• Correlation Coefficient
= Dimensionless Version of Covariance

$$\rho = E\left[\frac{(X - E[X])}{\sigma_x} \cdot \frac{(Y - E[Y])}{\sigma_y}\right]$$

$$= \frac{cov(X, Y)}{\sigma_x \sigma_y} \longrightarrow +\leq \rho \leq 1.$$

$$\begin{cases} |\rho| = 1 \iff (X - E[X]) = c(Y - E[Y]) \\ \quad\quad \hookrightarrow \text{ linearly Related.} \\ |\rho| = 0 \longrightarrow \text{ independent} \end{cases}$$

---

L.12 Iterated Expectations + Sum of Random Number
          of Random Variables

• Conditional Expectations

$$E[X | Y = y] = \sum_x x \, P_{X|Y}(x|y)$$

• Low of Iterated expectations

$$E[E[X|Y]] = \sum_y E[X|Y=y] \, P_Y(y) = E[X]$$

• $Var(X|Y=y) = E[(X - E[X|Y=y])^2 | Y = y)]$

↳ ⟨law of total variance⟩

$$Var(X) = E[Var(X|Y)] + Var(E[X|Y])$$

"X의 분산은    X의 Y조건부 분산의 Y평균
                    + X의 Y조건부 평균의 Y분산"

• ⟨Section Means⟩ & ⟨Variance⟩.

(ex) Y=1 (10 students) → mean score$|_{Y=1}$ = 90 = E[X|Y=1]
      Y=2 (20 students) → mean score$|_{Y=2}$ = 60 = E[X|Y=2]

then $E(X) = \frac{1}{30} \times (10 \times 90 + 20 \times 60) = 70.$

what happens? $E(X) = \sum_x \sum_y x \, P_Y(Y) P_X(X|Y)$

$$= \sum_y P_y(y) E(X|y)$$

$Var(E[X|Y]) = \sum_y P_y(y) (E[X|Y] - E[E[X|Y]])^2$

$$= \frac{1}{3}(90-70)^2 + \frac{2}{3}(60-70)^2 = 200.$$

$E(Var[X|Y]) = \sum_y P_y(y) Var[X|Y=y] = \frac{50}{3} \longrightarrow \fbox{Var(X)}$

↳ Section별 mean & variance &
각 section의 확률을 알면
total mean & Variance 구할수있다.

• Sum of random number of RN
                                      ↑ independent
↳ N : number of stores visited

$$\begin{cases} X_i : \text{money spent in store } i \\ (\text{assume i.i.d}) \\ \quad \hookrightarrow \text{independent & identically} \\ \quad\quad\quad\quad\quad\quad\quad \text{distributed} \end{cases}$$

↳ $Y = \sum_i^N X_i$

then $E[Y|N=n] = n \, E[X]$
         $E[Y|N] = N E[X]$
⟹ $E[Y] = E[E[Y|N]]$
                 $= E[N E[X]]$
                 $= \underline{E[N] E[X]}$
         ↑ expectation of          ↑ expectation of
         random number           i.i.d. RV X
         from 1 to N

• Variance of sum of random number of
                                           independent R.V
$$Var[Y] = E(Var(Y|N)) + Var(E(Y|N))$$
$$= E[N] Var[X] + E[X]^2 Var(N)$$

---

L.13 Bernoulli Process (Discrete) → P
L.14 Poisson Process (Continuous) → P(k, z)
L.15          ↳memoryless                $\frac{(\lambda z)^k e^{-\lambda z}}{k!}$

L.16~ Markov Chain
      ⟹ with memory / dependence
                          across time.

# (Lec16~18) Markov Processes

new state = $f$(oldstate, noise)

- Finite Markov Chain.
  - $X_n$: state after $n$ transitions
    belong to a finite set $\{1, \ldots, m\}$
    - $X_0$ is either given or random.
  - Markov Property / Assumption:
    - $P_{ij} = P(X_{n+1} = j \mid X_n = i)$
      $\underset{\uparrow}{=} P(X_{n+1} = j \mid X_n = i, X_{n-1}, \ldots X_0)$
      assumption that past doesn't matter.

    $\Rightarrow$ $\Bigg($ possible states
      transition
      probability for each transition

- n-step transition probability
  - $\gamma_{ij}(n) = P(X_n = j \mid X_0 = i)$
    $= \sum_{k=1}^{m} \gamma_{ik}(n-1) P_{kj}$ ; key recursion
  - With random initial state
    $P(X_n = j) = \sum_{i=1}^{m} P(X_0 = i) \gamma_{ij}(n)$

- Generic Convergence Questions?
  $\Rightarrow$ - Does $\gamma_{ij}(n)$ convergence to something?
    ; steady, oscillating, ....
    - Does the limit depend on initial state?

- Recurrent vs Transient States.

- Periodic states

- Steady-state probability

- Visit frequency interpretation

- The phone company problem.

  phone line
  number
  $\bigcirc$ ≡ $\textcircled{B}$

  - Calls originates a Poisson process $\lambda$ (rate)
    $\hookrightarrow$ call duration($\mu$) exponentially distributed.
    B lines available

  $\Rightarrow$ Discrete time intervals of small length $\delta$

  $i\mu\delta$

  - Balance equations; $\lambda \pi_{i-1} = i\mu\pi_i$
    (Finding steady states by $n \to \infty$ for key recursion)
    $\Rightarrow \pi_i = \pi_0 \dfrac{\lambda^i}{\mu^i i!}$  $\pi_0 = \dfrac{1}{\sum_{i=0}^{B} \dfrac{\lambda^i}{\mu^i i!}}$

- Mean First Passage and Recurrence Times
  - chain with one recurrent class;
  $\Rightarrow$ fix $\textcircled{s}$ recurrent
  - Mean first time passage time from $i$ to $s$:
    $t_i = E[\min\{n \geq 0 \text{ such that } X_n = s\} \mid X_0 = i]$
    $\Rightarrow \begin{cases} t_s = 0 \\ t_i = 1 + \sum_j P_{ij} t_j \end{cases}$
  - Mean recurrence time of $s$:
    $t_s^* = E[\min\{n \geq 1 \text{ s.t. } X_n = s\} \mid X_0 = s]$
    $\begin{cases} t_s^* = 1 + \sum_j P_{sj} t_j \end{cases}$

# Lect. 19) Limit Theorem

- Chebyshev's Theorem.

  r.v. $X$ w/ $(\mu, \sigma^2)$

  $\Rightarrow \sigma^2 = \int (x-\mu)^2 f_X(x)dx$

  $\geq \int_{-\infty}^{-c}(x-\mu)^2 f_X(x)dx + \int_{c}^{\infty}(x-\mu)^2 f_X(x)dx$

  $\geq c^2 P(|X-\mu| \geq c)$

  $\Rightarrow P(|X-\mu| \geq c) \leq \dfrac{\sigma^2}{c^2}$

  $P(|X-\mu| \geq k\sigma) \leq \dfrac{1}{k^2}$.

- Convergence in probability

  For every $\epsilon > 0$, $\lim\limits_{n \to \infty} P(|Y_n - a| \geq \epsilon) = 0$.

  example $Y_n = 1 - \frac{1}{n}$.

- Convergence of sample mean.

  $X_1, X_2, \ldots$ iid finite mean $\mu$ & variance $\sigma^2$

  $\Rightarrow M_n = \dfrac{X_1 + \cdots X_n}{n}$

  $E[M_n] = \dfrac{E[X_1] + \cdots + E[X_n]}{n} = \dfrac{N\mu}{N} = \mu$

  $Var[M_n] = \dfrac{n\sigma^2}{n^2} = \dfrac{\sigma^2}{n}$

  $\boxed{P(|M_n - \mu| < \epsilon) \leq \dfrac{Var(M_n)}{\epsilon^2} = \dfrac{\sigma^2}{n\epsilon^2}}$

- $M_n$ converges to $\mu$ in probability

  $\Downarrow$

  Weak Law of Large Numbers.

- the Central limit theorem

  $\hookrightarrow \forall c ; \; P(Z_n \leq c) \to P(Z < c)$

  $Z_n = \dfrac{S_n - E(S_n)}{\sigma S_n}$

  standard normal CDF

  (zero mean, unit variance), $S_n = X_1 + X_2 + \cdots X_n$;

---

# Lect 21. Bayesian Statistical Inference
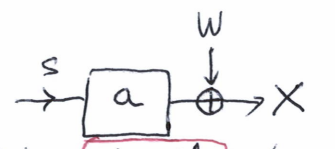


Reality — Model

Data

Application: polling, medical / pharmaceutical trials, netflix competition, finance signal processing (tracking detection, speaker identification)
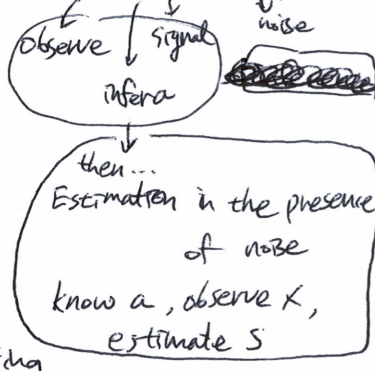
- Types of Inference models / approaches

  - Model Building vs Inferring unknown variables

    $(X = aS + W)$

    observe / signal / noise

    

    ① known  unknown infer   observe
    ② estimate  known   observe

    then...
    Estimation in the presence of noise
    know $a$, observe $X$, estimate $S$

  - Hypothesis testing

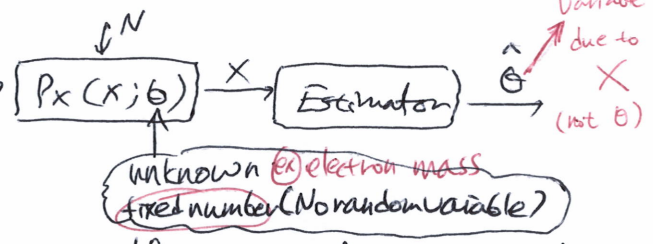    $\hookrightarrow$ unknown takes one of possible values & aim at small PR of incorrect decision

  - Estimation: aim at a small estimation error.

- Classical Statistics
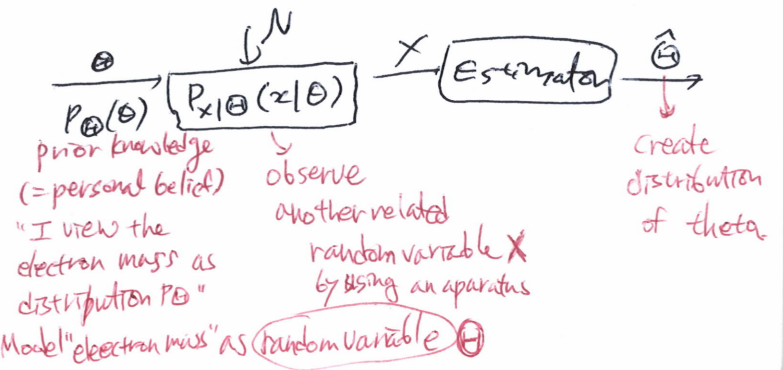
  different philosophical approach

  

  $\theta \to \boxed{P_X(X;\theta)} \xrightarrow{X} \boxed{\text{Estimator}} \xrightarrow{\hat{\theta}}$

  random variable due to $X$ (not $\theta$)

  unknown (ex. electron mass) fixed number (No random variable)

- Bayesian: Use priors & Bayes rule

  

  $\overset{\theta}{\underset{P_\Theta(\theta)}{\to}} \boxed{P_{X|\Theta}(x|\theta)} \xrightarrow{X} \boxed{\text{Estimator}} \xrightarrow{\hat{\Theta}}$

  prior knowledge (= personal belief) "I view the electron mass as distribution $P\theta$"

  observe another related random variable $X$ by using an apparatus

  create distribution of theta

  Model "electron mass" as random variable $\Theta$
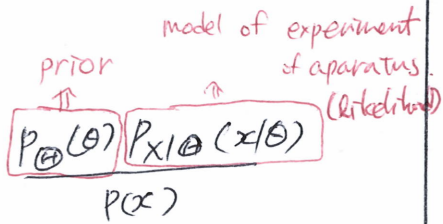
☆ Bayesian Inference : Use Bayes rule

- Hypothesis testing

  – discrete data

  posterior

  $$P_{\Theta|X}(\theta|x) = \frac{P_{\Theta}(\theta) \, P_{X|\Theta}(x|\theta)}{P(x)}$$

  prior — model of experiment of aparatus (likelihood)

  – continuous data

  $$P_{\Theta|X}(\theta|x) = \frac{P_{\Theta}(\theta) \, f_{X|\Theta}(x|\theta)}{f_X(x)}$$

- Estimation ; Continuous data

  $$f_{\Theta|X}(\theta|x) = \frac{f_{\Theta}(\theta) \, f_{X|\Theta}(x|\theta)}{f_X(x)}$$

  $Z_t = \Theta_0 + t\,\Theta_1 + t^2\Theta_2 \cdots$

  $X_t = Z_t + W_t$

  ⇒ Bayes rule gives

  $$f_{\Theta_0,\Theta_1,\Theta_2|X_1,X_2,\cdots X_n}(\theta_0,\theta_1,\theta_2 | x_1, \cdots x_n)$$
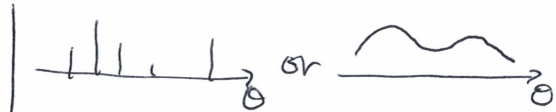
  ( See Monty Hall Problem )

- Estimation w/ ~~continuous~~ discrete data

  $$f_{\Theta|X}(\theta|x) = \frac{f_{\Theta}(\theta)\, P_{X|\Theta}(x|\theta)}{P_X(x)}$$

  $$P_X(x) = \int f_{\Theta}(\theta)\, P_{X|\Theta}(x|\theta)\, d\theta$$

- Output of ~~Bayesian Inference~~ Bayesian Inference

  ⇒ Posterior distribution

  

  ↳ If interested in a single answer,

    – Maximum a posteriori probability (MAP)

    ⇒ $P_{\Theta|X}(\theta^*|x) = \max_\theta P_{\Theta|X}(\theta|x)$

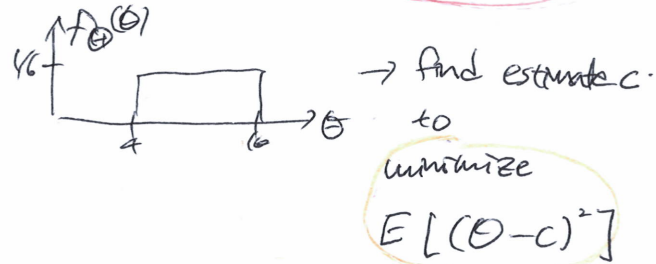    minimizes Pr of error; often used in hypothesis testing.

---

& $f_{\Theta|X}(\hat{\theta}|x) = \max_\theta f_{\Theta|X}(\theta|x)$

& conditional expectation

$$E[\Theta | X=y] = \int \theta \, f_{\Theta|X}(\theta|x)\, dx$$

& single answers can be misleading!

(ex1) Least Mean Square (Estimation)



→ find estimate $c$ to minimize

$$E[(\Theta - c)^2]$$

→ $c = E[\Theta]$
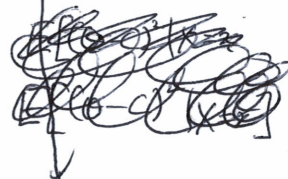
↳ optimal mean squared error

$$E[(\Theta - E[\Theta])^2] = Var(\Theta)$$

if we observe that $X=x$,

$E[(\Theta - c)^2 | X=x]$ is minimized by

$c = E[\Theta | X=x]$



$E[(\Theta - E[\Theta|X=x])^2 | X=x]$

$\le E[(\Theta - g(x))^2 | X=x]$
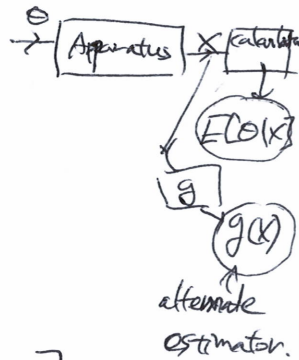
$E[(\Theta - E[\Theta|X])^2 | X] \le E[(\Theta - g(X))^2 | X]$

$E[(\Theta - E[\Theta|X])^2] \le E[(\Theta - g(X))^2]$

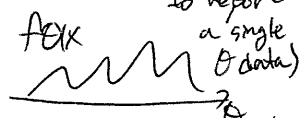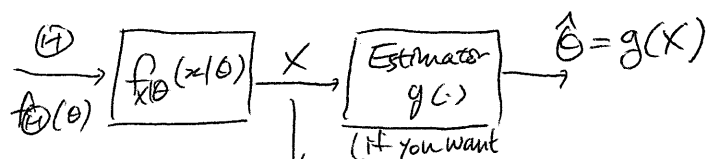$E[\Theta|X]$ minimizes $E[(\Theta - g(X))^2]$ over all estimators $g(\cdot)$

★ LMS Estimation w. several measurement

- unknown r.v $\Theta$
- Observe values of r.v.s $X_1, X_2, \cdots X_n$
- Best Estimator : $E[\Theta | X_1, \cdots X_n]$
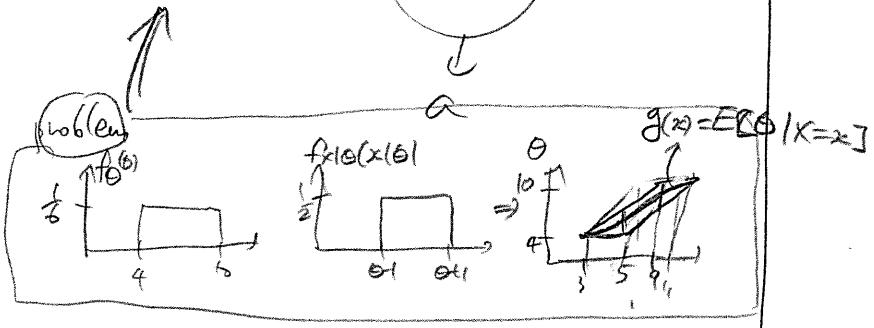
Bayesian

* Linear LMS estimation



$\xrightarrow[f_\Theta(\theta)]{(\theta)} \boxed{f_{X|\Theta}(x|\theta)} \xrightarrow{X} \boxed{\begin{array}{c} \text{Estimator} \\ g(\cdot) \end{array}} \xrightarrow{} \hat\Theta = g(X)$

(If you want to report a single $\theta$ data)

$f_{\Theta|X}$

$\to \hat\theta$

• MAP ($\hat\Theta_{MAP}$ = max $f_{\Theta|X}(\theta|x)$)

or

• LMS Estimation
$\hat\Theta = E[\Theta|X]$ minimize variation of $\hat\Theta$ overall $g(x)$

• Linear LMS

Consider estimator of $\Theta$, of the form
$\hat\Theta = aX + b$

⤷ minimize $E[(\Theta - aX - b)^2] = h(a,b)$ : quad
   $a, b$

⤷ Best choice of $a, b$; best linear estimator

$\hat\Theta_L = E[\Theta] + \boxed{\dfrac{\text{Cov}(X,\Theta)}{\text{Var}(X)}}(X - E[X])$



$g(x) = E[\Theta|X=x]$

• Linear LMS with multiple data, general estimator

$\hat\Theta = a_1 X_1 + \cdots a_n X_n + b$         $E[\Theta|X_1,\cdots X_n]$

⤷ Find best choices of $a_1 \sim a_n, b$

⤷ minimize $E[(a_1 X_1 + \cdots a_n X_n + b - \Theta)^2]$

⤷ Set derivative to zero. linear system rub and the $a_i$

⤷ only means, variances, covariances matter!!

---

* The cleanest linear LMS example

$X_i = \Theta + W_i$    $\Theta, W_1, \cdots W_n$ independent.

$\Theta \sim \mu, \sigma_0^2$   $W_i \sim 0, \sigma_i^2$

$\Rightarrow \hat\Theta_L = \dfrac{\mu/\sigma_0^2 + \left(\sum\limits_{i=1}^{n} X_i / \sigma_i^2\right)}{\left(\sum\limits_{i=0}^{n} 1/\sigma_i^2\right)}$

If all normal, $\hat\Theta_L = E[\Theta | X_1, X_2, \cdots X_n]$

---

Big picture

• Standard examples

  – $X_i$ uniform on $[0, \theta]$:
    ⤷ uniform prior $\theta$.
  – $X_i$ Bernoulli($p$)
    ⤷ uniform (on Beta) prior on $p$
  – $X_i$ normal w/ mean $\theta$, know variance $\sigma^2$
    ⤷ $X_i = \Theta + W_i$
       normal prior on $\theta$;

• Estimation Method

  ⎰ MAP
  ⎱ MSE
  ⎰ Linear MSE

# Lect.23 Classical Inference

• Maximum likelihood estimation

Model w/ unknown parameters $\Theta\Theta\Theta$

$X \sim P_X(x; \Theta)$

↳ Pick $\Theta$ that "makes data most likely"

$\hat{\Theta}_{ML} = \arg\max\limits_{\Theta} P_X(x; \Theta)$  ↗ likelihood ↓

↳ Compare to Bayesian MAP ?

$\hat{\Theta}_{MAP} = \arg\max\limits_{\Theta} P_{\Theta|X}(\Theta|x)$
↗ likelihood  ↗ prior

$\hat{\Theta}_{MAP} = \arg\max\limits_{\Theta} \dfrac{P_{X|\Theta}(x|\Theta) P_\Theta(\Theta)}{P_X(x)}$