

Statistical Data Analysis

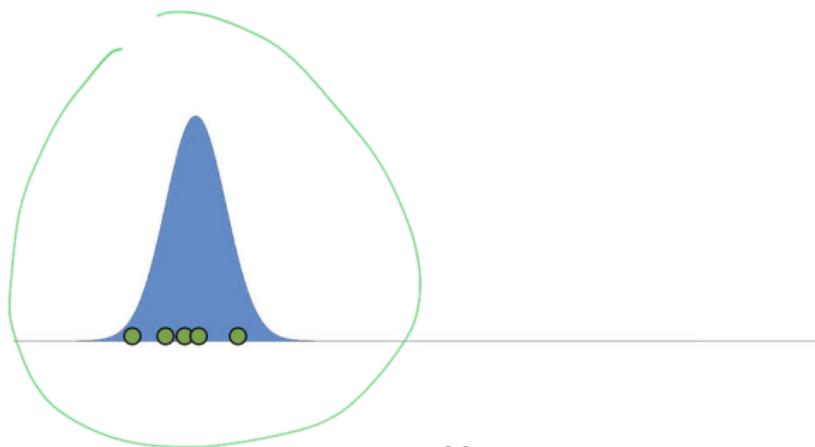
Dr. Jana de Wiljes

16. November 2021

Universität Potsdam

Lecture 12:30
next week
on tuesday

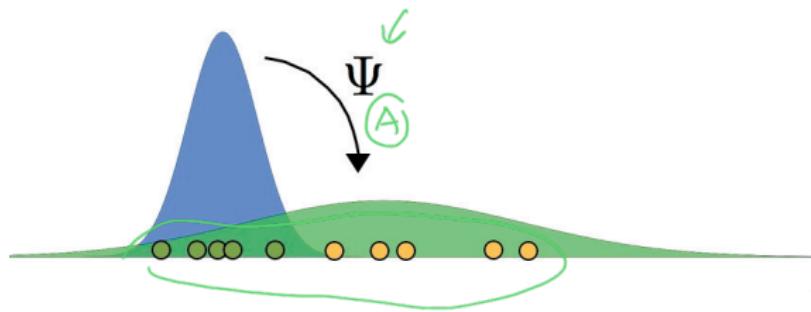
Ensemble Kalman filter



$$\mathcal{N}(\textcolor{blue}{m}_0, \textcolor{blue}{C}_0) \text{ with } \textcolor{blue}{m}_0 \approx \frac{1}{M} \sum_{i=1}^M \textcolor{blue}{z}_0^i$$

$$\textcolor{blue}{C}_0 \approx \frac{1}{M} \sum_{i=1}^M (\textcolor{blue}{z}_0^i - \textcolor{blue}{m}_0)(\textcolor{blue}{z}_0^i - \textcolor{blue}{m}_0)^\top$$

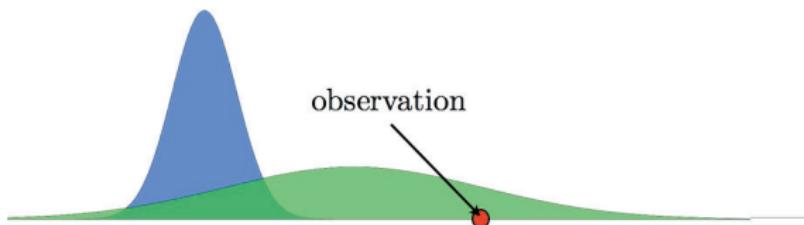
Ensemble Kalman filter



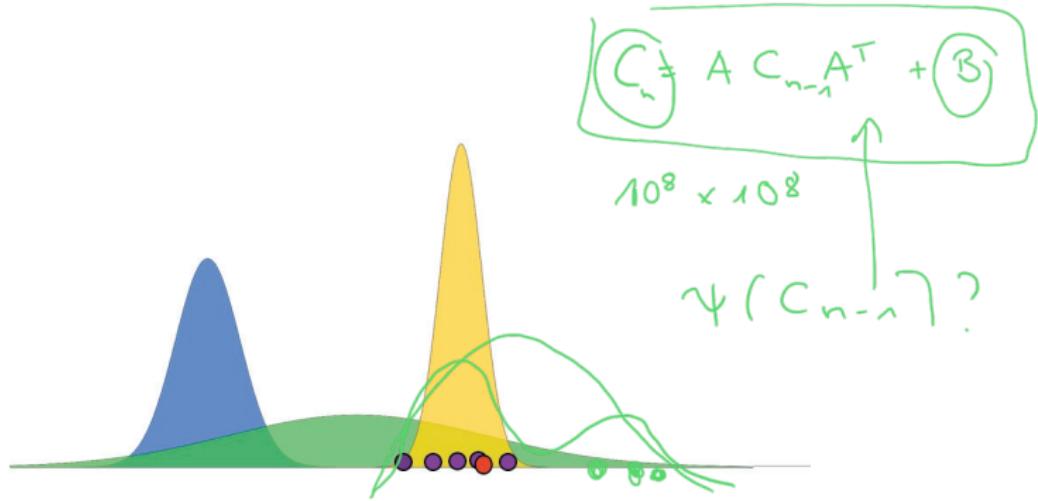
$$\mathcal{N}(\hat{\mathbf{m}}_1, \hat{\mathbf{C}}_1) \text{ with } \hat{\mathbf{m}}_1 \approx \frac{1}{M} \sum_{i=1}^M \hat{\mathbf{z}}_1^i = \frac{1}{M} \sum_{i=1}^M \Psi(\mathbf{z}_0^i)$$

$$\hat{\mathbf{C}}_1 \approx \frac{1}{M} \sum_{i=1}^M (\hat{\mathbf{z}}_1^i - \hat{\mathbf{m}}_1)(\hat{\mathbf{z}}_1^i - \hat{\mathbf{m}}_1)^\top$$

Ensemble Kalman filter



Ensemble Kalman filter



$$\mathcal{N}(\textcolor{blue}{m}_1, \textcolor{blue}{C}_1) \text{ with } \textcolor{blue}{m}_1 \approx \frac{1}{M} \sum_{i=1}^M \textcolor{blue}{z}_1^i$$

$$\textcolor{blue}{C}_1 \approx \frac{1}{M} \sum_{i=1}^M (\textcolor{blue}{z}_1^i - \textcolor{blue}{m}_1)(\textcolor{blue}{z}_1^i - \textcolor{blue}{m}_1)^\top$$

Ensemble Kalman filter

Goal: approximate $\pi(\textcolor{blue}{z}_n | \textcolor{blue}{y}_{1:n})$

Ensemble Kalman filter

Goal: approximate $\pi(\textcolor{blue}{z}_n | \textcolor{blue}{y}_{1:n})$

Approach: propagate samples $\hat{\textcolor{blue}{z}}_{n+1}^i$ with Kalman formula

$$\textcolor{blue}{z}_{n+1}^i = \hat{\textcolor{blue}{z}}_{n+1}^i - \textcolor{blue}{K}_{n+1}(H\hat{\textcolor{blue}{z}}_{n+1}^i - \tilde{\textcolor{red}{y}}_{n+1}^i)$$

Ensemble Kalman filter

Goal: approximate $\pi(z_n | y_{1:n})$

$F(x)$ want $F_n(x)$

Approach: propagate samples \hat{z}_{n+1}^i with Kalman formula

$$z_{n+1}^i = \hat{z}_{n+1}^i - K_{n+1}(H\hat{z}_{n+1}^i - \tilde{y}_{n+1}^i)$$

Need: perturbed observations

$$H z - y$$

$$+ \begin{cases} \sim N(0, R) \\ y = h(z) + \epsilon \end{cases}$$

$$\tilde{y}_{n+1}^i = y_{n+1} + \epsilon_{n+1}^i$$

with $\epsilon_{n+1}^i \sim \mathcal{N}(0, R)$ i.i.d. to get the correct mean and covariance
in the linear case for $M \rightarrow \infty$

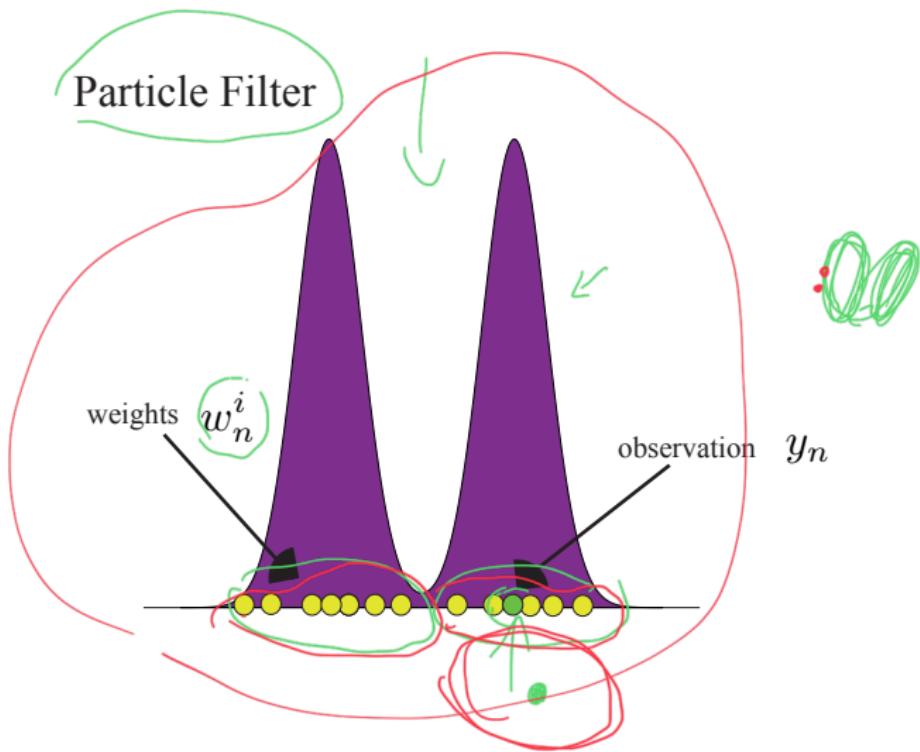
Ensemble Kalman filter

Works well in practice: e.g., EnKF is used for operational NWP
for z_n^i with dimension 10^8 only using $M = 100$

Yet: mathematical foundation largely missing

Recent study: accuracy results for EnKF for idealized setting:
 $H = Id$ and observational error small

Particle filter



Particle filter

Problem: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n})$ to approximate posterior via

$$\text{prior } \pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(z - (\mathbf{z}_n^i))$$

Idea: sampling from $\pi(\mathbf{z}_n | \mathbf{y}_{1:n-1})$ instead i.e.,

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \sum_{i=1}^M w_n^i \delta(z - \hat{\mathbf{z}}_n^i)$$

Bayes:

$$\pi(\mathbf{z}_{n+1} | \mathbf{y}_{1:n}) \propto \pi(\mathbf{y}_n | \mathbf{z}_n) \pi(\mathbf{z}_n | \mathbf{y}_{1:n-1}) \quad (1)$$

Weighting: unnormalized weights

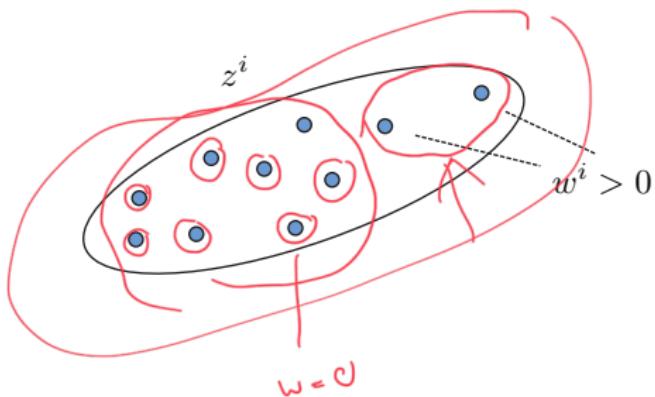
$$\tilde{w}_n^i = \pi(\mathbf{y}_n | \mathbf{z}_n^i) w_{n-1}^i \text{ with } w_0^i = \frac{1}{M}$$

and normalized weights

$$w_n^i = \frac{\tilde{w}_n^i}{\sum_{j=1}^M \tilde{w}_n^j}$$

Particle collapse

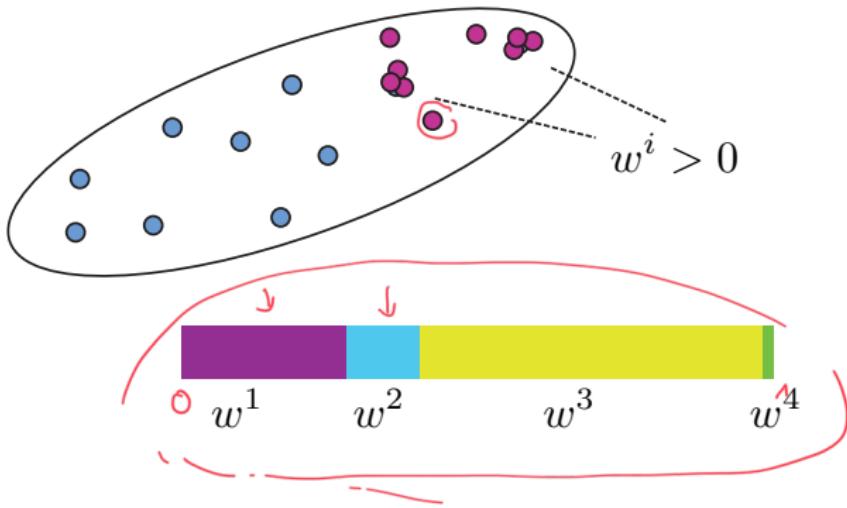
Degeneracy
of Particle Filter



Resampling

• y

Resampling



Resampling

Problem: weights w_n^i become very small

Ansatz: resampling

Input: w_n^i

For($k = 1 : M$)

1. Draw a number $u \in [0, 1]$ from the uniform distribution $U[0, 1]$
2. Compute $i^* \in \{1, \dots, M\}$ which satisfies

$$i^* = \arg \min_{i \geq 1} \sum_{j=1}^i w_j \geq u \quad (2)$$

3. Set $\xi_{i^*} = \xi_{i^*} + 1$

Return ξ_i

Still a lot of challenges....

Ansatz: approximative via empirical measure

$$\pi(\mathbf{z}_n | \mathbf{y}_{1:n}) = \frac{1}{M} \sum_{i=1}^M \delta(z - \mathbf{z}_n^i)$$

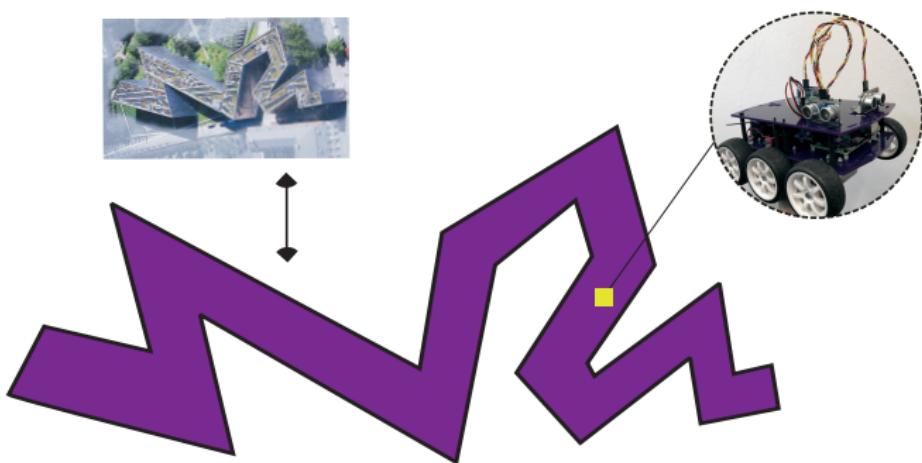
where

$$\mathbf{z}_n^i \sim \pi(\mathbf{z}_n | \mathbf{y}_n)$$

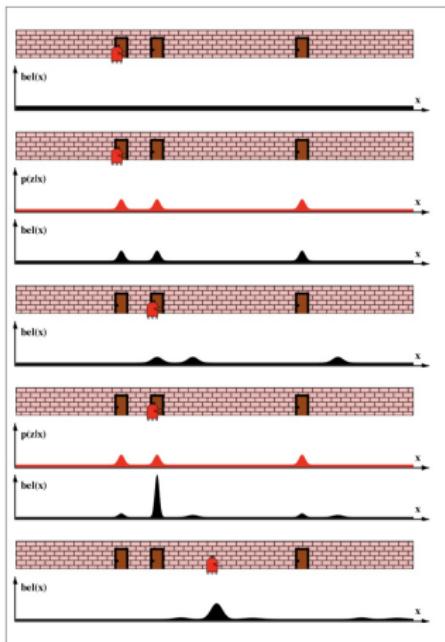
Monte Carlo approximation leads to a variety of filters e.g.,

- Particle filters (**curse of dimensionality**)
- Ensemble Kalman filter (**underlying Gaussian assumption**)

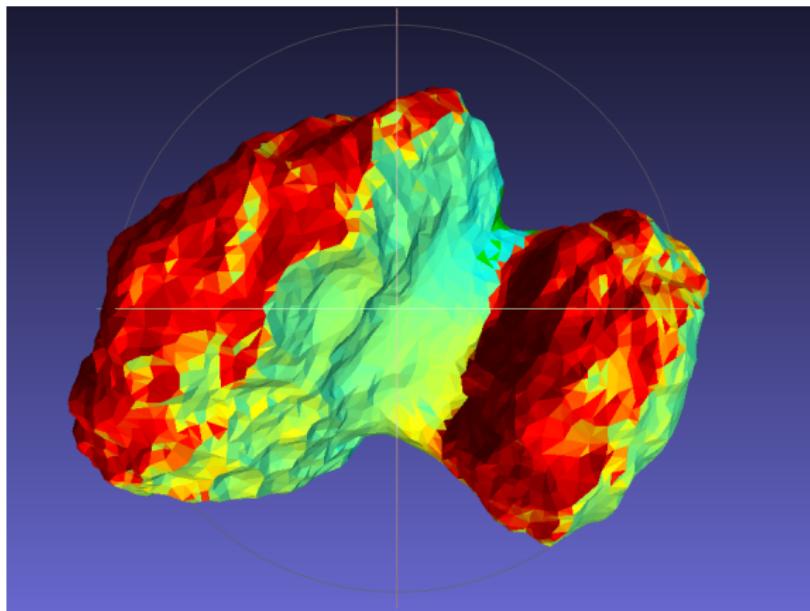
Same Math different Applications



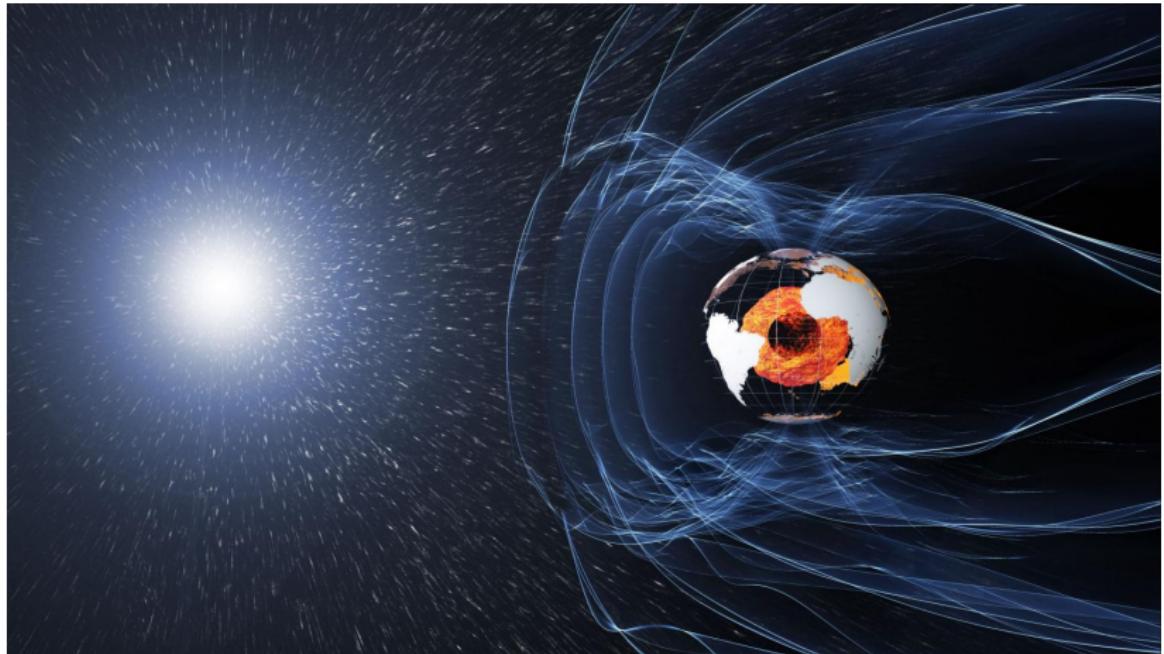
Simultaneous state and parameter estimation



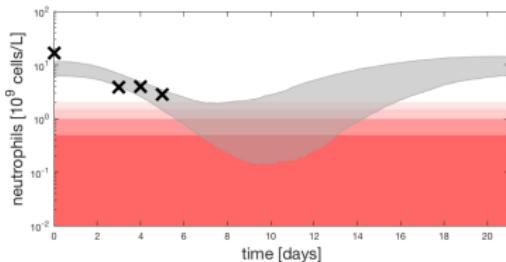
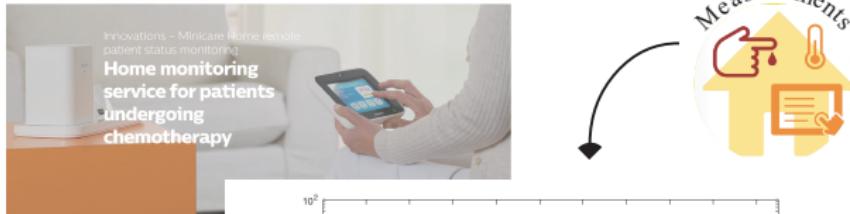
Thermophysical modeling



Space weather



Pharmacology

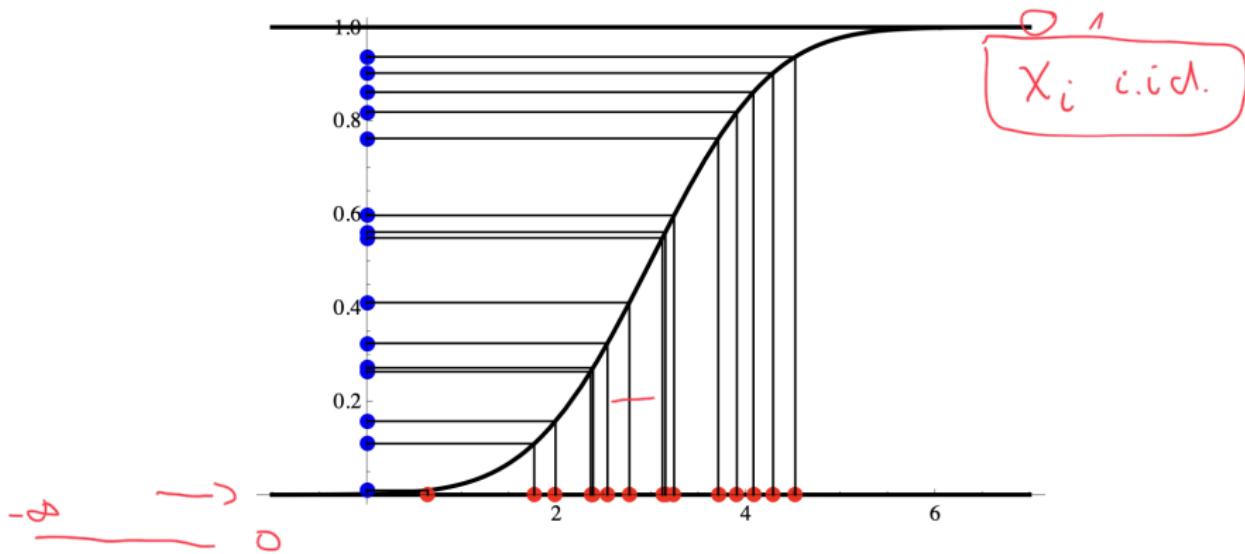


Personalised medicine

Maximum-spacing method

Maximum-spacing method

Lemma: Let the cdf F_θ be continuous and strictly monotone increasing. Under \mathbb{P}_θ the random variables $F_\theta(X_1), \dots, F_\theta(X_n)$ are independent and uniformly distributed on the $(0, 1)$ interval.



Proof

Prof: The independence follows from the independence of X_1, \dots, X_n

to show: $F_\theta(X_i)$ is uniformly distributed on $(0,1)$ $\forall i$

Note: $F_\theta(X_i)$ is assuming values in $(0,1)$

Let $x \in (0,1)$ then the following holds:

$$P_\theta[F_\theta(X_i) \leq x] = P_\theta[X_i \leq F_\theta^{-1}(x)] = F_\theta(F_\theta^{-1}(x)) = x$$

$\Rightarrow F_\theta(X_i)$ uniform



□

$$F_\theta(X_i): \Omega \rightarrow E = (0,1)$$

Proof

Maximum-spacing method

Lemma: Let $z_1, \dots, z_k \in [0, 1]$ be numbers that are subject to the condition $z_1 + \dots + z_k = 1$. Then

$$z_1 \cdot \dots \cdot z_k \leq \frac{1}{k^k}. \quad (3)$$

Equality is attained only if all the numbers are equal to $\frac{1}{k}$.

Beispiel

$k=2$

$$\begin{cases} z_1 = 0.9 & z_2 = 0.1 \\ z_1 = 0.5 & z_2 = 0.5 \end{cases}$$

$$z_1 \cdot z_2 = 0.09$$

$$\frac{1}{2^2} = \frac{1}{4} = 0.25$$

$$z_1 \cdot z_2 = 0.25$$

$$z_1 = \frac{1}{n}, z_2 = \frac{1}{n}, \dots, z_n = \frac{1}{n}, \dots$$

Example

Preparation for Max. Spacing Estimator. Given set sample (x_1, \dots, x_n)

Consider order statistic:

$$x_{(1)} \leq \dots \leq x_{(n)}$$



$$x_{(0)} := -\infty, x_{(n+1)} := +\infty$$

$\sim \theta$ is characterised by the so called spacings

unknown
 $D_i(\theta) = F_\theta(x_{(i)}) - F_\theta(x_{(i-1)})$ $i = 1, \dots, n+1$

Note: $D_i(\theta)$ is a realisation of the spacings of n RVs that are distributed on $(0,1)$



of the same size the more uniform $F_\theta(x_{(i)})$ all are spacings should be roughly

Maximum-spacing method

Lemma: The maximum-spacing method is defined via

$$\hat{\theta}_{MS} = \arg \max_{\theta \in \Theta} \prod_{i=1}^{n+1} (F_\theta(x_{(i)}) - F_\theta(x_{(i-1)})) \quad (4)$$

$\Pr_{\theta}[F_\theta(x_i) \leq x] = x$

$$F_\theta(F_\theta^{-1})$$

$$p = 0.2 \rightsquigarrow x_1, \dots, x_n$$

$$\hat{\theta} = \boxed{\tilde{p} = 0.8}$$

Example

Example

Unbiased estimators

Estimator

Def:

- An estimator is an arbitrary (Borel-measurable) function

$$\hat{\theta} : \mathcal{X} \rightarrow \Theta, \quad x \mapsto \hat{\theta}(x) \quad (5)$$

- An estimator $\hat{\theta}$ is called unbiased, if

$$\mathbb{E}_{\theta}[\hat{\theta}(X)] = \theta \quad (6)$$

for all $\theta \in \Theta$.

- The bias of an estimator $\hat{\theta}$ is

$$\text{Bias}_{\theta}(\hat{\theta}) = \mathbb{E}_{\theta}[\hat{\theta}(X)] - \theta \quad (7)$$

$$\frac{1}{N-1}$$

Note: $\text{Bias}_{\theta}(\hat{\theta})$ is a function in $\hat{\theta}$

$$\hat{\theta} - \theta = 0$$

Example

Example

Example

Mean square error

Def: Let $\Theta = (a, b) \subset \mathbb{R}$ be an interval. The mean square error (MSE) of an estimator $\hat{\theta} : \mathcal{X} \rightarrow \Theta$

$$\text{MSE}_\theta(\hat{\theta}) = \mathbb{E}_\theta[(\hat{\theta}(X) - \underbrace{\theta}_\text{green})^2] \quad (8)$$

Mean square error

Lemma: The relationship between the mean square error (MSE) of an estimator $\hat{\theta} : \mathcal{X} \rightarrow \Theta$ and the BIAS is given by

$$\text{MSE}_\theta(\hat{\theta}) = \text{Var}_\theta \hat{\theta} + (\text{Bias}_\theta(\hat{\theta}))^2 \quad (9)$$

$$\text{MSE}_\theta(\hat{\theta}) = \text{Var}_\theta \hat{\theta} + 0$$



Proof

Proof

Consistently better

Let $\hat{\Theta}_1$ and $\hat{\Theta}_2$ be two estimators. The estimator θ_1 is consistently better than θ_2 if,

$$MSE_{\theta}(\hat{\theta}_1) \leq MSE_{\theta}(\hat{\theta}_2) \quad \forall \theta \in \Theta \quad (10)$$

Minimum-variance unbiased estimator

Def: An unbiased estimator $\hat{\theta}$ is called minimum-variance unbiased estimator if all unbiased estimators $\tilde{\theta}$ the following inequality holds

$$\text{Var}_\theta \hat{\theta} \leq \text{Var}_\theta \tilde{\theta} \quad (11)$$

for all $\theta \in \Theta$.

Minimum-variance unbiased estimator

Lemma: Let $\hat{\theta}_1, \hat{\theta}_2 : \mathcal{X} \rightarrow \Theta$ are two minimum-variance unbiased estimator the

$$\hat{\theta}_1 = \hat{\theta}_2 \quad \text{almost surely under } \mathbb{P} \text{ for all } \theta \in \Theta \quad (12)$$

for all $\theta \in \Theta$.

Proof

Proof

Proof

Sufficient statistic

Def: A function $T : \mathcal{X} \rightarrow \mathbb{R}^r$ is called a sufficient statistic if the function

$$\theta \mapsto \mathbb{P}_\theta[X = x | T(X) = t] \quad (13)$$

is constant for all $x \in \mathcal{X}$ and for all $t \in \mathbb{R}^r$, i.e.,

$$\mathbb{P}_{\theta_1}[X = x | T(X) = t] \mathbb{P}_{\theta_2}[X = x | T(X) = t] \quad (14)$$

for all $t \in \mathbb{R}^r$ and all $\theta_1, \theta_2 \in \Theta$ with $\mathbb{P}_{\theta_1}[T(X) = t] \neq 0$ and
 $\mathbb{P}_{\theta_2}[T(X) = t] \neq 0$