

Announcements:

- 0.25 EC point for each lecture you attend for rest of the term.
- links.eecs16b.org/lecture-ec
- Lab Design Contest: See Ed post (EC opportunities!)
- Course Evaluations: 1.5 EC pts. with $\geq 80\%$ response rate.

Last time:

• PCA motivation/introduction

$$A_l = \sum_{i=1}^l \sigma_i \vec{u}_i \vec{v}_i^T \quad (\text{Note: } A = \sum_{j=1}^r \sigma_j \vec{u}_j \vec{v}_j^T)$$

"effective" rank of data matrix $A = l \ll n$ typically

- Young-Eckart Theorem:

$$\min \|A - B\|_F^2 *$$

s.t. B has rank l

is solved by $B = A_l = \sum_{i=1}^l \sigma_i \vec{u}_i \vec{v}_i^T$

* $(\|X\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n |X_{ij}|^2)$ e.g. $\| \begin{bmatrix} 1 & 3 & 2 \\ 2 & 0 & -1 \end{bmatrix} \|_F = 19$

Today:

- Wrapup PCA

- Derive $l=1$ case: $A_1 = \sigma_1 \vec{u}_1 \vec{v}_1^T$ is best Rank-1 approx. to A .

See Note 16

- Linearization for control
 - linear approx to non-linear function
 - Equilibrium/operating points
 - examples

User j : "sensitivity" or "affinity" for components, i.e. how much user j "likes" or "dislikes" the attributes or "features".

$$U_j : [s_{aj} \quad s_{bj} \quad s_{cj} \quad s_{dj}]$$

q_{ij} is the inner product between Movie i 's feature vector & User j 's sensitivity to features vector

$$q_{ij} = s_{aj} \cdot a_i + s_{bj} \cdot b_i + s_{cj} \cdot c_i + s_{dj} \cdot d_i$$

e.g. $q_{ij} = 80(0.2) + 0(0.1) + 77(0) + 20(0.7) = \boxed{30}$

$$\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{100} \end{bmatrix}$$

action scores of all videos $S_a = \in \mathbb{R}^{100}$
Similarly, for $\vec{b}, \vec{c}, \vec{d}$.

$$\begin{bmatrix} s_{a1} \\ s_{a2} \\ \vdots \\ s_{a1000} \end{bmatrix}$$

"action" sensitivity vector for users $\in \mathbb{R}^{1000}$
Similarly, for $\vec{s}_b, \vec{s}_c, \vec{s}_d$

Consider

$$\vec{a} \cdot \vec{s}_a^T$$

("outer product")

$$= \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{100} \end{bmatrix} \begin{bmatrix} s_{a1} & s_{a2} & \dots & s_{a1000} \end{bmatrix}$$

$$= \begin{bmatrix} a_1 s_{a1} & a_1 s_{a2} & \dots & a_1 s_{a1000} \\ a_2 s_{a1} & & & \vdots \\ \vdots & & & \\ a_{100} s_{a1} & \dots & \dots & a_{100} s_{a1000} \end{bmatrix}$$

$$Q = \vec{a} \cdot \vec{s}_a^T + \vec{b} \cdot \vec{s}_b^T + \vec{c} \cdot \vec{s}_c^T + \vec{d} \cdot \vec{s}_d^T$$

$$Q = U \Sigma V^T = \sum_{i=1}^r \sigma_i \vec{u}_i \vec{v}_i^T$$

r : rank of Q

The (k) principal components of Q

- along the columns: $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$
- along the rows: $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_k$

($m=100, n=1000$)

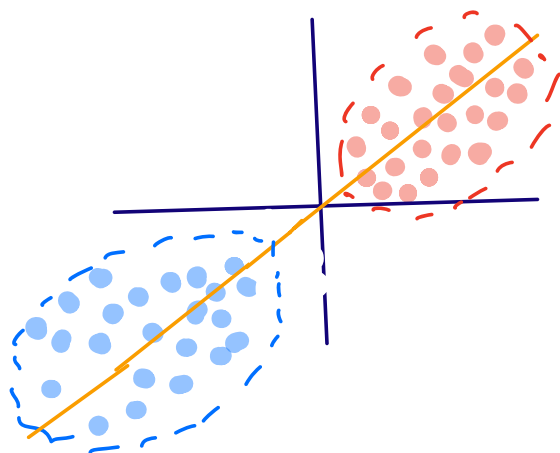
$$Q = \begin{matrix} & \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_n \\ \begin{matrix} v_1 \\ v_2 \\ \vdots \\ v_{100} \end{matrix} & \begin{bmatrix} | \\ | \\ | \\ | \end{bmatrix} & \begin{bmatrix} | \\ | \\ | \\ | \end{bmatrix} & \dots & \begin{bmatrix} | \\ | \\ | \\ | \end{bmatrix} \end{matrix} ; Q = \begin{matrix} & \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_n \\ \begin{matrix} v_1 \\ v_2 \\ \vdots \\ v_{100} \end{matrix} & \begin{bmatrix} \vec{m}_1 \\ \vec{m}_2 \\ \vdots \\ \vec{m}_m \end{bmatrix} \end{matrix}$$

- We assume data is organized by column:
 - i.e. each data point is a 100-dim. vector \vec{q}_j denoting U_j 's ratings for the 100 videos v_1, v_2, \dots, v_{100}

GOAL: Find the first principal component vector (direction) that is most informative about the data set.

The hope is that if we pick the right PC direction, the high-dim. data set will be projected to a much (low-dim.) representation that unveils distinct "clusters" in the dataset!

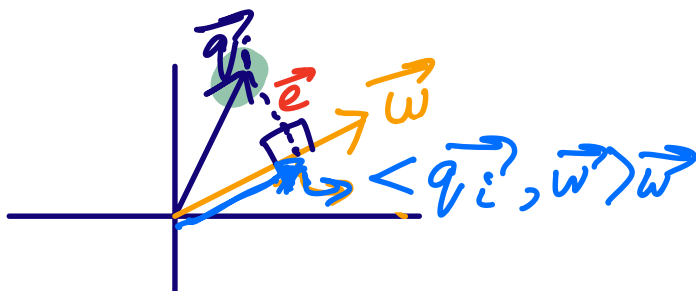
look in this direction



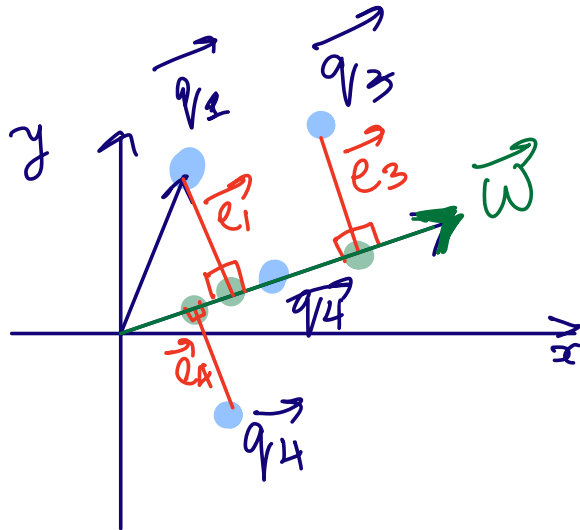
In typical apps, "clusters" will emerge from a low-dim. representation of the data using PCA.

→ We want a \vec{w} such that when q_i is projected onto \vec{w} , the error is minimized.

Let us normalize: $\|\vec{w}\| = 1$



Projection of \vec{q}_i onto \vec{w} : $\langle \vec{q}_i, \vec{w} \rangle \vec{w}$
 (since $\|\vec{w}\|=1$)



Squared-norm of error due to i^{th} data point } : $\|\vec{q}_i - \underbrace{\langle \vec{q}_i, \vec{w} \rangle \vec{w}}_{\vec{e}_i}\|^2$

Summing error over all points:

$$\min_{\vec{w}} \sum_{i=1}^n \|\vec{q}_i - \underbrace{\langle \vec{q}_i, \vec{w} \rangle \vec{w}}_{\vec{e}_i}\|^2$$

($\|\vec{w}\|=1$)

→ Find the \vec{w}^* that achieves the minimum error.

can also write as:

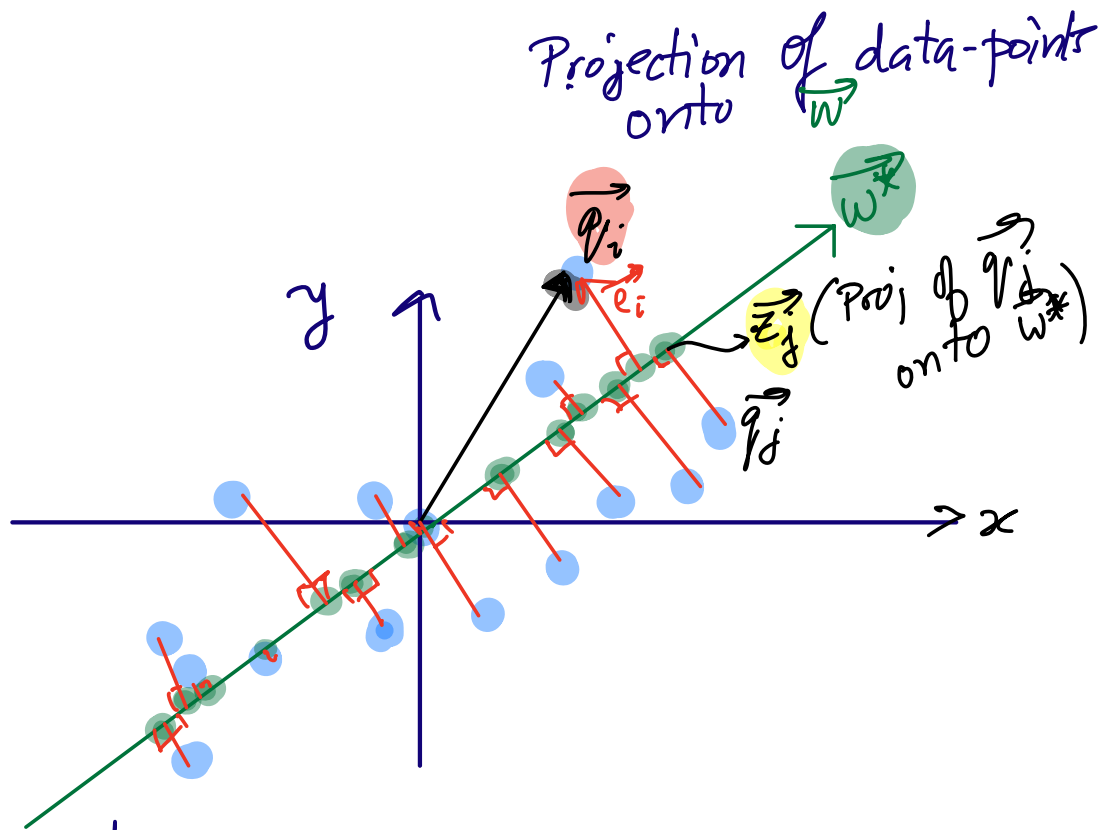
$$\vec{w}^* = \underset{\vec{w}}{\operatorname{argmin}} \sum_{i=1}^n \underbrace{\|\vec{q}_i - \langle \vec{q}_i, \vec{w} \rangle \vec{w}\|}_{\vec{e}_i}^2$$

s.t. $\|\vec{w}\| = 1$

$$\begin{aligned} \|\vec{e}_i\|^2 &= (\vec{q}_i - \langle \vec{q}_i, \vec{w} \rangle \vec{w})^T (\vec{q}_i - \langle \vec{q}_i, \vec{w} \rangle \vec{w}) \\ &= \vec{q}_i^T \vec{q}_i - \underbrace{\langle \vec{q}_i, \vec{w} \rangle \vec{w}^T \vec{q}_i} - \underbrace{\langle \vec{q}_i, \vec{w} \rangle \vec{q}_i^T \vec{w}} + \underbrace{\langle \vec{q}_i, \vec{w} \rangle^2 \vec{w}^T \vec{w}}_{\|\vec{w}\|^2=1} \\ &= \|\vec{q}_i\|^2 - 2 \langle \vec{q}_i, \vec{w} \rangle^2 + \langle \vec{q}_i, \vec{w} \rangle^2 \\ &= \|\vec{q}_i\|^2 - \langle \vec{q}_i, \vec{w} \rangle^2 \end{aligned}$$

$$\|\vec{e}_i\|^2 = \|\vec{q}_i\|^2 - \langle \vec{q}_i, \vec{w} \rangle^2$$

$$\vec{w}^* = \underset{\vec{w}}{\operatorname{argmin}} \left(\sum_{i=1}^n \underbrace{\|\vec{q}_i\|^2}_{\text{does not depend on } \vec{w}} - \sum_{i=1}^n \langle \vec{q}_i, \vec{w} \rangle^2 \right)$$



2 perspectives:

- \vec{w}^* minimizes the sum of projection errors of the data-points

$$\vec{w}^* = \underset{\vec{w}}{\operatorname{argmax}} \sum_{i=1}^n \langle \vec{q}_i, \vec{w} \rangle^2$$

($\|\vec{w}\|=1$)

- \vec{w}^* maximizes the "spread" of the data as measured by the sum of the squared magnitudes of the projection of the data pts

(Assumes the data-pts. are "mean-removed")

(a.k.a. the direction of maximum variance)

($\sum_{j=1}^n \|z_j\|^2$ is maximum)

↪ a probability concept
(out of scope for 16B)

$$\vec{w}^* = \underset{\substack{\vec{w} \\ (\text{s.t. } \|\vec{w}\|=1)}}{\text{argmax}} \sum_{i=1}^n \underbrace{\langle \vec{w}_i, \vec{q} \rangle \langle \vec{q}, \vec{w}_i \rangle}_{\langle \vec{w}_i, \vec{q} \rangle^2}$$

$$= \underset{\vec{w}}{\text{argmax}} \sum_{i=1}^n \vec{w}^T \vec{q}_i \vec{q}_i^T \vec{w}$$

outer product

$$\vec{w}^* = \underset{\vec{w}}{\text{argmax}} \left[\vec{w}^T \underbrace{Q Q^T}_S \vec{w} \right]$$

$$S = U \Lambda U^T$$

Recall:

$$\sum \vec{u}_i \vec{v}_i^T = U V^T$$

where $\Lambda = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_m^2)$
 and U is an $(m \times m)$ orthonormal matrix
 that is the same U in $Q = U \Sigma V^T$

$$\vec{w}^* = \underset{\substack{\vec{w} \\ \|\vec{w}\|=1}}{\text{argmax}} \left[\vec{w}^T S \vec{w} \right]$$

S has e-values σ_i^2
 and e-vectors \vec{u}_i :
 $S \vec{u}_i = \sigma_i^2 \vec{u}_i$

can be shown that $\vec{w}^* = \underset{\|\vec{w}\|=1}{\text{argmax}} \vec{w}^T S \vec{w}$ is $\vec{w} = \vec{u}_1$,
 where \vec{u}_1 is the maximum eigenvector of $S = Q Q^T$
 corresponding to $\lambda_1 = \sigma_1^2$, where the e-values of S
 are $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_m^2$.

"Proof":

Write $\vec{w} \in \mathbb{R}^m$ in eigenbasis $\{\vec{u}_i\}$ co-ordinate system:
 $\vec{w} = \alpha_1 \vec{u}_1 + \alpha_2 \vec{u}_2 + \dots + \alpha_m \vec{u}_m = \sum_{i=1}^m \alpha_i \vec{u}_i$

$$1 = \|\vec{w}\|^2 = \vec{w}^T \vec{w} = \left(\sum_{i=1}^m \alpha_i \vec{u}_i^T \right) \left(\sum_{j=1}^m \alpha_j \vec{u}_j \right)$$

$$= \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j \underbrace{\vec{u}_i^T \vec{u}_j}_{= \begin{cases} 1 & \text{if } i=j \\ 0 & \text{else} \end{cases}}$$

$$\Rightarrow \boxed{\sum_{i=1}^m \alpha_i^2 = 1} \Rightarrow \underline{\alpha_1^2 + \alpha_2^2 + \dots + \alpha_m^2 = 1}$$

$$\begin{aligned} S\vec{u}_1 &= \lambda_1 \vec{u}_1 \\ S\vec{u}_2 &= \lambda_2 \vec{u}_2 \\ &\vdots \\ S\vec{u}_m &= \lambda_m \vec{u}_m \end{aligned}$$

eigenvalue equations for S
($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$)

So, we want to find the \vec{w} that maximizes $\vec{w}^T S \vec{w}$ s.t. $\vec{w} = \sum_{i=1}^m \alpha_i \vec{u}_i$

where $\boxed{\alpha_1^2 + \alpha_2^2 + \dots + \alpha_m^2 = 1}$

- If $\vec{w} = \vec{u}_1$, $\alpha_1 = 1, \alpha_2 = \alpha_3 = \dots = \alpha_m = 0$,
 $\vec{w}^T S \vec{w} = \vec{u}_1^T S \vec{u}_1 = \vec{u}_1^T \lambda_1 \vec{u}_1 = \lambda_1 \underbrace{\|\vec{u}_1\|^2}_1$
- If $\vec{w} = \vec{u}_2$, $\alpha_1 = 0, \alpha_2 = 1, \alpha_3 = \dots = \alpha_m = 0$
 $\vec{w}^T S \vec{w} = \vec{u}_2^T S \vec{u}_2 = \lambda_2$

Similarly, If $\vec{w} = \vec{u}_i$, $\vec{w}^T S \vec{w} = \vec{u}_i^T S \vec{u}_i = \lambda_i$

We need to pick a weighted-combination of the \vec{u}_i 's (since $\vec{w} = \sum_{i=1}^m \alpha_i \vec{u}_i$) where $\alpha_1^2 + \alpha_2^2 + \dots + \alpha_m^2 = 1$.

Q) What is the optimal choice for the α 's?

A) $\vec{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m) = (1, 0, 0, \dots, 0)$

$\Rightarrow \vec{w}^* = \underset{\|\vec{w}\|=1}{\operatorname{argmin}} \vec{w}^T Q Q^T \vec{w} = \vec{u}_1$

□

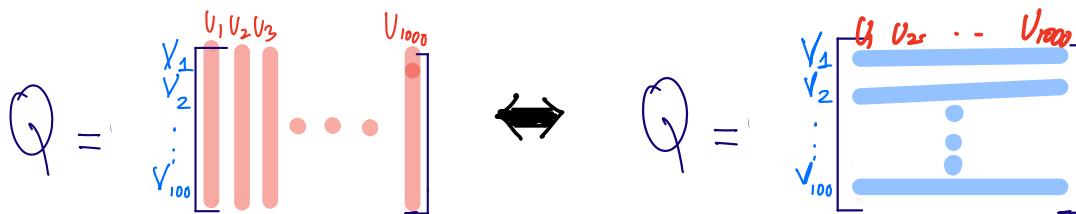
$$\underset{\vec{w}}{\text{argmin}} \sum_{i=1}^n \|\vec{q}_i - \langle \vec{q}_i, \vec{w} \rangle \vec{w}\|^2 = \vec{u}_1$$

Similarly, 2nd most important component? \vec{u}_2

More generally, if we want the l -dimensional low-rank approx. to Q , the solution will turn out to be $Q_l = \sum_{i=1}^l \sigma_i \vec{u}_i \vec{v}_i^T$

(Eckart-Young theorem)

If the data points of interest are organized by rows instead of by columns



$$\underset{\vec{w}}{\text{argmin}} \sum_{i=1}^n \|\vec{v}_i - \langle \vec{v}_i, \vec{w} \rangle \vec{w}\|^2 = \vec{v}_1$$

first right singular vector

- \vec{u}_i 's are principal components along columns (u_i 's)
- \vec{v}_i^T 's " " " " rows (v_i 's)

SUMMARY: Given data pts. $\{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$,
find the principal components of this data.

1) Arrange the data into a matrix X :

$$X = \begin{bmatrix} | & | & \dots & | \\ \vec{x}_1 & \vec{x}_2 & \dots & \vec{x}_n \\ | & | & \dots & | \end{bmatrix}$$

→ "Mean-removal": depending on application,
"center" each column around 0 by
subtracting the column means from each
column

2) Compute $X = U \Sigma V^T = \sum_{i=1}^n \sigma_i \vec{u}_i \vec{v}_i^T$

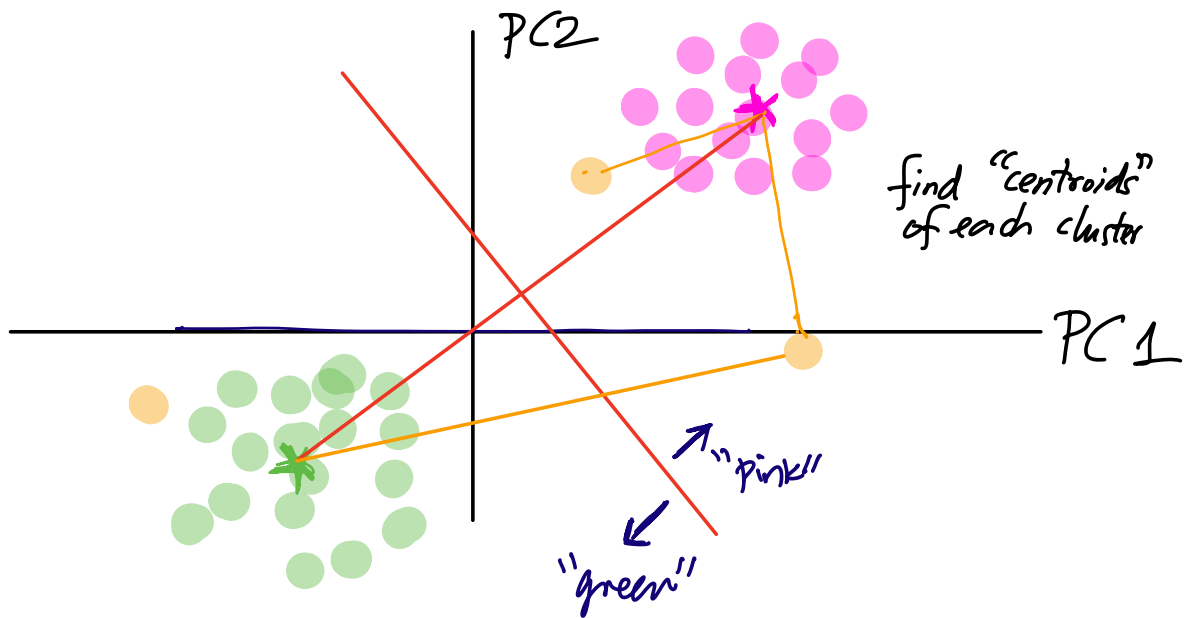
3) For k principal components along the columns, choose $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$.

For k principal components along the rows, choose $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_k$.

4) (OPTIONAL) Project data onto $\vec{u}_1, \vec{u}_2, \dots, \vec{u}_k$
to get low-dimensional structure
Use this to perform clustering/classification

Classification:

After projecting data onto top 2 Principal Components (200-dim. dataset)



Linearization for control

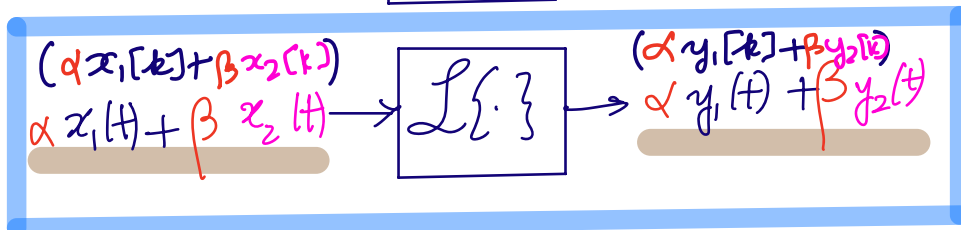
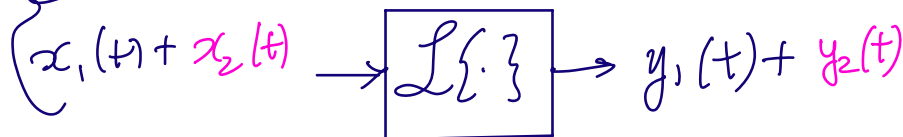
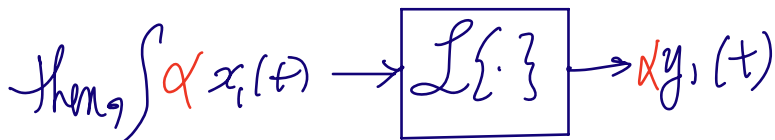
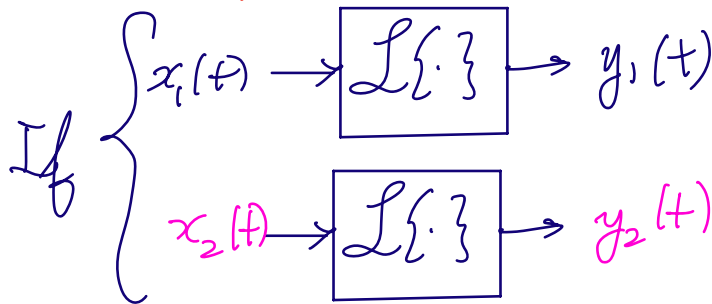
So far, in this course, we have focused on LINEAR control systems, where in discrete-time,

$$\vec{x}[k+1] = A \vec{x}[k] + B \vec{u}[k]$$

or in continuous-time,

$$\frac{d\vec{x}(t)}{dt} = A \vec{x}(t) + B \vec{u}(t)$$

RECAP First, what is a linear system?
(Holds for CT and DT systems)



Ex.: (i) $\vec{y} = A\vec{x}$ (matrix multiplication)
is a linear system.

Why? If $L\{\cdot\} = A$, and if $A\vec{x}_1 = \vec{y}_1$; $A\vec{x}_2 = \vec{y}_2$

then, $A[\alpha\vec{x}_1 + \beta\vec{x}_2] = \alpha A\vec{x}_1 + \beta A\vec{x}_2$. ✓

(ii) $y(t) = \frac{dx(t)}{dt}$ (differential equation)

$y_1(t) = \frac{dx_1(t)}{dt}$; $y_2(t) = \frac{dx_2(t)}{dt}$

$\frac{d}{dt} [\alpha x_1(t) + \beta x_2(t)] = \alpha \frac{dx_1(t)}{dt} + \beta \frac{dx_2(t)}{dt}$

$= \alpha y_1(t) + \beta y_2(t)$ ✓

However, in the real world, many systems exhibit NON-LINEAR behavior!

Ex.: — Transistors: not just "on-off" switches but work "continuously" with highly non-linear governing equations.

— Robotics & control: can have highly non-linear dynamics; e.g. applying feedback control in order to stabilize an inverted pendulum on a cart.

— Machine learning: can have highly non-linear systems. E.g. gradient descent can result in non-linear trajectories.

Generally, a nonlinear control system is characterized by:

Discrete-time: $\vec{x}[k+1] = \vec{f}(\vec{x}[k], \vec{u}[k])$

state vector arbitrary vector function input control vector

Continuous-time: $\frac{d\vec{x}(t)}{dt} = \vec{f}(\vec{x}(t), \vec{u}(t))$

Key concept: How to deal w/ NL systems by approx. them locally using linear models by a process called **LINEARIZATION** of the function f .

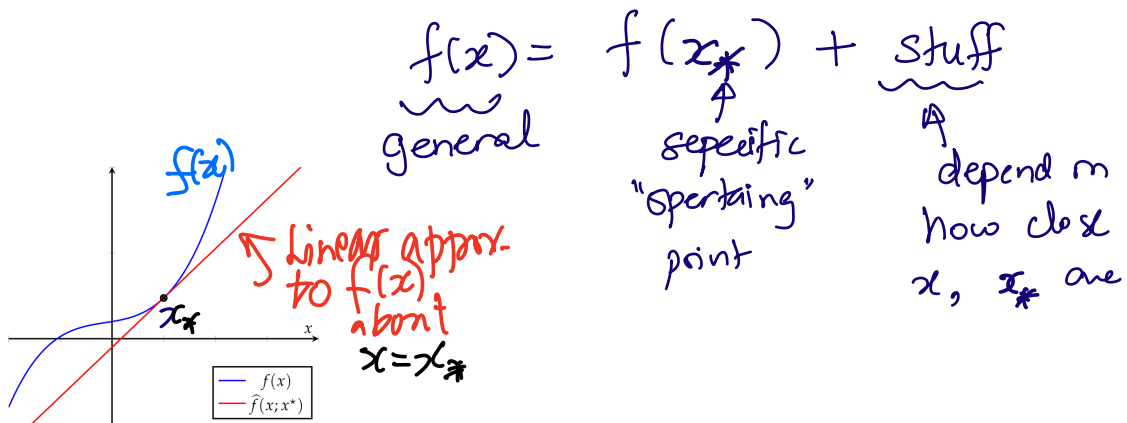
Ex:

Linearization.

$$\left. \begin{aligned} f(x) &= x^2 \\ f(x, u) &= x^2 u^3 \end{aligned} \right\}$$

$$\frac{d\vec{x}(t)}{dt} = f(\vec{x}, \vec{u}) = \underbrace{A\vec{x} + B\vec{u}}_{\text{Linear formulation}}$$

Scalar functions (One variable)



How does $f(x)$ change in the neighbourhood of x^* ?

Derivative of $f(x)$ @ x_* = $\left. \frac{df}{dx} \right|_{x=x_*}$

Taylor Expansion

$$f(x) = \underbrace{f(x_*)}_{\text{constant}} + \underbrace{\left. \frac{df}{dx} \right|_{x=x_*}}_{\text{linear}} (x-x_*) + \left. \frac{d^2f}{dx^2} \right|_{x=x_*} \frac{(x-x_*)^2}{2} \left. \vphantom{\frac{d^2f}{dx^2}} \right\} \text{quadratic} + \left. \frac{d^3f}{dx^3} \right|_{x=x_*} \frac{(x-x_*)^3}{6} \left. \vphantom{\frac{d^3f}{dx^3}} \right\} \text{cubic} + \dots$$

SCALAR FUNCTION CASE:

Linear approximation

$$f(x) \approx f(x_*) + f'(x_*) (x-x_*)$$

↑
(approx. equal)
