

The big picture ingredients in a problem

extra notes

involving uncertainty:

① θ , something

(15 Feb) ~~poster~~ ①

unknown α of interest to you (think of a vector of real numbers of length k) ② D , a dataset - You

judge relevant to decreasing your uncertainty about θ (think of a vector of real numbers of length n)

③ B , a (finite) collection of true.

False propositions, all regarded by you as true, (B_1 and B_2 and ... and B_n) capturing relevant background info.

probability (PT Cox) } $P(B|A)$ | T/F propositions ^②
 quantities your information about the truth of B
 gives that you regard A as T } can extend this

to cumulative distribution functions (CDFs) and
 densities for quantifying uncertainty about
 quantities of living ^{continuously} on the real line, by

considering $F_{(\theta|A)}(z) = P(\theta \leq z | A)$ (CDF)

and $\frac{d}{dz} F_{(\theta|A)}(z) = p_{(\theta|A)}(z)$ (density); i.e.

small abuse of notation people write $F(\theta|A)$ (CDF)
 and $p(\theta|A)$ (density)



total information

about θ

you have to specify & infer θ

③ $p(\theta | B)$, your info. about θ

external to D (your prior dist.)

for θ)

④ $p(D | \theta, B)$, your info.

about θ internal to D (your sampling dist. $\mathcal{G}(\theta)$)

More 2 are enough for inference & prediction

for decision, need to specify 2 more ingredients:

⑤ A = your action space of possible choices

⑥ $U(a, \theta)$, your utility function quantifying cost & benefits if you choose a & unknown θ

call $\{p(\theta|B), p(D|\theta B), a, u(a, \theta)\}$ your model for your uncertainty about θ (give

~~DRB~~)

(usually you're uncertain about θ but

you're also uncertain about how to specify

your uncertainty about θ through the model:

call this second type of uncertainty model

uncertainty) $\left\{ \begin{array}{l} \text{if we have so far pretended that} \\ \text{this type of uncertainty does} \end{array} \right.$

not exist)

poor off: with these 4 ingredients

specified, \mathcal{D}_n the optimal for yielding logically - \mathcal{D}
 internally - consistent answers is as follows:

① (inference) (a) write $L(\theta | \mathcal{D} \mathcal{B}) = c p(\mathcal{D} | \theta \mathcal{B})$
 (for arbitrary positive c) (b) compute:

$$p(\theta | \mathcal{D} \mathcal{B}) = c p(\theta | \mathcal{B}) L(\theta | \mathcal{D} \mathcal{B})$$

posterior
 dist.
 (full-
 info.
 dist.)

likelihood function

② (prediction) future \mathcal{D}^*

$$p(\mathcal{D}^* | \mathcal{D} \mathcal{B}) = \int p(\mathcal{D}^* | \theta \mathcal{B}) p(\theta | \mathcal{D} \mathcal{B}) d\theta$$

set of $\mathcal{B} \rightarrow$ possible values of θ

③ (Revision) find a^* that maximizes \mathcal{L}

E $U(a, \theta)$ (expected utility)
 $(\theta | D_B) \propto \int U(a, \theta) p(\theta | D_B) d\theta$

Challenges:
 ④ the
 specifications
 by D_B

⑤ computation: $p(\theta | D_B) = \int p(\theta | p) \mathcal{L}(a | D_B)$

$\theta = (\theta_1, \dots, \theta_k)$

hard to visualize

each of these \rightarrow is a k -dimensional probability dist.

work with

⑤: $p(\theta, | D_B) = \int \dots \int p(\theta, \theta_1, \dots, \theta_k)$

directly
 marginal \uparrow
 post. first for $\theta_1, \dots, \theta_k$

& here we $\textcircled{4}$ of here $(k-1)$ dim. integrals $\textcircled{7}$

$$(8) \text{ (prediction)} \quad p(\mathcal{D}^* | \mathcal{D} \mathcal{B}) = \int p(\mathcal{D}^* | \theta, \mathcal{D} \mathcal{B}) p(\theta | \mathcal{D} \mathcal{B}) d\theta$$

Here is another k -dimensional integral $\textcircled{8}$ (inference)

$$\text{No normalizing constant } c \text{ in } p(\theta | \mathcal{D} \mathcal{B}) = c p(\theta | \mathcal{B}) p(\mathcal{D} | \theta)$$

is another k -dimensional integral $\textcircled{9}$ (as expected)

utility $\int u(a, \theta) p(\theta | \mathcal{D} \mathcal{B}) d\theta$ is another

k -dimensional integral, and if \mathcal{D} is large you'll

have to compute or approximate this integral

a lot of times (ex. variable selection in regression.)

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im} + \varepsilon_i \leftarrow \text{IID } N(0, \sigma^2) \quad \textcircled{B}$$

which subset of (x_1, \dots, x_m) is "left" for prediction

Y? } there are 2^m subsets, so Q has 2^m actions
 to compare (ex: $m=100$; $2^m = 10^{30}$ actions
 to choose among) when k is small, sometimes

all these problems go away, if you can find a

conjugate prior to the sampling dist. is

ANTI $\textcircled{10}$

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m$$

your problem ex. $L(\theta | \mathcal{D}) = \theta^s (1-\theta)^{n-s}$, $s = \sum_{i=1}^n x_i$

$p(\theta | \mathcal{R}) = c \theta^{\alpha-1} (1-\theta)^{\beta-1}$ for $\alpha, \beta > 0$ (the family of Beta densities) has property that

$$p(\theta | \mathcal{D}_B) = c p(\theta | \mathcal{R}) L(\theta | \mathcal{D}_B)$$

Beta · Beta

but conjugate priors only

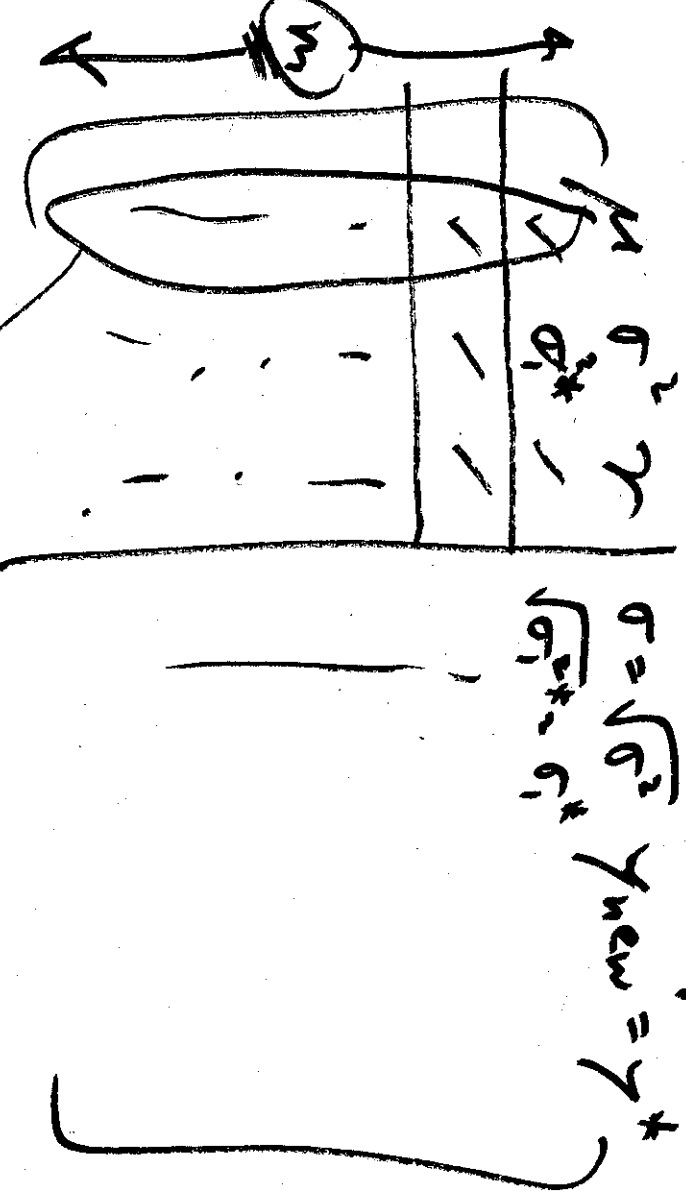
exist for a small subset of {real problems with small k }; need software more general

computing method: MCMC

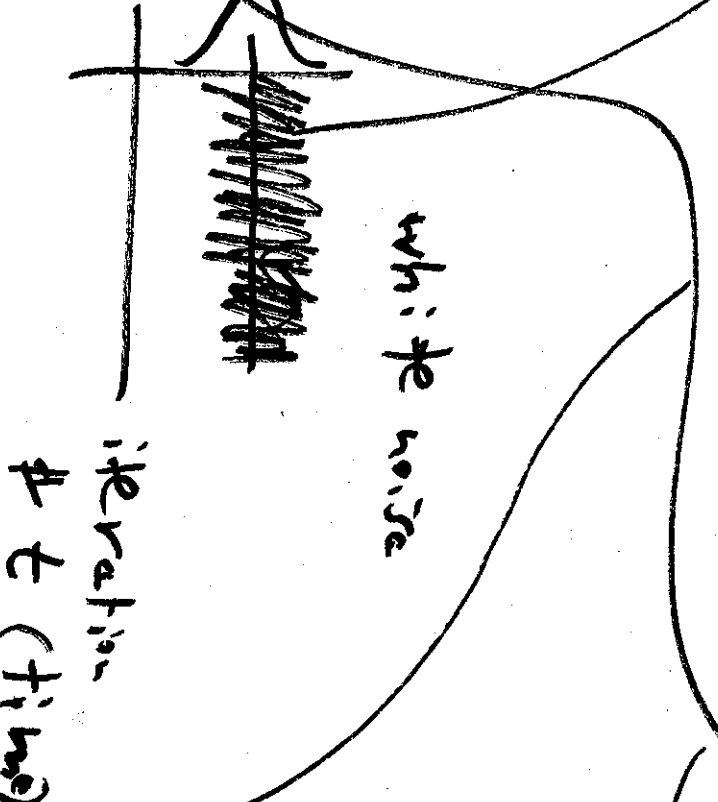
ex. 1 $(y_i | \mu, \sigma^2, \tau | \mathcal{D}) \stackrel{i.i.d.}{\sim}$ $E=3$

NR to data $(\mu, \sigma^2, \tau) \sim p(\mu, \sigma^2, \tau | \mathcal{R})$ $t_r(\mu, \sigma^2)$

Let us want to sample from $p(\mu, \sigma^2 | r | D^t)$ ⑩

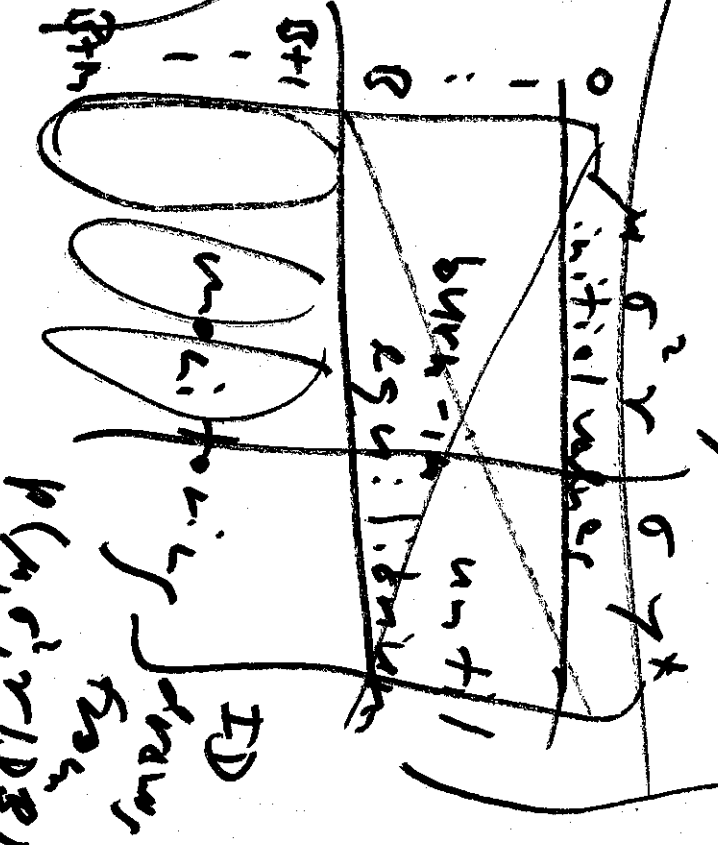


means μ^t
 σ^t, σ_{μ}^t
 density
 $p(\mu | D^t)$



if IID
 AC data set

if IID not possible, KCHC



McMC time series plot

poor mixing if $\#t$

strongly positively autocorrelated time series

(log autocorrelations)

$$\begin{pmatrix} \mu_{B+1}^* & \mu_{B+2}^* & \dots & \mu_{B+m}^* \\ \mu_{B+1}^* & \mu_{B+2}^* & \dots & \mu_{B+m}^* \\ \mu_{B+1}^* & \mu_{B+2}^* & \dots & \mu_{B+m}^* \\ \mu_{B+1}^* & \mu_{B+2}^* & \dots & \mu_{B+m}^* \end{pmatrix} \begin{matrix} r \\ r \\ r \\ r \end{matrix} \approx \begin{matrix} \rho_1 \\ \rho_1 \\ \rho_1 \\ \rho_1 \end{matrix}$$

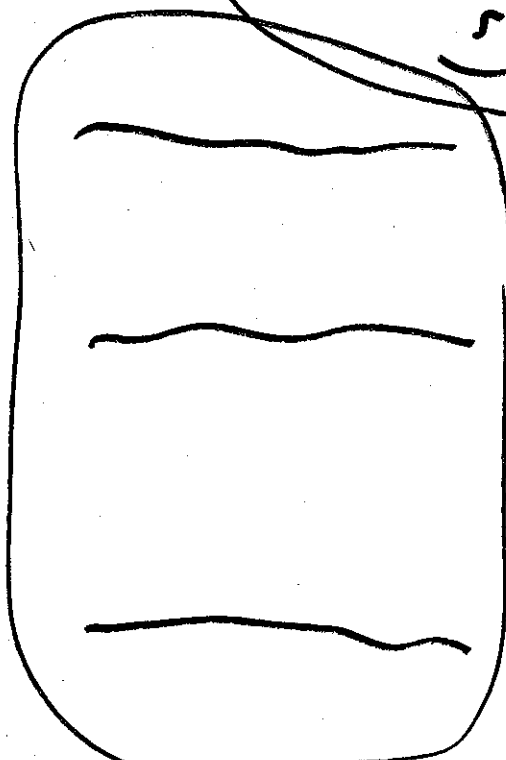
① brute-force parallelizing MCMC

node 1 node 2 ... node N

read # read 1 2 ... N

save init. values (or different)

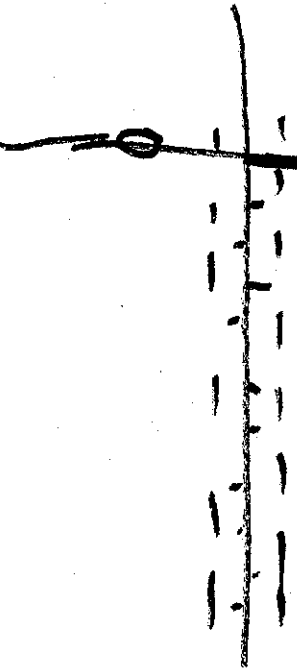
save down it



$\left[\begin{matrix} \mu_{B1}^* & \mu_{B3}^* \\ \mu_{B2}^* & \mu_{B3}^* \end{matrix} \right]$

low ρ_2

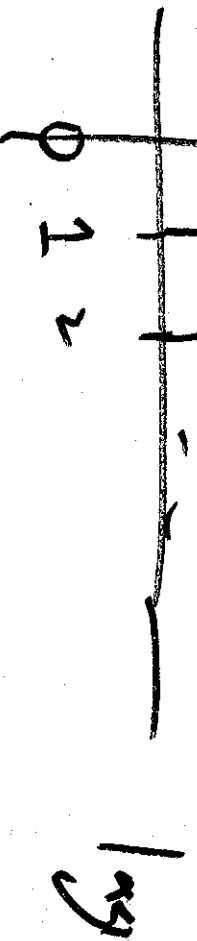
AcF white noise



...
 $\rho(\rho_1, \rho_2, \dots, \rho_k)$ AcF
 sample
 autocorrelation function

+1+

good mix
 (ρ_1, ρ_2, \dots)



poor mix
 geometric
 decay



poor mixing with
 large ρ_1

common class of time series models (Box-Jenkins) (13)
 time domain) is called autoregressive (AR),
 models of order (m) :

$$\theta_t = \alpha_1 \theta_{t-1} + \alpha_2 \theta_{t-2} + \dots + \alpha_m \theta_{t-m} + e_t$$

$(\alpha_1, \alpha_2, \dots, \alpha_m, e_t)$

Fact

many columns in matrix look
 like AR_1 (AR_1 order 1) time series
 with 1st order autocorrelation ρ_1

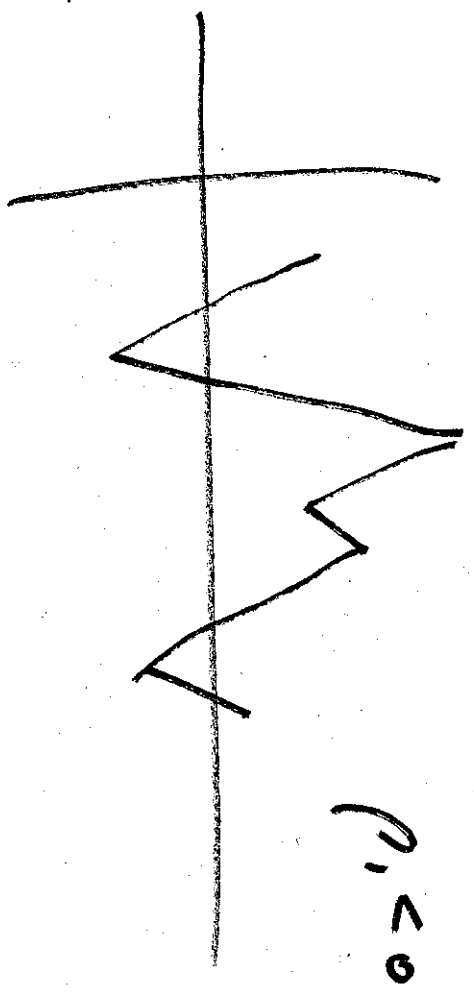
if $\rho_1 = 0$ + white noise (IID) (AR_0)

Gills sampling in the NB10 + worker } posterior (18)

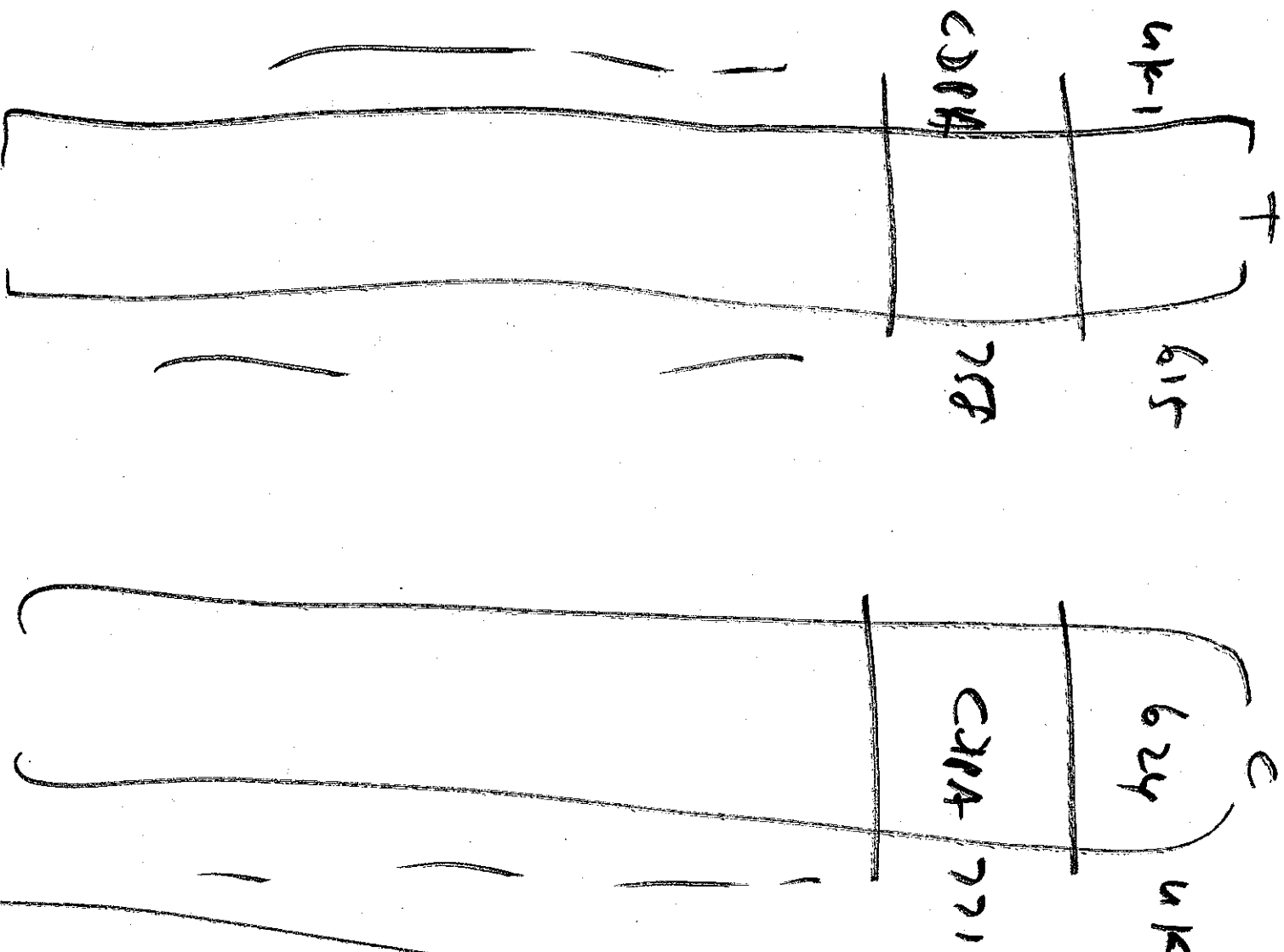
is $p(\mu, \sigma^2, r | D_B)$; full conditionals are

$$\left\{ \begin{array}{l} p(\mu | \sigma^2, r, D_B) \\ p(\sigma^2 | \mu, r, D_B) \\ p(r | \mu, \sigma^2, D_B) \end{array} \right\}$$

ex. of Markov
called
stochastic
relaxation



pooling (fixed effects analysis)
(chemilla idea)



it's my own study

$$SE(\bar{Y}_i) = SE(\hat{\beta}_{T_i})$$

$$\frac{1}{n_{T_i}} \sum y_{i,j} = \hat{\beta}_{T_i} = \sqrt{\frac{\hat{\beta}_{T_i}(1-\hat{\beta}_{T_i})}{n_{T_i}}}$$

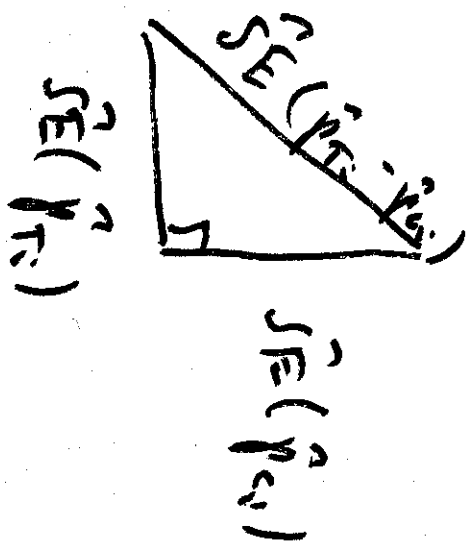
study
patient

$$SE(\bar{c}_i) = SE(\hat{\beta}_{c_i}) = \sqrt{\frac{\hat{\beta}_{c_i}(1-\hat{\beta}_{c_i})}{n_{c_i}}}$$

$$SE(\bar{T}_i - \bar{C}_i) = SE(\hat{\mu}_{T_i} - \hat{\mu}_{C_i}) =$$

(16)

4th indy.



$$= \sqrt{\left[SE(\hat{\mu}_{T_i}) \right]^2 + \left[SE(\hat{\mu}_{C_i}) \right]^2}$$

$$= \sqrt{\underbrace{\frac{\hat{\mu}_{T_i}(1-\hat{\mu}_{T_i})}{n_{T_i}} + \frac{\hat{\mu}_{C_i}(1-\hat{\mu}_{C_i})}{n_{C_i}}}}_{SE \text{ of diff } (\%)}$$

a straightforward way to make operational the $\text{Pr}(\mathcal{D})$
judgment of one's liability of unknown
quantities $\theta_1, \dots, \theta_k$ is to regard them as (conditionally)

ISD laws from the same Dist: e.g.

$$(\theta_{01}, \mu, \sigma^2) \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2)$$

data have no
info about

some quantity \rightarrow likelihood-flat \rightarrow

post $\hat{=}$ prior for that quantity

graphical representation of model (1), p. 4 of part 4 of PDF notes

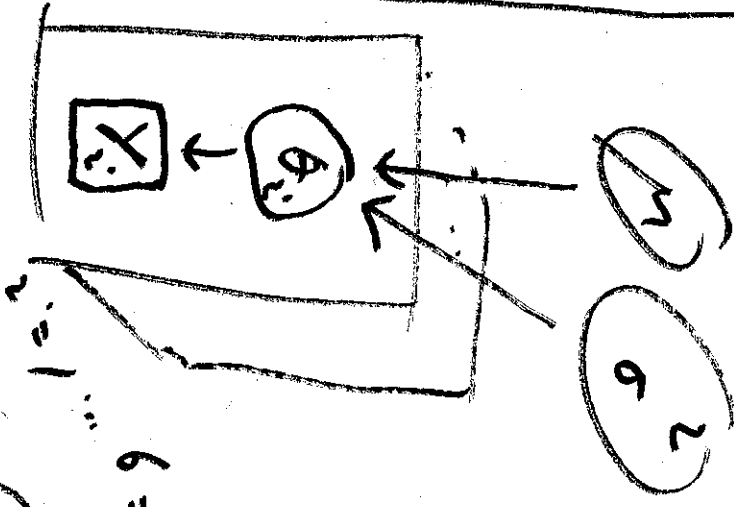
$$\begin{aligned} (\mu, \sigma^2) &\sim \gamma(\mu, \sigma^2) \\ (\theta_i | \mu, \sigma^2) &\stackrel{\text{ind}}{\sim} N(\mu, \sigma^2) \\ (y_i | \theta_i) &\stackrel{\text{ind}}{\sim} N(\theta_i, V_i) \end{aligned}$$

known

$$\theta = (\theta_1, \dots, \theta_k)$$

$$y = (y_1, \dots, y_n)$$

$$p(\mu, \sigma^2, \theta, y) =$$



$i = 1 \dots k$

μ, σ^2
 params. of θ_i
 θ_i part of y_i
 1 sheet (plate)
 for each θ_i

\circ unknown
 \square known

$(y_i | \mu, \sigma^2)$ conditionally indep. given θ

$$p(\mu) p(\sigma^2 | \mu) p(\theta | \mu, \sigma^2) p(y | \theta)$$

