

Approximate Bayesian Inference

$$\underbrace{p(\theta | x_{1:n})}_{\text{posterior}} \propto \underbrace{p(x_{1:n} | \theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}$$

→ models that have analytical expressions require that

$$\underbrace{p(x_{1:n})}_{\text{marginal likelihood}} = \int p(x_{1:n} | \theta) p(\theta) d\theta$$

$$x_i \sim \text{Bernoulli}(\pi)$$

$$\pi \sim \text{Beta}(\alpha_0, \beta_0)$$

$$p(x_i | \pi) = \pi^{x_i} (1-\pi)^{1-x_i}$$

$$p(\pi) = \frac{\pi^{\alpha_0-1} (1-\pi)^{\beta_0-1}}{B(\alpha_0, \beta_0)}$$

$$p(x_i) = \int_0^1 p(x_i | \pi) p(\pi) d\pi$$

$$= \int_0^1 \pi^{x_i} (1-\pi)^{1-x_i} \frac{\pi^{\alpha_0-1} (1-\pi)^{\beta_0-1}}{B(\alpha_0, \beta_0)} d\pi$$

$$= \frac{1}{B(\alpha_0, \beta_0)} \int_0^1 \pi^{\overbrace{x_i + \alpha_0 - 1}^{\alpha'}} (1-\pi)^{\overbrace{1 - x_i + \beta_0 - 1}^{\beta'}}$$

$$= \frac{1}{B(\alpha_0, \beta_0)} \int_0^1 \frac{B(\alpha', \beta')}{B(\alpha', \beta')} \pi^{\alpha' - 1} (1-\pi)^{\beta' - 1} d\pi$$

$$= \frac{B(\alpha', \beta')}{B(\alpha_0, \beta_0)}$$

But, what if we can't compute $p(x_{1:n})$?

(1) Parametric approximations

Find q closest as possible to $p(\theta | x_{1:n})$

Laplace Approximation

← most popular approach

$$\theta | X_{1:n} \sim N(\hat{\theta}_{\text{MAP}}, I_B^{-1}(\hat{\theta}_{\text{MAP}}))$$

Posterior mode:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta | X_{1:n})$$

$$I_B(\hat{\theta}_{\text{MAP}}) = -\mathbb{E} \left[\frac{\partial^2 \log p(\theta | X_{1:n})}{\partial \theta^2} \Big| \theta = \hat{\theta}_{\text{MAP}} \right]$$

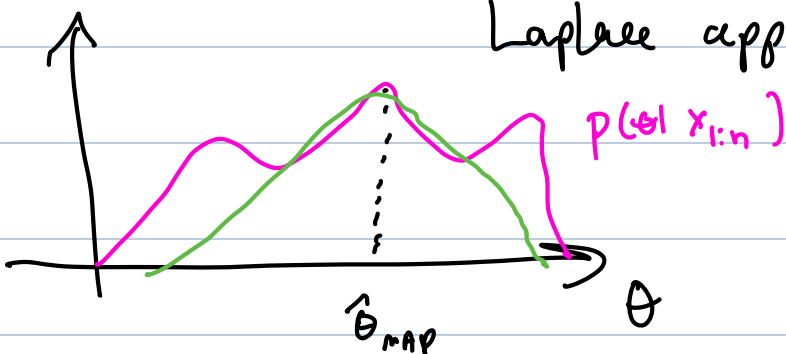
Limitations:

- If θ is multivariate, then you need the Hessian matrix of $\log p(\theta | X_{1:n})$

→ expensive to compute in high dimensions

- Asymptotic normality is for large n . If n is small, the Laplace approximation will be inaccurate

- If $p(\theta | X_{1:n})$ is multimodal, the Laplace approximation will fail



Monte Carlo Methods

→ backed by SLLN:

Given a random sample X_1, \dots, X_n with population $p(x)$, under the assumption $E[X] < \infty$ then

$$\begin{array}{l} X_1, \dots, X_n \\ \downarrow \\ Y_i = g(X_i), \dots, Y_n = g(X_n) \end{array} \begin{array}{l} \rightarrow E[X] < \infty \\ \rightarrow E[Y] < \infty \end{array} \quad \bar{X}_n \xrightarrow{\text{a.s.}} E[X]$$

→ Idea of Monte Carlo is to use samples from $p(\theta | X_{1:n})$ as a proxy for the distribution to compute any expectation values we want

- Standard Monte Carlo: Assumes we can sample from $\pi(\theta) \stackrel{\Delta}{=} p(\theta | X_{1:n})$

$$\theta^{(1)}, \dots, \theta^{(m)} \stackrel{\text{i.i.d.}}{\sim} \pi$$

$$\mathbb{E}_{\pi} [h(\theta)] \approx \frac{1}{M} \sum_{m=1}^M h(\theta^{(m)})$$

$$\xrightarrow{\text{a.s.}} \int h(\theta) \pi(\theta) d\theta$$

→ Standard Monte Carlo assumes we can sample from the posterior, which is generally not the case

Rejection Sampling: $\pi(\theta) \triangleq p(\theta | x_{1:n})$
 $\tilde{\pi}(\theta) \triangleq p(x_{1:n} | \theta) p(\theta)$

→ Idea: Sample θ from another distribution $q(\theta)$ and reject some of those samples with some probability

Algorithm: Requires $q(\theta)$, $M = \sup_{\theta} \frac{\tilde{\pi}(\theta)}{q(\theta)}$

$i=1$
While (keep-sampling):

1. Sample $\theta_* \sim q(\theta)$
2. Sample $u \sim \mathcal{U}(0, 1)$

$$3. \text{ If } u \leq \frac{\tilde{\pi}(\theta)}{Mq(\theta)} :$$

θ_* is accepted as a sample

$$\theta^{(i)} = \theta_*$$

$$i = i + 1$$

$$\{\theta^{(i)}\}_i$$

We can show that set of samples $\{\theta^{(i)}\}$ are samples from $\pi(\theta)$:

Let $\tilde{\theta}$ be a random generated from the algo.

$$F_{\tilde{\theta}}(a) = \mathbb{P}(\tilde{\theta} \leq a)$$

$$= \mathbb{P}(\theta_* \leq a \mid u \leq \frac{\tilde{\pi}(\theta_*)}{Mq(\theta_*)})$$

$$= \frac{\mathbb{P}(\theta_* \leq a, U \leq \frac{\tilde{\pi}(\theta_*)}{Mq(\theta_*)})}{\mathbb{P}(U \leq \frac{\tilde{\pi}(\theta_*)}{Mq(\theta_*)})}$$

$$\mathbb{P}(U \leq \frac{\tilde{\pi}(\theta_*)}{Mq(\theta_*)})$$

$$(1.) \quad P\left(\theta_* \leq a, U \leq \frac{\tilde{\pi}(\theta_*)}{M q(\theta_*)}\right) =$$

Law of total prob.

$$\int P\left(\theta_* \leq a, U \leq \frac{\tilde{\pi}(\theta_*)}{M q(\theta_*)} \mid \theta_* = \theta\right) q(\theta) d\theta$$

$$= \int_{-\infty}^a P\left(U \leq \frac{\tilde{\pi}(\theta_*)}{M q(\theta_*)} \mid \theta_* = \theta\right) q(\theta) d\theta$$

$$= \int_{-\infty}^a \frac{\tilde{\pi}(\theta)}{M q(\theta)} q(\theta) d\theta = \frac{1}{M} \int_{-\infty}^a \tilde{\pi}(\theta) d\theta$$

$$Z = \int_{-\infty}^{\infty} \tilde{\pi}(\theta) d\theta$$

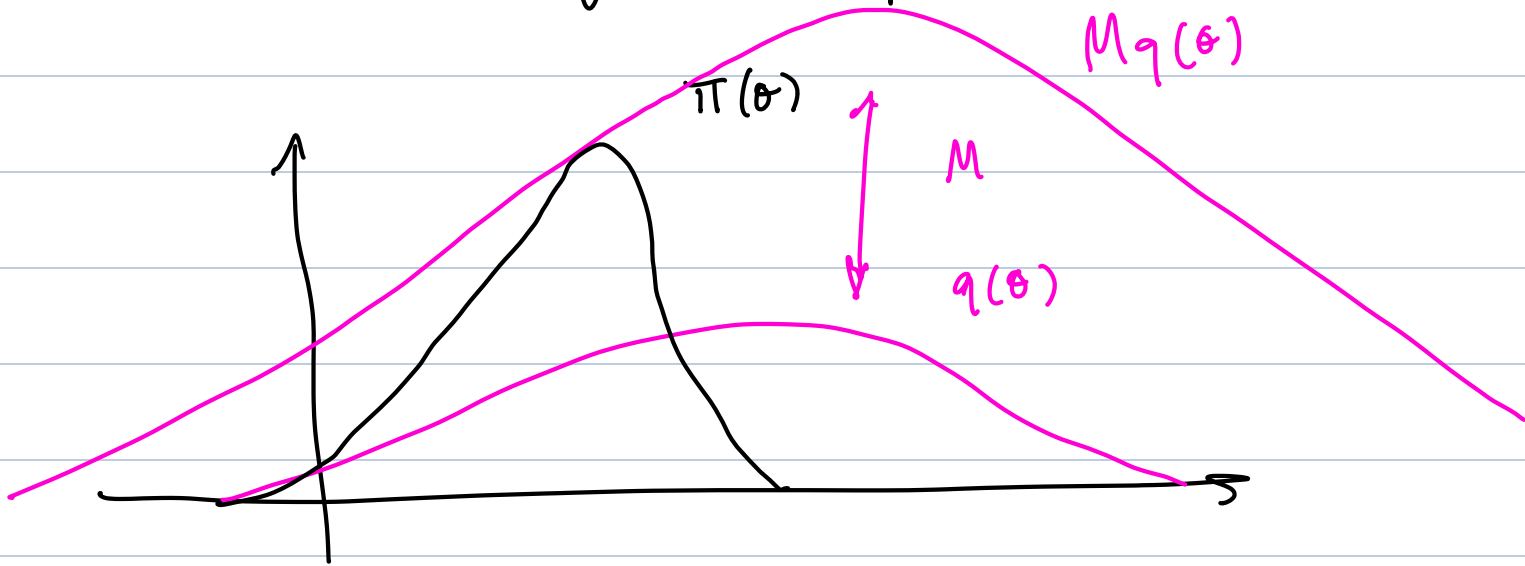
$$= \frac{1}{M} \int_{-\infty}^a \frac{Z}{Z} \tilde{\pi}(\theta) d\theta = \frac{Z}{M} F_{\theta}(a)$$

Using same approach, show (2.) = $\frac{Z}{M}$

$$\text{Thus, } F_{\tilde{\theta}}(a) = F_{\theta}(a) \implies \tilde{\theta} = \theta \quad \square$$

$$\bullet \text{ Rate of rejection: } \frac{Z}{M} \leftarrow P\left(U \leq \frac{\tilde{\pi}(\theta)}{M q(\theta)}\right)$$

- Limitations: If q is poorly chosen, we have to set M to be large and thus we reject too many samples



- Importance Sampling:

If we know $\pi(\theta)$:

$$\begin{aligned} \mathbb{E}_{\pi} [h(\theta)] &= \int h(\theta) \pi(\theta) d\theta \\ &= \int h(\theta) \frac{\pi(\theta)}{q(\theta)} q(\theta) d\theta \end{aligned}$$

$q(\theta) > 0$
whenever $\pi(\theta) > 0$

$$= \mathbb{E}_q [h(\theta) w(\theta)] \quad w(\theta) = \frac{\pi(\theta)}{q(\theta)}$$

Normalized Importance Sampling

$$\approx \frac{1}{M} \sum_{m=1}^M h(\theta^{(m)}) w(\theta^{(m)})$$

Importance weight
 $\theta^{(m)} \stackrel{iid}{\sim} q(\theta)$

a.s. $\rightarrow \mathbb{E}_\pi [h(\theta)]$ by SLLN

If we know $\tilde{\pi}(\theta)$:

$$\mathbb{E}_\pi [h(\theta)] = \int h(\theta) \frac{\tilde{\pi}(\theta)}{\int \tilde{\pi}(\theta) d\theta} d\theta$$

$\downarrow \hat{=} Z$
normalizing constant

$$= \frac{1}{Z} \int h(\theta) \frac{\tilde{\pi}(\theta)}{q(\theta)} q(\theta) d\theta = \frac{\tilde{\pi}(\theta)}{q(\theta)}$$

$$= \frac{1}{Z} \mathbb{E}_q [h(\theta) \tilde{w}(\theta)]$$

$$Z = \int \tilde{\pi}(\theta) d\theta = \int \tilde{\pi}(\theta) \frac{q(\theta)}{q(\theta)} d\theta$$

$$= \mathbb{E}_q [\tilde{w}(\theta)] \quad \hat{I}$$

$$\mathbb{E}_\pi [h(\theta)] = \frac{\mathbb{E}_q [h(\theta) \tilde{w}(\theta)]}{\mathbb{E}_q [\tilde{w}(\theta)]} \approx \frac{\frac{1}{n} \sum_{m=1}^n h(\theta^{(m)}) \tilde{w}(\theta^{(m)})}{\frac{1}{n} \sum_{m=1}^n \tilde{w}(\theta^{(m)})}$$

$$\xrightarrow{\text{a.s.}} E_{\tilde{\pi}} [h(\theta)]$$

→ numerator converges a.s. to $\int \tilde{\pi}(\theta) h(\theta) d\theta$

→ denominator converges a.s.

$$\int \tilde{\pi}(\theta) d\theta$$

Limitations : Variance of importance sampling estimator scales with the variance of $\tilde{w}(\theta)$

$$V[\hat{I}] \propto V[|h(\theta)| w(\theta)]$$

→ q should be close to $|h(\theta)| \pi(\theta)$ to reduce this variance

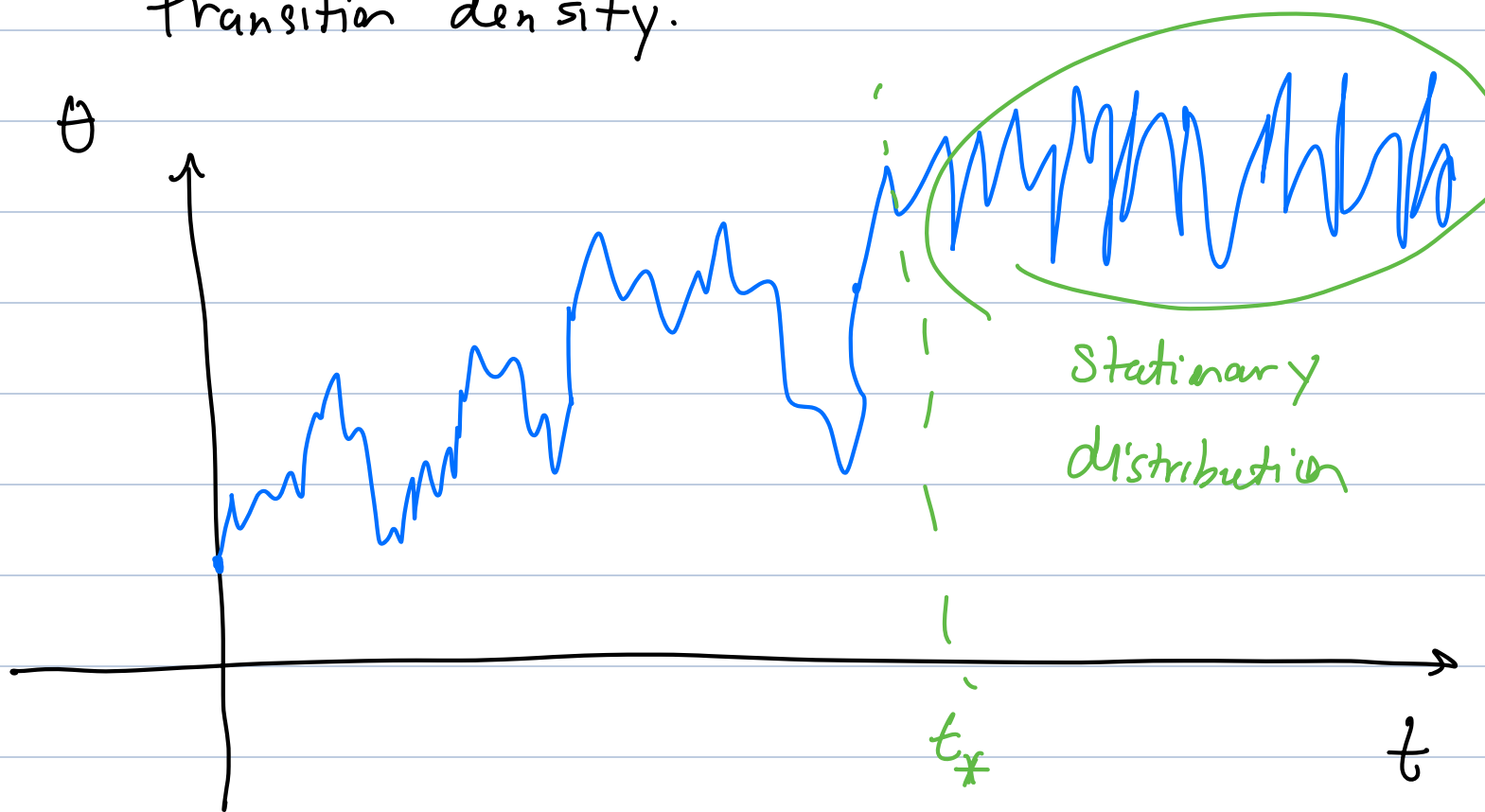
Markov chain Monte Carlo (MCMC)

A stochastic process is an indexed collection of RVs $\{\theta_s^{(i)}\}_S$ where S is the index set and it is assumed to be uncountable.

Def. Markov chain is a stochastic process $\{\theta^{(t)}\}$ that satisfies the Markov property

$$p(\theta^{(t+1)} | \theta^{(0)}, \dots, \theta^{(t)}) = p(\theta^{(t+1)} | \theta^{(t)})$$

The distribution $p(\theta^{(t+1)} | \theta^{(t)})$ is called transition density.



Def. A stationary distribution for the Markov chain $\{\theta^{(t)}\}$ with transition density $T(\theta^{(t+1)} | \theta^{(t)})$ has the property

$$p(\theta_*) = \int p(\theta) \pi(\theta | \theta_*) d\theta$$

MCMC idea: Design an algorithm that constructs a Markov chain with stationary distribution that is $\pi(\theta)$

↳ Existence of stationary distribution

↳ Markov chain should be ergodic

↳ aperiodic chain

↳ positive recurrence $E[T_{\theta_*}] < \infty$

expected time to revisit a state θ_* is finite

Algorithm: Initialize $\theta^{(0)}$

for $t = 1, \dots, \infty$:

1. Sample $\theta_* \sim q(\theta | \theta^{(i-1)})$

ex: $N(\theta^{(i-1)}, \sigma^2)$

2. Sample $u \sim U(0, 1)$

$\alpha = \frac{\tilde{\pi}(\theta_*) q(\theta^{(i-1)}|\theta_*)}{\tilde{\pi}(\theta^{(i-1)}) q(\theta_*|\theta^{(i-1)})}$ 3. Check if $u \leq \min(\alpha, 1)$:
if true:

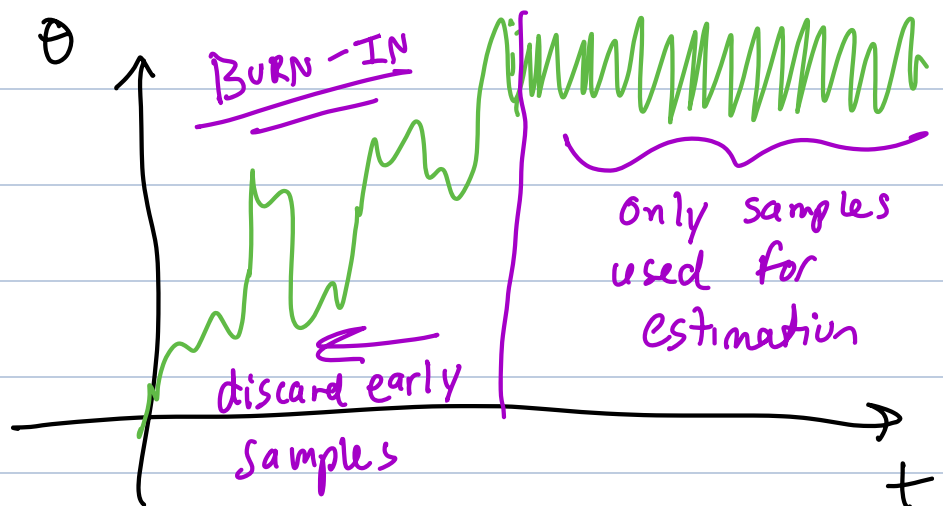
$$\theta^{(i)} = \theta_*$$

Otherwise

$$\theta^{(i)} = \theta^{(i-1)}$$

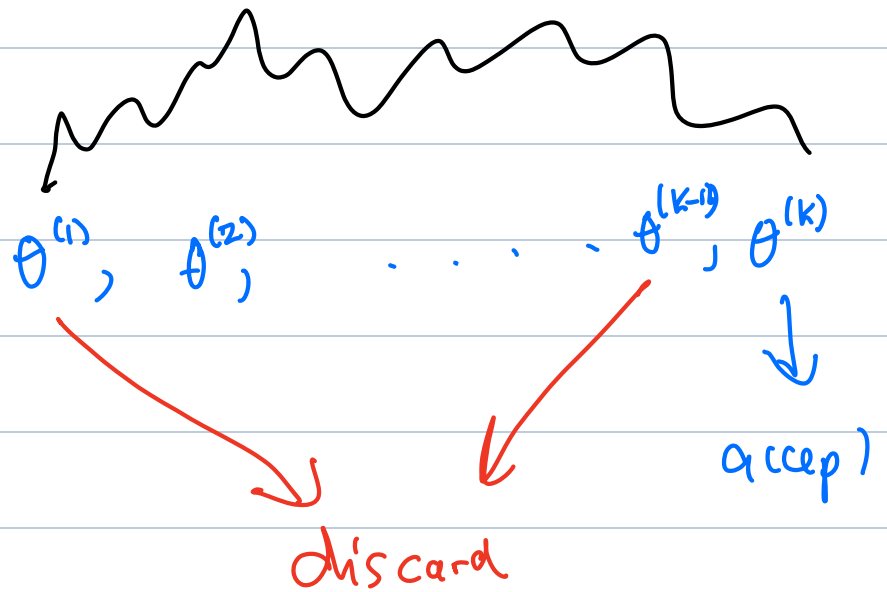
We can show that this algorithm's constructed Markov chain has stationary distribution = to $\pi(\theta)$

Important aspects: (1) Stationarity is not reached immediately, it will be achieved as $t \rightarrow \infty$



(2) Chains can be auto-correlated in practice:

→ Thinning is applied to Markov chains in practice where only every K^{th} sample is accepted



Linear Models (Chapter 4 - Kay Estimation Theory)

$$\underbrace{x}_{\mathbb{R}^n} = \underbrace{H}_{\mathbb{R}^{n \times d}} \underbrace{\theta}_{\mathbb{R}^d} + \underbrace{w}_{\mathbb{R}^n}$$

→ \underline{H} is an observation matrix and is assumed to be known

→ \underline{w} is zero-mean white noise

$$\mathbb{E}[\underline{w}] = 0$$

$$\mathbb{E}[\underline{w}\underline{w}^T] = \sigma^2 \mathbf{I}$$

identity matrix
 $n \times n$

Least Squares: Find a point estimator $\hat{\theta}$ that minimizes the squared error between our model's prediction and the observed data

$$\hat{\theta} = \arg \min_{\theta} \mathbb{E} \left[\left\| x - (H\hat{\theta} + \underline{v}) \right\|_2^2 \right]$$

$$L(\theta) = \left\| x - H\hat{\theta} \right\|_2^2 \xrightarrow{\text{take gradient}} \frac{\partial L}{\partial \theta} = -2H^T(x - H\hat{\theta}) = 0$$

$$H^T x - H^T H \hat{\theta} = 0 \rightarrow H^T x = H^T H \hat{\theta}$$

$$\hat{\theta} = (H^T H)^{-1} H^T x$$

• What is $\mathbb{E}[\hat{\theta}]$?

$$\mathbb{E} \left[(H^T H)^{-1} H^T x \right] = (H^T H)^{-1} H^T \mathbb{E}[x]$$

$$= \underbrace{(H^T H)^{-1}} \underbrace{H^T H} \theta$$

$$= \theta \quad \text{Therefore, unbiased estimator}$$

• What is $C_{\hat{\theta}} = \mathbb{E} [(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T]$

$$= \mathbb{E} [\hat{\theta} \hat{\theta}^T] - \underbrace{\mathbb{E} [\hat{\theta} \theta^T]}_{\theta \theta^T} - \underbrace{\mathbb{E} [\theta \hat{\theta}^T]}_{\theta \theta^T} + \underbrace{\mathbb{E} [\theta \theta^T]}_{\theta \theta^T}$$

$$= \mathbb{E} [\hat{\theta} \hat{\theta}^T] - \theta \theta^T$$

$$= \mathbb{E} \left[(H^T H)^{-1} H \underbrace{x x^T}_{H\theta + w} H^T (H^T H)^{-1} \right]$$

$$= \underbrace{\mathbb{E} \left[(H^T H)^{-1} H (H\theta + w)(\theta^T H^T + w^T) H^T (H^T H)^{-1} \right]}_{\sigma^2 (H^T H)^{-1} + \theta \theta^T}$$

$$= \sigma^2 (H^T H)^{-1}$$

←

In the case $p(w)$ is Gaussian, this is the MVUE and achieves the CRLB

Alternative derivation for linear model

$$w \sim \mathcal{N}(\underline{0}, \sigma^2 \mathbf{I})$$
$$x = H\theta + w$$

Then,

$$p(x; \theta) = \mathcal{N}(x; H\theta, \sigma^2 \mathbf{I})$$
$$= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} (x - H\theta)(x - H\theta)^T\right)$$

$$\ell(\theta; x) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x - H\theta\|_2^2$$

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta; x) \implies \arg \min_{\theta} \|x - H\theta\|_2^2$$

$$\hat{\theta} = (H^T H)^{-1} H^T x \quad \text{for both cases!}$$

Non-Isotropic Covariance: $w \sim N(\underline{0}, \underline{C})$

$$\underline{C}^{-1} = D^T D \quad \leftarrow$$

assume since \underline{C} and \underline{C}^{-1} by definition are positive definite, they can be factored via Cholesky

$$\begin{aligned} \mathbb{E}[(Dw)(Dw)^T] &= DC D^T \\ &= DD^{-1}(D^{-1})^T D^T = I \end{aligned}$$

$$X = H\theta + w, \quad \mathbb{E}[ww^T] = C$$

$$DX = DH\theta + Dw$$

\rightsquigarrow

$$\tilde{X} = \tilde{H}\theta + \tilde{w}$$

$$\left. \begin{aligned} \tilde{X} &= DX \\ \tilde{H} &= DH \\ \tilde{w} &= Dw \end{aligned} \right\}$$

$$\hat{\theta} = (\tilde{H}^T \tilde{H})^{-1} \tilde{H}^T \tilde{X}$$

$$= (H^T D^T D H)^{-1} H^T D X$$