

model with math

# Mathematical Modeling

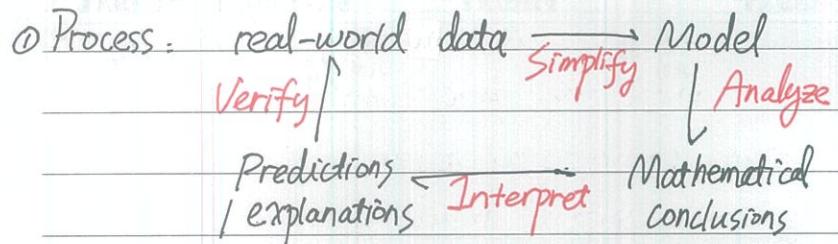
## MAT 3300 Notebook

Youthy WANG

Date: 2021 Sep 9.

## MAT 3300 Mathematical Modeling

## ● Introduction



## ② Classification

	Empirical	Mechanistic
Deterministic	regression relationship	(differential equations)
Stochastic	variance	(probabilistic equations)

## ③ Construction (process)

- (i) Identify the problem
- (ii) Make assumptions { identify & classify the variables  
determine interrelationships between variables
- (iii) Solve the model
- (iv) Verify / Interpret the model { Address problem?  
Make common sense?  
Correspond with real-world data?
- (v) Implement the model
- (vi) Maintain the model (some change)

Date: 2021 Sep 10

- ④ A modeling example (classic).
- ★ "The Secretary Problem" (Socrates' opinion of love)
- Q: a single position /  $n$  applicants, the value is known / can be ranked once seen / interview in random order / either accept or reject after each interview / decision only based on previous known rank!

- Aim: highest probability to select the best one
- Strategy chosen (choose a path & go farthest)
  - ↳ reject first  $r-1$  (as ranking sample), accept  $r^{th}$  better than  $r-1$  one. ⇒ choose a good  $r$  for highest  $P(r)$

$$\begin{aligned}
 P(r) &= \sum_{i=r}^n P(\text{applicant } i \text{ is selected and the best}) \\
 &= \sum_{i=r}^n P(i \text{ is selected} \mid i \text{ is the best}) P(i \text{ is the best}) \\
 &= \sum_{i=r}^n P(\text{the best in first } i-1 \text{ is in } r-1 \mid i \text{ is the best}) P(i \text{ is the best}) \\
 &= \sum_{i=r}^n \frac{r-1}{i-1} \times \frac{1}{n} = \sum_{i=r}^n \frac{\frac{r-1}{n} - \frac{1}{n}}{\frac{i-1}{n} - \frac{1}{n}} \times \frac{1}{n} \\
 \text{let } x = \frac{i}{n} &\quad \text{when } n \text{ large} \quad x \int_0^1 \frac{1}{t} dt \quad (\text{give up } \frac{1}{n} \text{ in numerator}) = -x \ln x \\
 &\quad P'(x) = 0 \Rightarrow x = \frac{1}{e} \approx 36.8\%. \quad P\left(\frac{1}{e}\right) = \frac{1}{e} \approx 36.8\%
 \end{aligned}$$

Date: 2021 Sep 10

## PART I. Differential Equations (ODE)

- ① Difference Equation: relates change ( $n^{\text{th}}$  order difference) to past values in a (discrete) sequence.

$$\Delta a_n := a_{n+1} - a_n = f(a_0, a_1, \dots, a_n)$$

another form (dynamic system)  $a_{n+1} = g(a_0, a_1, \dots, a_n)$

- ② Differential Equations: relates continuous change (approximation to difference equation)

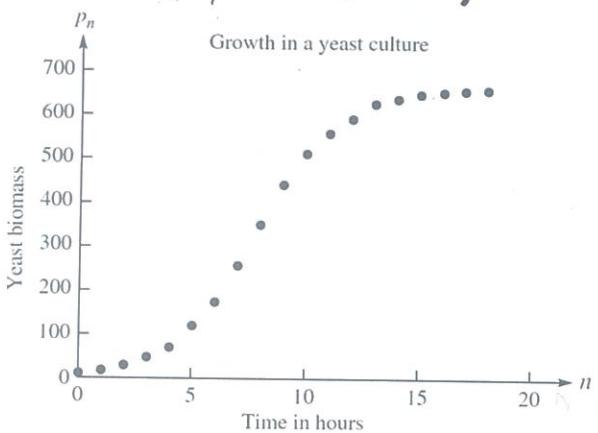
$$a'(t) = f(a(t)) \quad (\text{ODE})$$

Remark:  $f$  can be got from either theoretical derivation or model fitting.

e.g.: (i) Growth of a Yeast Cultivation (Yeast's "S" shape growth)

Analyze given data (by experiment)  
 ⇒ propose a difference equation

Time in hours $n$	Yeast biomass $p_n$	Change/ hour $p_{n+1} - p_n$
0	9.6	8.7
1	18.3	10.7
2	29.0	18.2
3	47.2	23.9
4	71.1	48.0
5	119.1	55.5
6	174.6	82.7
7	257.3	93.4
8	350.7	90.3
9	441.0	72.3
10	513.3	46.4
11	559.7	35.1
12	594.8	34.6
13	629.4	11.4
14	640.8	10.3
15	651.1	4.8
16	655.9	3.7
17	659.6	2.2
18	661.8	



Date: 2021 Sep 12

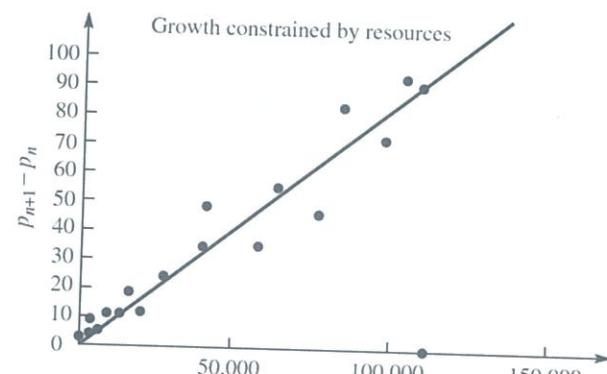
Carrying capacity = 665 due to resource limit

Model:  $\Delta p_n := p_{n+1} - p_n = k p_n (665 - p_n)$

difference equation → estimate  $k$  using a linear regression (with  $\Delta p_n$ ;  $p_n(665 - p_n)$ )

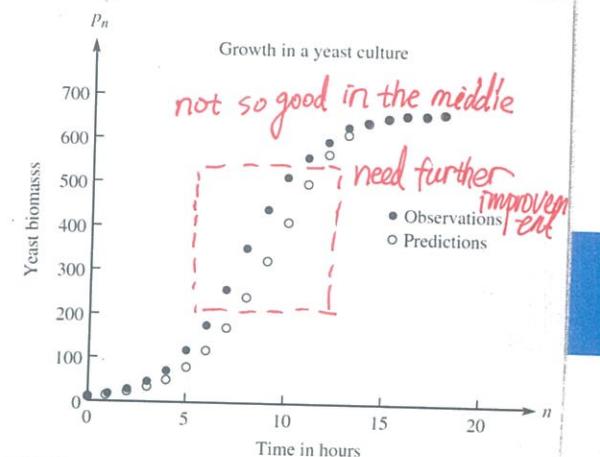
get  $k \approx 0.00082$  with no analytical sol.

$p_{n+1} - p_n$	$p_n (665 - p_n)$
8.7	6291.84
10.7	11834.61
18.2	18444.00
23.9	29160.16
48.0	42226.29
55.5	65016.69
82.7	85623.84
93.4	104901.21
90.3	110225.01
72.3	98784.00
46.4	77867.61
35.1	58936.41
34.6	41754.96
11.4	22406.64
10.3	15507.36
4.8	9050.29
3.7	5968.69
2.2	3561.84



Test the model - prediction with real data

Time in hours	Observation	Prediction
0	9.6	9.6
1	18.3	14.8
2	29.0	22.6
3	47.2	34.5
4	71.1	52.4
5	119.1	78.7
6	174.6	116.6
7	257.3	169.0
8	350.7	237.8
9	441.0	321.1
10	513.3	411.6
11	559.7	497.1
12	594.8	565.6
13	629.4	611.7
14	640.8	638.4
15	651.1	652.3
16	655.9	659.1
17	659.6	662.3
18	661.8	663.8



Date: 2021 Sep 12

Modification. (estimated by differential equation)

\* **Logistic Model**,  $\frac{dp(t)}{dt} = rP(t)(M-P(t))$ ,  $M$  - maximum possible population

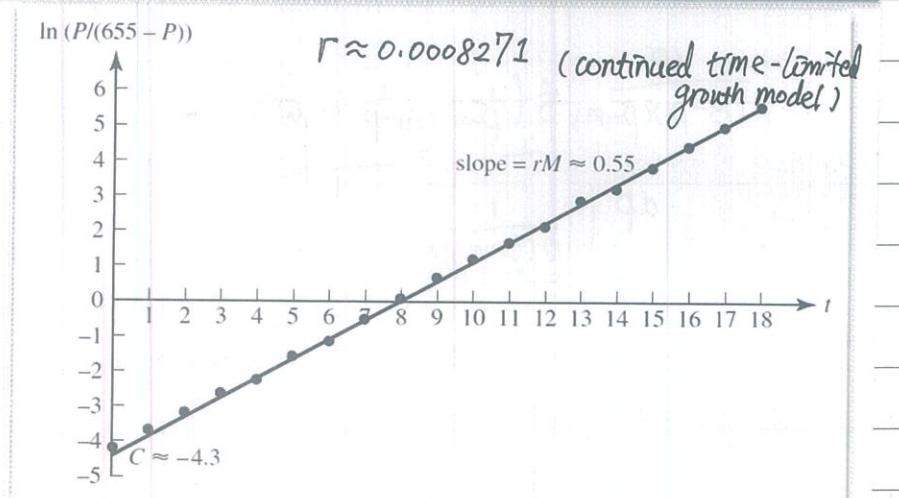
analysis:

$$\frac{dp(t)}{P(t)(M-P(t))} = r dt \Rightarrow \frac{1}{M} \left( \frac{dP(t)}{P(t)} - \frac{d(M-P(t))}{M-P(t)} \right) = r dt$$

$$\Rightarrow \ln \frac{P(t)}{P(0)} - \ln \frac{M-P(t)}{M-P(0)} = rMt$$

$$\Rightarrow P(t) = \frac{M P(0)}{P(0) + (M-P(0)) e^{-rMt}} \quad t \rightarrow \infty, P(t) \rightarrow M$$

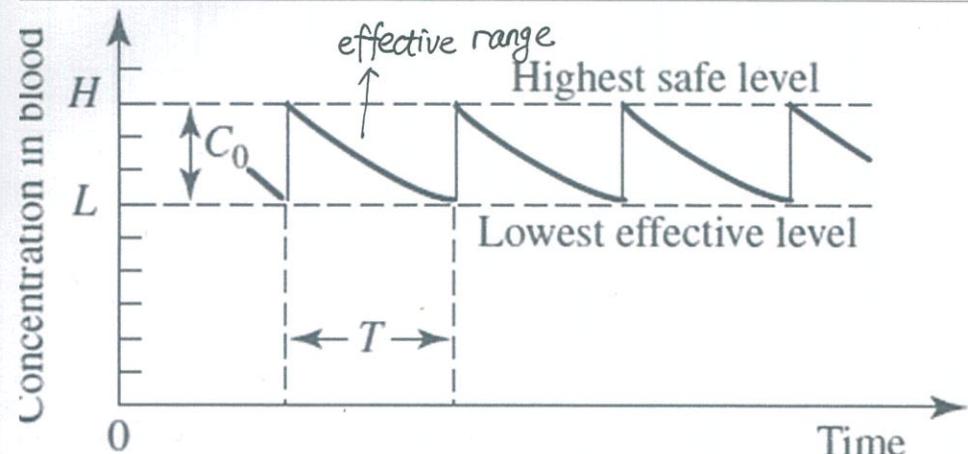
getting "r": linear regression  $\ln \frac{P(t)}{M-P(t)} = \ln \frac{P(0)}{M-P(0)} + rMt$



## (ii) Prescribing Drug Dosage:

How can doses & time between doses can be adjusted to maintain a safe but effective concentration?

Date: 2021 Sep 12



$$C(t) = f(\text{decay rate, assimilation rate, dosage amount, time})$$

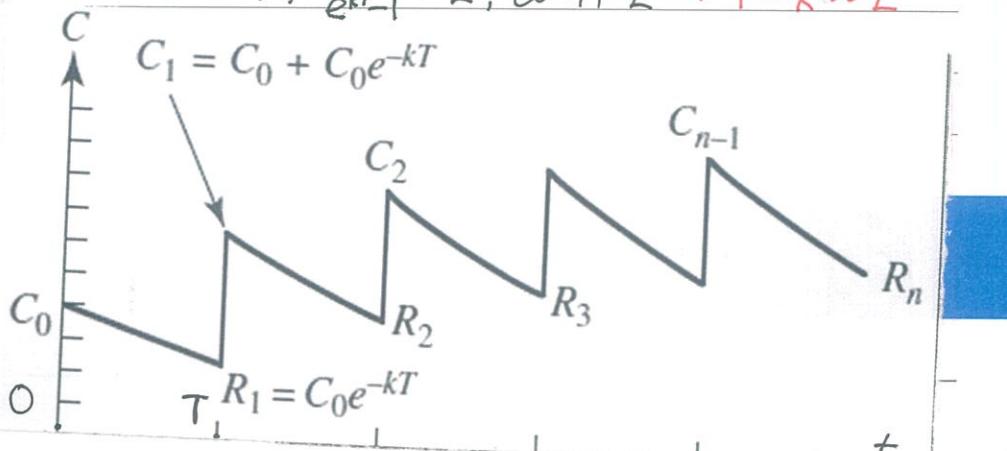
simplify assume that body weight & blood volume are constants

$$\Rightarrow \frac{dC(t)}{dt} = -kC(t) \quad (\text{just about total concentrate})$$

$$C(t) = C_0 \cdot e^{-kt} \quad \xrightarrow{\text{after interval } T} C_0 e^{-kT} \xrightarrow{\text{base}} C_0 e^{-kT} + C_0 \dots$$

$$\text{After } n^{\text{th}} \text{ dosage: } R_n = C_0 + C_0 e^{-kT} + \dots + C_0 e^{-n kT} \xrightarrow{n \rightarrow \infty} \frac{C_0}{e^{kT-1}}$$

$$\Rightarrow \frac{C_0}{e^{kT-1}} = L, C_0 = H - L \Rightarrow T = \frac{k}{L} \ln \frac{H}{L}$$



Date: 2021 Sep 16

### ③ Solving some ODEs (methods)

- { separation of variables (calculus I)
- first-order linear equation
- second-order (or higher) linear equation with constant coefficients (Linear Algebra)

Others with analytical solution = website "Eqworld"

What if not analytically solvable? { Graphically solve  
Numerically solve

### ④ Graphical Solutions

i) phase lines ( $\frac{dy}{dx} = g(x, y)$ , each tangent line at  $(x_0, y_0)$ )

ii) Autonomous Differential Equations,  $\frac{dy}{dx} = g(y)$   
(in calculus I) (equilibrium values / rest points + derivative analyze)

iii) System of ODEs

Some concepts: autonomous system of ODEs:  $\frac{dx}{dt} = f(x, y); \frac{dy}{dt} = g(x, y)$

Solutions = pairs of functions  $(x(t), y(t))$

phase plane = the  $x-y$  plane in this case satisfies "

trajectory / path / orbit; a curve consists of pts  $(x(t), y(t)) \forall t \geq 0$

rest point / equilibrium point

A point s.t.  $\frac{dx}{dt} = 0 = \frac{dy}{dt}$

Date: 2021 Sep 16

**Stable** = to describe the rest pts. if any trajectory starting close to it stays close to it for all future time  $t$

**Asymptotically stable** = ~ (same as stable) if any trajectory starting close to it approaches to it when  $t \rightarrow \infty$

**Unstable** (points not stable)

e.g.: (i) a linear autonomous system ( $\frac{dx}{dt} = -x+y; \frac{dy}{dt} = -x-y$ )

(ii) a non-linear autonomous system ( $\frac{dx}{dt} = xy; \frac{dy}{dt} = x^2$ )  
 $\hookrightarrow d(y^2 - x^2) = 2y dy - 2x dx = 0$ .

Properties of trajectories (omitted)

egs (i) **Competitive Hunter Model**

Imagine a small pond, mature enough to support wildlife.

Stock the pond with game fish, say trout & bass.

Let  $x(t)$  denote the population of the trout at time  $t$ , where  $y(t)$  denote the bass.

Q: How sensitive is the final solution of population levels to the initial stockage levels & external perturbations?

Modeling: Assume 2 species depend on an unlimited common food supply.

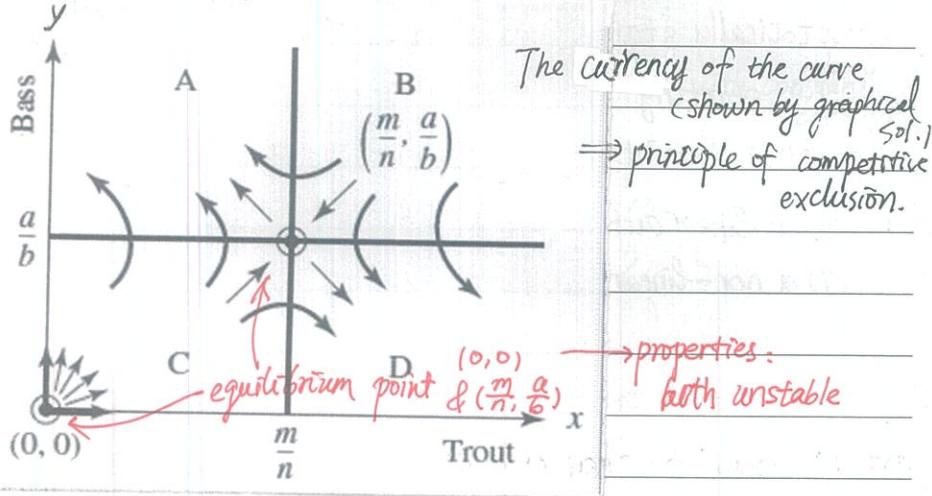
$$\frac{dx}{dt} = (a - by)x; \quad \frac{dy}{dt} = (m - nx)y$$

↑ compete: each species affects another



Date: 2021 Sep 21

**Principle of competitive exclusion:** mutual coexistence of the species is highly improbable.



Modification:  $\begin{cases} \frac{dx}{dt} = (a - ex - by)x \\ \frac{dy}{dt} = (m - ny - nx)y \end{cases}$

more realistic  
(if  $x \rightarrow \infty$ ,  $\frac{dx}{dt} \rightarrow 0$ )  
↑ environment limits  
it depends whether stable or not

equilibrium points:  
 $(x, y) = (0, 0); (0, \frac{m}{n}); (\frac{a}{e}, 0); (\frac{mb - an}{nb - en}, \frac{an - me}{nb - en})$

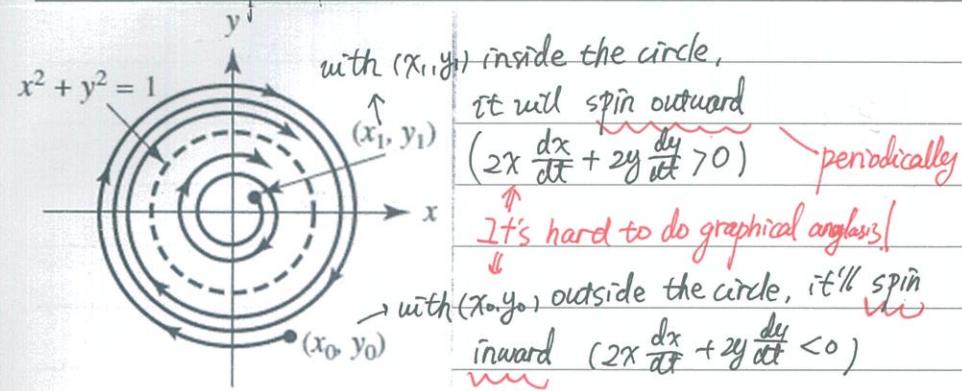
all unstable!

### ★ limitations of Graphical Analysis.

Consider  $\begin{cases} \frac{dx}{dt} = y + x - x(x^2 + y^2) \\ \frac{dy}{dt} = -x + y - y(x^2 + y^2) \end{cases}$  (the only equilibrium pt. (0, 0))

Analyze: when  $x^2 + y^2 = 1$ ,  $2x \frac{dx}{dt} + 2y \frac{dy}{dt} = 0$ .  
limit circle.  $2(x^2 + y^2)(1 - x^2 - y^2)$

Date: 2021 Sep 27



### egs(ii) Predator-Prey Model

Whale and krill live in the ocean, with former the predator, latter the prey.

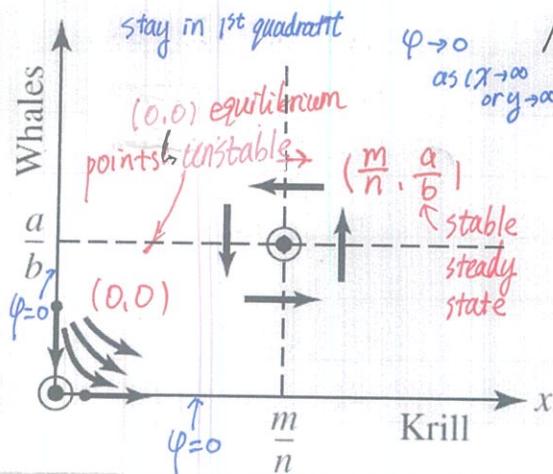
Q: Does the cycle of whales and krill population increasing or decreasing continue indefinitely or one of them die out?

Modeling: assume the ocean can support unlimited number of krill, so  $\dot{x} = ax$ , with  $a > 0$  without whales.

Similarly, without krill, whale decreases,  $\dot{y} = -my$  with  $m > 0$ .

By considering their interactions, we get

$\frac{dx}{dt} = ax - bxy = x(a - by)$  decreases in y  
 $\frac{dy}{dt} = -my + nxy = y(-m + nx)$  increases in x  
(correspondence with reality)



Date: 2021 Oct 1

Analysis:  $\frac{dy}{dx} = \frac{(-m+nx)y}{(a-by)x}$

$$\Rightarrow \frac{y^a}{e^{by}} = K e^{nx}, \text{ with } K > 0.$$

Let  $f(y) = \frac{y^a}{e^{by}}$ ,  $g(x) = \frac{x^m}{e^{nx}}$

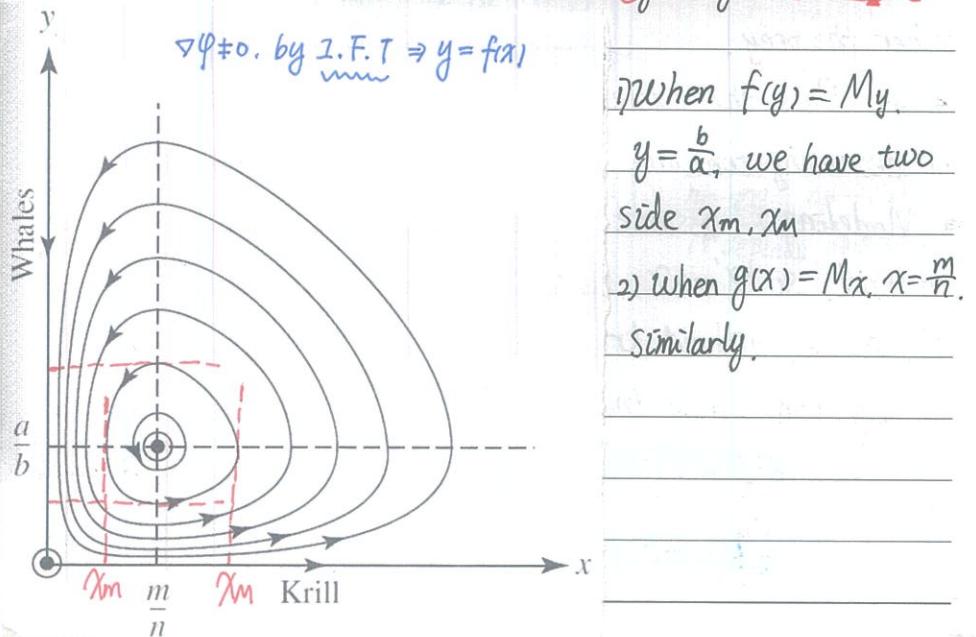
Equation  $f(y) = g(x) \Leftrightarrow f(x,y) = k$

By graphs of  $f$  &  $g$ , we know they're bdd by  $M_x, M_y$  respectively.

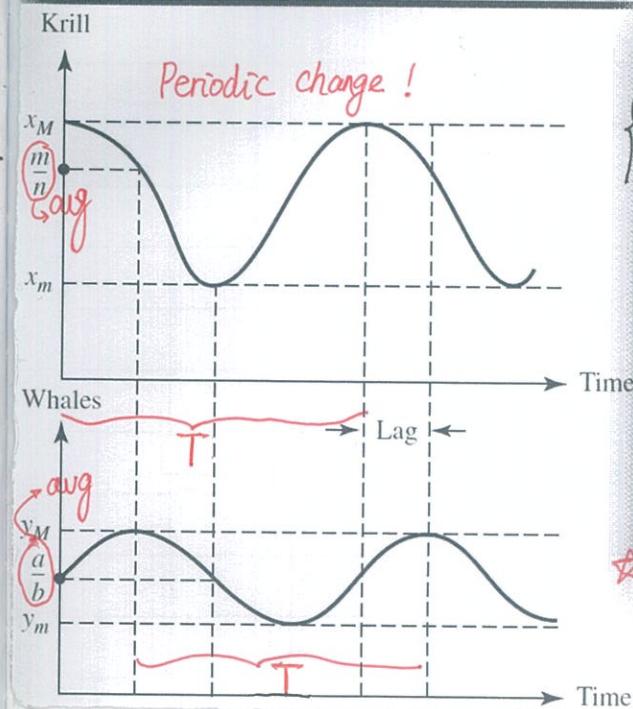
$\therefore K \leq M_x M_y$ . If  $K = s M_y$  with  $s < M_x$ ,

The equation  $g(x) = s$  has exactly two solutions when  $f(y) = M_y$ .

↳ indicates that the trajectory is a "circle".



Date: 2021 Oct 1



Calculate Average:

$$\bar{x} = \frac{1}{T} \int_0^T x(t) dt$$

$$\bar{y} = \frac{1}{T} \int_0^T y(t) dt$$

We know that  $\frac{1}{T} \int_0^T \frac{d}{dt} x(t) dt = \frac{1}{T} \int_0^T (a-by) dt$  ( $\frac{d}{dt} \frac{x(t)}{t} = a-by$ )

$$LHS = \ln \frac{x(T)}{x(0)} = 0$$

$$RHS = aT - bTy$$

★  $\left\{ \begin{array}{l} \bar{y} = \frac{a}{b} \\ \text{Similarly, } \bar{x} = \frac{m}{n} \end{array} \right.$

► Modification: (the effect of harvesting)

Harvest krill:  $\begin{cases} \frac{dx}{dt} = (a-by)x - rx = ((a-r)-by)x \\ \frac{dy}{dt} = (-m+nx)y - ry = (-m+(r-n)x)y \end{cases}$

make both animals decrease in rate

$\therefore$  The average change to  $\bar{y} = \frac{a-r}{b}$ ,  $\bar{x} = \frac{m+r}{n}$  ( $r < a$ )

Moderate amount of harvesting  $\leftarrow$  Volterra's Principle  
preys actually increases the average

Level of preys, while decreasing the avg level of predators.  
Fishing gives the opposite results.

Date: 2021 Oct 1

## egs (iii) Two military.

Situation of Combat between two homogeneous forces.

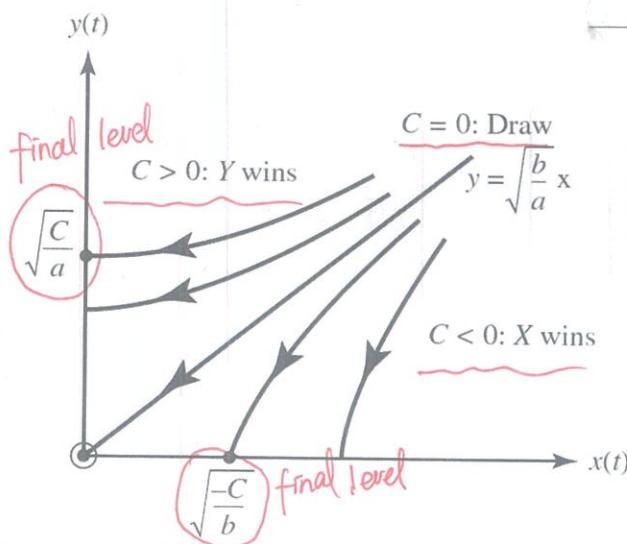
- Q: Will one force win / end in draw? How do force levels decrease? How many survivors? How long will it last? How does initial force level affect the outcome?

- Modeling: Decreasing rates of  $X$  /  $Y$  proportional to strength of their opponents.
- $$\begin{cases} \frac{dx}{dt} = -ay, \quad a > 0 \\ \frac{dy}{dt} = -bx, \quad b > 0 \end{cases}$$

Analysis:  $-\frac{dx}{dt} + bxy = 0$

$$\therefore -bx^2 + ay^2 = C \text{ for some } C \in \mathbb{R}$$

Lanchester square law model.



Two cases:

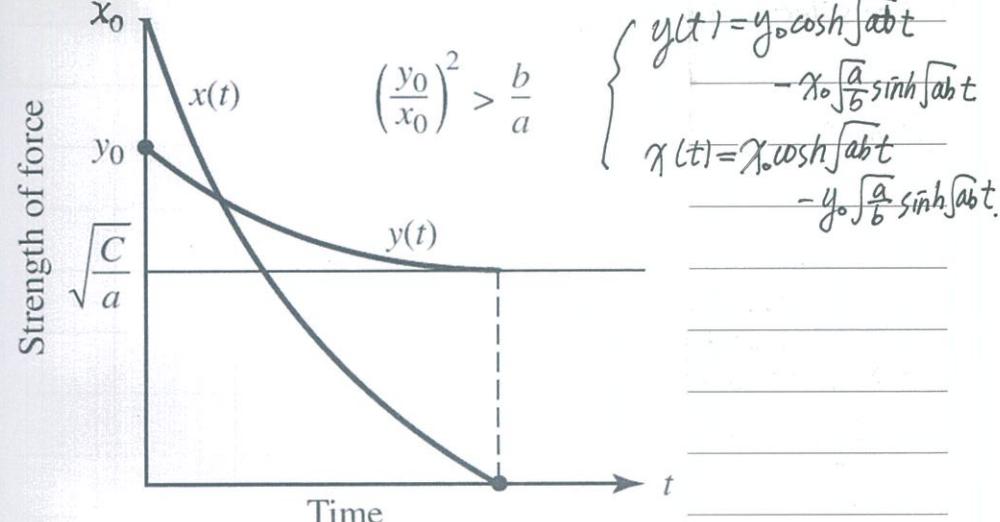
$$-bx_0^2 + ay_0^2 > 0, \quad Y \text{ wins.}$$

$$-bx_0^2 + ay_0^2 < 0, \quad X \text{ wins.}$$

Date: 2021. Oct. 1

Solve for equations (consider the case  $Y$  wins)

$$y'' = \frac{d^2y}{dt^2} = aby \quad \text{Solving it. } (x'' = abx, \text{ characteristic equation})$$



- Modification: Let the Lanchester attrition-rate coefficients depend on # of targets, refine it as

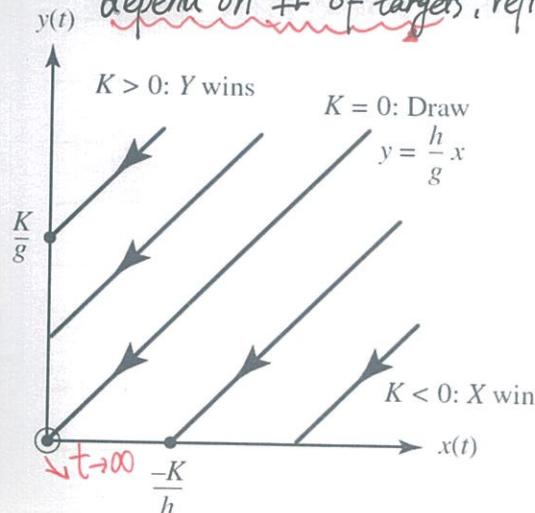
$$\begin{cases} x' = -gxy \\ y' = -hxy \end{cases}$$

$$\therefore \frac{dy}{dx} = \frac{h}{g} \text{ we have}$$

$$gy - hx = K$$

Find  $x(t)$  as

$$x(t) = \begin{cases} x_0 \left[ \frac{h^2 - gy_0}{h^2 - gy_0} e^{-\frac{h}{g}(t-t_0)} \right] & \text{for } h^2 \neq gy_0 \\ \frac{x_0}{1 + h^2 g t}, & \text{for } h^2 = gy_0 \end{cases}$$



Date: 2021 Oct. 1

## ⑤ Numerical Solutions: Euler's Method

### (i) Euler's Method:

Consider a system of ODEs  $\begin{cases} \frac{dx}{dt} = f(t, x, y) \\ \frac{dy}{dt} = g(t, x, y) \end{cases}$

Step-1 (Take time discretization)

$$t_i = t_0 + i\Delta t, \quad i \in \mathbb{N}$$

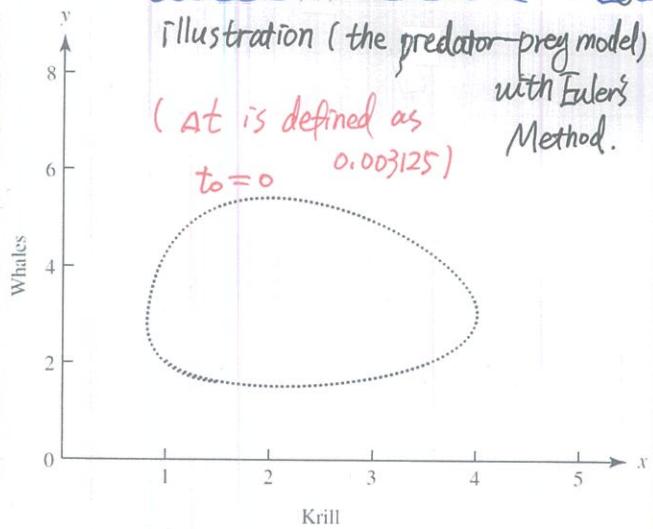
Step-2 (Approximation, iteration)

$$\frac{x_{i+1} - x_i}{\Delta t} = f(t_i, x_i, y_i); \quad \frac{y_{i+1} - y_i}{\Delta t} = g(t_i, x_i, y_i)$$

$$\Rightarrow \begin{cases} x_1 = x_0 + f(t_0, x_0, y_0)\Delta t \\ y_1 = y_0 + g(t_0, x_0, y_0)\Delta t \end{cases}; \quad \begin{cases} x_2 = x_1 + f(t_1, x_1, y_1)\Delta t \\ y_2 = y_1 + g(t_1, x_1, y_1)\Delta t \end{cases}; \quad \dots$$

Illustration (the predator-prey model)

with Euler's Method.



(ii) Examples.

Date: 2021 Oct. 1

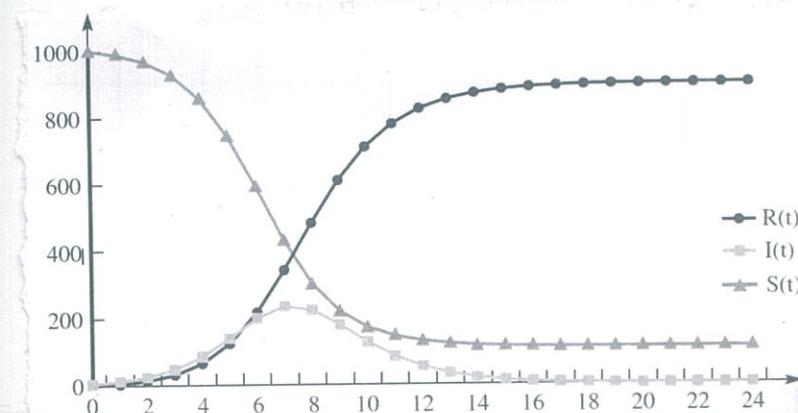
### egs ① Continuous SIR Model for epidemics

Consider a disease spreading through out the US. CDC is interested in knowing about & experimenting with a model for it.

► Modeling: Divide people into 3 categories.

- every person is either susceptible S / infected I / removed R
- no one leaves / enters, no other contacts with outside community.
- Someone got the flu can never be infected again.
- Time period = per week; Time needed to recover:  $\frac{5}{3}$  weeks avg.

$$\begin{cases} \frac{dR}{dt} = 0.6I \quad \text{per week 0.6 of them (avg) recovered.} \\ \frac{dI}{dt} = -0.6I + 0.001407IS \quad \text{density-dependent / mass action} \\ \frac{dS}{dt} = -0.001407IS \quad \text{got by data analyses.} \end{cases}$$



(data in the next page)

Date: 2021. Oct 12

Week	$R(t)$	$I(t)$	$S(t)$
0	0	5	995
1	4.29623	13.0107	982.693
2	13.6301	25.0714	961.299
3	31.503	47.0769	921.42
4	64.6671	84.2088	851.124
5	122.74	138.163	739.097
6	214.755	197.109	588.136
7	339.747	231.908	428.346
8	478.985	221.114	299.902
9	605.748	176.888	217.364
10	704.076	125.809	170.115
11	772.755	83.435	143.81
12	817.839	53.1578	129.003
13	846.397	33.0964	120.506
14	864.118	20.332	115.55
15	874.983	12.3924	112.624
16	881.598	7.51803	110.884
17	885.608	4.54821	109.844
18	888.033	2.74696	109.22
19	889.497	1.6574	108.845
20	890.381	0.99941	108.62
21	890.913	0.60242	108.484
22	891.234	0.36505	108.403
23	891.428	0.21876	108.354
24	891.544	0.13181	108.324

## egs ② Risk Aversion & Risk Premium

In finance (FIN 2010), we defined utility func & certainty-equivalent (CE). Assume  $U(x)$  is the utility func.

$$\therefore X_{CE} := U'(E(U(x))) \quad (\text{risk-aversion})$$

$\Rightarrow E(U(x)) < U(E(x))$  ( $U(x), U'(x) \uparrow$ ), by Jensen's equality.

$U(x)$  is concave, concavity  $\rightarrow$  extent of risk-aversion

$$\text{risk-premium} = E(x) - X_{CE}; (\pi_A)$$

$$\text{relative risk-premium} = \frac{\pi_A}{E(x)} = 1 - \frac{X_{CE}}{E(x)}; (\pi_R)$$

$$(E(x) \triangleq \bar{x}, \text{Var}(x) \triangleq \delta_x^2)$$

$$\text{Expand } U(x) \text{ around } \bar{x}: U(x) \approx U(\bar{x}) + \frac{1}{1!}U'(\bar{x})(x-\bar{x}) + \frac{1}{2!}U''(\bar{x})(x-\bar{x})^2$$

$$\text{when } x = X_{CE}, U(X_{CE}) \approx U(\bar{x}) + U'(\bar{x})(X_{CE}-\bar{x})$$

$$E(U(x)) \approx U(\bar{x}) + \frac{1}{2!}U''(\bar{x})\delta_x^2$$

Date: 2021. Oct. 12

Because  $E(U(x)) = U(X_{CE})$ , we have

$$U'(\bar{x})(X_{CE} - \bar{x}) \approx \frac{1}{2}U''(\bar{x})\delta_x^2 \quad -\frac{U''(\bar{x})}{U'(\bar{x})} =: A(x)$$

We can approximate  $\pi_A = \bar{x} - X_{CE} \approx -\frac{1}{2}\frac{U''(\bar{x})}{U'(\bar{x})}\delta_x^2$  Absolute Risk-aversion

$$\therefore \pi_A \approx \frac{1}{2}A(x)\delta_x^2; \quad -\frac{U''(\bar{x}) \cdot \bar{x}}{U'(\bar{x})} =: R(x)$$

$$\text{Similarly, } \pi_R = \frac{\pi_A}{\bar{x}} = \left[ -\frac{1}{2}\frac{U''(\bar{x}) \cdot \bar{x}}{U'(\bar{x})} \right] \cdot \frac{\delta_x^2}{\bar{x}} \quad \text{Relative Risk-aversion}$$

$$\therefore \pi_R = \frac{1}{2}R(x)\delta_x^2$$

▲ Special cases:

Constant Absolute Risk-Aversion (CARA)

$$A(x) = -\frac{U''(x)}{U'(x)} = r \Rightarrow U(x) = -ae^{-rx} + b$$

Constant Relative Risk-Aversion (CRRA)

$$R(x) = -\frac{U''(x) \cdot x}{U'(x)} = r \Rightarrow U(x) = ax^{1-r} + b, (0 < r < 1) \quad \text{or } a \ln x + b, (r=1)$$

Date: 2021 Oct 13

## PART II Differential Equations (PDE)

① Partial differential equations: often used to model the evolution over time & space.

② Derivation of PDEs.

(i) Wave Equation.

$$\frac{dm}{dx} = \frac{T_2 - T_1}{T_1} \approx \frac{T_2 - T_1}{T_1} \tan\theta_1 \quad (y = y(t, x))$$

$$F_y = T_2 \sin\theta_2 - T_1 \sin\theta_1 \approx T_2 \frac{\partial y}{\partial x}(t, x_2) - T_1 \frac{\partial y}{\partial x}(t, x_1) \approx T \frac{\partial^2 y}{\partial x^2}(t, x)(x_2 - x_1)$$

infinitesimal string

$\frac{dm}{dx}$

$y_{tt}(t, x)$

$\frac{\partial^2 y}{\partial x^2}$

linear density

$$\therefore y_{tt} = \frac{1}{c^2} y_{xx} \quad (\text{wave equation})$$

Solution of wave equation: (Characteristics)

$$\text{Let } c^2 = \frac{T}{\rho}, \quad y_{tt} = c^2 y_{xx}.$$

Let  $u = x - ct$ ,  $v = x + ct$ , then (changing variables)

$$\begin{cases} \frac{\partial}{\partial t} = & \left( \frac{\partial u}{\partial t} \frac{\partial}{\partial u} + \frac{\partial v}{\partial t} \frac{\partial}{\partial v} \right) = -c \frac{\partial}{\partial u} + c \frac{\partial}{\partial v} \\ \frac{\partial}{\partial x} = & \left( \frac{\partial u}{\partial x} \frac{\partial}{\partial u} + \frac{\partial v}{\partial x} \frac{\partial}{\partial v} \right) = \frac{\partial}{\partial u} + \frac{\partial}{\partial v} \end{cases}$$

$$\therefore \left( \frac{\partial}{\partial u} + \frac{\partial}{\partial v} \right)^2 y = \frac{1}{c^2} y_{tt} = \frac{1}{c^2} \cdot c^2 \left( \frac{\partial}{\partial u} - \frac{\partial}{\partial v} \right)^2 y$$

$$y_{xx} \Rightarrow y_{uv} = 0$$

which means  $\begin{cases} y_{uv} = f(u) + g(v) = f(x-ct) + g(x+ct) \\ y(x, t) \end{cases}$

Date: 2021 Oct 15

characteristic line: on line  $\begin{cases} u = x - ct \\ v = x + ct \end{cases}$  on the  $u-v$  plane,  
 $\frac{\partial^2 y}{\partial u \partial v} = 0$ .

### (ii) One-dim Conservation Law

Space-time variation of quantity density in a very thin tube

\*  $u(t, x)$  — denotes the quantity density at time  $t$ , location  $x$

⇒ Total amount of quantity in  $[x, x+h]$

$$A \int_x^{x+h} u(t, y) dy$$

(like volume)

\*  $\Phi(t, x)$  — denotes the flux density at time  $t$ , location  $x$

⇒ Net rate the quantity flows into  $[x, x+h]$  is

$$A [\Phi(t, x) - \Phi(t, x+h)]$$

\*  $f(t, x, \Phi(t, x))$  — denotes the source func (e.g. chemical reactions)

⇒ Rate of quantity produced in  $[x, x+h]$  by sources.

$$A \int_x^{x+h} f(t, y, u(t, y)) dy$$

The conservation law: rate of change amount = rate flows in + rate produced

$$\Rightarrow \frac{\partial}{\partial t} A \int_x^{x+h} u(t, y) dy = A [\Phi(t, x) - \Phi(t, x+h)] + A \int_x^{x+h} f(t, y, u(t, y)) dy$$

Divide both sides by  $h$ . take  $h \rightarrow 0$ , we have

$$u_t(t, x) = -\partial_x \Phi(t, x) + f(t, x, u(t, x))$$

i.e.  $\Phi_x(t, x)$

Date: 2021 Oct 15

### (iii) Transport Equations

(suppose) The quantity of interest moves together with the surrounding medium with a given velocity,  $v(t, x)$ . &  $f=0$  (no source)

Then  $\dot{\phi}(t, x) = u(t, x)v(t, x)$  (defn of  $\dot{\phi}$ )

$$\partial_t \dot{\phi}(t, x) = u_x(t, x) \cdot v(t, x) + v_x(t, x)u(t, x)$$

Use (ii) (a special case of 1-dim conservation law)

$$u_t(t, x) = -v u_x(t, x)$$

### (iv) Diffusion / Heat Equations

(suppose) No sources, the 1-dim conservation law:  $u_t + \dot{\phi}_x = 0$

Because heat moves from high temperature to low one

we have  $\dot{\phi}(t, x) = -D u_x(t, x)$  concentration

$\therefore$  1-dim diffusion equation  $\{u_t = D u_{xx}\}$

### (v) Initial & Boundary Conditions

PDEs need the initial & boundary conditions to be unique.

\* Initial conditions  $u(0, x) = \psi(x)$  form

\* Boundary conditions.

i) Dirichlet  $u(t, x) = g(t, x), \forall x \in \Omega$

Date: 2021 Oct 20

(e.g. heat equation, submerging in large container with melting ice.)

ii) Neumann  $\partial_x u(t, x) = g(t, x), \forall x \in \partial\Omega$

(e.g. heat equation, head is confined in the body.)

iii) Robin  $\partial_x u(t, x) + \alpha u(t, x) = g(t, x), \forall x \in \partial\Omega$

(e.g. heat equation, body immersed in a reservoir with known temperature)

### (3) Modeling Examples.

#### Eg 1. (Traffic Flow Modeling)

Different scales: (microscopic) a specific car & driver's respond  
(macroscopic) approx. as gas / liquid (continuum)

#### Modeling.

i) Continuum limit: suppose at every position  $x$ ,  $\exists$  a car with velocity  $v(t, x)$  at time  $t$ .

$$\begin{cases} \frac{dx}{dt} = v(t, x(t)), \\ x(t_0) = x_0 \end{cases}$$

(each car)

ii) Traffic density  $\rho(t, x)$ : avg # vehicles per unit length of road.

Suppose length of car =  $L$  & equal distance  $d$ .

$$\rho(t, x) = \frac{1}{L+d} \leq \frac{1}{L} \quad (d \geq 0)$$

Traffic flux  $q(t, x) = \rho(t, x) v(t, x)$

iii) Fix an arbitrary interval  $[A, B]$

By the conservation law  $\frac{d}{dt} \int_A^B \rho(t, x) dx = q(t, A) - q(t, B)$



Take  $B \rightarrow A = T$ , we have

Date: 2021 Oct 20

$$\frac{\partial P}{\partial t} + \frac{\partial P}{\partial x} = 0, \text{ with } g = PV, \text{ we have}$$

$$\frac{\partial g}{\partial t} + \frac{\partial (vP)}{\partial x} = 0$$

(Remark: let  $x_1(t), x_2(t)$  be positions of two cars (without any other cars / with some cars))

Then  $\frac{d}{dt} \int_{x_1(t)}^{x_2(t)} P(t, x) dx = x_2'(t)P(t, x_2(t)) - x_1'(t)P(t, x_1(t)) + \int_{x_1(t)}^{x_2(t)} \frac{\partial P}{\partial t} dx$

$\downarrow$   
no cars move  
in/out of the interval

$$0 = \int_{x_1(t)}^{x_2(t)} \frac{\partial (vP)}{\partial x} dx$$

$v(t, x)$  can be specified as  $V(P(t, x))$  (reasonable with reality)

& Let  $F(P) = PV(P)$ , we have  $\frac{\partial P}{\partial t} + F'(P) \frac{\partial P}{\partial x} = 0$

⇒ Applications: one of the possible  $V-g$  relations

$$V(P) = V_{max}(1 - P/P_{max})$$

$$\Rightarrow F'(P) = V_{max}(1 - 2P/P_{max}) \approx V_{max}(1 - 2P/P_{max}) := V_0$$

( $P \approx P_0$  which is a small deviation from a constant density)

Then, above model of PDE  $\frac{\partial P}{\partial t} + V_0 \frac{\partial P}{\partial x} = 0$

Use method of characteristics  $P(t, x) = f(x - V_0 t)$

with  $V_0$  - speed of linear traffic wave.

solve for  $P$ : parametrize line (characteristic) by  $(x(s), t(s))$

Consider  $\frac{dp}{ds} = \frac{dx}{ds} \frac{\partial P}{\partial x} + \frac{dt}{ds} \frac{\partial P}{\partial t}$

when  $\frac{dx}{ds} = V_0, \frac{dt}{ds} = 1, \frac{dp}{ds} = 0$ .

which means if  $\begin{cases} x = x_0 + V_0 s \\ t = t_0 + s \end{cases}$   $P(t, x) = P(t(s), x(s)) = P(t_0, x_0)$

choose different parametrization,

$(x_0, t_0) / (x_0 - V_0 t_0)$  we can get all pts  $(x, t)$

$$= P(0, x - V_0 t) \stackrel{?}{=} f(x - V_0 t)$$

Date: 2021 Oct 20

## Eg 2. (Population Growth Modeling)

Previously, we talk about population growth model without age structure. (e.g.  $N'(t) = rN(t)(M - N(t))$  logistic)

Modeling: population growth with age structure

parametrize  $t(s) = t$  (time);  $a(s) = a$  (age structure)

$$\Rightarrow n(t, a) = n(t(s), a(s)) = n(s).$$

i) Total population at time  $t$ :  $n(t) = \int_0^\infty n(t, a) da$

Let  $\mu(a)$  be the death rate at age  $a$   
collected by data analyses

$$\frac{dn(t, a)}{dt} = \frac{\partial n}{\partial t} dt + \frac{\partial n}{\partial a} da = -\mu(a)n(t, a)dt$$

Because  $\frac{da}{dt} = 1$  (defn),  $\Rightarrow \frac{\partial n}{\partial t} + \frac{\partial n}{\partial a} = -\mu(a)n(t, a)$

Von Foerster Equation

ii) Initial distribution age structure:  $n(0, a) = f(a), a > 0$

iii) Let  $b(a)$  be the birth rate at age  $a$  (age  $a$  people give birth to age 0)

then,  $n(t, 0) = \int_0^\infty b(a)n(t, a) da$

$(b(a) \rightarrow 0 \text{ when } a \rightarrow 0/\infty)$

Solve for the model:

By parametrize above, we want

$$-\mu(a(s))n(s) = n'(s) = \frac{\partial n}{\partial t} t'(s) + \frac{\partial n}{\partial a} b'(s)$$

Want  $\frac{\partial n}{\partial t} t'(s) \stackrel{?}{=} 1$  set to 1

Date: 2021 Oct 21

We then have  $n(s) = n_0 \exp\left(-\int_{a_0}^{a_0+s} \mu(s) ds\right)$   
(with  $t = t_0 + s$ ,  $a = a_0 + s$ )

For  $a > t$ , set  $t_0 = 0$ , then  $n(t, a) = n(0, a_0) \exp\left(-\int_{a_0}^a \mu(s) ds\right)$   
 $\uparrow$  In this case,  $n_0 = n(0, a_0)$

$$n(t, a) = f(a-t) \exp\left(-\int_{a-t}^a \mu(s) ds\right), a > t$$

For  $a < t$ , set  $a_0 = 0$ .  $\Rightarrow a = t - t_0 = s$

$$n(t, a) = n(t_0, a_0) \exp\left(-\int_0^a \mu(s) ds\right) = n(t-a, 0) \exp\left(-\int_0^a \mu(s) ds\right)$$

$$= \int_0^t b(a) n(t-a) da + \int_t^\infty b(a) n(t-a) da = \int_0^t b(a) n(t-a) da$$

$$\exp\left(-\int_0^a \mu(s) ds\right) da + \int_t^\infty b(a) f(t-a) \exp\left(-\int_{a-t}^a \mu(s) ds\right) da$$

2 terms  
boundary conditions

### Long-term asymptotic solutions

On the long run, the initial age structure is not important.

Seek for the solution in form  $n(t, a) = e^{rt} r(a)$

(satisfying  $\frac{n(t, a)}{\int_0^\infty n(t, a) da} = \frac{r(a)}{\int_0^\infty r(a) da}$ )

$$\text{Then } \frac{\partial n}{\partial t} + \frac{\partial n}{\partial a} = re^{rt}r'(a) + e^{rt}r'(a) = -\mu(a)e^{rt}r(a)$$

$$\Rightarrow r'(a) = -(\mu(a) + r)r(a)$$

$$\therefore r(a) = r(0) \exp\left(-ra - \int_0^a \mu(s) ds\right)$$

Date: 2021 Oct 21

By the boundary condition,  $n(t, 0) = e^{rt} r(0) = \int_0^\infty b(a) e^{rt} r(a) \exp(-ra - \int_0^a \mu(s) ds) da$   
 $r$  determined by  $1 = \int_0^\infty b(a) \exp(-ra - \int_0^a \mu(s) ds) da$ .

### Eg 3. (Financial Derivatives Pricing – Black-Scholes Model.)

Option = getting the right but no obligation of buying a stock at \$100 per year ( $S_0$ )

Aim: give a fair option price (do option pricing)

Pricing by replication:

consider the over-simplified case: the value of stock traded at  $\$100$ , and can become  $uS_0$ ,  $dS_0$  ( $u=1 > d$ ) with probability  $p$ ,  $1-p$ . The risk-free asset have the rate of return  $r$ . (No friction in trading)

IDEA: If there is a portfolio perfectly replicating the option payoff, its value = the option price.

(Otherwise, arbitrage – e.g. if option price > portfolio value,

$\Rightarrow$  sell the option & buy the portfolio)

$\Delta$  shares  $\{ b(1+r) + \Delta S_0 = V_i^u \}$  option payoff case I  
&  $b$  risk-free assets  $\{ b(1+r) + \Delta d S_0 = V_i^d = 0 \}$  option payoff case II

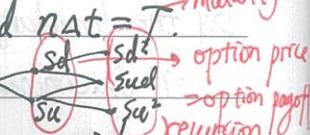
$$\therefore V_i = b + \Delta S_0 = \frac{1}{1+r} (\tilde{p} V_i^u + (1-\tilde{p}) V_i^d)$$

$$= \frac{1}{1+r} E^a(V_i) \text{ with } \tilde{p} = \frac{1+r-d}{d-r} \quad d < 1+r < u$$



Date: 2021 Oct 25

From one period, we consider multi-period  $\eta \Delta t = T$  maturity with one period interest rate  $r \Delta t$ .



After  $k$  step,  $S_k(w) = S_0 d^{k-j} u^j$  with  $j$  ups,  $k-j$  downs

$$d = 1 + r \Delta t - \delta \sqrt{\Delta t} < 1 + r \Delta t < 1 + r \Delta t + \delta \sqrt{\Delta t} = u, \text{ with } \tilde{p} = \frac{1}{2}$$

### Modeling:

Suppose the option price is  $C(t, S) = C$ , where  $0 \leq S < \infty$ ,  $t \leq T$ ,  $C$  is smooth. like

$$\text{Then } C^+ := C(t_0 + \Delta t, uS_0) = C(t_0 + \Delta t, S_0 + (u-1)S_0) \quad (V_i^u)$$

$$C^- := C(t_0 + \Delta t, dS_0) = C(t_0 + \Delta t, S_0 + (d-1)S_0) \quad (V_i^d)$$

$$\begin{aligned} \text{By Taylor's theorem, } C^+ &= C(t_0, S_0) + \frac{\partial C}{\partial S}(t_0, S_0)(\delta \sqrt{\Delta t} + r \Delta t)S_0 \\ &\quad + \frac{\partial C}{\partial t}(t_0, S_0)\Delta t + \frac{1}{2} \frac{\partial^2 C}{\partial S^2}(t_0, S_0)(\delta \sqrt{\Delta t} + r \Delta t)^2 S_0^2 + O(\sqrt{\Delta t}^3) \\ C^- &= C(t_0, S_0) + \frac{\partial C}{\partial S}(t_0, S_0)(-\delta \sqrt{\Delta t} + r \Delta t)S_0 + \frac{\partial C}{\partial t}(t_0, S_0)\Delta t + \frac{1}{2} \frac{\partial^2 C}{\partial S^2}(t_0, S_0)(-\delta \sqrt{\Delta t} + r \Delta t)^2 S_0^2 + O(\sqrt{\Delta t}^3) \end{aligned}$$

According to the previous conclusion, we have

$$(1 + r \Delta t) C(t_0, S_0) = \frac{1}{2}(C^+ + C^-) + O(\Delta t^2)$$

Up to  $O(\sqrt{\Delta t}^3)$ ,  $\Delta t \rightarrow 0$ , recursion (Sell option at price)

$$\begin{aligned} \text{we get } (1 + r \Delta t) C(t_0, S_0) &= C(t_0, S_0) + \frac{\partial C}{\partial S}(t_0, S_0) r \Delta t S_0 + \frac{\partial C}{\partial t}(t_0, S_0) \Delta t \\ &\quad + \frac{1}{2} \frac{\partial^2 C}{\partial S^2}(t_0, S_0) S_0^2 \cdot \delta^2 \Delta t \end{aligned}$$

$$\Rightarrow r C(t_0, S_0) = \frac{\partial C}{\partial S}(t_0, S_0) r S_0 + \frac{\partial C}{\partial t}(t_0, S_0) + \frac{1}{2} \frac{\partial^2 C}{\partial S^2}(t_0, S_0) \delta^2 S_0^2$$

Suppose at maturity the stock price is  $K$ .

$$\text{we have initial condition } C(T, S) = (S - K)^+$$

Date: 2021 Oct 25

(if  $S \leq K$ , get 0, otherwise  $S - K$ )

$\Rightarrow$  Black-Scholes PDE

$$\left. \begin{cases} C_t + r S C_S + \frac{1}{2} \delta^2 S^2 C_{SS} - r C = 0 \\ C(T, S) = (S - K)^+ \end{cases} \right\}$$

Solve it by changing variables  $C(t, S) = e^{-r(T-t)} G(S, X)$

$$X = \ln S + (r - \frac{\delta^2}{2})(T-t), S = e^{X + (2r/\delta^2 - 1)t}$$

$$T-t = \frac{2S}{\delta^2}, S = e^{X - (2r/\delta^2 - 1)t}$$

$$\Rightarrow \begin{cases} G_S = G_{XX}, X \in \mathbb{R}, S > 0 \quad (\text{Heat Equation}) \\ G(0, X) = (e^X - K)^+ \end{cases}$$

By Fourier Transform  $(\widehat{G}(S, \xi))$  need more explanations

$$\star F\left(\frac{\partial G}{\partial S}\right) = \frac{\partial F(G)}{\partial S} = \widehat{G}_S$$

$$F\left(\frac{\partial^2 G}{\partial X^2}\right) = (i\xi)^2 F(G) = -\xi^2 \widehat{G} \quad g(x) = G(0, x)$$

$$\Rightarrow \partial_x \widehat{G} + \xi^2 \widehat{G} = 0. \quad \widehat{G} = \widehat{\varphi}(\xi) e^{-\xi^2 x} \quad (\text{let } \widehat{G}(0, x) = \widehat{\varphi}(\xi))$$

$$\begin{aligned} \therefore G &= F^{-1}(\widehat{G}) = \int_{-\infty}^{\infty} (e^y - K)^+ \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-y)^2}{4s}} dy \\ &= e^{xt+s} N\left(\frac{x - \ln K + 2s}{\sqrt{2s}}\right) - KN\left(\frac{x - \ln K}{\sqrt{2s}}\right) \end{aligned}$$

$$\text{where } N(\cdot) = \int_{-\infty}^{\cdot} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx.$$

$$\Rightarrow Y(t, S) = S N(d+) - e^{-r(T-t)} N(d-)$$

$$d_{\pm} = \frac{\ln(S/K) + (r \pm \frac{\delta^2}{2})(T-t)}{\delta \sqrt{T-t}}$$

Date: 2021 Oct 29

## • PART III Optimization - Linear Programming

### ① Theoretical Knowledge:

Defn; Geometric solutions; the simplex method;

sensitivity analysis (in MAT 3007 course)

### ② Modeling Examples:

#### e.g. 1) Foreign-Currency Trading

- Assumptions. Big FX market with over \$1 trillion daily trading  
One day's currency rate as follows

To:	US Dollar	Pound	FFranc	D-Mark	Yen
From: US Dollar		0.6390	5.3712	1.5712	98.8901
Pound	1.5648		8.4304	2.4590	154.7733
FFranc	0.1856	0.1186		0.2921	18.4122
D-Mark	0.6361	0.4063	3.4233		62.9400
Yen	0.01011	0.00645	0.05431	0.01588	

- Aim. No arbitrage! (currency transaction creates profit without initiating funds called arbitrage)

- Modeling. objective - maximize final amount of one currency  
constraints - nonnegative final net amount (no investment)

$i = 1, 2, \dots, 5$  represent the currencies US dollar, British pound, French franc, German D-mark & Japanese Yen,  
respectively.  $X_{ij}$  - amount of currency from  $i$  to  $j$  ( $i \neq j$ )

Date: 2021 Oct 29.

Let  $f_k$  be final net amount of currency  $k$ .

$$\text{we have } f_1 = 1.5648X_{61} + 0.1856X_{31} + 0.6361X_{41} + 0.01011X_{51}$$

$$- \sum_{i \neq 1} X_{ii}$$

$$f_2 = 0.6390X_{12} + 0.1186X_{22} + 0.4063X_{42} + 0.00645X_{52} - \sum_{i \neq 2} X_{ii}$$

$f_3, f_4, f_5$  can be written in the same way.

The Model =

(LP model)

$$\max_{\sum_i X_{ij} = 1} f_i$$

$$\text{subject to } f_i, X_{ij} \geq 0, \forall i \neq j, i, j = 1, \dots, 5$$

$$f_i \leq 1$$

$$\Rightarrow \max_{\sum_i X_{ij} = 1} f_i \leq 1$$

Once arbitrage emerges  
 $f_i \rightarrow \infty$  (in case)  
normalize

#### e.g. 2) Option Price in Incomplete Market

- Assumptions.  $n$  states,  $m$  assets, with  $m \times n$  payoff

matrix  $X = (X_{ij})$  means  $i$  asset at  $j$  state (price)

The asset pricing vector  $p = (p_i) \in \mathbb{R}^m$

- Modeling. If there is an arbitrage opportunity,

$\exists$  amount vector  $a \in \mathbb{R}^m$  s.t.  $a^T X \geq 0$  &  $a^T p = 0$

or  $a^T X = 0$  &  $a^T p < 0$

$\Leftrightarrow$  an optimization problem as

$$\begin{aligned} & \min a^T p \\ & \text{subject to } a^T X \geq 0 \end{aligned}$$



Date: 2021 Oct 29

By alternating theorem (Farka's Lemma)

It is equivalent to another expression.

(The First Fundamental Theorem of Asset Pricing - FTAP)

The market admits NO arbitrage exactly when  $\exists$  a

positive vector  $b \in \mathbb{R}^n$  s.t.  $p = Xb$

( $b$  works as a probability vector)

### Application.

(i) Binomial Model. in Black-Schole PDE part. we can use FTAP to get  $E^Q(U) \cdot \frac{1+r}{1+r}$ , by  $b = \begin{bmatrix} \frac{1+r}{1+r} \\ \frac{a-(1+r)}{1+r} \\ \frac{a-a}{1+r} \end{bmatrix}$ .

(ii) Incomplete Markets:

Consider the case of  $n$  possible values.  $S_1 = u_1 S_0, u_2 S_0, \dots, u_n S_0$

No arbitrage price is not unique (underdetermined system)

$$V_0 = [V_0^-, V_0^+] = \begin{bmatrix} \min_{S_i^- b = S_0} V_i^T b & \max_{S_i^+ b = S_0} V_i^T b \\ \text{s.t. } S_i^- b = S_0 & \text{s.t. } S_i^+ b = S_0 \\ 1^T b = 1+r & 1^T b = 1+r \end{bmatrix}$$

Dual problem of  $V_0^-$  — super replication ( $\min y^T p$ )

Dual problem of  $V_0^+$  — sub-replication ( $\max_{\text{subject to } X^T y \geq V_0^+} (-y)^T p$ )

$$\text{subject to } X^T(-y) \geq V_0^+$$

### e.g. 3) Estimating & Managing Risk

#### Estimating Value-at-Risk (VaR)

Denote random variable  $L$  as the loss of a portfolio

Date: 2021 Oct 30

in the future. The VaR is defined as

$$\text{VaR}_\alpha = \sup \{x \mid F_L(x) \leq \alpha\} = F_L^{-1}(\alpha)$$

where  $F_L(x)$  can be seen as the cumulative distribution func. (CDF)

$$F_L(x) = P(L \leq x)$$

Estimate VaR (Monte Carlo Simulation)

(i) Generate  $n$  i.i.d. (independent, identically distributed)

samples of  $L$ , sorted as  $L_1 \leq L_2 \leq \dots \leq L_n$

(ii) Approx.  $F$  as

$$\star F_L^{(n)}(x) = \begin{cases} 0, & x \leq L_1 \\ \frac{i-1}{n-1} + \frac{1}{n-1} \cdot \frac{x-L_i}{L_{i+1}-L_i}, & L_i \leq x \leq L_{i+1}, i=1, \dots, n-1 \\ 1, & x \geq L_n \end{cases} \quad (\text{piece-wise linear})$$

$$(iii) \text{ If } \alpha \in [\frac{i-1}{n-1}, \frac{i}{n-1}], \text{ VaR}_\alpha^{(n)}(L) = L_i + (L_{i+1} - L_i)(\alpha(n-1) - i+1)$$

#### Conditional Value-at-Risk (CVaR)

$$\text{Defined as } \text{CVaR}_\alpha^{(n)} = \frac{\mathbb{E}(L | L > \text{VaR}_\alpha)}{\text{Pr}(L > \text{VaR}_\alpha)}$$

Estimate CVaR

(i) estimate VaR as above shows use  $n$  samples to estimate

$$(ii) \text{ Approx. CVaR as } \text{CVaR}_\alpha^{(n)} = \frac{\frac{1}{n} \sum_{i=1}^n L_i \mathbb{I}_{\{L_i > \text{VaR}_\alpha^{(n)}\}}}{\frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{L_i > \text{VaR}_\alpha^{(n)}\}}} \rightarrow \text{indicator function}$$

(By Law of Large Numbers)



A numerical example:

$$F_L = 1 - e^{-x} \text{ (suppose)} \quad \text{CVaR}_\alpha = \frac{1}{1-\alpha} \int_{\text{VaR}_\alpha}^{\infty} xe^{-x} dx \stackrel{\alpha=0.95}{=} 3.996$$

Date: 2021 Oct 30

The result of estimation:

Stochastic Optimization:

$$\min_x \mathbb{E}(f(x, \xi))$$

subject to  $x \in X$

can be approximate by (Monte Carlo Simulation)

Take a sample  $\xi_1, \xi_2, \dots, \xi_N$ .

$$\min_x \frac{1}{N} \sum_{i=1}^N f(x, \xi_i)$$

subject to  $x \in X$

Modeling:

Investment problem (simple)

$$\min_{x_i} CVaR_\alpha(-\sum_{i=1}^n x_i R_i)$$

loss

minimize loss risk

$$\text{subject to } \mathbb{E}(\sum_{i=1}^n x_i R_i) \geq r_0$$

$$\sum_{i=1}^n x_i = 1$$

promise certain profit

$n$  assets  $x_1, \dots, x_n$  with random rate of return  $R_i, i=1, \dots, n$ .

An alternative formulation of CVaR

$$CVaR_\alpha(L) = \min_{\xi \in R} [\xi + \frac{1}{1-\alpha} \mathbb{E}[\max\{L - \xi, 0\}]]$$

$$\text{Reason: } g(\xi) = \xi + \frac{1}{1-\alpha} \mathbb{E}[(L-\xi)^+] = \xi + \frac{1}{1-\alpha} \int_{-\infty}^{\infty} (x-\xi) f_L(x) dx$$

$$\Rightarrow g'(\xi) = \frac{1}{1-\alpha} (F_L(\xi) - \alpha). \text{ only candidate is } \xi = VaR_\alpha$$

Date: 2021 Oct 30

$$\therefore CVaR_\alpha(L) = \min_{\xi \in R} g(\xi)$$

We can write the initial problem as

$$\min_{x_i, \xi} \xi + \frac{1}{1-\alpha} \mathbb{E}[\max\{\sum_{i=1}^n x_i R_i - \xi, 0\}]$$

$$\text{subject to } \mathbb{E}(\sum_{i=1}^n x_i R_i) \geq r_0$$

$$\sum_{i=1}^n x_i = 1.$$

Use simulation to estimate  $R_i$  with realizations  $R_{ij}$  & probabilities  $P_j = \frac{1}{S}$

$$\mathbb{E}(\sum_{i=1}^n x_i R_i) = \sum_{i=1}^n x_i \mathbb{E}(R_i) = \sum_{i=1}^n x_i \frac{1}{S} \sum_{j=1}^S R_{ij} \triangleq \sum_{i=1}^n x_i \bar{R}_i$$

$$\mathbb{E}[\max\{\sum_{i=1}^n x_i R_i - \xi, 0\}] = \mathbb{E}[\max\{\sum_{i=1}^n \sum_{j=1}^S x_i R_{ij} - \xi, 0\}]$$

$$= \mathbb{E}[\max\{\sum_{j=1}^S \sum_{i=1}^n x_i R_{ij} - \xi, 0\}] = \frac{1}{S} \sum_{j=1}^S u_j$$

denoted as  $u_j, \sum_{j=1}^S u_j$

the estimation LP:

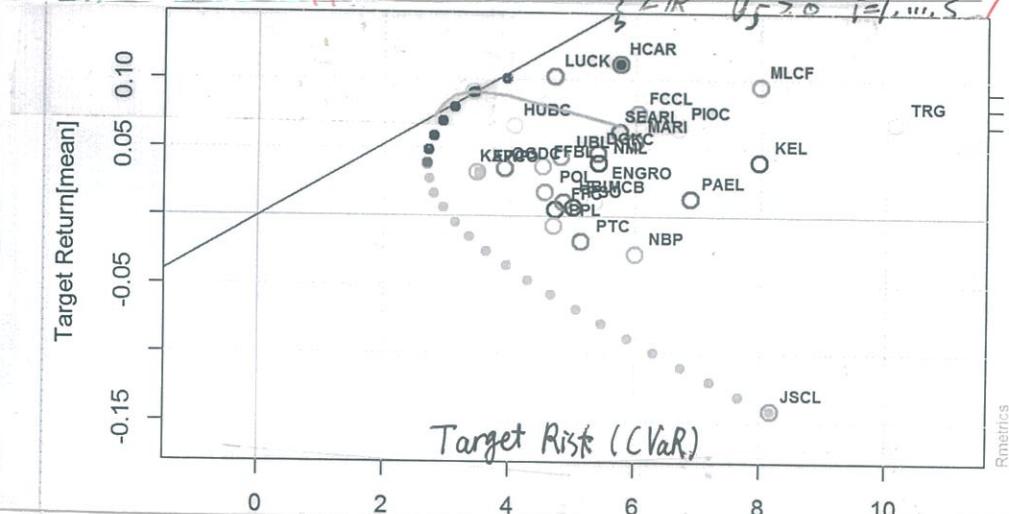
$$\begin{cases} \text{minimize } \xi + \frac{1}{(1-\alpha)S} \sum_{j=1}^S u_j \\ \xi, u_j, x_i \end{cases}$$

$$\text{subject to } u_j \geq -\sum_{i=1}^n x_i R_{ij} - \xi, j=1, \dots, S$$

$$\sum_{i=1}^n x_i R_i \geq r_0, \sum_{i=1}^n x_i = 1$$

$$\xi \leq L, u_j \geq 0, j=1, \dots, S$$

Efficient Frontier



Date: 2021 Nov 6

## ● PART IV – Optimization – Non-linear Programming

### ① Theoretical (part 1)

- (i) Existence of solutions; solving NLP by software (in MATLAB)
- (ii) Examples using NLP modeling

#### e.g. 1 (Regression)

Given a sequence of data points  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  with  $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{R}$ . want to find a mapping  $f(x; \theta)$  with  $\theta$  as the parameter. (e.g.  $f(x) = a + b x$  with  $\theta = (a, b)$ )

formulated as  $\min_{\theta} \frac{1}{n} \sum_{i=1}^n |y_i - f(x_i; \theta)|^r$  (or L<sub>r</sub> norm  $\|y - f(x; \theta)\|_r$ )

#### e.g. 2 (Likelihood maximization Estimation)

Suppose  $\exists$  a sequence of independent random observations  $\{x_1, x_2, \dots, x_n\}$  of a RV  $X$ . Want to infer a parametric distribution  $f(x; \theta)$  of  $X$  from the observations. (e.g. Gaussian  $f(x; \theta) = \frac{1}{\sqrt{2\pi^n |\Sigma|}} \exp^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$ )

where  $\Sigma$  is the covariance matrix  $\Sigma$

formulated as  $\max_{\theta} \prod_{i=1}^n f(x_i; \theta)$  (Because  $P(X=x_1 \text{ or } x_2 \text{ or } \dots \text{ or } x_n) \sim$  (density)  $\prod_{i=1}^n f(x_i; \theta)$ )

#### e.g. 3 (Portfolio Optimization)

An investor has  $n$  available assets where she allocates her money. RV  $R_i$  represents the (random) return of asset  $i$

Date: 2021 Nov 6

Naive formulation  $\max_{\mathbf{x}} \sum_{i=1}^n R_i x_i$ , s.t.  $\sum_{i=1}^n x_i = 1$   
(NOT well-defined)

Further improvement  $\max_{\mathbf{x}} \sum_{i=1}^n \mu_i x_i$ , s.t.  $\sum_{i=1}^n x_i = 1$   
where  $\mu_i = \mathbb{E}[R_i]$  (NOT well-posed  $\Rightarrow$  only all  $\mu_i$  equals) <sup>not unbdd</sup>

#### ★ Mean-Variance Framework (Harry Markowitz)

$$\max_{\mathbf{x}} \mu^T \mathbf{x} - \frac{1}{2} \gamma \mathbf{x}^T \Sigma \mathbf{x} \quad \text{s.t. } \mathbf{e}^T \mathbf{x} = 1$$

where  $\mu = \mathbb{E}[R]$ ,  $\Sigma = \text{Var}[R]$  and  $\gamma$  - risk-aversion coefficient / covariance matrix penalty

(\*)

$$(\mathbf{x}^T \Sigma \mathbf{x} = \text{Var}[\mathbf{x}^T \mathbf{R}] \quad (\text{Var}(\sum_{i=1}^n R_i x_i) = \mathbf{x}^T \text{Var}[R] \mathbf{x}))$$

the property of Var

#### Alternative Formulations

(same Lagrangian)  $\max_{\mathbf{x}} \mu^T \mathbf{x}$ , s.t.  $\mathbf{x}^T \Sigma \mathbf{x} \leq \sigma^2$ ,  $\mathbf{e}^T \mathbf{x} = 1$  "efficient frontier"

$$\min_{\mathbf{x}} \mathbf{x}^T \Sigma \mathbf{x}, \text{ s.t. } \mu^T \mathbf{x} \geq z, \mathbf{e}^T \mathbf{x} = 1$$

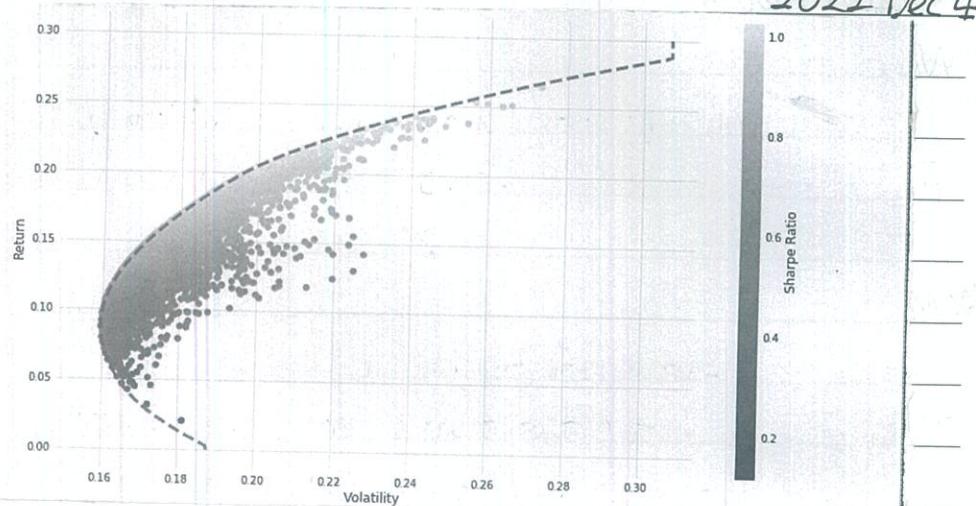
(Giver Return  $\Rightarrow$  minimum risk)

(Estimations of  $\mu$ ,  $\Sigma$  = factor analysis, Black-Litterman shrunk estimation, ... )

By KKT conditions, we can find analytical solution for (\*)

$$\mathbf{x} = \frac{1}{\gamma} \Sigma^{-1} \mu - \frac{\lambda \mathbf{e}}{\gamma} \Sigma^{-1} \mathbf{e}, \text{ where } \lambda = \frac{\mu^T \Sigma^{-1} \mathbf{e} - r}{\mathbf{e}^T \Sigma^{-1} \mathbf{e}}$$

The corresponding efficient frontier.  
(next page)



## ② Theoretical (Part II)

(i) Optimality conditions, convex optimizations, KKT conditions  
(in MAT3007 Notes)

(ii) Non-differentiable Convex Optimization & sub-gradient

The subgradient of  $f$ , denoted as  $\partial f$ , is given by

$$\partial f(x) := \{g_x | f(y) \geq f(x) + g_x^T(y-x)\}, \forall y$$

Given  $f$  is convex,  $x^*$  is a global minimizer iff  $0 \in \partial f(x^*)$ .

e.g.  $f = \|x\|_2$  (Euclidean norm).  $\partial f(0) = \{x | \|x\|_\infty \leq 1\}$

(However,  $f$  is not differentiable at  $x=0$ )

\* Application - LASSO (l1 regularized linear regression)

For the model  $y = X\beta$ , LASSO estimator is defined as

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 + \rho \|\beta\|_1$$

(Generalized LASSO  $\frac{1}{2} \|X\beta - y\|_2^2 + \rho \|F\beta\|_1$ , where  $F$  is

Date: 2021 Dec 4

an arbitrary linear transformation.)

$$(sol.) \text{ ADMM form: } \begin{cases} \beta^{k+1} := (X^T X + \lambda I)^{-1} (X^T y + \lambda (\beta^k - u^k)) \\ \text{Add constrain: } z^{k+1} := S_{\lambda}^{\rho}(\beta^{k+1} + u^k) \\ \beta = z / (\beta - z = 0) \\ u^{k+1} := u^k + \beta^{k+1} - z^{k+1} \end{cases} \quad \text{ridge regression}$$

(iii) Traveling Salesman Problem (TSP) & Simulated Annealing (SA)

TSP: find the shortest tour towards set of given locations.

s.t. 1° each location is passed once; 2° end at the starting point.

\* Simulated Annealing Algorithm (famous together with NN, GA)

- 1° Choose a random state  $s$ , initial Temperature  $T$  (large)
- & the cooling rate  $\beta$ . ( $\beta \in (0, 1)$ )
- 2° Create a new state  $s'$ .
- 3° Compute the cost function of 2 states & compare:  
If  $\{ C(s') < C(s) \} \left( \delta = \frac{C(s') - C(s)}{C(s)} < 0 \right)$  Accept  $s'$   
Otherwise ( $\delta \geq 0$ ) Accept  $s'$  with probability  $e^{-\delta/T}$   
& cool  $T$  as  $\beta T$  ( $T = \beta T$ )
- 4° Stopping criterion:  
Small temperature / reject enough times new  $s'$

Traits of SA: heuristic (not necessary get optimal)

Have probability go out of local optima

Date: 2021 Dec 4

### ③ Dynamic Programming

(i) Core IDEA: break complex problem into simple ones  
(Divide & conquer, recursion)

(ii) Examples:

e.g. 1 Chance = draw a die at most 3 times. Earnings: face value of a die. Option: stop after each roll.

Strategy: (DP) consider the third row  $E(3^{rd}) = \frac{\sum_{i=1}^6 i}{6} = 3.5$

→ consider the second row } Stop if  $X_2 > 3.5$  or  
proceed to 3<sup>rd</sup>

$$E(2^{nd}) = \frac{3}{6} \times \left( \frac{4+5+6}{3} \right) + \frac{3}{6} \times 3.5 = 4.25$$

→ consider the first row } Stop if  $X_1 > 4.25$  or  
proceed to 2<sup>nd</sup>

$$E = E(1^{st}) = \frac{2}{6} \times \left( \frac{5+6}{2} \right) + \frac{4}{6} \times 4.25 = \frac{14}{3}$$

### e.g. 2 (The Knapsack Problem)

A set of  $n$  items, each with weight  $w_i$  & value  $v_i$ ,  
find maximum total value, given the capacity of knapsack  $W$ .

$$\text{maximize } \sum_{i=1}^n v_i x_i$$

$$\text{subject to } \sum_{i=1}^n w_i x_i \leq W$$

$$x_i \in \{0, 1\}, i=1, \dots, n$$

Integer Programming

Date: 2021 Dec 5

Steps: split to multi-stages; identify state & action;  
identify the opt. cost-to-go; recursive relationship

Sol.  $m(i, w)$  denotes maximum total value can be obtained using the first  $i$  items with capacity limit  $w$ .

$$m(0, w) = 0, \forall 0 \leq w \leq W$$

$$m(i, w) = \max \{ m(i-1, w-w_i) + v_i + m(i-1, w) \}, \forall 1 \leq i \leq n, 0 < w \leq W$$

(The complexity:  $O(nW)$  given that  $w_i \in \mathbb{Z}$ )

→ Use form to solve:  $\frac{0}{2} / \dots / \frac{0}{2}$

### e.g. 3 (Matrix Product Parenthesization)

Suppose we need to multiply a series of matrices:  $A_1, \dots, A_n$   
Given the dimensions, find the best way to add parenthesis.

$$A_i \in \mathbb{R}^{P_{i-1} \times P_i}$$

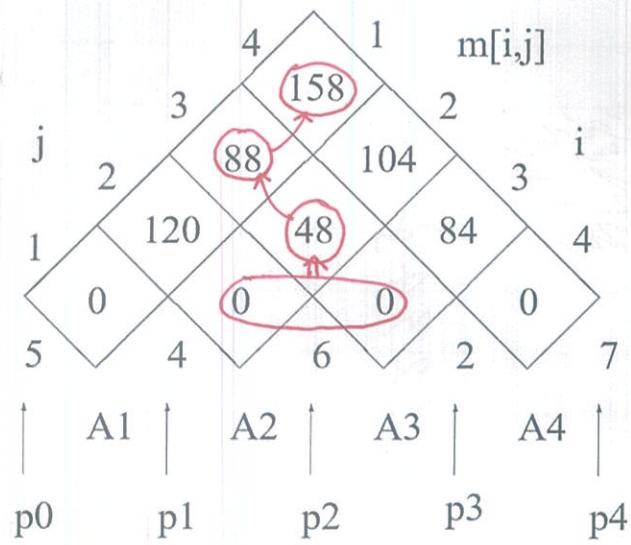
Sol. Let  $m(i, j)$  be the minimum number of operations to compute  $A_i \dots A_j$ . then

$$m(i, j) = \begin{cases} 0, & \text{if } i=j \\ \min_{1 \leq k < j} \{ m(i, k) + m(k+1, j) + p_{i-1} p_k p_j \} \end{cases}$$

Step I  $\Rightarrow$  Step II  $\Rightarrow$  Step n-1  
 $m(1, 2), m(3, 4), \dots, m(n-1, n) \quad m(1, 3), m(2, 4), \dots \quad m(1, n) = \text{target}$

Complexity:  $O(n^2)$

An example:  $A_1, A_2, A_3, A_4$  with  $p_0 = 5, p_1 = 4, p_2 = 6, p_3 = 2, p_4 = 7$



Date: 2021 Dec 5

The optimal steps  
& all  $m(i,j)$  in  
the graph.

#### e.g.4. (Gene/DNA sequence alignment)

Given two strings:  $x = x_1 x_2 \dots x_m$ ,  $y = y_1 y_2 \dots y_n$ . Alignment - assignment of gaps to positions  $0, \dots, M$  in  $x$  &  $0, \dots, N$  in  $y$ , so as to line up each letter in one sequence with either a letter or a gap in another sequence. (like  $\begin{matrix} AGGCTAGTT \\ AGCGAAGTT \end{matrix} \Rightarrow \begin{matrix} AG\text{--}G-T-A-G-T-T \\ A-G-C-G-A-A-G-T-T \end{matrix}$ )

Sol. Define scoring function (match + m, gap - d, mismatch - s)

$$\text{scoring } \leftarrow F(i, j) = \max / F(i-1, j-1) + S(x_i, y_i), F(i-1, j) - d, F(i, j-1) - d$$

Recursive DP procedure (consider the last letter of the 2)  
Complexity =  $O(MN)$

Date: 2021 Dec 5

match = 1      mismatch = -1      gap = -1

	G	C	A	T	G	C	U	
G	0	-1	-2	-3	-4	-5	-6	-7
A	-1	0	0	1	0	-1	-2	-3
T	-2	0	0	1	0	-1	-2	-3
T	-3	-1	-1	0	2	1	0	-1
A	-4	-2	-2	-1	1	1	0	-1
A	-5	-3	-3	-1	0	0	0	-1
C	-6	-4	-2	-2	-1	-1	1	0
A	-7	-5	-3	-1	-2	-2	0	0

\* 2 assumptions for DP:

i) Memoryless: Decisions made "sequentially" in states, depending on (only) states.

ii) Additive Cost/reward function (backward separable & forward monotonic)

#### e.g.5 (Expected Utility Maximization / EUM)

Discrete time market model:  $R_f = 1 + r$  (risk-free)

(stock)  $R_t$  i.i.d overtime. Goal = maximize  $E(U(W_T))$

where  $W_{t+1} = W_t (\phi_t R_{t+1} + (1-\phi_t) R_f)$ ,  $t=0, 1, \dots, T-1$

( $\phi_t$  — decisions to make)

Sol.  $J_0 = E(U(W_T))$ ,  $J_t = E(J_{t+1} | W_t)$ ,  $J_T = U(W_T)$

Let  $J_t^*(W_t) = \max_{\phi_t, \dots, \phi_T} E(U(W_T) | W_t)$

Bellman Equation  $J_t^*(W_t) = \max_{\phi_t} E[J_{t+1}^*(W_t (\phi_t R_{t+1} + (1-\phi_t) R_f)) | W_t]$

(Reason:  $E(J_{t+1}(W_{t+1}) | W_t) = E(E(J_{t+2} | W_{t+2}) | W_t) = E(J_{t+2} | W_t)$   
 $= \dots = E(U(W_T) | W_t)$ )

(APP) Consider  $S_{t+1} = \{u^s_t, d^s_t\}$ , with  $p$  & CRRA ( $U(W_t) = \frac{1}{r} W_t^{1-r}$ )  
 $\Rightarrow$  (backward induction).  $J_t^* = f_t W_t^{1-r}$ ,  $\phi_t = \frac{1-r}{r \sigma^2}$   
with  $u = 1 + \mu u t + \sigma \sqrt{u} t$ ,  $d = 1 + \mu u t - \sigma \sqrt{u} t$ ,  $R_f = 1 + r u t$ ,  $P = \frac{1}{2}$ .

#### ④ Graph Model.

(i) Vertex, Edge, adjacent, degree... (see in CSC3100 course)

(ii) Examples.

#### e.g. 1 (Maximum-Flow Problems.)

Descriptions & optimization formulation (see in MAT3007)

$$\text{maximize } Z = \sum_j x_{sj}$$

$$\text{subject to } \sum_i x_{ij} = \sum_k x_{jk}, \forall j \in V(G) \setminus \{s, t\}$$

$$0 \leq x_{ij} \leq u_{ij}, \forall i \in A(G)$$

Sol. The Ford-Fulkerson Algorithm:

Step 1: Initialize  $f_c = 0$

Step 2: Find a directed path from  $s$  to  $t$ , if no such path,

stop.  $f_{\max} = \text{current } f_c$

Step 3: compute  $U_{\min} = \text{minimum capacity of all arcs in}$   
the current path.

Step 4: update the residual capacity.  $U_{ij} = U_{ij} - U_{\min}$

(If  $U_{ij} \geq 0$ , remove  $U_{ij}$ )

$$U_{ji} = U_{ji} + U_{\min}$$

Step 5: update  $f_c = f_c + U_{\min}$ .  $\leftarrow$  iterations

#### e.g. 2 (The shortest path problem)

Given a graph  $G = (V(G), E(G))$  & vertex  $s, t$ . length  $c_{ij}$ .

if  $(i, j) \in E(G)$ . Find the shortest path from  $s$  to  $t$  along the edges

Sol. Dijkstra's Algorithm: (idea of DP)

Step 1: Initialize —  $L(s) = 0$ ,  $L(i) = \infty, \forall i \in V(G) \setminus \{s\}$

Step 2: Find the vertex with smallest temporary label

(randomly pick when tie/draw) Make it permanent.

Step 3: Update the temporary label, based on permanent ones

$$\text{as } L_{\text{temp}}(j) = \min_i L(i) + c_{ij}$$



Date: 2021 Dec 5

### e.g.3 (Vertex Cover Problem)

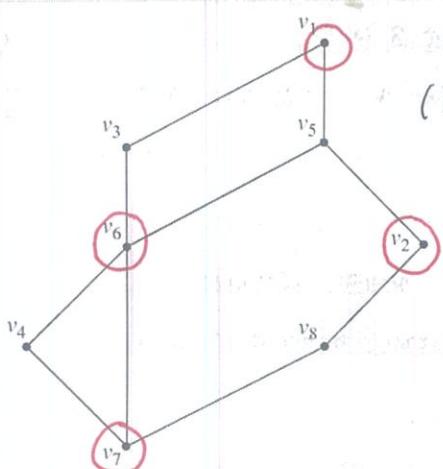
Given  $G = (V(G), E(G))$ , Find  $SC \subseteq V(G)$  as the smallest size s.t. every edge in  $E(G)$  can find its at least 1 end point in  $S$ .

$$\text{minimize } S = \sum_{i \in V(G)} x_i$$

$$\text{subject to } x_i + x_j \geq 1, \forall i, j \in E(G)$$

$$x_i \in \{0, 1\}, \forall i \in V(G)$$

(NP-complete problem!)



### e.g.4 (Euler's seven bridges Problem)

Can a graph  $(V(G), E(G))$  be drawn without repetition in one draw?

(Or find a closed walk traversing every edge exactly once.)

### e.g.5 (Four color problem)

Using only 4 colors to color the vertex of a graph s.t. no 2 vertexes with an edge connected have the same color.

Date: 2021 Dec 14

### • PART V – Statistics Models

#### ① Empirical Model

(i) Defn: infer a function relationship  $y = f(x; \theta)$  based on observations  $(x_i, y_i), i = 1, 2, \dots, n$

traits: { NOT care about how  $x$  affects  $y$

{ NOT imply causal effect of  $x$  on  $y$

{ provide information of  $y$  for  $x$  in the data set

(ii) Error sources: model error (misspecification, missing factors, & data error (truncation, round-off, measurement))

(iii) Model Fitting Criteria: solving minimize distances

$$\min_{\theta \in \Theta} d(\{y_i\}_{i=1}^m, \{f(x_i; \theta)\}_{i=1}^m) \quad (\text{distances between 2 sets})$$

i) Chebyshev Approximation ( $L_\infty$  norm)

$$d(\{y_i\}_{i=1}^m, \{f(x_i; \theta)\}_{i=1}^m) = \max_{1 \leq i \leq m} |y_i - f(x_i; \theta)|$$

ii) Sum of Absolute Deviation ( $L_1$  norm)

$$d(\{y_i\}_{i=1}^m, \{f(x_i; \theta)\}_{i=1}^m) = \sum_{i=1}^m |y_i - f(x_i; \theta)|$$

iii) Least square criterion ( $L_2$  norm)

$$d(\{y_i\}_{i=1}^m, \{f(x_i; \theta)\}_{i=1}^m) = \sum_{i=1}^m (y_i - f(x_i; \theta))^2$$

Relating criteria:  $c_i = |y_i - f_1(x_i)|$ ,  $C_{\max} = \max_{1 \leq i \leq m} c_i$

$$d_C = |y_i - f_2(x_i)|, \quad d_{\max} = \max_{1 \leq i \leq m} d_C$$

$$D = \sqrt{\frac{d_1^2 + \dots + d_m^2}{m}} \leq \sqrt{\frac{c_1^2 + \dots + c_m^2}{m}} \leq C_{\max} \leq d_{\max}$$



• If  $|D - C_{\max}| < |C_{\max} - d_{\max}|$  Chebyshev!

Date: 2021 Dec 14

Linear & generalized linear model ( $\hat{\theta} = (G^T G)^{-1} G^T y$ )

Transformed least-square fit e.g.  $f(x_i, \theta) = a e^{bx} \Rightarrow \sum_{i=1}^m (\ln y_i - (\ln a + b x_i))^2$   
(the transformed minimizer is not necessarily the original one)

(iv) Choose the best model.

\* Quantitative indicators: max deviation;  $\sum |\text{deviation}|$ ;  $\sum \text{deviation}^2$

\* Qualitative indicators: # parameters (less is better); obvious characteristics of the data (trend); patterns of deviations ( $e(x_i)$ )  
good model = no explicit pattern on  $e(x)$ .

(v) High-order polynomial models

Lagrangian form of the polynomial  $P(x) = \sum_{i=1}^m y_i L_i(x)$   
where  $L_k(x) = \prod_{i \neq k} \frac{x - x_i}{x_k - x_i}$ . (Unique polynomial  $P(x)$  of at most degree  $m-1$ )

Advantages: perfect fitting, easily integrated & differentiated.

Disadvantages: too oscillatory

{ possible wrong trends }  $\rightarrow$  in reality, not used  
sensitive to measurement errors

Smoothing - the procedure to use lower order polynomials to mitigate the disadvantages.

Divided Difference Table

Data	First divided difference	Second divided difference
$x_1 \quad y_1$	$y_2 - y_1$	$\frac{y_3 - y_2}{x_3 - x_2} - \frac{y_2 - y_1}{x_2 - x_1}$
$x_2 \quad y_2$	$\frac{y_3 - y_1}{x_3 - x_1}$	
...		

Date: 2021 Dec 14

choose: constant, if 0; linear, if first  $\neq 0$ , second = 0  
& n-order if 0, ..., n  $\neq 0$ ,  $n+1 = 0$ .

(If all not zero, polynomial model may not fit.)

Interpolation (linear, cubic spline)

(Natural spline,  $S_1''(x_1) = S_2''(x_3) = 0$ ; Clamped spline  $S_1'(x_1) = f'(x_1), S_2'(x_3) = f'(x_3)$ )

② Big Data Analysis

(i) Revisiting generalized regression models:

Infer function  $f$  in  $Y = f(X) + \epsilon \rightarrow$  with 0 mean  
 $\Rightarrow f(x) = E(Y|X), Y \in R, X \in R^p$

(Linear) regression assumes:  $f(x) = \beta_0 + \sum_{j=1}^p X_j \beta_j$

(sources of  $X$ : quantitative inputs, basis expansions ( $X_2 = x_1^2, X_3 = x_1^3$ ), dummy variables (e.g. (1,0), (0,1)), interactions ( $X_3 = x_1 x_2$ ))

(ii) The Gauss-Markov Theorem: the least square estimates of parameter  $\beta$  have the smallest variance among all linear unbiased estimates (BLUE property - best linear unbiased estimate)

i.e. If a, suppose  $\hat{\theta} = C^T y$  is unbiased for  $\alpha^T \beta$

$$\text{Var}[\alpha^T \hat{\theta}] \leq \text{Var}[C^T y] \quad E[\hat{\theta}] = \alpha^T \beta$$

Estimator  $\hat{\theta}$  for  $\theta$   $\text{MSE}[\hat{\theta}] = E[(\hat{\theta} - \theta)^2]$

$$\text{mean squared error} = \text{Var}[\hat{\theta}] + (E[\hat{\theta}] - \theta)^2 = \text{Var}[\hat{\theta}] + \text{Bias}$$

Trade-off between variance & bias



Date: 2021 Dec 15

dim of  $\Theta$  ↗ - var ↑, bias ↓ (overfitting)  
 ↘ - var ↓, bias ↑ (underfitting)

## (iii) Model Selection:

\* **Forward-Stepwise**: starts with the intercept ( $\beta_0$ ), then sequentially adding " $\beta_i$ " into model to improve the fitting. ( $\beta_0 \rightarrow \beta_1 \dots \beta_p$  fit n times)

\* **Backward-Stepwise**: starts with the full model, sequentially delete the predictor with the smallest Z-score ( $= \frac{\text{coefficient}}{\text{s.t.d error}}$ )

↳ Only work for  $p < n$  (otherwise  $X^T X$  not pos. def.)

\* **Forward-stragewise**: starts with forward-stepwise with  $\beta_0 = \bar{y}$ , adds the most correlated predictor every time (for some variables). Stop when all variables NOT correlated with residuals  $\Rightarrow$  Fit for  $p < n$

Other popular ways

1) **Ridge Regression**  $\min_{\beta} \|y - X\beta\|^2 \text{ s.t. } \beta^T \beta - \beta_0^2 \leq t$

(w.l.o.g. assume  $\bar{y} = 0$ ,  $\bar{x} = 0$ ,  $\beta \in \mathbb{R}^p$ ,  $\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$ )

Degree of freedom  $df(\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j + \lambda}$ ,  $d_j$ 's are eigenvalues of  $X^T X$ .

2) **LASSO**  $\min_{\beta} \|y - X\beta\|^2 \text{ s.t. } \|\beta\|_1 - |\beta_0| \leq s$  more strict than ridge regression

(solved with proximal gradient descent)

Date: 2021 Dec 15

→ usually 5~10

3) **Cross Validation**: K-fold cross validation:

Every time (break into  $\frac{n}{K}$  size), train the model with  $\frac{n}{K}-1$  data groups / variables, validation & checking errors with the chosen 1 (choose the least test data (mean squared error))

## (iv) Logistic Regression (classification)

$\Rightarrow$  Aim: infer  $E(Y=0|X)$ ,  $E(Y=1|X)$  (linear models not good)

A non-linear transformation  $P(Y=1|X) = \frac{1}{1+exp(X\beta)}$

$\Rightarrow$  Fitting method maximize log-likelihood

$\max \ln \prod_{i=1}^n \left( \frac{1}{1+exp(x_i^T \beta)} \right)^{1-y_i} \left( \frac{exp(x_i^T \beta)}{1+exp(x_i^T \beta)} \right)^{y_i}$  (Gradient Descent Method)

$\Rightarrow$  Decision Boundary  $\ln \left( \frac{P(Y=1|X)}{P(Y=0|X)} \right) = 0 \Leftrightarrow X^T \beta = 0$  (hyperplane)

\* For multiple cases, need K vectors  $\beta_1, \dots, \beta_K$   $P(Y=\beta_k|X) = \frac{exp(X\beta_k)}{1+\sum_{k=1}^K exp(X\beta_k)}$

Decision Boundary:  $\hat{c}(x) = \arg \max_{k \in \{0, \dots, K\}} X^T \hat{\beta}_k$  (with  $\hat{\beta}_0 = 0$ )

\* **L1-regularized Logistic Regression**:  $\max_{\beta} \left( \sum_{i=1}^n y_i x_i^T \beta - \ln(1+exp(x_i^T \beta)) \right) - \lambda \sum_{j=1}^p |\beta_j|$  (LASSO)

## (v) Data Compression &amp; PCA (principal component analysis)

PCA: Give  $x_1, \dots, x_n \in \mathbb{R}^d$  (large  $N$  &  $d$ ), we want

a close representation in  $\mathbb{R}^{d'}$  with  $d' \ll d$ .

i.e.  $\mu + \sum_{j=1}^{d'} \lambda_j v_j = \mu + V_d^T \lambda$  ( $\lambda \in \mathbb{R}^d$ ,  $V_d \in \mathbb{R}^{d \times d'}$ ,  $V_d^T V_d = I$ )

Date: 2021 Dec 15

minimize the reconstruction error  $\min_{\mu, \Phi_k, V_d} \sum_{i=1}^n \|x_i - \mu - V_d \Phi_k(x_i)\|^2$

Another approach (SVD)  $X = UDV^T$ , with  $U_1 > U_2 > \dots$

Optimal  $V_d$  consists of first  $d'$  column of  $V$

$$(X^T = (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x})) \quad X \in \mathbb{R}^{d' \times d}$$

(vi) Data Clustering - partition into groups, pairwise dissimilarities between same cluster < different ~s.

\* Combinatorial Algorithm (no probability model)

\* Mixture Modeling (suppose all data is i.i.d.)

### K-Mean clustering:

1° For a given assignment (cluster)  $C$ , minimize the total cluster variance w.r.t.  $(m_1, \dots, m_K)$  yield means of currently assigned clusters.

$$\bar{x}_s = \arg \min_{m_s} \sum_{i \in S} \|x_i - m_s\|^2$$

2° Given a current set of means  $(m_1, \dots, m_K)$  assign each observation  $C(i) = \arg \min_{1 \leq k \leq K} \|x_i - m_k\|^2$

3° Repeat 1° & 2°, stop when the assignment does NOT change

### EM Algorithm (expectation & maximize)

Density for  $K$  i.i.d classes (Gaussian Distribution)

$$\Phi_k(x) = \frac{1}{(2\pi)^{\frac{D}{2}} \sqrt{|Z_k|}} \exp(-\frac{1}{2} (x - \mu_k)^T Z_k^{-1} (x - \mu_k))$$

1° Take initial guesses on  $\mu_k, Z_k, \pi_k$  (proportion of data in  $K$ )

Date: 2021 Dec 15

2° Expectation step: (Baye's Theorem)

$$\hat{\gamma}_{i,k} = \frac{\pi_k \Phi_k(x_i)}{\sum_{k=1}^K \pi_k \Phi_k(x_i)} \rightarrow \text{likelihood}$$

3° Maximization step:  $\hat{\mu}_k = \frac{\sum_{i=1}^N \hat{\gamma}_{i,k} x_i}{\sum_{i=1}^N \hat{\gamma}_{i,k}}$  weighted average

$$\hat{Z}_k = \frac{\sum_{i=1}^N \hat{\gamma}_{i,k} (x_i - \bar{x})(x_i - \bar{x})^T}{\sum_{i=1}^N \hat{\gamma}_{i,k}}$$

$$\hat{\pi}_k = \frac{\sum_{i=1}^N \hat{\gamma}_{i,k}}{N}$$

weighted covariance estimation

4° Repeat 2° & 3°. Stop when convergence.

MAT3300 Mathematical Modeling: Midterm Exam  
March 25, 2019

Answer the questions in the answer book.

Question:	1	2	3	4	5	6	Total
Points:	10	15	20	20	20	15	100
Score:							

Name: \_\_\_\_\_ Student ID: \_\_\_\_\_ Time limit: 2 hours

1. (10 points) Consider the multiple linear regression model  $\mathbf{y} = \mathbf{ax} + \mathbf{b}$  where  $\mathbf{y}, \mathbf{a}, \mathbf{b} \in \mathbb{R}^d$  and  $x \in \mathbb{R}$ . We regress this model to a data set  $(x_i, \mathbf{y}_i)$ ,  $i = 1, 2, \dots, n$  by the least square criteria, i.e., minimizing  $\sum_{i=1}^n \|\mathbf{y}_i - \mathbf{ax}_i - \mathbf{b}\|_2^2$  where  $\|\mathbf{z}\|_2 = \sqrt{\sum_{j=1}^d z_j^2}$  for  $\mathbf{z} \in \mathbb{R}^n$ . Find the optimal  $\mathbf{a}, \mathbf{b}$  analytically.

To minimize the function  $f(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^n (\mathbf{y}_i - \mathbf{ax}_i - \mathbf{b})^T (\mathbf{y}_i - \mathbf{ax}_i - \mathbf{b})$ , we have to find the station point of the function, which is

$$\begin{aligned} \nabla_{\mathbf{a}} f(\mathbf{a}, \mathbf{b}) &= 0, \quad \Rightarrow \quad \sum_{i=1}^n -2x_i(\mathbf{y}_i - \mathbf{ax}_i - \mathbf{b}) = 0, \\ \nabla_{\mathbf{b}} f(\mathbf{a}, \mathbf{b}) &= 0, \quad \Rightarrow \quad \sum_{i=1}^n -2(\mathbf{y}_i - \mathbf{ax}_i - \mathbf{b}) = 0. \end{aligned}$$

Solve the equations, we have

$$\begin{aligned} \mathbf{a} &= \frac{\sum_{i=1}^n (x_i - \bar{x})(\mathbf{y}_i - \bar{\mathbf{y}})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \\ \mathbf{b} &= \bar{\mathbf{y}} - \mathbf{a}\bar{x}, \end{aligned}$$

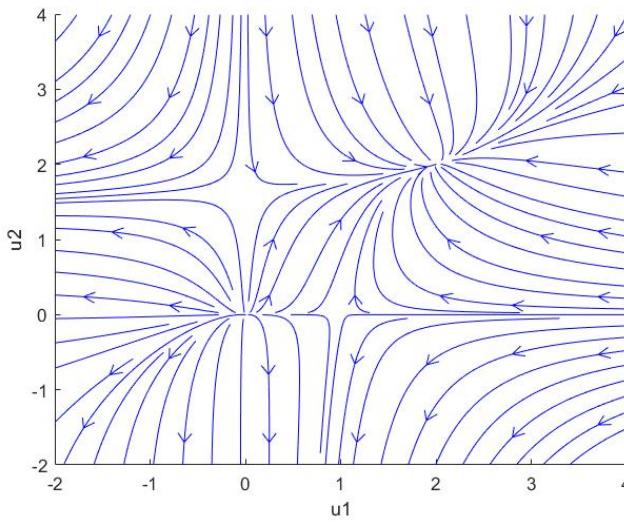
where  $\bar{x} = \sum_{i=1}^n x_i/n$ ,  $\bar{\mathbf{y}} = \sum_{i=1}^n \mathbf{y}_i/n$ .

2. Consider the following population model for two species.

$$\begin{aligned} u'_1(t) &= u_1(t)(1 - u_1(t) + 0.5u_2(t)), \\ u'_2(t) &= u_2(t)(2.5 - 1.5u_2(t) + 0.25u_1(t)). \end{aligned}$$

Here  $u_1(t)$  and  $u_2(t)$  are the populations (in suitable units, say thousands of animals) at time  $t$  of two species that share the same environment.

- (a) (5 points) Find all the equilibrium points of the system.
- (b) (5 points) Identify the stability properties of these equilibrium.
- (c) (5 points) Is it the case that, for all initial populations with  $u_1(0) > 0$  and  $u_2(0) > 0$ , the trajectories all converge to the same equilibrium point? Explain by sketching phase lines.



(a) Let

$$\begin{aligned} u_1(t)(1 - u_1(t) + 0.5u_2(t)) &= 0, \\ u_2(t)(2.5 - 1.5u_2(t) + 0.25u_1(t)) &= 0. \end{aligned}$$

Solve the equation system, we can obtain the equilibrium points are

$$\begin{array}{ll} \textcircled{1} \left\{ \begin{array}{l} u_1(t) = 0, \\ u_2(t) = 0, \end{array} \right. & \textcircled{2} \left\{ \begin{array}{l} u_1(t) = 0, \\ u_2(t) = \frac{5}{3}, \end{array} \right. \quad \textcircled{3} \left\{ \begin{array}{l} u_1(t) = 1, \\ u_2(t) = 0, \end{array} \right. \quad \textcircled{4} \left\{ \begin{array}{l} u_1(t) = 2, \\ u_2(t) = 2, \end{array} \right. \end{array}$$

(b) From the phase diagram, we can easily see the result. The point  $\textcircled{1}\textcircled{2}\textcircled{3}$  are unstable nodes. The point  $\textcircled{4}$  is a stable node.

(c) YES. From the phase diagram, we can see that when  $u_1(0) > 0$  and  $u_2(0) > 0$ , there is only one equilibrium point, which is a stable node. Therefore, the trajectories all converge to the same equilibrium point.

3. Consider the following PDE for  $u(t, x)$ .

$$\begin{cases} \partial_t u = \partial_{xx} u + \gamma \partial_x u, & t > 0, -\infty < x < \infty, \\ u(0, x) = f(x). \end{cases}$$

- (a) (5 points) Find the PDE satisfied by  $v(t, x) = u(t, x)/g(x)$  where  $g(x)$  is a known function.
- (b) (5 points) Find an appropriate  $g(x)$  such that the PDE of  $v(t, x)$  does not involve  $\partial_x v$ .
- (c) (10 points) Solve for  $v(t, x)$  based on the  $g(x)$  you find in (b) and then find  $u(t, x)$  for the original PDE.

(a) Substitute  $u(t, x) = g(x)v(t, x)$  in the original function, we can get

$$\begin{cases} g(x)\partial_t v(t, x) = (\partial_{xx}g(x) + \gamma\partial_xg(x))v(t, x) + (2\partial_xg(x) + \gamma g(x))\partial_xv(t, x) + g(x)\partial_{xx}v(t, x), \\ g(x)v(0, x) = f(x). \end{cases}$$

(b) From the condition, it's easy to see that,

$$2\partial_xg(x) + \gamma g(x) = 0.$$

The solution is  $g(x) = ce^{-\frac{\gamma x}{2}}$ .

(c) From the solution of (b), we have

$$\partial_{xx}g(x) + \gamma\partial_xg(x) = -\frac{\gamma^2}{4}g(x).$$

Then the original PDE can be transformed to

$$\begin{cases} \partial_t v(t, x) = -\frac{\gamma^2}{4}v(t, x) + \partial_{xx}v(t, x), \\ v(0, x) = \frac{1}{c}e^{\frac{\gamma x}{2}}f(x) \end{cases}$$

Here, we do not consider the trivial solution  $g(x) = 0$ . Let  $\hat{v}(t, \omega) = \mathcal{F}[v(t, x)](\omega)$ . Taking Fourier transform on both sides,

$$\begin{cases} \partial_t \hat{v}(t, \omega) = -\frac{\gamma^2}{4}\hat{v}(t, \omega) - \omega^2\hat{v}(t, \omega), \\ v(0, x) = \frac{1}{c}\hat{f}(\omega + i\gamma) \end{cases}$$

Solve the ODE, we have

$$\hat{v}(t, \omega) = \frac{1}{c}\hat{f}(\omega + i\gamma)e^{-(\frac{\gamma^2}{4} + \omega^2)t}$$

By using the IFT,

$$v(t, x) = \frac{1}{2c\sqrt{\pi t}}e^{-\frac{\gamma^2 t}{4}} \int_{-\infty}^{\infty} f(\xi)e^{-(x-\xi)^2/(4t)}d\xi$$

With the solution in (b), we can find that

$$u(t, x) = \frac{1}{2\sqrt{\pi t}}e^{-\frac{\gamma^2 t + 2\gamma x}{4}} \int_{-\infty}^{\infty} f(\xi)e^{-(x-\xi)^2/(4t)}d\xi$$

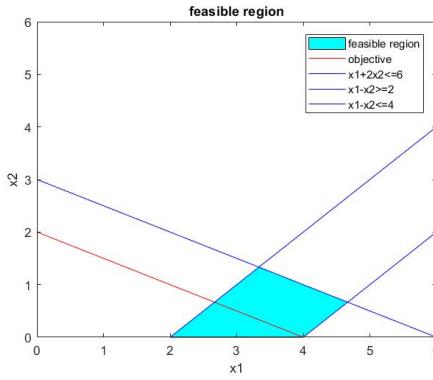
4. Consider the following linear programming (LP) problem.

Minimize  $z = 3x_1 + 6x_2$ , subject to  $x_1 + 2x_2 \leq 6$ ,  $x_1 - x_2 \geq 2$ ,  $x_1 - x_2 \leq 4$ ,  $x_1, x_2 \geq 0$ .

(a) (5 points) Sketch the feasible region of this LP problem and solve it graphically.

(b) (5 points) Enumerate all extreme points of the feasible region and determine the optimal solution.

(c) (10 points) Transform the LP problem into the standard form and solve it with the simplex method.



- (a) As shown in the figure.  
(b) There are four extreme points, which are

$$\begin{cases} x_1 = 2, \\ x_2 = 0, \end{cases} \quad \begin{cases} x_1 = 4, \\ x_2 = 0, \end{cases} \quad \begin{cases} x_1 = \frac{10}{3}, \\ x_2 = \frac{4}{3}, \end{cases} \quad \begin{cases} x_1 = \frac{14}{3}, \\ x_2 = \frac{2}{3}, \end{cases}$$

From the picture, we can see that the first point is the solution,  $z = 6$ .

(c) The standard form is

$$\begin{aligned} \max \quad & z' = -3x_1 - 6x_2 \\ \text{s.t.} \quad & \begin{cases} x_1 + 2x_2 + s_1 = 6, \\ x_1 - x_2 - s_2 = 2, \\ x_1 - x_2 + s_3 = 4, \\ x_1, x_2, s_1, s_2, s_3 \geq 0. \end{cases} \end{aligned}$$

To remedy the predicament, artificial variables are created. For  $\geq$  or  $=$  constraint, add artificial variables. Add sign restriction  $s_i \geq 0$ . The problem can be converted as

$$\begin{aligned} \max \quad & z' = -3x_1 - 6x_2 - Ms_4 \\ \text{s.t.} \quad & \begin{cases} x_1 + 2x_2 + s_1 = 6, \\ x_1 - x_2 - s_2 + s_4 = 2, \\ x_1 - x_2 + s_3 = 4, \\ x_1, x_2, s_1, s_2, s_3, s_4 \geq 0. \end{cases} \end{aligned}$$

where  $M$  denotes a very large positive number.

Build the tableau:

		$x_1$	$x_2$	$s_1$	$s_2$	$s_3$	$s_4$	$b$
		-3	-6	0	0	0	-M	
0	$s_1$	1	2	1	0	0	0	6
$-M$	$s_4$	[1]	-1	0	-1	0	1	2
0	$s_3$	1	-1	0	0	1	0	4
$\sigma$		$-3 + M$	$-6 - M$	0	$-M$	0	0	
0	$s_1$	0	3	1	1	0	-1	4
-3	$x_1$	1	-1	0	-1	0	1	2
0	$s_3$	0	0	0	1	1	-1	2
$\sigma$		0	-9	0	-3	0	$-M + 3$	

From the table, we can see that,  $x_1 = 2$ ,  $s_3 = 2$ ,  $s_1 = 4$ . Hence,  $x_2 = 0$ ,  $s_2 = 0$ . The solution is  $\max z' = -6$ ,  $\min z = 6$ .

5. You have to design a 3-D block as a water storage. The volume of the block has to be at least 9 cubic meters. The base area of the block is at most 6 square meters, while the height of the block is at least 1 and at most 2 meters. Your task is to design the block that satisfies the above conditions and the difference between the base area and the height of the block is maximal.
- (5 points) Formulate the above problem as a nonlinear programming (NLP).
  - (5 points) Show that the height is not 1 at an optimal solution.
  - (5 points) Write down the KKT condition of the NLP in (a).
  - (5 points) Check if the height=1.5 meters, base-area=6 square meters corresponds to an optimal solution.

(a)

$$\begin{aligned} \max \quad & f = s - h \\ \text{s.t.} \quad & \begin{cases} sh \geq 9, \\ s \leq 6, \\ 1 \leq h \leq 2, \\ s \geq 0. \end{cases} \end{aligned}$$

where  $s = ld$ .

(b) If  $h = 1$ , then from the first condition  $ld \geq 9$  which contradicts with the second condition. Hence it is not an optimal solution.

(c)

$$\begin{aligned} \min \quad & f(h, s) = h - s \\ \text{s.t.} \quad & \begin{cases} g_1(h, s) = -sh + 9 \leq 0, \\ g_2(h, s) = s - 6 \leq 0, \\ g_3(h, s) = h - 2 \leq 0, \\ g_4(h, s) = -h + 1 \leq 0, \\ g_5(h, s) = s \geq 0. \end{cases} \end{aligned}$$

The KKT condition is

$$\begin{aligned} \nabla_h f(h, s) + \sum_{i=1}^5 \mu_i \nabla_h g_i(h, s) &= 1 + \mu_1 s + \mu_3 - \mu_4 = 0, \\ \nabla_s f(h, s) + \sum_{i=1}^5 \mu_i \nabla_s g_i(h, s) &= -1 + \mu_1 h + \mu_2 - \mu_5 = 0, \\ \mu_j^* \geq 0, \quad j &= 1, \dots, 5, \\ \mu_j^* g_j(h^*, s^*) &= 0, \quad j = 1, \dots, 5. \end{aligned}$$

**WRONG answer!**

where  $\mu^*$  is the value of  $\mu_1, \dots, \mu_5$  at the optimal,  $g_j(x^*)$  is the condition of the function.

(d) Substitute  $h = 1.5, s = 6$  to the KKT condition, we have,

$$\begin{aligned} \mu_3^* g_3(1.5, 6) = -0.5\mu_3^* &= 0, \quad \Rightarrow \quad \mu_3^* = 0, \\ \mu_4^* g_4(1.5, 6) = -0.5\mu_4^* &= 0, \quad \Rightarrow \quad \mu_4^* = 0. \end{aligned}$$

**WRONG answer!**

Then, from the first KKT condition,  $\mu_1 = -1/6$  which contradicts with  $\mu_1^* \geq 0$ . Hence, it is not an optimal solution.

6. (15 points) Solve the following knapsack problem with dynamic programming. There are four items and the  $i$ th item has value  $v_i$ , size  $w_i$  listed in the table. The knapsack size limit is  $W$ . Find the maximum total value of items that can be put into the knapsack and the corresponding optimal strategy for any positive real  $W$ .

$i$	1	2	3	4
$v_i$	10	40	30	50
$w_i$	5	4	6	3

The problem can be written as

$$\begin{aligned} \max V &= \sum_{i=1}^4 v_i x_i, \\ \text{s.t. } &\sum_{i=1}^4 x_i w_i \leq W, \quad x_i = 0, 1. \end{aligned}$$

Let  $m(i, w)$  denotes the maximum total value can be obtained using the first  $i$  items with capacity limit  $w$ .

$V(i, w)$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$i = 0$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	10	10	10	10	10	10	10	10	10	10	10	10	10	10
2	0	0	0	0	40	40	40	40	40	50	50	50	50	50	50	50	50	50	50
3	0	0	0	0	40	40	40	40	40	50	70	70	70	70	70	80	80	80	80
4	0	0	0	50	50	50	90	90	90	90	100	120	120	120	120	120	120	130	

From the table, we can conclude that

$$V(i, w) = \begin{cases} 0, & 0 \leq w \leq 2, \quad x_i = 0, i = 1, 2, 3, 4, \\ 50, & 3 \leq w \leq 6, \quad x_4 = 1, \\ 90, & 7 \leq w \leq 11, \quad x_4 = x_2 = 1, \\ 100, & w = 12, \quad x_4 = x_2 = x_1 = 1, \\ 120, & 13 \leq w \leq 17, \quad x_4 = x_2 = x_3 = 1, \\ 130, & w \geq 18, \quad x_i = 1, i = 1, 2, 3, 4. \end{cases}$$

$\gamma(x, y, t)$	$\Delta_t = \Delta_x^2 = \frac{1}{4}$	$\Delta_t = \Delta_x^2 = \frac{1}{16}$	$\Delta_t = \Delta_x^2 = \frac{1}{64}$	$\Delta_t = \Delta_x^2 = \frac{1}{256}$
$\sin(xy t + \frac{2\pi}{5})$	$2.4286 \times 10^{-2}$	$7.8620 \times 10^{-3}$	$2.0681 \times 10^{-3}$	$5.3845 \times 10^{-4}$
$\cos(xy t + \frac{1}{100})$	$2.3408 \times 10^{-2}$	$7.4865 \times 10^{-3}$	$1.9660 \times 10^{-3}$	$5.0486 \times 10^{-4}$
$\frac{e^{xyt} - \sin(xy t)}{10}$	$1.2208 \times 10^{-2}$	$3.9143 \times 10^{-3}$	$1.0309 \times 10^{-3}$	$3.4173 \times 10^{-4}$
$\frac{e^{xyt} + \cos(xy t)}{20}$	$1.1952 \times 10^{-2}$	$3.8370 \times 10^{-3}$	$1.0112 \times 10^{-3}$	$3.3466 \times 10^{-4}$
$\frac{e^{xyt} - xy t}{8}$	$1.1091 \times 10^{-2}$	$3.5259 \times 10^{-3}$	$9.2618 \times 10^{-4}$	$3.0723 \times 10^{-4}$
$\frac{e^{xyt} - (xy t)^3}{12}$	$1.2447 \times 10^{-2}$	$4.0124 \times 10^{-3}$	$1.0597 \times 10^{-3}$	$3.4946 \times 10^{-4}$
$e^{xyt - 2.5}$	$1.2492 \times 10^{-2}$	$4.0195 \times 10^{-3}$	$1.0606 \times 10^{-3}$	$3.5039 \times 10^{-4}$
$e^{-xyt - 1.8}$	$1.0361 \times 10^{-2}$	$3.2696 \times 10^{-3}$	$8.5603 \times 10^{-4}$	$2.8827 \times 10^{-4}$
$\frac{\sqrt{xyt} + 1}{15}$	$1.2664 \times 10^{-2}$	$4.1117 \times 10^{-3}$	$1.0883 \times 10^{-3}$	$3.6031 \times 10^{-4}$
$\frac{1 + (xy t)^5}{9}$	$1.1781 \times 10^{-2}$	$3.7736 \times 10^{-3}$	$9.9398 \times 10^{-4}$	$3.2941 \times 10^{-4}$

# Numerical Analysis

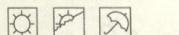
## MAT 4001 Notebook

Youthy WANG

## 通訊錄

## Address List

Name/姓名	<input type="text"/>	<input type="text"/>
<input type="text"/>	Fax	
<input type="text"/>	E-Mail	
<input type="text"/>	QQ	
Name/姓名	<input type="text"/>	<input type="text"/>
<input type="text"/>	Fax	
<input type="text"/>	E-Mail	
<input type="text"/>	QQ	
Name/姓名	<input type="text"/>	<input type="text"/>
<input type="text"/>	Fax	
<input type="text"/>	E-Mail	
<input type="text"/>	QQ	
Name/姓名	<input type="text"/>	<input type="text"/>
<input type="text"/>	Fax	
<input type="text"/>	E-Mail	
<input type="text"/>	QQ	
Name/姓名	<input type="text"/>	<input type="text"/>
<input type="text"/>	Fax	
<input type="text"/>	E-Mail	
<input type="text"/>	QQ	
Name/姓名	<input type="text"/>	<input type="text"/>
<input type="text"/>	Fax	
<input type="text"/>	E-Mail	
<input type="text"/>	QQ	



Mo Tu We Th Fr Sa Su

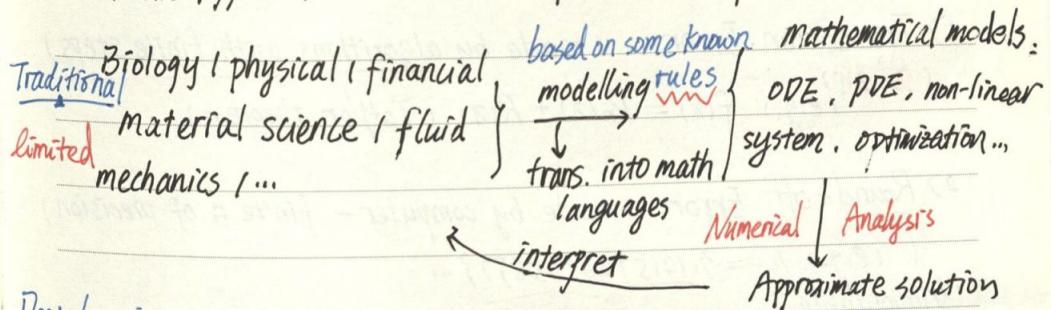
Memo No. 01

Date 2022/Sep/06

## MAT 4001 Numerical Analysis

## • Introduction &amp; Motivations

From 1990's, with the development of computer.



Deep Learning - with Deep Neural Networks (AI for science)

\* 1. Some concrete topics below:

## 1) Numerically Solving Non-linear Equations

(e.g.) Find root for  $f(x) = xe^x - 1$  on  $[-1, 1]$

(e.g.) LambertW function: inverse func of  $f(w) = we^w$ , find values of  $w(z)$ .

## 2) Interpolation

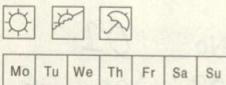


(different interpolations)

## 3) Numerical Integration:

(e.g.)  $\int_0^2 e^{-\frac{t^3}{2}} dt$  } Taylor expansion / make sure convergence  
 $\sum_{n=1}^{\infty} e^{-\frac{n^3}{2}}$

4) Linear System:  $Ax = b$  (or  $Ax = \lambda x$  - eigenvalue prob)



Memo No. 02  
Date 2022 / Sep / 06

## 5) Least Square Approximations:

### \* 2. Computational Errors

1) Truncation Error: (made by algorithms with finite steps)  
(截断)  
e.g.,  $f(x) = P_n(x) + R(x)$  (Taylor theorem)

2) Round-off Error: (made by computer - finite # of precision)  
(e.g.)  $\pi \approx 3.14159265358979\dots$

Unavoidable!

3) Computer Arithmetic: binary system; 64-bit system - double precision.

$\xrightarrow{(-1)^s 2^{c-1023} (1+f)}$   
gives floating point number ( $s-1$  digit,  $c-11$  digits  
 $f-52$  digits)

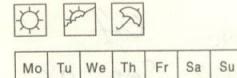
Smallest normalized machine number:  $s=0, c=1;$   
 $(-1)^0 2^{1-1023} (1+0) = 2^{-1022}$  (Any number  $\leq$  that  $2^{-1022}$

positive  
will cause underflow, neglected to be zero

Largest normalized machine number:  $s=0, c=\sum_{i=0}^{52} 2^i, f=\sum_{i=1}^{52} 2^{-i}$

we then get  $2^{1023} (2 - 2^{-52}) \approx 2^{1024}$

4) (Defn) The error that results from replacing by its floating-point form is called round-off error. (regardless of whether chopping/rounding is applied)



Memo No. 03  
Date 2022 / Sep / 08

(Defn) Suppose  $p^*$  is an approximation of  $p$ , we then call  $p-p^*$  actual error,  $|p-p^*|$  absolute error &  $\frac{|p-p^*|}{p}$  relative error provided that  $p \neq 0$ .

(Defn) The number  $p^*$  is said to approximate  $p$  to  $t$  significant digits if  $\frac{|p-p^*|}{p} \leq 5 \times 10^{-t}$  but  $\frac{|p-p^*|}{p} > 5 \times 10^{-t-1}$

Q: How to avoid the loss of accuracy due to rounding/chopping?

Ans. (i) reformulate the problem (without subtraction of nearly equal numbers)

(ii) Horner's Method (Horner's Algorithm)

→ express polynomials into nested forms minimized the total number of calculation

(iii) Avoid large number eat small numbers  
when doing  $f(\sum_{i=1}^n f_i x^i)$ , the way computer calculates, add small numbers firstly.

Replace using floating numbers

### \* 3. Numerical Algorithm

1) Algorithms & Pseudo codes: (skipped)

2) Convergence & Stability:

(i) Rate of Convergence: (Defn) Suppose  $\{B_n\}_{n=1}^{\infty} \rightarrow 0$ , while  $\{a_n\}_{n=1}^{\infty} \rightarrow a \in \mathbb{R}$ , if  $\exists K > 0$ , with  $|a_n - a| \leq K |B_n|$  (For large  $n$ ) we say  $\{a_n\}_{n=1}^{\infty}$  converges to  $a$  with the rate/order of convergence  $O(B_n)$ , written as  $a_n = a + O(B_n)$  (Usually  $B_n = 1/n^p$ )



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 04  
Date 2022 / Sep / 08

Similarly, if  $\lim_{h \rightarrow 0} F(h) = L$ , with  $|F(h) - L| \leq kh^p$  with small positive  $h$ , written as  $F(h) = L + O(h^p)$ .

- (ii) Stability: "small" changes in initial data  $\xrightarrow{\text{stable}}$  "small" changes in results
- (iii) Relationship: convergence - theoretical errors  
stability - computation errors  $\xrightarrow{\text{cannot imply}}$

## • Solve Non-linear Equations Numerically

### \* 1. Solns of Eqns with single variable (i.e. $f(x) = 0$ )

i) Bi-section Method: "concrete algorithm" - see in MAT3007 optimization Notes

#### \* Convergence Analysis:

(i) Suppose that  $f \in C[a,b]$  &  $f(a)f(b) < 0$ . Bisection Method gives sequence  $\{p_n\}_{n=1}^{\infty}$ , approximating zero  $p$  of  $f$ .

We then have that  $|p_n - p| \leq \frac{b-a}{2^n}$   $\xrightarrow{\text{rate of convergence}} \text{linear } O(\frac{1}{2^n})$

(PROOF.)  $a_n - b_n = \frac{b-a}{2^{n-1}}$ , by the algorithm,  $n \in \mathbb{N}$   
Because  $p_n = \frac{a_n + b_n}{2}$  &  $p \in [a_n, b_n]$ , we get  $|p_n - p| \leq \frac{1}{2}(a_n - b_n) \leq \frac{b-a}{2^n}$

(ii) Conservative Error Bound: the actual error may be quite small, compared with quite conservative error bound.

bisection  
only a bound

Good for finding initial point (converge slowly!)



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 05  
Date 2022 / Sep / 13

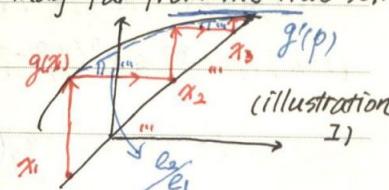
2) Fixed Point Iteration: want to solve  $f(p) = 0$ , to achieve it, construct an auxiliary function  $g(x)$ , s.t.  $[p = g(p)]$  exactly when  $f(p) = 0$ . Then the question becomes to find " $p$ " s.t.  $p = g(p)$ .

(e.g., consider  $f(x) = x - \cos x$ . then  $g(x) = \cos x$   
(NOT unique construction))

(i) (Thm) Existence of fixed point: if  $g \in C[a,b]$ ,  $g(x) \in [a,b] \forall x \in [a,b]$ , then  $\exists c \in [a,b]$  s.t.  $g(c) = c$ .  
(proof is omitted)

(ii) (Thm) Uniqueness of fixed point: Provided that the above holds, with an addition that  $|g'(x)| \leq k < 1$ .  $\forall x \in [a,b]$ . Then the fixed point is unique on  $[a,b]$ . strong assumption, carefully find  $g$ !

\* Iteration Steps for [fixed-point algorithm]:  
Every step let  $x_{n+1} = g(x_n)$ , stop if  $|x_{n+1} - x_n| < \epsilon$   
(or  $|\frac{x_{n+1} - x_n}{x_{n+1}}| < \epsilon$ ). Note that in bi-section method,  $|f(x_n)| < \epsilon$  is NOT a good stop criterion, since  $x_n$  may far from the true soln.



(e.g.) Suppose  $x^3 + 4x^2 - 10 = 0$ , some choice of  $g(x) = x - x^3 - 4x^2 + 10$ ; start from 1.5 complex root  $\sqrt[3]{\frac{10}{x} - 4x}$ ;  $\frac{1}{2}\sqrt{10 - x^3}$ ;  $\sqrt{\frac{10}{x+1}}$   
 $x = \frac{x^3 + 4x^2 - 10}{3x^2 + 8x}$  (Newton's Method)  $\xrightarrow{\text{different convergence rate!}}$



Mo Tu We Th Fr Sa Su

Memo No. 06  
Date 2022 / Sep / 15

\* Convergence Analysis: (from experiment,  $\frac{f_n}{f_{n-1}} \rightarrow f(p)$ ) a sense of convergence analysis

(iii) (Thm) Let  $g \in C[a,b]$  with  $g(x) \in [a,b], \forall x \in [a,b]$

Let  $g(x)$  exist on  $[a,b]$  with  $|g'(x)| \leq k < 1, \forall x \in [a,b]$

At point  $p_0 \in [a,b]$ ,  $\{p_n\}_{n=0}^{\infty}$  is defined by  $p_n = g(p_{n-1}), \forall n \geq 1, n \in \mathbb{N}$   
will converge to the unique fixed point  $p \in [a,b]$ .

(proof.) Consider  $p_{n+1} - p_n = g(p_n) - g(p_{n-1}) \stackrel{\text{MVT}}{=} g'(f_n)[p_n - p_{n-1}]$

$$\therefore |p_{n+1} - p_n| = |g'(f_n)| \cdot |p_n - p_{n-1}| \leq k |p_n - p_{n-1}| \quad (\text{w.l.o.g. with } f_n \in [p_n, p_{n-1}])$$

Thus  $|p_{n+1} - p_n| \leq k^n |p_1 - p_0| = k^n |g(p_0) - p_0| \leq k^n (b-a) \rightarrow 0$

which means  $\{p_n\}_{n=0}^{\infty}$  converges (C.C. or comparison test) (say, to  $p$ )

$$\therefore p = \lim_{n \rightarrow \infty} p_{n+1} = \lim_{n \rightarrow \infty} g(p_n) \stackrel{\text{cts}}{\Rightarrow} g(\lim_{n \rightarrow \infty} p_n) = g(p). \quad \square$$

(iv) (Corollary - convergence rate) With above hypothesis,

$$|p_n - p| \leq \frac{(k^n)}{1-k} |p_1 - p_0|. \quad \begin{array}{l} \text{can be faster than bi-section with good } k \\ \text{only a bound!} \end{array}$$

(proof.)  $|p_n - p| = |g(p_n) - g(p)| \stackrel{\text{M.V.T.}}{=} |g'(f_n)| |p_n - p| \leq k |p_n - p|$

$$\therefore |p_n - p| \leq k^n |p_0 - p| \quad \& \quad |p_n - p| \leq k^{n-1} |p_1 - p|$$

$$(1-k) |p_n - p| \leq k^n (|p_0 - p| - |p_1 - p|) \leq k^n |p_0 - p_1|. \quad \square$$

Back to the previous e.g.  $|\left(\frac{1}{2}\sqrt{10-x^3}\right)'| = \left|\frac{3x^2}{4\sqrt{10-x^3}}\right| \leq \frac{3}{\sqrt{2}}$ , but get "1.7",  
it is true.

$$\left|\left(\frac{\sqrt{10}}{4+x}\right)'\right| = \left|\frac{\sqrt{10}}{2}(4+x)^{-\frac{3}{2}}\right| \leq \frac{\sqrt{2}}{10} \quad \& \quad \left|\left(x - \frac{x^3+4x^2-10}{3x^2+8x}\right)'\right| = \left|\frac{(x^3+4x^2-10)(6x+8)}{(3x^2+8x)^2}\right|$$

$\leq 0.58$ . can be used.

Note that  $g'(p)=0$ !!  
Fastest!



Mo Tu We Th Fr Sa Su

Memo No. 07

Date 2022 / Sep / 15

Note that  $x - x^3 - 4x^2 + 10 \notin [1,2] \& \text{No interval s.t. } |\sqrt{\frac{10}{x}} - 4x|' < 1$ .  
(diverges!)

### 3) Newton's Method:

\* Derivation: Suppose  $f \in C^2[a,b]$ , let  $p_0 \in [a,b]$  be an initial approximation to  $p$  s.t.  $f(p_0) \neq 0$  &  $|p_0 - p|$  is small enough.

Consider  $f(p) = f(p_0) + f'(p_0)(p - p_0) + \frac{f''(\tilde{p})}{2!} (p - p_0)^2$  (Taylor's thm)

where  $\tilde{p}$  between  $p$  &  $p_0$  since  $f(p)=0$ , we get  $0 = f(p_0) + f'(p_0)(p - p_0) + \frac{f''(\tilde{p})}{2!} (p - p_0)^2 \Rightarrow p = p_0 - \frac{f(p_0)}{f'(p_0)}$   
"linearization drop" (Or graphically, use tangent lines to estimate zeros.)

(i) Iteration steps: similar to "bi-section" & "fixed point" with  $p_{n+1} = p_n - \frac{f(p_n)}{f'(p_n)}$ . (Fastest among all fixed-point functions!)

### ii) Convergence Analysis:

(Thm) Let  $f \in C^2[a,b]$ , if  $p \in (a,b)$  is such that  $f(p)=0$  &  $f'(p) \neq 0$   
Then there exists a ( $\delta > 0$ ) s.t. Newton's method generates a sequence  $\{p_n\}_{n=1}^{\infty}$  defined by  $p_n = p_{n-1} - \frac{f(p_{n-1})}{f'(p_{n-1})}$  converges to  $p$  with  $p_0 \in [p-\delta, p+\delta]$ . weak

(proof.) Write  $g(x) = x - \frac{f(x)}{f'(x)}$  (idea - use the results from "fixed point")  
Let  $k \in [0, 1)$ , since  $f'$  is continuous &  $f'(p) \neq 0$ ,  $\exists \delta_1 > 0$  s.t.  $f'(x) \neq 0, \forall x \in V_{\delta_1}(p)$ . Thus,  $g$  is well-defined & continuous on  $V_{\delta_1}(p)$ .

$$g'(x) = \frac{-f(x)f''(x)}{(f'(x))^2}, \quad \forall x \in V_{\delta_1}(p). \quad g'(x) \text{ is continuous, for}$$



Mo Tu We Th Fr Sa Su

Memo No. 08

Date 2022 / Sep / 20

$f \in C^2[a, b]$ .  $\exists \delta_2 > 0$  s.t.  $|g'(x)| \leq k$ ,  $\forall x \in V_{\delta_2}(p)$ , since  $g'(p) = 0$  ( $\delta_2 < \delta_1$ , obviously). For  $x_0 \in [p - \delta_2, p + \delta_2]$ , MVT  $\Rightarrow |g(x_0) - g(p)| \leq |g'(\xi)| |x_0 - p| \leq k \cdot \delta_2 < \delta_2$ , thus  $|g(x_0) - g(p)| = |g(x_0) - p| < \delta_2$ .  
 $\therefore g(x_0) \in V_{\delta_2}(p)$ . Thus  $g(x) \in [p - \delta_2, p + \delta_2]$ . Let  $\delta_2 = \delta$ , with  
 thm for fixed-point iteration,  $\{p_n\}_{n=1}^{\infty} \rightarrow p$ .

\* Remark: seldom used in practice, for we don't know  $\delta$  generally.  
 An initial approx. should be carefully selected.

(iii) Defects with Newton's Method:

Need to know  $f'$  at every approximation  
 $\rightarrow$  generally,  $f'$  is hard & cost many computations

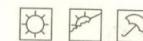
To approximate  $f'(p_n)$  by  $\frac{f(p_{n+1}) - f(p_n)}{p_{n+1} - p_n}$  to get.  
 called Secant Methods  $\rightarrow$  only one func. evaluation  $f(p)$

is required every iteration!

Convergence of secant method is much faster than ... but slight slower than Newton's method.

(iv) The Method of False Position: (盈不足法)  
 (Regula Falsi)

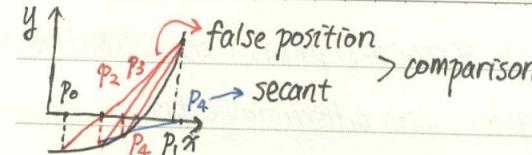
Choose  $p_0$  &  $p_1$  s.t.  $f(p_0)f(p_1) \Rightarrow p_2$  - intersection of secant crossing  $(p_0, f(p_0)), (p_1, f(p_1))$ ;  $(p_2, f(p_2)) \Rightarrow \begin{cases} \text{sgn } f(p_2) \cdot \text{sgn } f(p_1) < 0 \text{ use } p_1, p_2 \Rightarrow p_3 \\ \text{otherwise use } p_0, p_2 \Rightarrow p_3 \end{cases}$



Mo Tu We Th Fr Sa Su

Memo No. 09

Date 2022 / Sep / 20



\* Convergence Rate:

(i) Suppose  $\{p_n\}_{n=1}^{\infty}$  converges to  $p$ , with  $p_n \neq p$ ,  $\forall n \in \mathbb{N}$ . If  $\exists$  constants  $\lambda > 0$  &  $\alpha > 0$  s.t.  $\left[ \lim_{n \rightarrow \infty} \frac{|p_{n+1} - p|}{|p_n - p|^{\alpha}} = \lambda \right]$  generally,  $\lambda < 1$ . Then  $\{p_n\}_{n=1}^{\infty}$  converges to  $p$  of order  $\alpha$  with asymptotic error constant  $\lambda$ . ( $\alpha = 1$ , linearly;  $\alpha = 2$ , quadratically)

(ii) (Defn) A sequence  $\{p_n\}_{n=1}^{\infty}$  is said to be super-linearly convergent to  $p$  if  $\left[ \lim_{n \rightarrow \infty} \frac{|p_{n+1} - p|}{|p_n - p|} = 0 \right]$ .

(Thm) Fixed point iteration converges linearly (obvious to prove)  
 $\left( \lim_{n \rightarrow \infty} \frac{|p_{n+1} - p|}{|p_n - p|} = \lim_{n \rightarrow \infty} |g'(\xi)| = |g'(p)| \right)$ , linear convergence

\* Remark: high-order convergence for fixed-point must have  $g'(p) = 0$ .

(Thm') Let  $p$  is a solution s.t.  $x = g(x)$ . Suppose  $g'(p) = 0$  &  $g''$  continuous,  $|g''(x)| < M$  on an open interval  $I$ ;  $\exists \delta > 0$  s.t.  $\forall p_n \in [p - \delta, p + \delta]$

The sequence given by  $p_{n+1} = g(p_n)$ ,  $n \in \mathbb{N}$  converges at least quadratically, to  $p$ . Moreover,  $\exists N \in \mathbb{N}$  s.t.  $|p_{n+1} - p| < \frac{M}{2} |p_n - p|^2$

(proof.) Taylor's thm:  $p_{n+1} = p + g(p)(p_n - p) + \frac{g''(\xi_n)}{2!} (p_n - p)^2$ ,  $\forall n \geq N$ . Thus  $\frac{|p_{n+1} - p|}{|p_n - p|^2} \leq \frac{|g''(\xi_n)|}{2!} < \frac{M}{2}$   $\xrightarrow{\xi_n \rightarrow p} p$  ( $p_n \rightarrow p$ )



Mo Tu We Th Fr Sa Su

Memo No. 10

Date 2022 / Sep / 22

(Because we can set  $\delta$  s.t.  $|g'(x)| < 1$ ,  $\forall x \in [p-\delta, p+\delta]$  & also  $|g(x)| \in [p-\delta, p+\delta]$ .)

Construction:  $g(x) = x - \phi(x_1 f(x))$ , with differentiable  $\phi(x)$ .

$$\text{need } g'(x) = 1 - (\phi' f + f' \phi). \quad g'(p) = 0 = f(p) \Rightarrow \phi(p) = \frac{1}{f'(p)}$$

[Corollary] With  $f(p) = 0$ ,  $f'(p) \neq 0$  &  $f \in C^2$ . For starting value sufficiently close to  $p$ ,  $\star$  Newton's method converges at least quadratically.

(Thm") Secant Methods converge super-linearly!

$$(proof.) \quad P_{n+1} = p_n - \frac{f(p_n)(p_n - p_{n-1})}{f(p_n) - f(p_{n-1})}. \quad \text{Let } e_n = p_n - p$$

$$\text{Then } e_{n+1} = e_n - \frac{f(p_n)(e_n - e_{n-1})}{f(p_n) - f(p_{n-1})} = \frac{e_n f(p_n) - e_n f(p_{n-1})}{f(p_n) - f(p_{n-1})} \stackrel{\text{MVT}}{=} \frac{e_n}{e_n - e_{n-1}} \frac{f'(e_n) - f'(e_{n-1})}{f'(p_n) - f'(p_{n-1})} \stackrel{\text{MVT}}{=} \frac{e_n}{e_n - e_{n-1}} \frac{f'(e_n) - f'(e_{n-1})}{f'(p_n) - f'(p_{n-1})} \stackrel{\text{L'Hopital Rule}}{\rightarrow} \frac{f''(p)}{f'(p)} \text{ & } \frac{f'(e_n)}{f'(p)} \rightarrow f'(p). \quad \text{Thus, } e_{n+1} \propto e_n e_{n-1}$$

Suppose  $e_{n+1} \propto e_n^p \Rightarrow e_{n+1} \propto e_n^p \Rightarrow p = \frac{1}{p} + 1 \Rightarrow p = \frac{1 + \sqrt{5}}{2} \approx 1.618$ .

## \*2. Solns for Systems:

Bi-section method is hard to apply for systems, thus we generalize Newton's method.

1) For  $2 \times 2$  systems: Let  $F(\vec{x}) = \begin{bmatrix} f(x, y) \\ g(x, y) \end{bmatrix}$ ,  $\vec{x} = \begin{bmatrix} x \\ y \end{bmatrix}$ .

Jacobian matrix  $D_F = \begin{bmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{bmatrix}$ , update  $\vec{x}_{k+1} = \vec{x}_k - D_F^{-1}(\vec{x}_k) F(\vec{x}_k)$

2) For  $n \times n$  systems:  $\vec{x}_{k+1} = \vec{x}_k - D_F^{-1}(\vec{x}_k) F(\vec{x}_k)$ .

Need some further techniques for  $D_F(\vec{x}_k) \Delta \vec{x}_k = F(\vec{x}_k)$ .



Mo Tu We Th Fr Sa Su

Memo No. 11

Date 2022 / Sep / 27

Another way in proofs of convergence of secant method:

$$e_{n+1} = - \left[ \frac{P_n - P_{n-1}}{f(P_n) - f(P_{n-1})} \right] \left[ \frac{f(P_n)e_n - f(P_{n-1})e_{n-1}}{P_n - P_{n-1}} \right] e_n e_{n-1} \quad \text{(use Taylor's thm)} \quad \text{simplify}$$

## \*3. Ways to get complex roots.

1) Müller's Method: (for polynomials with real coefficients, it can converge to a complex root)

(i) Connection with secant method: with  $(p_0, f(p_0)), (p_1, f(p_1))$  — linear approximation (secant method); with  $(p_0, f(p_0)), (p_1, f(p_1)), (p_2, f(p_2))$  — quadratic approximation (Müller's method)

(ii) Derivation: Lagrange's Interpolation Formula

The quadratic curve:  $\sum_{i=0}^2 f(p_i) \frac{\prod_{j \neq i} (x - p_j)}{\prod_{j \neq i} (p_i - p_j)}$ , send it to zero, we get  $p_3 = r_1, r_2$ , we need to choose  $r_1$  or  $r_2$  closer to  $p_2$

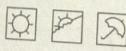
(iii) Algorithms: in order to save computation.

write the curve as  $a(x - p_2)^2 + b(x - p_2) + c$ ; then (to reduce the round-off error)  $p_3 = p_2 + \frac{-2c}{b + \text{sgn}(b)\sqrt{b^2 - 4ac}}$

$$\text{write } a = \frac{f(p_2) - f(p_1)}{p_2 - p_1} - \frac{(p_1 - p_0)}{p_1 - p_0} \quad \& \quad b = \frac{f(p_2) - f(p_1)}{p_2 - p_1} + (p_2 - p_1)a$$

$\star$  Müller's method gives complex roots, for roots of a quadratic curve can be complex.

(iv) Convergence rate: comes from  $p^3 - p^2 - p - 1 = 0$ ,  $p \approx 1.84$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 12  
Date 2022 / Sep / 27

Summary: slow but reliable: bi-section, fixed point  
fast but may NOT converge: Newton's, secant, Müller's

2) Multiple roots:

(i) Drawback of Newton's Method:  $f'(p) \neq 0$  must hold.  
cannot work for  $f(x) = (x-p)^m g(x)$ , with multiplicity  $m$

(ii) Modifications: define  $\mu(x) = \frac{f(x)}{f'(x)}$  ( $f'(p) \neq 0$ )

write  $f = (x-p)^m g(x)$ , then  $\mu(x) = (x-p)^{-m} g(x) + (x-p)^{-1} g'(x)$

Thus,  $\mu(x)$  has a simple zero at  $p$ !

⇒ Do Newton's method for  $\mu(x) \Rightarrow g(x) = x - \frac{f(x) \cdot f'(x)}{f'(x)^2 - f(x) \cdot f''(x)}$

(★ Drawback: require  $f''$  calculations!)

(iii) Analysis: If  $g$  has the required continuity conditions,  
fixed-point iteration applied to  $g$  will be quadratically convergent  
In practice, multiple roots can cause serious round-off error!

[Thm] The function  $f \in C^1[a, b]$  has a simple zero iff  $f(p)=0$   
&  $f'(p) \neq 0$ .

(proof.) Let  $f(x) = (x-p)h(x)$ . then  $f'(p) = h(p) + (p-p)h'(p)$   
 $h(p) \neq 0 \Leftrightarrow f'(p) \neq 0$ .

Conversely,  $f(x) = (x-p)f'(\xi_x)$ .  $f'(p) = f' \lim_{x \rightarrow p} \xi_x = \lim_{x \rightarrow p} f'(\xi_x)$   
 $\neq 0$ . Then, let  $g(x) = f'(\xi_x)$  with  $g(p) \neq 0$ .  
[Corollary] multiplying "m" zero  
⇒  $f = f' = \dots = f^{(m-1)} \text{ at } p = 0$  but  $f^{(m)}(p) \neq 0$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 13  
Date 2022 / Sep / 29

## • Interpolation & Polynomial Approximation

### \*1. Lagrange interpolating polynomials:

(i) Algebraic Polynomials:  $P_n(x) = \sum_{i=0}^n a_i x^i$  with  $n \in \mathbb{N}$  &  $a_i \in \mathbb{R}$

(ii) Algebraic polynomials are important.

[Weierstrass Approximation Thm] (see in MAT2006 Notes)

Suppose that  $f \in C[a, b]$ ,  $\forall \varepsilon > 0$ ,  $\exists$  a polynomial  $P(x)$ ,  
with the property that  $|f(x) - P(x)| < \varepsilon$ .  $\forall x \in [a, b]$ . good approximation

(PROOF.) See later.

### (iii) Degree of n Lagrangian Polynomial:

Addition of simple curves: let  $L_{n,m} = \frac{\prod_{i \neq m} (x-x_i)}{\prod_{i \neq m} (x_m-x_i)}$  (crossing  $(x_i, 0)$  &  
 $(x_m, 1)$ ) Thus, the polynomial is attained as  $[P(x) = \sum_{m=0}^n f(x_m) \cdot L_{n,m}(x)]$

### (iv) Error Analysis:

(Theoretical error bound) Suppose  $x_0, \dots, x_n$  are distinct numbers  
in  $[a, b]$  &  $f \in C^{n+1}[a, b]$ . Then, for every  $x \in [a, b]$ , a number  $\xi(x)$   
between  $x_0, x_1, \dots, x_n$  in  $(a, b)$  exists with

$$[f(x) = P(x) + \frac{f^{(n+1)}(\xi(x))}{(n+1)!} (x-x_0)(x-x_1) \dots (x-x_n)]$$

(proof.) If  $x \in \{x_0, x_1, \dots, x_n\}$  done. Otherwise, let  $g(t) = P(t) - f(t) + [f(x) - f(x_i)] \prod_{i=0}^n \frac{(t-x_i)}{(x-x_i)}$  with fixed  $x_0$ . Then  $g(x_i) = 0$  &  
 $g(x) = 0$ . Thus, by Rolle's thm,  $\exists \xi(x)$  s.t.  $f^{(n+1)}(\xi(x)) = 0$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 14

Date 2022 / Sep / 29

Since we go from  $n+2$  zeros for  $f$  to 1 zero for  $f^{(n+1)}$ , use Rolle's thm  $n+1$  times.  $\therefore 0 = -f^{(n+1)}(\xi) + [p(x) - f(x)](n+1)! \underset{\substack{\text{if } \xi \\ \in (a, b)}}{\underset{\text{finished}}{\approx}} \frac{1}{n+1} \underset{\substack{\text{if } \xi \\ \in (a, b)}}{\underset{\text{finished}}{\approx}} \frac{1}{n+1} (x-a)^{n+1}$

### ★ Proofs of the "WAT"

#### Step 1° Approximate the Absolute value Functions

(i) Taylor expansion with Cauchy's Remainder for  $\sqrt{1-x}$  around  $x=0$

$$\sqrt{1-x} = 1 - \frac{1}{2}x + \sum_{n=2}^N \frac{-(2n-3)!!}{(2n)!!} x^n + \frac{f^{(N+1)}(c)}{N!} (x-c)^N \cdot x$$

Take  $g(x)$  as above, we get on  $[-1, 1]$  (with  $x$  replaced by  $-x^2$ )

$$|1-x| = |\sqrt{1-(1-x^2)} - g(x)| = \left| \frac{f^{(N+1)}(c)}{N!} (1-x^2-c)^N (1-x^2) \right|$$

where  $c \in [0, 1-x^2]$ . RHS  $= \frac{(2N-3)!!}{(2N)!!} \cdot (1-c)^{-\frac{N+1}{2}} (1-x^2-c)^N (1-x^2) \underset{N \rightarrow \infty}{\rightarrow} |E_N(x)|$

$$\lim_{N \rightarrow \infty} \left| \frac{E_{N+1}(x)}{E_N(x)} \right| = \lim_{N \rightarrow \infty} \left| \frac{2N-1}{2N+2} \cdot \frac{(1-x^2-c)^N}{(1-c)^{N+1}} \right| = 1 - \frac{x^2}{1-c} < 1, \forall x \in [-1, 1] \text{ (so)}$$

which means  $\sum_{N=1}^{\infty} E_N(x)$  converges on  $[-1, 0] \cup [0, 1]$ , thus  $E_N(x) \xrightarrow{N \rightarrow \infty} 0$

When  $x=0$ , the question becomes the convergence of  $\sqrt{1-x} = \sum_{n=0}^{\infty} a_n x^n$  (at  $x=1$ ) with  $a_0$  as coefficients above. Cauchy Remainder fails here,

but  $0 = \sum_{n=0}^{\infty} a_n$  also holds, since  $a_k = (-1)^k \left( \frac{1}{k} \right)$ ,  $\forall k \in \mathbb{N}$

( $\sum_{k=0}^{\infty} a_k = (1-1)^{\frac{1}{2}} = 0$ ). Thus  $E_N(x) \xrightarrow{N \rightarrow \infty} 0$ , which means  $\exists$  some  $M \in \mathbb{N}$

s.t.  $|E_N(x)| < \varepsilon, \forall x \in [-1, 1]$ . Then, we can find such  $g(x)$ .

(ii) To generalize  $|1-x| - g(x) | < \varepsilon$  from  $[-1, 1]$  to  $[a, b]$ , we have  $\exists \tilde{g}(x)$  s.t.  $| \frac{x-a}{b-a} t - \tilde{g}(\frac{x-a}{b-a}) | < \frac{\varepsilon}{b-a}$ . Since  $t = \frac{x-a}{b-a} \in [-1, 1]$  Let  $h(x) = \tilde{g}(\frac{x-a}{b-a})$  still a polynomial, thus  $|1-x| - h(x) | < \varepsilon$ .

Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 15

Date 2022 / Sep / 30

Do translation:  $|1-x| - h(x+a) | = |1-x| - g(x) | < \varepsilon, \forall x \in [a, b]$

Step 2° WAT is true on  $[-1, 1]$  for polygons!

(i) Consider ha(x) =  $\frac{1}{2}(1-x+a)+(x-a)$ , with fixed  $a \in [-1, 1]$

ha(x) can be uniformly approximated on  $[-1, 1]$  since  $\exists \tilde{g}(x)$  s.t.

$|1-x| - \tilde{g}(x) | < \frac{1}{n}, \forall x \in [-1, 1]$  (from above). Then, let  $g_n(x) = \frac{1}{2}(1-x-a) + \tilde{g}_n(x)$ , we have  $|ha(x) - g_n(x)| = \frac{1}{2}|1-x-a - \tilde{g}_n(x)| < \frac{1}{2n}$ . Thus,

$\{g_n(x)\}_{n=1}^{\infty} \rightarrow ha(x)$  uniformly. (uniform approximation)

(ii) Generalize to any polygonal  $\phi$ : linear on every subinterval of  $-1 = a_0 < a_1 < \dots < a_n = 1$ . Consider  $\phi(x) - \phi(-1)$  on  $[a_0, a_1]$ , it must be a linear function with zero  $a_0$  & slope  $b_0 \in \mathbb{R}$ . then  $b_0 \cdot ha_0(x)$  can express it wonderfully. Next step consider  $\phi(x) - \phi(-1) - b_0 \cdot ha_0(x)$

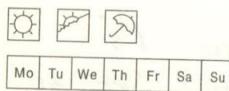
Then, zero becomes  $a_1$ , on  $[a_1, a_2]$ , similarly we let  $b_1, ha_1(x)$  be this

function on  $[a_1, a_2]$ . Sequentially do this. we have.  $[\phi(x) = \phi(-1) + \sum_{j=0}^{n-1} b_j \cdot ha_j(x)]$  general polygonal.

Then without any doubt we can find a polynomial  $g(x) = \sum_{j=0}^{n-1} g_j(x) + \phi(-1)$ , making  $|\phi(x) - g(x)| < \varepsilon, \forall x \in [-1, 1]$ .

#### Step 3° Approximate uniformly continuous function by polygons

Since  $f \in C[a, b]$ ,  $f$  is uniformly continuous on  $[a, b]$ . Fix  $\varepsilon > 0$   $\exists \delta > 0$  s.t.  $|f(x) - f(y)| < \varepsilon$ .  $\forall |x-y| < \delta$ . Then do partitions of  $[a, b]$   $a_0 = a < a_1 < \dots < a_n = b$ , with  $a_{i+1} - a_i \leq \frac{\delta}{2}, \forall i \in \{0, \dots, n-1\}$ .



Memo No. 16  
Date 2022 / Sep / 30

$$\text{Let } \phi(x) = \frac{f(a_i) - f(a_i)}{a_{i+1} - a_i}(x - a_i), \forall x \in [a_i, a_{i+1}], i \in \{0, 1, \dots, N-1\}.$$

Then  $|f(x) - \phi(x)| = |f(x) - (\lambda f(a_i) + (1-\lambda)f(a_{i+1}))| < \varepsilon, \forall x \in [a_i, a_{i+1}], i \in \{0, \dots, N-1\}$ . Thus,  $\forall x \in [a, b], |f(x) - \phi(x)| < \varepsilon$ .

#### Step 4° Proof of WAT.

From above, let  $f(x) \in C[-1, 1]$ , then  $\exists$  such  $\phi_\varepsilon(x) = \phi_\varepsilon(x)$  ( $f(x) \neq 0$ ) s.t.  $|f(x) - \phi_\varepsilon(x)| < \frac{\varepsilon}{2}, \forall x \in [-1, 1]$ . Since  $\phi_\varepsilon(x)$  is a polygonal, thus satisfying WAT,  $\exists g(x) \in R[x]$  (polynomials) s.t.  $|\phi_\varepsilon(x) - g(x)| < \frac{\varepsilon}{2}$ .  $\forall x \in [-1, 1] \Rightarrow |f(x) - g(x)| < \varepsilon, \forall x \in [-1, 1]$ . Then, go to  $[a, b]$  like 1°.  $\square$

(v) Dangerous Error Bound: "Runge's Phenomenon"  $f(x) = (1+25x^2)^{-1}$   
we can find that  $\lim_{n \rightarrow \infty} \left( \max_{-1 \leq x \leq 1} |f(x) - P_n(x)| \right) = \infty$  (worse & worse performance as  $n$  increases.)

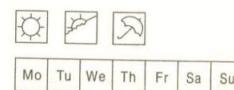
Other: Gibbs' phenomenon in Fourier analysis

## \*2. Improvements of Lagrangian Methods.

### (i) Drawbacks of Lagrange Interpolation:

Difficult to apply the error term (NOT known degree  $n$  until computations done). Waste of calculations (repeat).

(ii) Let  $f$  be defined at  $x_0, \dots, x_n$ , & suppose  $m_1, \dots, m_k$  are  $k$  distinct integers, with  $0 \leq m_i \leq n, \forall i$ .



Memo No. 17  
Date 2022 / Oct / 11

Denote (save, in real-word computations) the Lagrangian agreeing with  $f(x)$  at  $k$  points  $x_{m_1}, \dots, x_{m_k}$  by  $P_{m_1, \dots, m_k}(x)$

$$\text{Then, recursively } P_{m_1, \dots, m_n}(x) = \frac{(x - x_j) P_{m_1, \dots, m_{j-1}}(x)}{x_i - x_j} - (x - x_i) P_{m_1, \dots, m_{i+1}}(x)$$

with  $i \neq j \in \{0, \dots, n\}$ .

(Since  $P_{0, \dots, n}(x)$  passes all  $(x_i, f(x_i))$ ,  $i=0, \dots, n$ , easily verified.)

### (iii) Neville's Method:

Recursion can be applied for finding  $P_n$  with table

$x_0$	$p_0$			
$x_1$	$p_1$	$p_{0,1}$		
$\vdots$	$\vdots$			"Lower Triangular"
$x_n$	$p_n$	$p_{m,n}$	$\dots$	$p_{0,n}$

Const/data  
 $p_0$  linear  
 $p_1 \rightarrow p_{0,1}$  quad  
 $p_2 \rightarrow p_{1,2} \rightarrow p_{0,2} \dots O(n^2)$   
Derivation with table

$$Q_{ij} = \frac{(x - x_{i-1}) Q_{i-1,j-1} - (x - x_i) Q_{i-1,j}}{x_i - x_{i-1}}, \text{ for } i=1, \dots, n \text{ & } j=1, \dots, i.$$

$p_{i,j}$  position of the form/table  
iterated

★ Stopping Criterion:  $|Q_{i,i} - Q_{i-1,i-1}| \leq \varepsilon$

### (iv) A new algebraic representation:

Suppose  $P_n(x)$  is the  $n^{\text{th}}$ -Lagrangian polynomial with  $x_0, \dots, x_n$ . The divide differences of  $f$  w.r.t.  $\{x_i\}_{i=0}^n$  is of the form

$$P_n(x) = a_0 + a_1(x - x_1) + \dots + a_n(x - x_1) \dots (x - x_n)$$

### ★ Aitken's $\Delta^2$ Notation for divided-difference

$$f[x_i] := f(x_i) \quad \text{&} \quad f[x_i, \dots, x_{i+n}] = \frac{f[x_{i+1}, \dots, x_{i+n}] - f[x_i, \dots, x_{i+n-1}]}{x_{i+n} - x_i}$$

$\forall i < n < i+n \in \mathbb{N}, n \in \mathbb{N}$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

First divided-diff.

Memo No. 18

Date 2022/Oct/13

$$\begin{aligned} x_0 &> f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0} \\ x_1 &> f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_1} \\ x_2 &> f[x_1, x_2] = \frac{f[x_2] - f[x_1]}{x_2 - x_1} \\ \vdots & \quad \quad \quad \vdots \end{aligned}$$

Second divided-diff.

(Thm)  $P_n(x)$  can be written as  $P_n(x) = f(x_0) + \sum_{k=1}^n f[x_0, \dots, x_k] \prod_{i=1}^{k-1} (x - x_i)$ (proof.) By defn,  $f[x, x_0, \dots, x_k] = f[x_0, \dots, x_k] + f[x, x_0, \dots, x_k](x - x_k)$  $\forall k \in \{1, \dots, n\}$  &  $f(x) = f(x_0) + f[x, x_0]$ . Thus  $f(x) = P_n(x) + R_n(x)$ where  $R_n(x) = f[x, x_0, \dots, x_n] \prod_{i=1}^n (x - x_i)$ . Then construct  $g(t) =$ 

$$f(t) - P_n(t) - [f(x) - P_n(x)] \frac{\pi(t - x_i)}{\pi(x - x_i)}, \quad g(\pi(x)) = 0 \quad (\text{Rolle's, let } f \in C^m) \\ \Rightarrow f[x, \dots, x_n] = \frac{f^{(m)}(x)}{(m+1)!} \quad \square$$

### \*3. Hermite Interpolation (Osculating Polynomials)

Given  $x_0, \dots, x_n$  distinct in  $[a, b]$  &  $m_0, m_1, \dots, m_n \in \mathbb{Z}^+ \cup \{0\}$ .

$m = \max_{i=0}^n \{m_i\}$ . The Osculating polynomials approximate  $f \in C^m[a, b]$  for  $\{x_i\}_{i=0}^m$ , and has the least degree with the same values as  $f$  and [all its derivatives up to  $m_i$  at  $x_i$ ,  $\forall i$ ].

(i) [Defn] Let  $x_0, \dots, x_n, m_0, \dots, m_n$  &  $m$  be given as above.Suppose  $f \in C^m[a, b]$ , the osculating polynomial approximating  $f$  is  $P(x)$  of least degree s.t.  $\frac{d^k P}{dx^k}(x_i) = \frac{d^k f}{dx^k}(x_i), \forall i = 0, \dots, n$  &  $k = 0, \dots, m_i$ .Note that the degree of  $P$  is at most  $\sum_{i=0}^n m_i + n$ 

Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 19

Date 2022/Oct/13

(Remark) The osculating poly... is the  $m_0^{\text{th}}$  Taylor poly. for  $f$  at  $x_0$  & the  $n^{\text{th}}$  Lagrangian poly. interpolating  $f$  on  $x_0, \dots, x_n$  when  $m_i = 0, \forall i > 0$ . (Generalization of the two)

### (ii) Explicit Form for Hermite Polynomials

when  $m_i = 1, \forall i \in \{0, 1, 2, \dots, n\}$ . the osculating polynomials become Hermite polynomials.(Thm) Let  $f \in C^1[a, b]$ ,  $x_0, \dots, x_n \in [a, b]$  are distinct.The unique polynomial of least degree agreeing with  $f$  &  $f'$  at  $x_0, \dots, x_n$  is the Hermite polynomial of degree at most  $2n+1$ 

$$\text{given by } H_{2n+1}(x) = \sum_{j=0}^n f(x_j) H_{n,j}(x) + \sum_{j=0}^n f'(x_j) \hat{H}_{n,j}(x)$$

$$\text{where } H_{n,j} = [1 - 2(x - x_j) L'_{n,j}(x_j)] L_{n,j}^2(x) \quad \&$$

$$\hat{H}_{n,j} = (x - x_j) L_{n,j}^2(x) \quad (L_{n,j} = \frac{\prod_{i \neq j} (x - x_i)}{\prod_{i \neq j} (x_j - x_i)})$$

$$\text{Moreover, if } f \in C^{2n+2}[a, b], \quad f(x) = H_{2n+1}(x) + \frac{f^{(2n+2)}(x)}{(2n+2)!} \prod_{i=0}^n (x - x_i)^2$$

For some  $g(x) \in [a, b]$ .(proof.) Firstly  $H_{2n+1}(x_k) = f(x_k)$ , &  $H'_{2n+1}(x_k) = f'(x_k), \forall k \in \{0, \dots, n\}$ Suppose  $f \in C^{2n+2}[a, b]$ , then let  $g(t) = f(t) - H_{2n+1}(t) - [f(x) - H_{2n+1}(x)] \cdot \frac{\pi(t - x_i)}{\pi(x - x_i)}$ ,  $g'(t)$  has  $n+1$  number of zeros. (Rolle's  $x_0, \dots, x_n$ )Thus, continuously using Rolle's thm,  $g^{(2n+2)}(x) = 0$  ( $\exists g(x) \in [a, b]$ )  $\Rightarrow \dots$



Mo Tu We Th Fr Sa Su

Memo No. 20

Date 2022 / Oct / 13

(iii) Determine Hermite polynomials by Newton divided-difference formula: setting  $Z_{2i+1} = Z_{2i} = x_i$ ,  $H_i \in [0, \dots, n]$ . with

$$f[Z_{2i+1}] = f[Z_{2i}] = f(x_i) \text{ & } [f[Z_{2i}, Z_{2i+1}] = f'(x_i)]$$

$$\begin{array}{ll} Z & f(Z) \\ Z_0 & f(Z_0) = f(x_0) \\ Z_1 & f(Z_1) = f(x_1) \\ Z_2 & f(Z_2) = f(x_2) \\ Z_3 & f(Z_3) = f(x_3) \\ \vdots & \vdots \\ Z_n & f(Z_n) = f(x_n) \end{array}$$

$\begin{array}{l} 1^{\text{st diff.}} \\ 2^{\text{nd diff.}} \end{array}$

$$f[Z_0, Z_1] = \frac{f(Z_1) - f(Z_0)}{Z_1 - Z_0}$$

$$f[Z_0, Z_1, Z_2] = \frac{f(Z_2) - f(Z_0)}{Z_2 - Z_0}$$

$$\vdots$$

$$\text{Then, } H_{2n+1}(x) = f[Z_0] + \sum_{k=1}^n f[Z_0, \dots, Z_k] \prod_{i=0}^{k-1} (x - Z_i)$$

\* Interpolation Algorithm:

$$\begin{aligned} Z_{2i+1} &= x_i, \quad Q_{2i,0} = f(x_i), \quad Q_{2i+1,0} = f'(x_i), \quad Q_{2i+1,1} = f''(x_i) \\ (\neq 0: \quad Q_{2i+1,1} &= \frac{Q_{2i,0} - Q_{2i+1,0}}{Z_{2i} - Z_{2i+1}} \quad \& \quad i = 2, \dots, 2n+1 \quad j = 2, \dots, n \quad Q_{ij} = \frac{Q_{ij-1} - Q_{i-1,j-1}}{Z_i - Z_{i-1}} \end{aligned}$$

#### \* 4. Cubic-Spline Interpolation

(i) locally, (piece-wise polynomials) approximation

cubic-spline: a special (order 3) piece-wise polynomial approx.  
containing "derivatives" information

(ii) Requirements: cross points  $x_0, x_1, x_2, \dots, x_n$  &  $S'_1(x_0) = S'_n(x_n)$ ,  
 $\dots, S'_{n-1}(x_{n-1}) = S'_n(x_{n-1})$  Totally  $2n + n - 1 = 3n - 1$  equations  
but unknown coefficients  $\geq 3n$ . (for 3 quadratic polys.)

Boundary Conditions:  $S'_1(x_0)$  &  $S'_n(x_n)$  (need at least 2 for symmetry.) — Quadratic NOT good enough.



Mo Tu We Th Fr Sa Su

Memo No. 21

Date 2022 / Oct / 18

(iii) Quadratic = similarly with  $s$  &  $s'$  conditions. with additional  $s''$  conditions (end-points). Totally,  $2n + 2(n-1) = 4n - 2$  equations with 2 more boundary conditions  $\Rightarrow 4n$  unknowns symmetric

{ choice I natural/free boundary  $S''_1(x_0) = S''_n(x_n) = 0$   
choice II clamped boundary  $S'_1(x_0) = f'(x_0)$  &  $S'_n(x_n) = f'(x_n)$

(iv) Construction: (cubic spline)

$$S_j(x) = a_j + b_j(x - x_j) + c_j(x - x_j)^2 + d_j(x - x_j)^3, \quad j = 0, 1, \dots, n-1$$

$$\& \quad S_j(x_j) = f(x_j), \quad S_j(x_{j+1}) = f(x_{j+1}), \quad j = 0, \dots, n-1 \quad (2n \text{ eqns})$$

$$S'_j(x_{j+1}) = S'_{j+1}(x_{j+1}), \quad j = 0, \dots, n-2 \quad (n-1 \text{ eqns})$$

$$S''_j(x_{j+1}) = S''_{j+1}(x_{j+1}), \quad j = 0, \dots, n-2 \quad (n-1 \text{ eqns})$$

Boundary conditions (2 eqns)

Procedures: get  $a_j \Rightarrow a_{j+1} = a_j + b_j h_j + c_j h_j^2 + d_j h_j^3, \quad j = 0, \dots, n-1$

$$\Rightarrow b_{j+1} = b_j + 2c_j h_j + 3d_j h_j^2 \quad \text{back} \quad \Rightarrow c_j = \frac{g_{j+1} - g_j}{3h_j}, \quad j = 0, \dots, n-2, n-1$$

$$(C_n = C_0 = 0) \rightarrow B.C. \quad (j = 0, \dots, n-2) \quad \text{with } h_j = x_{j+1} - x_j, \quad (a_n = f(x_n))$$

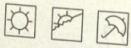
$$\therefore \frac{1}{3}h_j c_{j+1} + \frac{2}{3}(h_{j+1} + h_j) g_j + \frac{1}{3}h_j c_{j+1} = \frac{1}{h_j} (a_{j+1} - a_j) - \frac{1}{h_{j+1}} (a_j - a_{j+1})$$

Let  $d_j = \frac{3}{h_j} (a_{j+1} - a_j) - \frac{3}{h_{j+1}} (a_j - a_{j+1})$ , we get for  $j = 1, 2, \dots, n-1$

$$\begin{bmatrix} 1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 & 0 \\ 0 & 2(h_{j+1}) & h_1 & 0 & \dots & 0 & \dots & 0 & 0 \\ 0 & h_1 & 2(h_{j+2}) & h_2 & \dots & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & C_n \end{bmatrix} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} \in R^{n+1} \quad a_0 = C_0 = 0 = a_n = C_n$$

(Invoke algorithms for

| arr? |  $\sum_{i \neq j} | a_{ij}|$  |  $\rightarrow$  Diagonally Dominant matrix (invertible) solving tridiagonal linear systems!  
(proof.  $A \neq 0 \wedge A^{-1} \neq 0$ )



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 22  
Date 2022 Oct 20

For clamped boundary condition:

(thm) If  $f$  is defined on  $a = x_0 < x_1 < \dots < x_n = b$  & differentiable at  $a$  &  $b$ . Then,  $f$  has a unique clamped spline interpolation on  $S$  with nodes  $x_0, x_1, \dots, x_n$ .

(proof.) Boundary conditions:  $2h_0 c_0 + h_0 c_1 = \frac{3}{h_0} (a_1 - a_0) - 3f(a_1) \stackrel{!}{=} x_0$

Others  $\frac{1}{3}h_j c_{j-1} + \frac{2}{3}(h_{j-1} + h_j)c_j + \frac{1}{3}h_j c_{j+1} = \frac{1}{3}x_j$ ,  $j = 1, \dots, n-1$ .

Define  $c_n = c_{n-1} + 3h_{n-1}d_{n-1}$ ,  $2c_n h_{n-1} + c_{n-1} h_{n-1} = 3f(b) - 3\frac{a_n - a_{n-1}}{h_{n-1}} \stackrel{!}{=} x_n$

Similarly  $\begin{bmatrix} 2h_0 & h_0 & 0 & \dots & 0 & 0 \\ h_0 & 2h_0 & 2h_1 & h_1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & h_{n-1} & 2h_{n-1} & \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ c_n \\ d_n \end{bmatrix} = \begin{bmatrix} a_0 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$ .

Diagonally Dominant Matrix (invertible)  $\Rightarrow$  uniqueness.  $\square$

## Numerical Differentiations & Integrations.

### 1. Numerical Differentiations.

(i) Recall  $f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0+h) - f(x_0)}{h}$ ; we can use  $\left( \frac{f(x_0+h) - f(x_0)}{h} \right)$  to estimate  $f'(x_0)$  numerically. (derive: Taylor's thm  $\downarrow$  Lagrange poly. with  $x_0, x_0+h$ )

Remark:  $\begin{cases} h > 0, (*) \text{ forward difference formula} \\ h < 0, (*) \text{ backward difference formula} \end{cases}$

(ii) General derivative approximation formulas with  $n+1$  points

Suppose  $x_0, \dots, x_n$  are distinct  $n+1$  points on  $[a, b]$  containing  $x_j, j \in \{0, \dots, n\}$   $f \in C^m[a, b]$ , then  $f(x) = P_{n+1}(x) + \frac{\pi(x-x_i)}{(n+1)!} f^{(n+1)}(\xi(x))$



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 23  
Date 2022 Oct 25

while  $P_{n+1}(x)$  is the Lagrangian polynomials with  $x_0, \dots, x_n$

$$\text{Then } f(x) = \sum_{j=0}^{n+1} f(x_j) L_{n+1}^j(x) + \sum_{k=0}^{n+1} \left[ \frac{\pi(x-x_k)}{(n+1)!} f^{(n+1)}(\xi(x)) \right] + \frac{\pi(x-x_k)}{(n+1)!} f^{(n+1)}(\xi(x))$$

$$\text{Take } x = x_j \text{ & we get } f(x_j) = \sum_{k=0}^{n+1} f(x_k) L_{n+1,k}^j(x_j) + \frac{\pi(x-x_j)}{(n+1)!} f^{(n+1)}(\xi(x))$$

(Remark: by Taylor's thm, also can get the formula) more points, less error!

Special cases: (differentiation formula)

$$\text{Three points } x_0, x_1, x_2 \quad \begin{cases} x_0, x_0+h, x_0+2h \text{ (end-point)} \\ f'(x_0) = \frac{1}{2h} (-3f(x_0) + 4f(x_0+h) - f(x_0+2h)) + \frac{h^2}{3} f'''(\xi_0) \end{cases}$$

$$\begin{cases} x_0-h, x_0, x_0+h \text{ (mid-point)} \\ f'(x_0) = \frac{1}{2h} (f(x_0+h) - f(x_0-h)) - \frac{h^2}{6} f'''(\xi_1) \end{cases} \quad \begin{matrix} \nearrow 2\text{nd order} \\ \searrow \text{error} \end{matrix}$$

(iii) High-order derivatives & their approximation formulas

For 2nd-order: Taylor's thm  $\Rightarrow f(x_0+h) = f(x_0) + hf'(x_0) + \frac{h^2}{2} f''(x_0) + O(h^3)$

&  $f(x_0-h) = f(x_0) - hf'(x_0) + \frac{h^2}{2} f''(x_0)$ . Thus,  $\left[ \frac{1}{h^2} [f(x_0+h) + f(x_0-h) - 2f(x_0)] \right]$

$$+ \frac{1}{h^2} O(h^2) = f''(x_0) \quad \text{estimation}$$

(Note that  $f''(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0+h) - f(x_0)}{h} = \lim_{h \rightarrow 0} \frac{f(x_0+h) - f(x_0)}{h} - \lim_{h \rightarrow 0} \frac{f(x_0+h) - f(x_0-h)}{h}$   
 $\uparrow \text{forward} \quad \downarrow \text{backward}$   
 $= \lim_{h \rightarrow 0} \frac{f(x_0+h) + f(x_0-h) - 2f(x_0)}{h^2}$ )

Richardson's Extrapolation (for higher orders)

can be applied whenever an approximation technique is known to have an error term with a predictable form. (e.g.  $O(h), O(h^2) \dots$ )

Suppose  $\exists N_1(h) \xleftarrow{\text{step-size}} s.t.$  the truncation error involved has the form  $M - N_1(h) = \sum_{i=1}^{\infty} k_i h^i \approx k_1 h$ .



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 24

Date 2022 / Oct / 27

The objective of extrapolation is to combine these rather inaccurate approximations  $\Rightarrow$  formula with a higher-order trunc. errors.  
 (e.g.  $O(h)$  to  $O(h^2)$ , construct  $N_2(h) = 2N_1(\frac{h}{2}) - N_1(h)$ .)  
 When truncation error has the form  $\sum_{j=1}^{m-1} k_j h^{aj} + O(h^{am})$  with  $a_1 < a_2 < \dots < a_m$  known (to get  $O(h^{am})$  error), it can be applied.  
 (3-point & 5-point methods from extrapolation.)

Additional Materials: use Bernstein function:  $\sum_{i=0}^n |i| x^i (1-x)^{n-i} / i!$  to prove WAT  
 Last time we use Polygons to estimate uniformly cts functions & then use polynomials to estimate polygons.

Here, a concise proof is given, with Bernstein function. (for  $x \in [0, 1]$ )

$$\left| \sum_{i=0}^n \binom{n}{i} x^i (1-x)^{n-i} f\left(\frac{i}{n}\right) - f(x) \right| \leq \sum_{i=0}^n \binom{n}{i} x^i (1-x)^{n-i} \left| f\left(\frac{i}{n}\right) - f(x) \right|$$

$$(\leq) \sum_{\substack{i: \frac{i}{n} < x \\ i \geq n^{\frac{2}{3}}}} \binom{n}{i} x^i (1-x)^{n-i} \left| f\left(\frac{i}{n}\right) - f(x) \right| + \sum_{\substack{i: \frac{i}{n} > x \\ i \leq n^{\frac{2}{3}}}} \binom{n}{i} x^i (1-x)^{n-i} \left| f\left(\frac{i}{n}\right) - f(x) \right|$$

For the 1<sup>st</sup> term, with uniformly continuous (compact set)  $< \frac{\epsilon}{2}$

For the 2<sup>nd</sup> term, Chebyshev Inequality comes to rescue.

$$\text{2nd term} = \sum_{i: |i-nx| \geq n^{\frac{2}{3}}} P(X=i) |f(\frac{i}{n}) - f(x)| \leq 2M P(|X-nx| \geq n^{\frac{2}{3}})$$

Chebyshev's  $\leq n^{-\frac{4}{3}} \cdot n \pi(1-\pi) \cdot 2M \leq \frac{M}{2} n^{-\frac{2}{3}} < \frac{\epsilon}{2}$ , with overlapping:

Let  $n \geq \max\left\{\left(\frac{M}{\epsilon}\right)^3, \left(\frac{1}{\delta}\right)^3\right\}$ , while  $|f(x) - f(y)| < \frac{\epsilon}{2}$  if  $|x-y| < \delta$ .  $\square$   
 (Generalize to  $[a, b]$ )



Mo	Tu	We	Th	Fr	Sa	Su
----	----	----	----	----	----	----

Memo No. 25

Date 2022 / Oct / 28

Taylor's Series & Its remainders: (w.l.o.g.)  $\rightarrow$  Maclaurin ( $x_0=0$ )

(i) Peano's Form:  $E_N(x) = O(x^{N+1})$

(ii) Lagrangian Form:  $E_N(x) = \frac{f^{(N+1)}(\xi)}{(N+1)!} x^{N+1}$ , where  $\xi$  between 0 &  $x$ .

(iii) Cauchy Form:  $E_N(x) = \frac{f^{(N+1)}(\xi)}{N!} (x-\xi)^N x$ , where  $\xi$  is between 0 &  $x$ .

(iv) Integral Form:  $E_N(x) = \int_0^x \frac{f^{(N+1)}(t)}{N!} (x-t)^N dt$

(proof.) From  $\int_a^b u v^{(n+1)} dx = [uv^{(n)} + (-1)u'v^{(n)} + u^{(n)}(-1)^n] \Big|_a^b + (-1)^{n+1} \int_a^b u^{(n)} v^{(n)} dx$

let  $u = f(t)$ ,  $v = \frac{(x-t)^N}{N!}$ ,  $n=N$ , we get Integral form (remove  $(-1)^{n+1}$ )  $\Rightarrow E_N(x) = \int_0^x \frac{f^{(N+1)}(t)}{N!} (x-t)^N dt$ .

From integral form  $\Rightarrow$  Lagrangian form: by MVT  $E_N(x) = \frac{f^{(N+1)}(\xi)}{N!} \int_0^x dt$

From integral form  $\Rightarrow$  Cauchy form: by MVT  $E_N(x) = \frac{f^{(N+1)}(\xi)}{N!} (x-\xi)^N \int_0^x dt$

Another approach:  $F(t) = f(x) - P_n(t; x)$ , with  $F'(t) = -\frac{f^{(n+1)}(t)}{n!} (x-t)^n$

Cauchy-MVT  $\Rightarrow F(x) - F(0) = \frac{F'(\xi)}{\phi'(\xi)} [\phi(x) - \phi(0)]$ , for some  $\phi \in C^1[0, x]$ .

$-E_N(x) = LHS = RHS = -\frac{f^{(N+1)}(\xi)}{\phi'(\xi) N!} (x-\xi)^N [\phi(x) - \phi(0)]$ . Take  $\phi(x) = x \Rightarrow$

Cauchy form, while taking  $\phi(t) = (x-t)^{N+1}$ .  $\square$

## \*2. Numerical Integrations.

(i) Idea: to approximate  $\int_a^b f(x) dx$ , use Lagrangian Polynomials

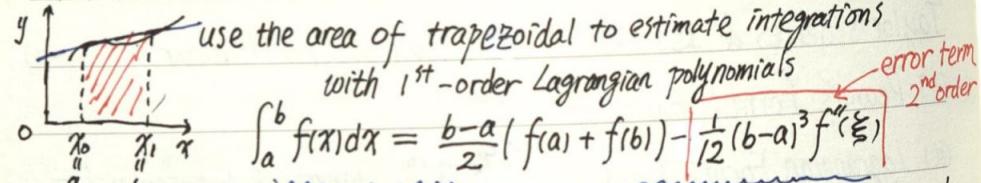
$\int_a^b f(x) dx \equiv \sum_i f(x_i) \int_a^b L_i(x) dx = \sum_i a_i f(x_i)$  Quadrature formula

(ii) The Trapezoidal & Simpson's Rule:

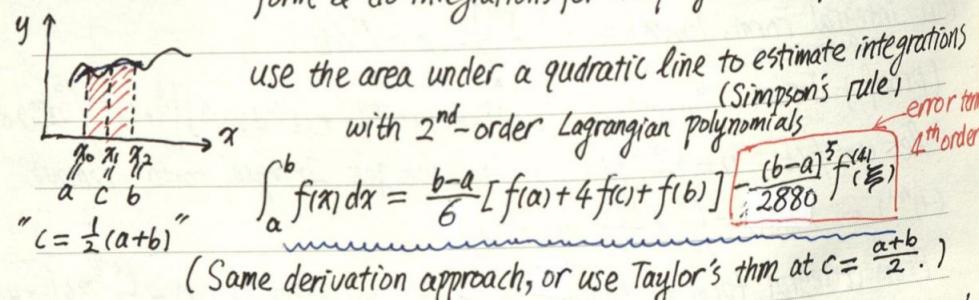


Mo Tu We Th Fr Sa Su

Memo No. 26  
Date 2022 / Nov / 1



(Derivation: Error term in Lagrangian polynomial, MVT in integral form & do integrations for the polynomials.)



(Same derivation approach, or use Taylor's thm at  $c = \frac{a+b}{2}$ .)

(R.m.k. Generalize to  $n+1$  points case,  $\{x_0 + kh\}_{k=0}^n$ ,  $a = x_0$  &  $b = x_0 + nh$ )

$$\text{Trapezoidal rule: } \int_a^b f(x) dx = \frac{h}{2} [(y_0 + y_n) + 2 \sum_{i=1}^{n-1} y_i] - \frac{1}{12} (b-a) h^2 f''(\xi)$$

$$\text{Simpson's rule: } \int_a^b f(x) dx = \frac{h}{3} [(y_0 + y_n) + 4(y_1 + y_3 + \dots + y_{n-1}) + 2(y_2 + \dots + y_{n-2})] - \frac{f^{(4)}(\xi)}{180} (b-a) h^4, \text{ when } n \text{ is even.}$$

(iii) Measuring Precision: the standard derivation of quadrature error formula is based on determining the class of polynomials for which these formulas produce exact solutions.

The degree of accuracy/precision: largest pos.  $n$  that error is 0 for  $\pi^n$

(e.g. Simpson's rule: deg = 3, trapezoidal: deg = 1)

Sun Cloud Moon  
Mo Tu We Th Fr Sa Su

Memo No. 27  
Date 2022 / Nov / 1

#### (iv) Newton-Cotes Formula:

( $n+1$ -pts closed formula:  $\int_a^b f(x) dx \approx \sum_{i=0}^n a_i f(x_i)$ , where

$$x_i = x_0 + ih, h = \frac{b-a}{n} \text{ & } a_i = \int_a^{x_i} L_{n,i}(x) dx$$

(Thm) Error terms: if  $\sum_{i=0}^n a_i f(x_i)$  denotes  $(n+1)$ -pt closed Newton-Cotes formula,  $\exists \xi \in (a, b)$  s.t.  $\int_a^b f(x) dx = \sum_{i=0}^n a_i f(x_i) + \frac{h^{n+3} f^{(n+2)}(\xi)}{(n+2)!} \int_0^n t^{n+2} (t-n)^{n+1} dt$

when  $n$  is even &

$$f \in C^{n+2}[a,b] \quad \int_a^b f(x) dx = \sum_{i=0}^n a_i f(x_i) + \frac{h^{n+2} f^{(n+1)}(\xi)}{(n+2)!} \int_0^n t(t-1) \dots (t-n) dt, \text{ due to "symmetry"!}$$

when  $n$  is odd &  
 $f \in C^{n+1}[a,b]$

(2) ( $n+1$ -pts open formula (excluding end points))

$$x_i = x_0 + ih, \text{ with } x_1 = a, x_{n+1} = b \text{ & } x_0 = a + ih, x_n = b - h, h = \frac{b-a}{n+2}$$

Still  $\int_a^b \approx \sum_{i=0}^n$ , with  $a_i = \int_a^{x_i} L_{n,i}(x) dx$  without  $x_1, x_{n+1}$ .

(Thm) Error terms: if  $\sum_{i=0}^n a_i f(x_i)$  denotes  $(n+1)$ -pt open Newton-Cotes formula,  $\exists \xi \in (a, b)$  s.t.  $\int_a^b f(x) dx = \sum_{i=0}^n a_i f(x_i) + \frac{h^{n+3} f^{(n+2)}(\xi)}{(n+2)!} \int_{-1}^{n+1} t^2 (t-1) \dots (t-n) dt$

$$\text{when } n \text{ is even &} \int_a^b f(x) dx = \sum_{i=0}^n a_i f(x_i) + \frac{h^{n+2} f^{(n+1)}(\xi)}{(n+2)!} \int_{-1}^{n+1} t(t-1) \dots (t-n) dt$$

\* Special Cases:   
closed -  $n=1$  (Trapezoidal),  $n=2$  (Simpson)  
open -  $n=1$  (mid-point)

(R.m.k Trapezoidal with periodic functions) (i.e.  $f(a) = f(b)$ ) can achieve spectral accuracy (error  $\sim e^{-kh}$ ) fast (F.F.T. good for periodic funcs)

$n$ -points mid-point rule:  $h = \frac{|a-b|}{n+2}, x_j = a + (j+1)h, \forall j = -1, \dots, n+1$

$$(n \text{ is even}) \quad \int_a^b f(x) dx \approx 2h \sum_{j=0}^{\frac{n}{2}} f(x_j) + \frac{b-a}{6} h^2 f''(\mu), \mu \in (a, b)$$



Memo No. 28  
Date 2022 / Nov / 3

### (v) Romberg Integrations.

(e.g.) Composite Trapezoidal Rule:  $\int_a^b f = \frac{h}{2} [f(a) + 2\sum_{i=1}^{n-1} f(x_i) + f(b)] + K_1 h^2 + K_2 h^4 + \dots$

where  $K_i$  is a constant depending on  $f^{(2i-1)}(a)$  &  $f^{(2i-1)}(b)$

Apply Richardson's extrapolation:

$$O(h^{2j+2}) = R_{k,j+1} = R_{k,j} + \frac{1}{4^{j+1}} (R_{k,j} - R_{k-1,j}), \quad \forall k=j+1, \dots$$

with "2<sup>j</sup> points" (2 times more points  $\Rightarrow (\frac{h}{2})$  in the error term)

Thus, we get a more precise approximation with error  $O(h^{2j+2})$  for trapezoidal Rule.

(Note: recursive formulas for  $R_{k,1}$ :  $R_{k,1} = \frac{1}{2} R_{k-1,1} + \text{some new points}$ )

\* Algorithm:  $\begin{matrix} x & O(h^2) & O(h^4) & \dots & O(h^{2n}) \\ 1 & R_{1,1} & - & \dots & - \\ 2 & R_{2,1} & R_{2,2} & \dots & - \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n & R_{n,1} & R_{n,2} & \dots & R_{nn} \end{matrix}$  reduce repetitive computations

(vi) Gaussian Quadrature: choose the points (partitions)  $x_0, x_1, \dots, x_n$  in the optimal way rather than uniformly distributed way.

(change to optimization)  $\max_{x_i, c_i} \deg \text{pref} [\int_a^b f - \sum_i c_i f(x_i)]$   $\rightarrow$  maximize degree of precision with  $2n$  parameters to choose

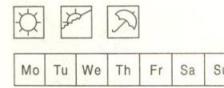
\* Legendre Polynomials:

a collection of polynomials  $\{P_0, P_1, \dots\}$  with properties

(i)  $P_n$  is a monic polynomial of degree  $n$ .

(ii)  $\int_{-1}^1 P(x) P_n(x) dx = 0, \forall p \in \mathbb{R}[x]$  with  $\deg(p) \leq n-1$ .

(Note that  $\langle P_m, P_n \rangle_{L^2} = 0$ , orthogonal basis)



Memo No. 29  
Date 2022 / Nov / 8

By Gram-Schmidt,  $y_{n+1} = x_n - \sum_i \langle x_n, p_i \rangle_{L^2} p_i, P_{n+1} = \frac{y_{n+1}}{\|y_{n+1}\|_{L^2}}$   $\leftarrow$  monic ...

$\therefore \exists P_n(x)$  s.t.  $\deg(P_n) \leq 2n-1, \exists Q, R \in \mathbb{R}[x]$ . s.t.  $Q, R$  has degree less than  $n$ , s.t.  $P_n(x) = Q(x) P_n(x) + R(x) \int_{-1}^1 R(x) = 0$

Choose  $x_i$  s.t.  $P_n(x_i) = 0$ , then  $\sum_i c_i Q_n(x_i) P_n(x_i) = 0$ , then  $\sum_i c_i R(x_i)$  with fixed  $x_i$  — only need to choose  $c_i$ .

$n$  real roots in  $(-1, 1)$

[Theorem] Suppose  $x_1, \dots, x_n$  are the roots of  $n$ th Legendre polynomial

$P_n(x) \& \forall i, c_i = \int_{-1}^1 \prod_{j \neq i} \frac{x-x_j}{x_i-x_j} dx$ . Then, if  $P(x) \in \mathbb{R}[x]$ ,

$\deg(P) < 2n, \int_{-1}^1 P(x) dx = \sum_i c_i P(x_i)$

(proof.) If  $\deg(P) < n, P(x) = \sum_i P(x_i) \prod_{j \neq i} \frac{x-x_j}{x_i-x_j}$  is exact.

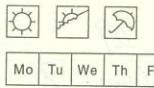
$\int_{-1}^1 P = \sum_i P(x_i) \int_{-1}^1 \prod_{j \neq i} \frac{x-x_j}{x_i-x_j} dx = \sum_i c_i P(x_i)$ . When  $\deg(P) < 2n$

$\exists Q, R$  as above s.t.  $P = Q P_n + R$ . Then  $\int_{-1}^1 P = \sum_i R(x_i) \int_{-1}^1 \dots = \sum_i c_i R(x_i)$

(since  $\deg(R) < n$ ),  $RHS = \sum_i c_i P(x_i)$  since  $P(x_i) = R(x_i)$ .  $\square$

(PACKAGE: Chebfun - MATLAB; FFT (trapezoidal rule for periodic func); NUFFT (Gaussian Quadrature ...).)

Generalized to arbitrary  $[a, b]$ :  $t = \frac{(x-a)+(x-b)}{b-a}$  ( $x = \frac{1}{2}(b-a)(t+a+b)$ )



Mo Tu We Th Fr Sa Su

Memo No. 30  
Date 2022 / Nov / 10

## • Numerical Linear Algebra

### \* 1. Solving Linear Systems ( $Ax = b$ )

(i) When  $A$  is a square matrix

"Gaussian Elimination & Backward Substitution":  $O(n^3)$

\* Complexity of basic operations:  $flop$  unit (floating point operations)

we are interested in rough "flops". of floating point numbers

① "Vec-Vec" operations:  $\vec{a}, \vec{b} \in \mathbb{R}^n$  |  $\vec{a} + \vec{b}$  n flops

② "Mat-Vec" product:  $A \in M_{m \times n}(\mathbb{R})$  |  $C \cdot \vec{a}$  (c  $\in \mathbb{R}$ ) n flops

In general  $2mn$  flops. For special  $A$  | s-sparse 2s flops  
k-banded ( $M_{m \times n}(\mathbb{R})$ )  $2nk$  flops  
permutation  $\sum_{i=1}^r v_i u_i t$   $2r(m+n)$  flops

③ "Mat-Mat" product:  $0$  flops  
 $A \in M_{m \times n}(\mathbb{R}), B \in M_{n \times k}(\mathbb{R})$  (in general)  $2mnk$  flops

s-sparse  $A$ :  $2sp$  flops (less if  $B$  also sparse)

Pseudo Codes: (for Gaussian-Elimination)

$\forall k=1, \dots, n-1 \quad \forall i=k+1, \dots, n : a_{ik} = \frac{a_{ik}}{a_{kk}}$ ;  $\forall j=k+1, \dots, n$   
 $a_{ij} = a_{ij} - a_{ik}a_{kj}$  |  $\rightarrow 3$  loops  
NOT 0!

Pivoting: | Partial ~:  $|a_{\mu, k}| = \max_{k \leq i \leq n} |a_{ik}|$  (the 1<sup>st</sup>  $\mu$ )  
then interchange  $\mu$ <sup>th</sup> &  $k$ <sup>th</sup> rows  
total ~:  $|a_{\mu, \nu}| = \max_{k \leq i, j \leq n} |a_{ij}|$  (the 1<sup>st</sup>  $\mu, \nu$ ) for several,  
then interchange the  $\mu$ <sup>th</sup>,  $k$ <sup>th</sup> rows &  $\nu$ <sup>th</sup>,  $k$ <sup>th</sup> columns  
Keep track of index vectors !!



Mo Tu We Th Fr Sa Su

Memo No. 31  
Date 2022 / Nov / 11

Thus, with forward / backward substitutions  $\Rightarrow$  all  $b$ 's.

(ii) "Matrix Factorizations"

Factorize  $A = \prod_{i=1}^r A_i$ , then compute  $x = A^{-1} A^{-1} \cdots A^{-1} b$ .

(r.m.k ① Usually  $k=2$  or  $3$ ; ② Computing factorization is expensive

③ Applying  $A_i^{-1} \cdots A_k^{-1}$  is cheaper, usually  $n^2$  flops (  $n^3$  flops )  
↳ orthogonal / triangular / diagonal / permutation ....

"LU"-Factorization  $\rightarrow$  equivalent to "Gaussian Elimination"

↳ pivoting using permutation matrix  $P$ :  $A = PLU$  for partial ~

(Note that  $P^T P = I$ ,  $P^T A Q^T = LU$ )  $A = P_1 L U P_2$  for total ~

\* Measuring "Numerical Effort" by flops counting.  
(another criterion — only count ":", "/" since amount of alg.  
time others like "+", "-", loops, etc. is roughly prop. ( $\propto$ ) to flops).

Special structure: positive definite matrices

Cholesky Factorization  $A = LL^T$  | lower triangular

(existence of Cholesky factorization)

Induction:  $M_{K \times K}(\mathbb{R})$  holds, that is  $A_K = L_K L_K^T$ . For  $M_{K+1 \times K+1}(\mathbb{R})$

$A_{K+1} = \left[ \begin{array}{c|c} A_K & b_K \\ \hline b_K^T & a_{KK} \end{array} \right]$  pos. def., construct  $L_{K+1} = \left[ \begin{array}{c|c} L_K & 0 \\ \hline b_K^T L_K^T & \sqrt{a_{KK}} \end{array} \right]$ ,  $L_K^T$  exists  
since  $\forall x \in \mathbb{R}^K$ ,  $L_K^T x \neq 0$  (invertible).  $\rightarrow r_K = -b_K^T A_K^{-1} b_K + a_{KK}$   
 $\Rightarrow x_{K+1}^T A_{K+1} x_{K+1} = a_{KK} - b_K^T A_K^{-1} b_K > 0$ . since let  $x_{K+1} = \left[ \begin{array}{c|c} -A_K^T b_K \\ 1 \end{array} \right]$  positive



Mo Tu We Th Fr Sa Su

Memo No. 32  
Date 2022 / Nov / 15

"Algorithm" to find  $L$ : col 1 -  $l_{11} = \sqrt{a_{11}}$ ,  $l_{i1} = \frac{a_{i1}}{\sqrt{a_{11}}} \quad (\forall i \geq 2)$   
 col 2 -  $l_{22} = \sqrt{a_{22} - l_{11}^2}$ ,  $l_{i2} = \frac{a_{i2} - l_{11}l_{21}}{l_{22}} \quad (\forall i \geq 2)$   
 Generally:  $l_{ik} = \sqrt{a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2}$ ,  $l_{ik} = \frac{a_{ik} - \sum_{j=1}^{k-1} l_{ij}l_{kj}}{l_{kk}} \quad (i \geq k)$   
 (From left  $\rightarrow$  right, up  $\rightarrow$  bottom)  $\frac{6n^3 + \frac{1}{2}n^2 - \frac{2}{3}n}{6}$

\* Numerical Stability, theoretically  $|l_{ij}| \leq \sqrt{a_{ii}}$  bdd  
 (No need to do pivoting!!) numerically stable!

Complete fill-in: Gaussian-elimination change (may) sparse matrix to be dense!  
 (e.g.)  $\begin{bmatrix} x & x & x & x & x \\ x & x & 0 & 0 & 0 \\ x & 0 & x & 0 & 0 \\ x & 0 & 0 & x & 0 \\ x & 0 & 0 & 0 & x \end{bmatrix} = A$ . but with  $A^T$ , fine.) solve by pivoting  
 the tip of large ice-hall.

Banded-matrix:

[Defn] 2 integers  $1 < p, q < n$  s.t.  $a_{ij} = 0$ ,  $\forall i+j \geq p$  or  $j-i \leq -q$   
 band length =  $p+q-1$



(r.m.k special case  $p=q=2$ ) LU factorization:  $T = L_U$   $\xrightarrow{\text{2-banded}}$  matrix

"tridiagonal" matrices) "tridiagonal"

"Algorithm":  $a_{ii-1} = l_{i,i-1}$   $\xrightarrow{\text{Crout factorization}}$   
 $a_{ii} = l_{i,i} \cdot u_{i-1,i} + l_{i,i}$   
 $a_{i,i+1} = l_{i,i} \cdot u_{i,i+1}$

(numerically solving  $u'' = f$  at distinct points  $u''_i = \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2}$  "tridiagonal"  
 $O(n)$   $\rightarrow$  tridiagonal matrix  $f''_i$  matrix)



Mo Tu We Th Fr Sa Su

Memo No. 33  
Date 2022 / Nov / 17

Product property:  $\|Ax\|_p \leq \|A\|_{p,q} \|x\|_q$   
 $\|AB\|_{pq} \leq \|A\|_{p,r} \|B\|_{r,q}$

The infinity norm:  $\|A\|_\infty = \max_{\|x\|_\infty=1} \|Ax\|_\infty = \max_{i=1,\dots,n} \sum_{j=1}^m |a_{ij}|$   
 (In contrast,  $\|A\|_{1,1} = \max_{j=1,\dots,n} \sum_{i=1}^m |a_{ij}|$ ).

Spectral Radius:  $\|A\|_{2,2} = \sqrt{\rho(A^TA)}$   
 $\rho(B) := \max_i |\lambda_i(B)|$   $\xrightarrow{\text{spectral radius of square matrix}}$   
 maximum of all absolute of eigenvalues of  $B$

If  $A$  is symmetric,  $\|A\|_{2,2} = \rho(A)$ .

[Thm]  $\forall \varepsilon > 0$ .  $\exists$  an induced norm  $\|\cdot\|_{p,q}$  s.t.  $\|A\|_{p,q} < \rho(A) + \varepsilon$ ,  $\forall$  matrix  $A$   
 &  $\rho(A) \leq \|A\|_{p,q}$   $\xrightarrow{\rho(A) \text{ is the greatest lower bound of all induced norm}}$   $\rho(A) = \inf_{p,q \in \mathbb{Z}_{\geq 1}^{1 \times m}} \|A\|_{p,q}$

## (iv) Iterative Methods:

Direct methods: Gaussian elimination, "no truncation error" but  
 "expensive computation"

Iterative methods: "effective storage & computation"

① Jacobi's Method:  $x^{(K)} = + \frac{1}{a_{ii}} \left[ - \sum_{j=1, j \neq i}^n (a_{ij} x_j^{(K-1)}) + b_i \right]$   $\xrightarrow{\text{last}}$   
 next  $x_i^{(K)} = + \frac{1}{a_{ii}} \left[ - \sum_{j=1, j \neq i}^n (a_{ij} x_j^{(K-1)}) + b_i \right] \quad \forall K \geq 1$   
 (when  $K = K-1$ , exact)

Matrix form:  $A = D - L - U$   $\xrightarrow{\text{upper}}$   
 $\xrightarrow{\text{lower}}$  with 0 in diagonal

If  $D^{-1}$  exists  $x^{(K)} = D^{-1}(L+U)x^{(K-1)} + D^{-1}b$

Pseudo Codes: initialize  $x^{(0)}$ ,  $T$ ,  $C$   
 $T_j$  for  $j^{\text{th}}$  row  $C_j$  for  $j^{\text{th}}$  row



Mo Tu We Th Fr Sa Su

Memo No. 34  
Date 2022 / Nov / 22

while  $\frac{\|x^{(k)} - x^{(k-1)}\|_\infty}{\|x^{(k)}\|_\infty} \geq \text{tolerance}$

for  $j=1, \dots, n$   $\{ x_j^{(k)} = T_j x^{(k-1)} + c_j \}$  end

end

Another type: Storage savingInitialize  $x^{(0)}$ ,  $T$ ,  $C$ ,  $y = \mathbf{1}$  (auxiliary)while  $\frac{\|x - y\|_\infty}{\|x\|_\infty} \geq \text{tolerance}$  $y = x$ :for  $j=1, \dots, n$   $\{ x_j = T_j \bar{x} + c_j \}$  end

end

② Improvement: the Gauss-Seidel Method

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[ - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^n (a_{ij} x_j^{(k-1)}) + b_i \right]$$

matrix form:  $[x^{(k)} = (D-L)^{-1} U x^{(k-1)} + (D-L)^{-1} b]$ 

general iteration method:

$$x^{(k)} = T x^{(k-1)} + C \quad (T = T_j/T_g, C = c_j/c_g)$$

[Thm] (Convergence result)  $\forall x^{(0)} \in \mathbb{R}^n$ ,  $\{x^{(k)}\}_{k=0}^\infty$  defined by $x^{(k)} = T x^{(k-1)} + C$  converges to a unique solution  $x = T x + C$ iff [spectral radius  $P(T) < 1$ ](proof.) ( $\Leftarrow$ )  $x^{(k)} = T^k x^{(0)} + \left( \sum_{i=0}^{k-1} T^i \right) C$ , take limit:  $x = (I-T)^{-1} C$ .( $\Rightarrow$ ) arbitrary  $\mathbf{z}$  &  $x$  unique s.t.  $x = T x + C$ ,  $x^{(0)} \triangleq x - \mathbf{z}$ Since by iteration  $x - x^{(k)} = T(x - x^{(k-1)}) = \dots = T^k \mathbf{z}$ , take limits $LHS = 0$ , that is  $T^k \mathbf{z} \rightarrow 0$ ,  $\forall \mathbf{z}$ . By equivalent conditions below, done.

Mo Tu We Th Fr Sa Su

Memo No. 35  
Date 2022 / Nov / 24[Defn] Convergent matrix: if  $\lim_{k \rightarrow \infty} (A^k)_{ij} = 0 \quad \forall i, j \leq n$ 

(r.m.k. Equivalent conditions of convergent matrix):

①  $\lim_{k \rightarrow \infty} \|A^k\| = 0$  for some/all induced norms;②  $P(A) < 1$ ; ③  $\lim_{k \rightarrow \infty} A^k x = 0, \forall x \in \mathbb{R}^n$ . )

trivial

(proof.) convergent matrix  $\Rightarrow \| \cdot \|_\infty = 0$  (defn), some  $\Rightarrow$  all special  $x_i = e_i$ .  $\|P(A)\|$  is g.l.b. (equivalence of norms)  $\Rightarrow$  ③ ( $\|A^k x\| \leq \|A^k\| \|x\|$ )  $\Rightarrow$  ② (contradiction)No  $\lambda_i$  s.t.  $\lambda_i \geq 1$ . otherwise  $\lambda_i^k \rightarrow 1$  or  $\infty$ .  $\square$ (Lemma) If  $P(T) < 1$ ,  $(I-T)^{-1} = I + T + T^2 + \dots$  converges.(proof.)  $(I-T)x = (1-\lambda)x$ , i.e.,  $\lambda$  is an ei-value of  $T$ , iff  $1-\lambda$  is an ei-value of  $I-T$ . & 0 is NOT an ei-value of  $I-T$ , ( $I-T$  invertible)

$$S_n = \sum_{k=0}^n T^n \Rightarrow S_n = (I-T)^{-1} (I-T^{n+1})$$

[Coro.] Jacobi & Gauss-Seidel methods for strictly diagonally dominant matrix  $A$  converges to unique  $x$  s.t.  $Ax = b$ .(proof.) Check  $P(T_j)$ ,  $P(T_g)$ 

(v) Relaxation Techniques:

$$\text{(error bound)} \|x - x^{(k)}\| \leq \frac{\|T^k\|}{1 - \|T\|} \|x^{(1)} - x^{(0)}\|$$

construct  $T$ :  $P(T)$  small enough[Defn] (residual vectors)  $\tilde{x}$  is approx. soln to  $Ax = b$ . Then the residual vector  $[r := b - A\tilde{x}]$  for  $r$ .



Mo Tu We Th Fr Sa Su

$$\vec{r}_i^{(k)} = (r_{i,1}^{(k)}, r_{i,2}^{(k)}, \dots, r_{i,n}^{(k)})$$

Memo No. 36

Date 2022 / Nov / 29

residual vector for Gauss-Seidel,  $\vec{x}_i^{(k)} = x_i^{(k-1)} + \frac{r_{ii}^{(k)}}{a_{ii}} \vec{r}_{i+1}^{(k)}$  choose  $x_i^{(k)}$  in  
 $x_i^{(k)}$  s.t.  $r_{i,i+1}^{(k)} = 0$ , (since  $r_{i,i+1}^{(k)} = b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k-1)} - \sum_{j=i+2}^n a_{ij}x_j^{(k-1)}$ )

"makes  $\|\vec{r}_{i+1}^{(k)}\|$  less" better  
 $\star$  relaxation techs:  $x_i^{(k)} = x_i^{(k-1)} + w \frac{r_{ii}^{(k)}}{a_{ii}}$  }  $w \in (0, 1)$  under-relaxation  
 Used to accelerate the convergence of system  
 that are convergent with G-S.

$\|\vec{r}_{i+1}^{(k)}\| = \|\vec{r}_i^{(k)} - w \frac{r_{ii}^{(k)}}{a_{ii}} A_{i+1}\|$ , where  $A_i$  means the  $i^{\text{th}}$  column of  $A$ ,

varying  $w$ ,  $\|\vec{r}_{i+1}^{(k)}\|$  can be reduced (more than when  $w=1$ )

Update method:  $x_i^{(k)} = (1-w)x_i^{(k-1)} + \frac{w}{a_{ii}} [b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k-1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)}]$

(matrix form)  $\vec{x}^{(k)} = T_w \vec{x}^{(k-1)} + C_w$  }  $T_w = (D - wL)^{-1} [(1-w)D + wU]$   
 $C_w = w(D - wL)^{-1} b$   
 (D, L, U defined as previous)

choose optimal  $w$ :

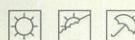
[Thm] (Kahan) If  $a_{ii} \neq 0 \quad \forall i=1, \dots, n$ .  $p(T_w) \geq |w-1|$ .

(Thus, only when  $w < 2$ , SOR method converges.)

(proof.) Since  $T_w = (I - wL)^{-1} [D + w(U - D)]$ ,  $\det(T_w) = \frac{1}{\det(I - wL)} \det(D + w(U - D)) = \frac{1}{\det D} (1-w)^n \cdot \det D = (1-w)^n$ .

( $\det D \neq 0$  since  $a_{ii} \neq 0 \quad \forall i$ ). Then, by Pigeonhole principle,

$p(T_w) = \max_{1 \leq i \leq n} |\lambda_i| \geq \sqrt[n]{(1-w)^n} = |w-1|$ , since  $\prod \lambda_i = \det(T_w) = (1-w)^n$ .  $\square$



Mo Tu We Th Fr Sa Su

symmetric

Memo No. 37

Date 2022 / Nov / 29

[Thm'] If  $A$  is pos. def. &  $0 < w < 2$ . The SOR method converges

& initial point  $\vec{x}^{(0)}$ . (Ostrowski-Reich)

(proof.) Firstly,  $A = \frac{1}{w}(D - wL) - \frac{1}{w}[(1-w)I + wU]$  is positive definite. Thus,  $[A - T_w^T A T_w] = [w(D - wL)^{-1} A]^T [\frac{1}{w}(D - wL) + \frac{1}{w}(D^T - wL^T) - A][w(D - wL)^{-1} A]$  is positive definite since the middle one is  $\frac{2-w}{w}D$ . Thus, let any eigenvector  $\vec{x}$  of  $T_w$   $\Rightarrow \vec{x}^T A \vec{x} - \lambda^2 \vec{x}^T A \vec{x} > 0 \Rightarrow |\lambda| < 1$  eigenvalue of  $T_w$ .

[Thm''] If  $A$  is positive definite & tri-diagonal. Then  $p(T_w) = [p(T_j)]^2 < 1$ . The optimal choice of  $w$  for SOR is  $w = \frac{2}{1 + \sqrt{1 - p(T_j)^2}}$ . With this  $w$ ,  $p(T_w) = w-1$ .

(proof.) sketch - firstly show  $A$  is consistently ordered & 2-cyclic. Under this condition,  $p(T_j^2) = p(T_{j-1}) = p(T_g)$

Moreover  $\lambda_{T_w}(\lambda) = (\lambda + w - 1)^p \prod_{i=1}^p [(\lambda + w - 1)^2 - \lambda w^2 \mu_i^2]$ , while  $\lambda_{T_g}(\mu) = \mu^p \prod_{i=1}^p (\mu^2 - \mu_i^2)$ . (relations between eigenvalues of  $T_w$  &  $T_g$ )

Optimal case (draw curves for analysis):  $\mu_i^2 w^2 - 4w + 4 = 0 \dots$

### (vi) Condition number of a matrix

[Defn] The condition number of an invertible matrix  $A$  is (w.r.t.

$\|\cdot\|$ ),  $K(A) := \|A\| \cdot \|A^{-1}\|$  }  $\begin{cases} K(A) \text{ closed to 1} : \text{well-conditioned} \\ K(A) \text{ far away from 1} : \text{ill-conditioned} \\ (K(A) \gg 1) \end{cases}$



Mo Tu We Th Fr Sa Su

Memo No. 38  
Date 2022 / Dec / 1

\* Estimation of "cond(.)": LAPACK & MATLAB. (never use  $A^{-1}$  sharp to approximate  $\text{cond}(A)$ ). With  $p(A) \leq \|A\| \Rightarrow \text{cond}(A) \geq \frac{\lambda_{\max}}{\lambda_{\min}} \geq 1$

| For orthogonal matrix  $Q$  ( $Q^T Q = I$ ),  $\text{cond}(Q) = 1$  ← as well conditioned  
 ... singular matrix  $X$  -  $\text{cond}(\cdot) \rightarrow \infty$  can never get reliable soln with as possible  
 numerical difficulties iterated method (large relative error)

$$\frac{\|\vec{r}\|}{\|A\| \|A^{-1}\|} \leq \frac{\|\vec{x} - \vec{x}'\|}{\|\vec{x}\|} \leq \|A\| \|A^{-1}\| \frac{\|\vec{r}\|}{\|\vec{b}\|}$$

lower bd      relative error      upper bound

$$K(A)^{-1} = \text{cond}(A)^{-1}$$

$$K(A) = \text{cond}(A)$$

The key to handling ill-conditioning is to avoid it.

(e.g.) Hilbert matrix  $(H)_{ij} = \frac{1}{i+j-1}$ . pos-def. ill-conditioned

(can be derived from polynomial estimation of  $f$  ( $\min_{\alpha_1, \dots, \alpha_n} \|f - \sum_{i=1}^n \alpha_i x^{i-1}\|_2$ )

$\Rightarrow [H_n \cdot \vec{a} = \vec{b}]$ . Transform to well-conditioned use Legendre's polynomial

$$\alpha_k = \sum_{i=1}^n \frac{\langle f, p_i \rangle_{L_2}}{\langle p_i, p_i \rangle_{L_2}} = \frac{\langle f, p_k \rangle_{L_2}}{\langle p_k, p_k \rangle_{L_2}}$$

[Thm] Suppose  $A$  is invertible &  $\|SA\| < \frac{1}{\|A^{-1}\|}$ .

The solution  $\vec{x}$  to  $(A + SA)\vec{x} = \vec{b} + S\vec{b}$  approximates the soln  $\vec{x}$  of  $A\vec{x} = \vec{b}$  with error bdd: perturbation

$$\left[ \frac{\|\vec{x} - \vec{x}'\|}{\|\vec{x}\|} \leq \frac{K(A) \|A\|}{\|A\| - K(A) \|SA\|} \left( \frac{\|S\vec{b}\|}{\|\vec{b}\|} + \frac{\|SA\|}{\|A\|} \right) \right]$$

(proof.)  $\vec{x}' = \vec{x} + S\vec{x}$ ,  $\Rightarrow (I + A^{-1}SA)(\vec{x} + S\vec{x}) = A^{-1}(\vec{b} + S\vec{b})$

$\Rightarrow (I + A^{-1}SA)\vec{x}' = A^{-1}\vec{b} - A^{-1}S\vec{A}\vec{x}$ , thus,  $\|\vec{x}'\| \leq \|(I + A^{-1}SA)^{-1}\| \|A^{-1}(\vec{b} - S\vec{A}\vec{x})\| \leq \|(I + A^{-1}SA)^{-1}\| K(A) (\|S\vec{A}\| \|\vec{x}\| + \|\vec{b}\|) \leq \frac{K(A)}{1 - \|A^{-1}\| \|SA\|} \left( \frac{\|S\vec{A}\|}{\|A\|} + \frac{\|\vec{b}\|}{\|A\|} \right) \|\vec{x}\|$ . Since  $\|A^{-1}S\vec{A}\| \leq \|A^{-1}\| \|SA\| < 1$

$\therefore \|(I + A^{-1}SA)^{-1}\| \leq (1 - \|A^{-1}\| \|SA\|)^{-1}$ . ( $\|(I \pm B)^{-1}\| \leq (1 - \|B\|)^{-1}$ )

need check!



Mo Tu We Th Fr Sa Su

Memo No. 39  
Date 2022 / Dec / 1

## \* 2. Eigen-value Problems

### (i) Q-R factorization:

Least-square revisit = { normal egn  $A^T A \vec{x} = \vec{b} \Rightarrow A \vec{x} - \vec{b} \perp \text{Col}(A)$   
 If  $A^T A \succ 0$ ,  $R(A^T A) = R^2(A)$

[Defn] Let  $A \in M_{m \times n}(\mathbb{R})$ , QR-factorization writes  $A$  as

$[A = QR]$  where  $Q \in M_{m \times m}(\mathbb{R})$ , orthogonal,  $R \in M_{m \times n}(\mathbb{R})$  "upper-triangular"  $\rightarrow$  if  $m \geq n$ ,  $R = \begin{bmatrix} U_{n \times n} \\ 0_{m-n \times n} \end{bmatrix}$  upper-triangular matrix zero-matrix

(e.g. Apply QR-factorization to least-square.

$$\|A\vec{x} - \vec{b}\|_2^2 = \|Q^T(A\vec{x} - \vec{b})\|_2^2 = \|R\vec{x} - Q^T\vec{b}\|_2^2 \Rightarrow \text{becomes } \min_{\vec{x}} \|R\vec{x} - Q^T\vec{b}\|_2^2,$$

where  $R_1 = U_{n \times n}$ ,  $\vec{b}_1 = \vec{b}(1:n)$ , only need to let  $R_1 \vec{x} = Q^T \vec{b}_1$  substitution

### (ii) Computation QR-factorization:

[Defn] Householder Reflection: a matrix  $H_u$  w.r.t.  $U$  s.t.  $[H_u = I - 2\vec{u}\vec{u}^T]$ ,

where  $\|\vec{u}\|_2 = 1$ .

Properties of  $H_u$ :  $H_u$  is symmetric, orthogonal ei-values either 1 or -1

$H_u \vec{v}$  reflects vector  $\vec{v}$  through the hyper-plane orthogonal to  $\vec{u}$

Bzeroing a column:  $A = QR \Leftrightarrow Q^T A = R$  r-make 1<sup>st</sup> col.  $\vec{R}_1 = -\vec{u}$

of  $R$  only with the first entry. Let  $\vec{d}_1 = \vec{A}_1 + \alpha_1 \vec{e}_1$ ,  $H_{\vec{u}_1}$ , where  $\vec{u}_1 = \frac{\vec{d}_1}{\|\vec{d}_1\|_2}$

$$\Rightarrow H_{\vec{u}_1} \cdot \vec{A}_1 = \vec{A}_1 - 2 \frac{\vec{u}_1^T \vec{A}_1}{\|\vec{u}_1\|_2^2} \cdot \vec{u}_1 = \left[ 1 - \frac{2(\vec{A}_1 + \alpha_1 \vec{e}_1)^T \vec{A}_1}{\vec{A}_1^T \vec{A}_1 + 2\alpha_1 \vec{A}_{11} + \alpha_1^2} \right] \vec{A}_1 - \frac{2(\vec{A}_1 + \alpha_1 \vec{e}_1)^T \vec{A}_{11}}{\vec{A}_1^T \vec{A}_1 + 2\alpha_1 \vec{A}_{11} + \alpha_1^2} \vec{e}_1$$

Want:  $H_{\vec{u}_1} \cdot \vec{A}_1 = r \vec{e}_1 \Rightarrow \vec{A}_1^T \vec{A}_1 = \alpha_1^2 \Rightarrow \alpha_1 = \pm \|\vec{A}_1\|_2$  choice  $\alpha_1 = \text{sgn}(\vec{A}_{11}) \|\vec{A}_1\|_2$

Finally,  $\vec{u}_1 = \frac{\vec{A}_1 + \alpha_1 \vec{e}_1}{\sqrt{2\alpha_1(\alpha_1 + \vec{A}_{11})}}$ .

For the 2<sup>nd</sup> - let  $H_{\vec{u}_2}$  be  $(n-1) \times (n-1)$  householder reflection,  $H_{\vec{u}_2} = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 \\ 0 & H_{\vec{u}_2} \end{bmatrix} \dots$



Mo Tu We Th Fr Sa Su

Memo No. 40  
Date 2022 / Dec / 6

Let  $Q = H_1 \cdot H_2 \cdots H_n$ , thus get  $Q^T A = R$ . (Note  $H_i \cdot A \sim O(n^2)$ )

(ii) Find all ei-values of a square matrix

Finding ei-values by computing  $\chi_A(\lambda)$  & finding roots is ill-conditioned.  
Nonetheless - the opposite way works well  
(normalize a polynomial  $p \rightarrow$  construct companion matrix  $A$   
(e.g. "roots" in MATLAB)  $\xrightarrow{(-1)^n \text{ for } x^n}$  compute ei-values of  $A$ )

\* The [power method]: (compute ei-values/ei-vectors with QR method)

Idea: Start with a random vector & keep multiplying it with  $A$

Suppose  $A$  has  $n$  ei-values, denoted as  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$

[1<sup>st</sup>-version] every time multiplying with  $A$ , amplifies the dominant components

Assume ei-vectors of  $A$  can be a basis of  $\mathbb{R}^n$ , denoted as  $\vec{x}_1, \dots, \vec{x}_n$

$\vec{q}_0 = \sum_{i=1}^n \alpha_i \vec{x}_i \Rightarrow q^{(k)} = A^k \vec{q}_0 = \sum_{i=1}^n \alpha_i \lambda_i^k \vec{x}_i$   $\xrightarrow{\text{as } k \uparrow, \lambda_1^k \text{ dominates}}$   
If  $|\lambda_1| > 1$ , tends to  $\infty$

[Normalization] Still with a random  $\vec{q}_0$ , but  $\vec{z}^{(k)} = A \vec{q}^{(k)}$  &  $\vec{q}^{(k)} = \frac{\vec{z}^{(k)}}{\|\vec{z}^{(k)}\|}$   
(often  $\|\cdot\| = \|\cdot\|_\infty$ )

Thus,  $\vec{q}^{(k)}$  converges to ei-vector corresponding to  $\lambda_1$

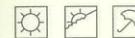
Find  $\lambda$  w.r.t. ei-vector  $\vec{q}$ : least square  $\left[ \min \|A\vec{q} - \lambda \vec{q}\|_2^2 \right]$   $\xrightarrow{\text{root of an}}$  2<sup>nd</sup> order egn

also Rayleigh Quotient  $\lambda = \frac{\vec{q}^T A \vec{q}}{\vec{q}^T \vec{q}}$  (symmetric  $A$ );  $\frac{\vec{q}^T (A + A^T) \vec{q}}{2 \vec{q}^T \vec{q}}$  (generally)

[Issue:  $|\lambda_1| = |\lambda_2| = \dots$ , diff largest ei-value with the same modulus and different ei-vectors]

CASE I:  $\lambda_1 + \lambda_2 = 0$  become oscillations  $\vec{q}^{(k)} = \frac{\lambda_1^k (\alpha_1 \vec{x}_1 + (-1)^k \alpha_2 \vec{x}_2)}{\|\lambda_1^k (\alpha_1 \vec{x}_1 + (-1)^k \alpha_2 \vec{x}_2)\|}$

do  $\vec{q}^{(2k)}, \vec{q}^{(2k+1)}$  converge to diff vectors  $\Rightarrow$  can get  $\vec{x}_1, \vec{x}_2$  (with multiplicity)



Mo Tu We Th Fr Sa Su

Memo No. 41

Date 2022 / Dec / 8

Case II:  $\lambda_1 = \bar{\lambda}_2$ , still can get  $\vec{x}_1, \vec{x}_2$  ...

Note that "floating arithmetic", rounding-off error makes  $\lambda_1$  component non-zero  $\Rightarrow$  whatever initial point, converges to the dominant ( $\lambda_1$ )

[Shift the origin] Let  $B = A - \mu I$ , converges to  $\lambda$ , which is (how to choose?) other ei-values the dominant ei-value of  $B$   
 $\therefore$  get  $\mu + \lambda$   $\rightarrow$  ei-value of  $A$ .

[Inverse Power Method] Find the smallest (absolute value sense) ei-value (3<sup>rd</sup> version)

$$\begin{aligned} \vec{z}^{(k+1)} &= A^{-1} \vec{q}^{(k)} \quad (\text{by solving } A \vec{z}^{(k+1)} = \vec{q}^{(k)}) \\ \vec{q}^{(k+1)} &= \frac{\vec{z}^{(k+1)}}{\|\vec{z}^{(k+1)}\|} \quad (\text{normalization}) \end{aligned}$$