# Lecture 1

- Topics:

  - Regression

  - CART, adaboost, XGB

  - kernel regression, Gaussian Process, SVM

  - Deep learning: MLP, SGD, CNN, RNN, Transformer,
    GPT, BERT, generative: GAN, VAE, Diffusion

  - Reinforcement Learning: MDP, policy, value, Alpha Go,
    policy gradient, Q-learning,
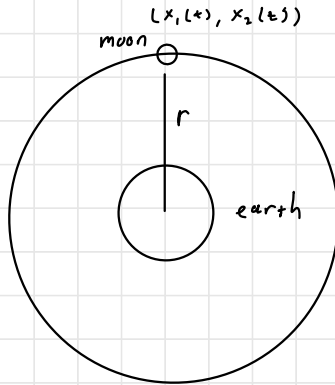    Decision Transformer.

- Course work:

  - bi weekly hw: coding & theoretical
  - coding Python / Py Torch
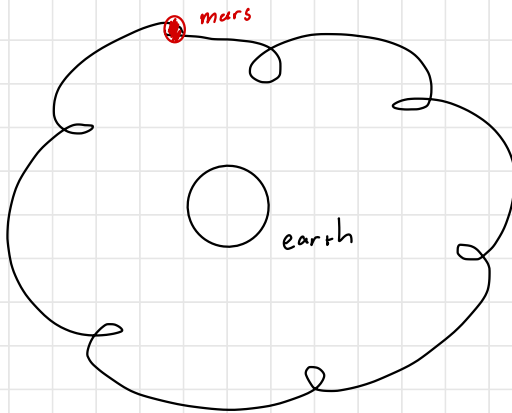
- Machine Learning in ancient time:

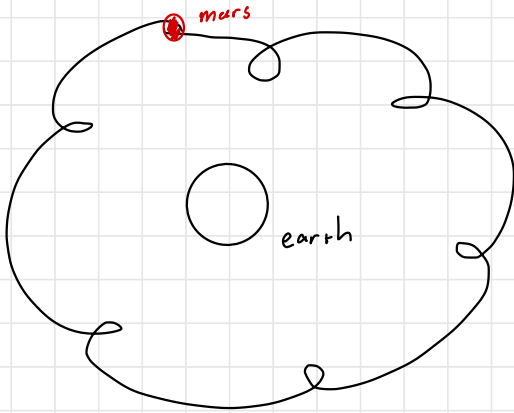  - Astronomy : Observe positions of planets → Predict motion
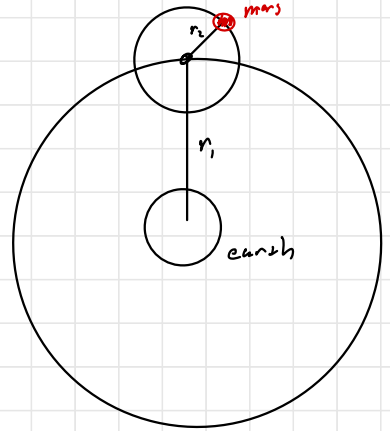
  - P$_{tolemy}$ Epicycle :

$(x_1(t), x_2(t))$

moon

r

earth

$$X(t) = x_1(t) + i x_2(t)$$

$$= r e^{i \omega t}$$

mars

earth

mars

earth

$$Ptolemy \longmapsto$$

$r_2$  mars

$r_1$

earth

$$X(t) = r_1 e^{i w_1 t} + r_2 e^{i w_2 t}$$

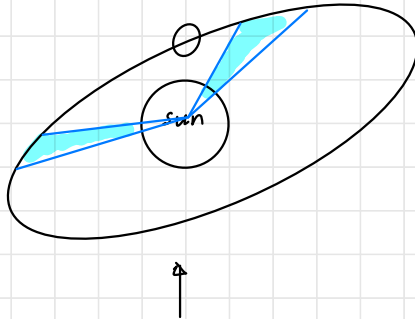- Then for $r_1 > r_2 > \cdots$ we have

$$X(t) = r_1 e^{i w_1 t} + r_2 e^{i w_2 t} + r_3 e^{i w_3 t} + \cdots$$
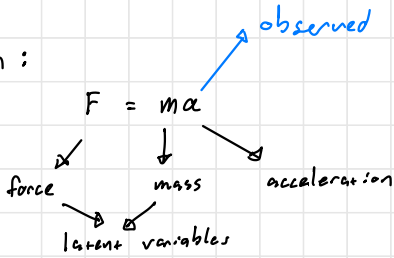
}  boosting

  - We can draw an analogy with the MLP,
    i.e. adding a perceptron on top of a perceptron

- Kepler :
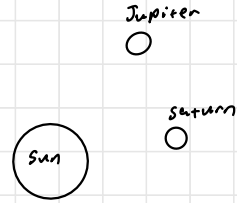


sun

↑

Simpler model under this re-representation
But Ptolemy's model is much more
general

- **Newton:**

  observed

  $$F = ma$$

  force    mass    acceleration

  latent variables

  $$F = G \frac{m_1 m_2}{r^2}$$

  This is more general.
  Can explain 3-Body
  interactions.

  action at a distance, as real as epicycle
  (unreal)

  Jupiter

  Saturn

  Sun

- **Einstein:**

  apple

  ground

  x

  space    actually a straight line,
  but earth as a mass
  curves space time.

  t

  mercury

  Sun

  Einstein can
  model this
  better than
  Newton

  } But just as Ptolemy's epicycles are imaginary
  so is Einstein's notion of space-time

○ Quantum


electron
nucleus
**Wrong**

• Schrödinger :

$$V_{t+\Delta t} = (I + A\Delta t) V_t$$   **: Linear RNN**
(quantised) **(hidden layer)**

$$A = -i H / \hbar$$   Hamiltonian, discrete eigenvalues

even more unreal than epicycle

**embedding**
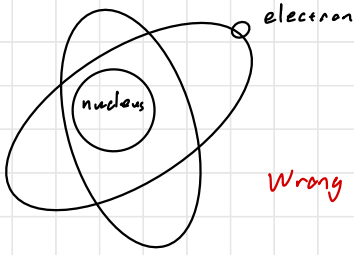(state vector)
thought vector

query vector   Key vector

○ Born :   $$P_t(x) \propto |\langle u(x), V_t \rangle|^2$$   **: emission**
(output layer)

**observed state**

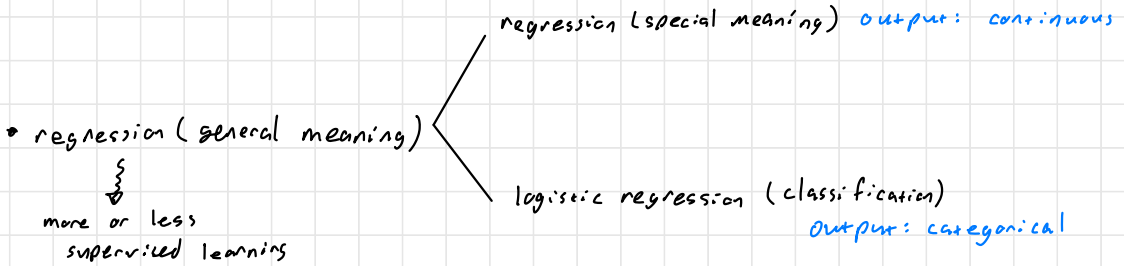$\phi_t(x)$ : wave function, not a physical wave

• Bohr, Heisenberg. Pauli : Copenhagen interpretation

observer outside system, collect data from system
Quantum mechanics is not responsible for explaining "reality"
beyond observed data

- Machine Learning in Modern Time

  - Jan De leeuw : everything is regression

- regression (general meaning)
  $\rightleftharpoons$
  more or less
  supervised learning

  regression (special meaning) output: continuous

  logistic regression (classification)
  output: categorical

|   | input | output |
|---|-------|--------|
| 1 |       |        |
| : |       |        |
| : |       |        |
| : |       |        |
| i | $x_i$ | $y_i$  |
| : |       |        |
| : |       |        |
| : |       |        |
| n |       |        |

- Gauss Paradigm :
  - invented Least Squares Method : Used to predict Ceres
  - Gauss distribution : maximum likelihood
  - Gauss - Markov $\rightarrow$ optimality : LS is the best linear
    unbiased estimator

- Chat - GPT :

  - text :      Prompt

    word

    $x_0, x_1, x_2, \ldots, x_{t-1}, x_z$ ← tokens (50 k of them)

  - learn generative model $P(x_t \mid x_{<t})$ , auto-regressive model.
    
    ↑ output    ↑ input

  - $\max_\theta \; \mathbb{E}_{data} \left[ \log P_\theta (x_t \mid x_{<t}) \right]$    maximum likelihood

    $\underbrace{\hspace{7cm}}$

    The negative of this gives us
    cross-entropy loss.

    Similar to a super parrot ⟶ Oracle

- Diffusion model:

$$X_0 \longrightarrow \cdots \longrightarrow X_t \longrightarrow X_{t+\Delta t} \longrightarrow \cdots \longrightarrow X_t$$

$$X_{t+\Delta t} = X_t + \sigma \sqrt{\Delta t}\, \varepsilon_t \;,\quad \varepsilon_t \sim N(0, I)$$

learn $P(X_t | X_{t+\Delta t}) \sim N(f(X_{t+\Delta t}, t), \sigma^2 \Delta t\, I)$

$$\max_\theta \; \mathbb{E}_{data}\left[\log P_\theta(X_t | X_{t+\Delta t})\right]$$

$\updownarrow$

negative of likelihood gives us
least squares loss.

$$P_\theta(X_t | X_{t+\Delta t}, \text{ text input}) \longrightarrow \text{generating arts}$$