm}} \\ \\ \end{matrix} \left. \begin{matrix} \\ \\ \\ \end{matrix} \right\} \Rightarrow A^T R^{-1} A = [1 \dots 1] \begin{bmatrix} 1/\sigma^2 & & \\ & \ddots & \\ & & 1/\sigma^2 \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

$A \triangleq \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad n \times 1$

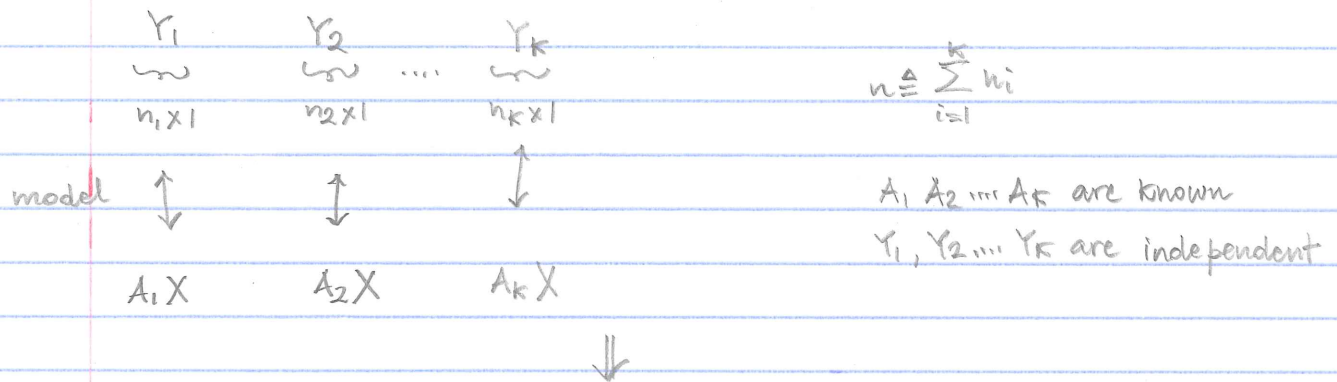
$= \sum_{i=1}^n \frac{1}{\sigma^2} = \frac{n}{\sigma^2}$

$\Rightarrow P = (A^T R^{-1} A)^{-1} = \sigma^2/n$ - The larger n , the lower the covariance of the error

Note: The estimation of P and \hat{X} proceeds in "batch" mode (\Rightarrow All measurements must be available in Y) and assumes that the A is constant \Rightarrow What if A changes in time?

6

Assume that there are consecutive batches of measurements



$$E = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_k \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_k \end{bmatrix} X \quad \text{X - Block partition}$$

If we consider the first batch, we can write:

$$P_1^{-1} = A_1^T R_1^{-1} A_1 \quad \text{where } R_1 \text{ is the covariance matrix of the error computed for the first batch}$$

$$\hat{X}_1 = P_1 A_1^T R_1^{-1} Y_1$$

If we consider the first two batches, we have:

$$P_2^{-1} = \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} R_1^{-1} & 0 \\ 0 & R_2^{-1} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$$

$$= A_1^T R_1^{-1} A_1 + A_2^T R_2^{-1} A_2 = P_1^{-1} + A_2^T R_2^{-1} A_2 \quad (a)$$

$$\hat{X}_2 = P_2 \begin{bmatrix} A_1^T & A_2^T \end{bmatrix} \begin{bmatrix} R_1^{-1} & 0 \\ 0 & R_2^{-1} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} =$$

$$\begin{aligned}
 &= P_2 \left(A_1^T R_1^{-1} Y_1 + A_2^T R_2^{-1} Y_2 \right) \\
 &= P_2 \left(P_1^{-1} \hat{X}_1 + A_2^T R_2^{-1} Y_2 \right) \\
 &= P_2 \left(P_2^{-1} \hat{X}_1 - A_2^T R_2^{-1} A_2 \hat{X}_1 + A_2^T R_2^{-1} Y_2 \right) \\
 &\uparrow \\
 \text{Replace } P_1^{-1} \text{ with (a)} &= \hat{X}_1 + \underbrace{P_2 A_2^T R_2^{-1}}_{\hat{= k}_2 \text{-correction factor}} \underbrace{\left(Y_2 - A_2 \hat{X}_1 \right)}_{\text{error (innovation)}}
 \end{aligned}$$

$\Rightarrow \hat{X}_2$ and P_2^{-1} are obtained recursively by using \hat{X} and P^{-1} at the previous step \Rightarrow The procedure can be generalized:

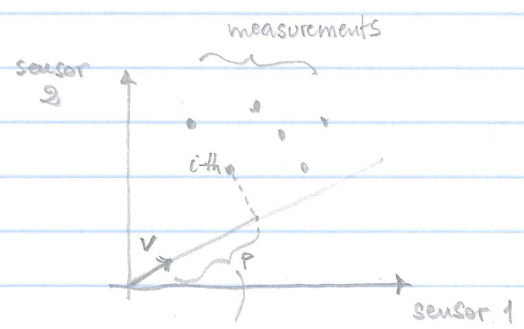
$$P_i^{-1} = P_{i-1}^{-1} + A_i^T R_i^{-1} A_i$$

$$K_i = P_i A_i^T R_i^{-1}$$

$$\hat{X}_i = X_{i-1} + K_i \left(Y_i - A_i \hat{X}_{i-1} \right)$$

2) Singular Value Decomposition

$$A = \left[\begin{array}{c} n \times m \\ \text{e.g., measurements} \end{array} \right] \left. \begin{array}{l} \text{e.g., sensors} \\ n \gg m \end{array} \right\}$$



$$p = \underbrace{a_i}_{\substack{\uparrow \\ \text{i-th row}}} v \Rightarrow \text{For all rows: } \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{bmatrix} = Av$$

projection of the i-th row of A on the direction of vector v

$v_1 \hat{=}$ first singular vector is the vector that maximizes $\|Av\|^2$

⑧

Analogously, one can define:

$v_2 \triangleq$ second singular vector is the vector that maximizes $\|Av\|^2$
conditioned to being perpendicular to v_1

The procedure can be iterated and I can find m orthogonal unit vectors that form a base for the space of rows of A

$$V \triangleq [v_1 | v_2 \dots | v_m]$$

For each v_i , let us call: $\sigma_i \triangleq \|Av_i\|$ - i -th singular value

Now, let us consider the transformation that A performs on the vectors v_i

$u_i \triangleq \frac{1}{\sigma_i} Av_i$ - $n \times 1$ vector \Rightarrow It is the i -th LEFT singular vector of A (it has length = 1)

$$U = [u_1 | u_2 | \dots | u_m]$$

Now note thws:

$$\|Av_i\|^2 = v_i^T A^T A v_i = \sigma_i^2 \Rightarrow v_i \text{ is an eigenvector of } A^T A$$

$$\|A^T u_i\|^2 = u_i^T A A^T u_i = \frac{1}{\sigma_i^2} v_i^T \underbrace{A A^T A}_{\sigma_i^2 v_i} v_i = \frac{1}{\sigma_i^2} \sigma_i^2 v_i^T \cdot \sigma_i^2 v_i = \sigma_i^2$$

because
of the definition
of u_i

$\Rightarrow u_i$ is an eigenvector of $A A^T$

Therefore, the following happens:

$$[\sigma_1 u_1 \quad \sigma_2 u_2 \quad \dots \quad \sigma_m u_m] = A V$$

and by defining $S \stackrel{\text{def}}{=} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_m \\ & & & & \dots \\ & & & & & & 0 \end{bmatrix}$, one can write:

diagonal matrix
($m \times m$)

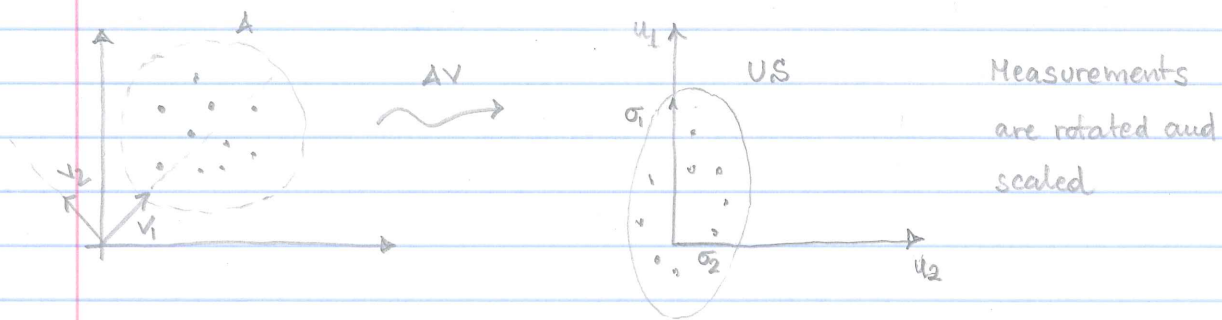
$US = AV \Rightarrow A = USV^T$ - Singular Value Decomposition of A

Note the construction process \Rightarrow Decomposition is not unique

$\Rightarrow V$ is a set of orthonormal eigenvectors of ATA (It always exists because ATA is symmetric and with real entries), i.e.:

$$V^T V = I_m$$

Example:



Note: AA^T is an $n \times n$ matrix \Rightarrow It may have n eigenvectors

Therefore, a typical SVD computes:

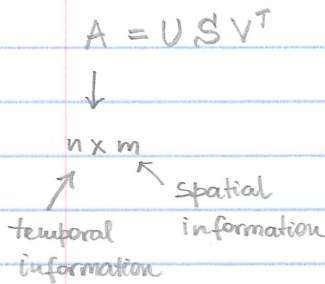
$$U: n \times n$$

$$V: m \times m$$

$$S = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_m \\ & & & & \dots \\ & & & & & & 0 \end{bmatrix} : n \times m \text{ with only up to } m \text{ nonnegative singular values}$$

10

Let us return to the case of A being a matrix of measurements:



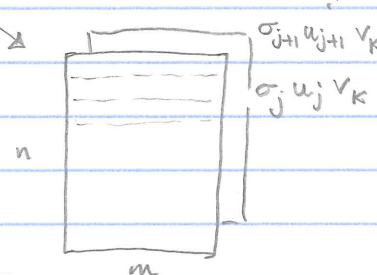
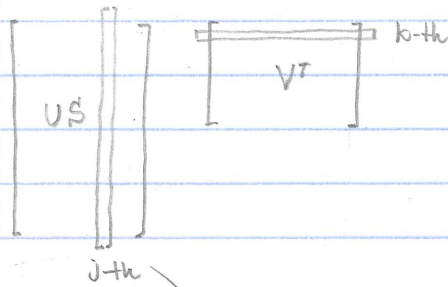
U is $n \times n \Rightarrow$ It is a collection of n temporal "modes"

V is $m \times m \Rightarrow$ It is a collection of m spatial "modes"

It is called "matrix of scores"

US is $n \times m \Rightarrow$ Only m temporal modes are retained and each one is weighted by the correspondent singular value

According to this interpretation, we are saying that the k -th right singular vector contributes to the measurements collected by every sensor but, at each time point, the contribution is different across temporal modes



Note: we are not adding or removing information. We are simply rearranging information via linear operations

What is the advantage of doing this?

Data compression: Since σ_i are ordered from the largest to the smallest, one can approximate A with the summation of just a few elements in the stack of matrices $\sigma_j u_j v_k$

A criterion to stop: Choose h such that $\frac{\sum_{i=1}^h \sigma_i}{\sum_{i=1}^m \sigma_i}$ is above a

given threshold and use the first h singular vectors \Rightarrow That ratio indicates the probability of the first h components accounting for the energy (or variance) of A

Data separation: Different spatial modes may correspond to different players (e.g., sources) \Rightarrow Data points mainly contributed by one mode are likely related to that source

Noise cancellation: The first singular value (if mean is not removed from the data) and the smallest singular values may be associated with non-physiologically-meaningful modes \Rightarrow Removal allows to increase the SNR of the data points of interest □

