International
Association
of Oil & Gas
Producers

virtual engagement session for OSRC 2021
fixed offshore structures

pre-meeting video 1
basic probability & statistics

# content – notation & terminology

— notation – typically compressed

— discreet variable  v  continuous variable

— probability mass function  v  probability density function

— random variable  v  deterministic variable

— parameters  v  variables

— independently and identically distributed ($iid$)

— addition rule,  multiplication rule  &  chain rule

— parent  v extreme distributions

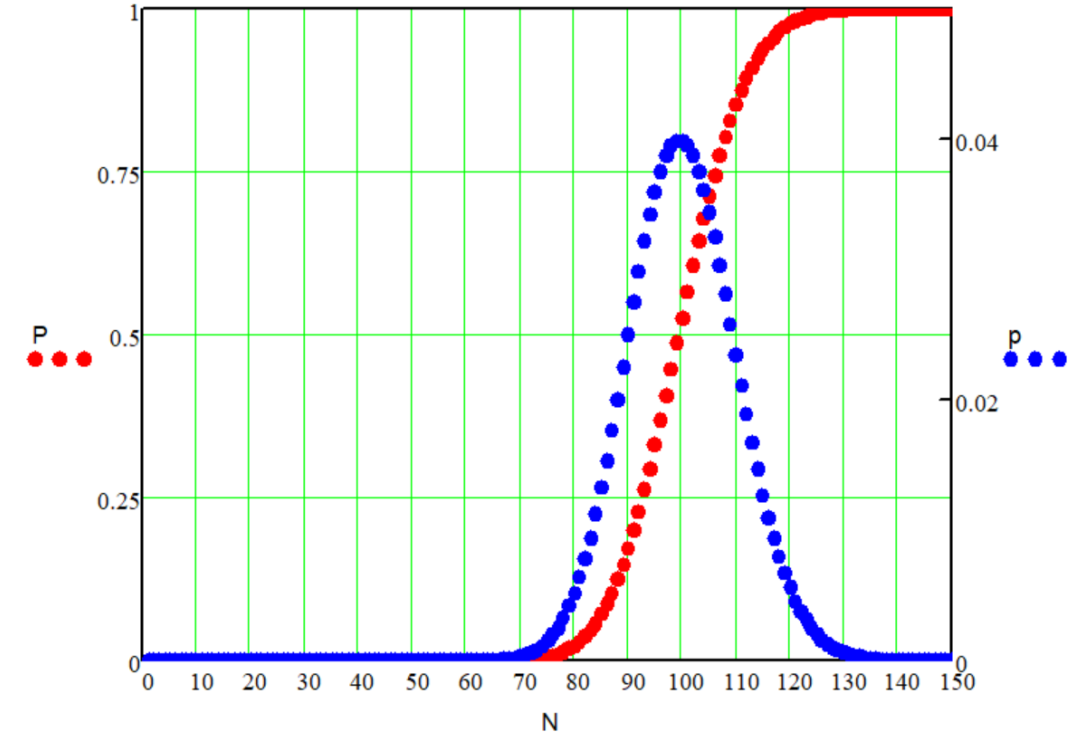— marginal probability,  conditional probability, Bayes' rule  &  Bayesian inference
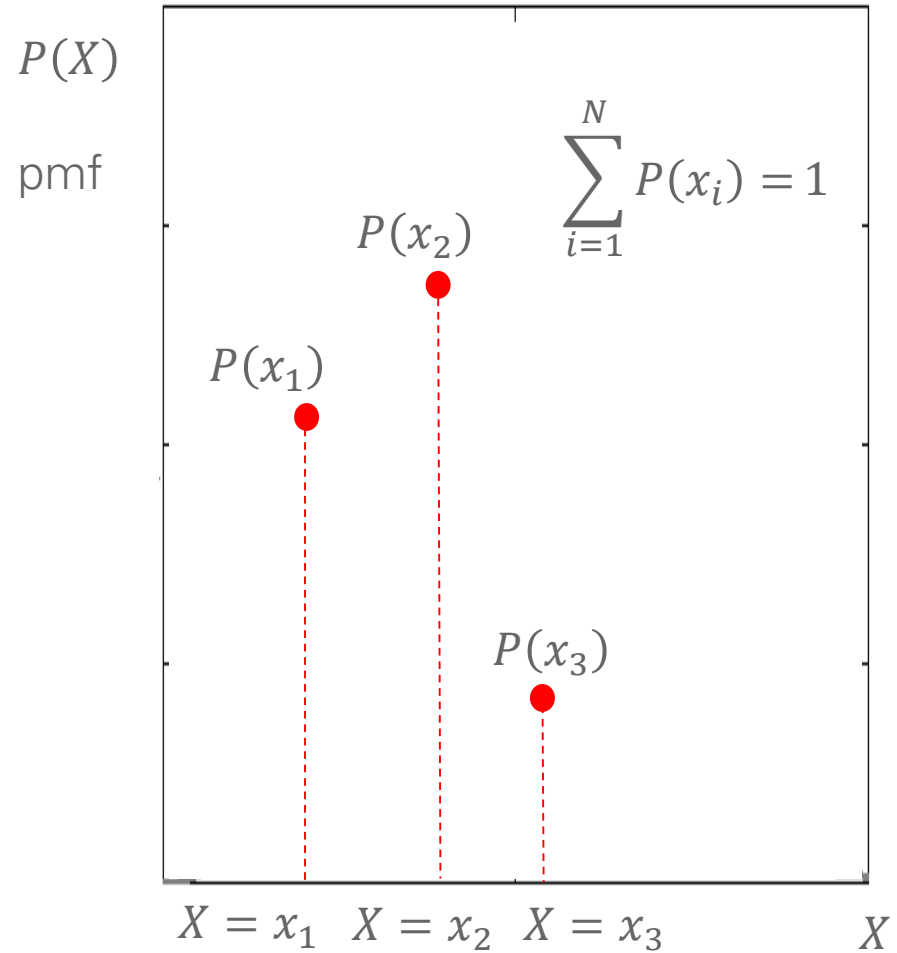
# probability – definition

probabilities, $P$ ,are numerical quantities…

— defined on a set of "outcomes" ……………eg $P(C \leq 20m), \ P(20 < C \leq 22m), \ P(C > 22m)$

— non-negative

— additive over mutually exclusive outcomes ……………eg $P(C \leq 20m) + P(20 < C \leq 22m)$

— sum to 1 over all possible mutually exclusive outcomes

probability of an event A $= P(A) = \dfrac{\text{number of ways event A can occur}}{\text{total number of possible outcomes}}$
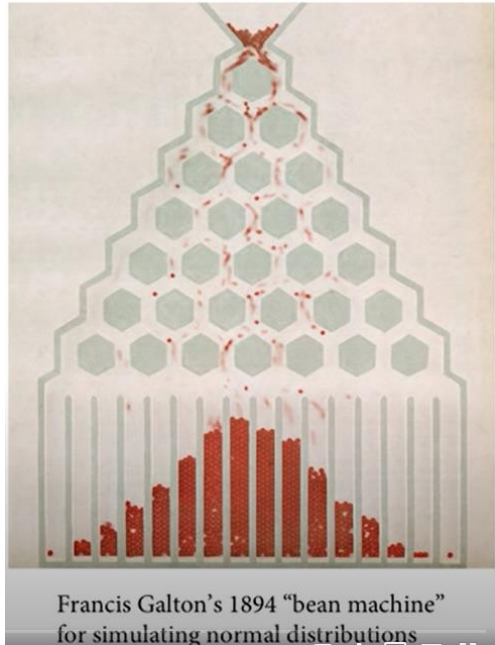
# probability mass function (pmf) & probability



$P(X)$

pmf

$$\sum_{i=1}^{N} P(x_i) = 1$$

$P(x_2)$

$P(x_1)$

$P(x_3)$

$X = x_1 \quad X = x_2 \quad X = x_3 \qquad X$

$X -$ the variable

$x_i -$ a specific value of the variable

# probability density function (PDF) – in 1 dimension



International
Association
of Oil & Gas
Producers

$X$ – the variable

$x$ – a specific value of the variable

Francis Galton's 1894 "bean machine"
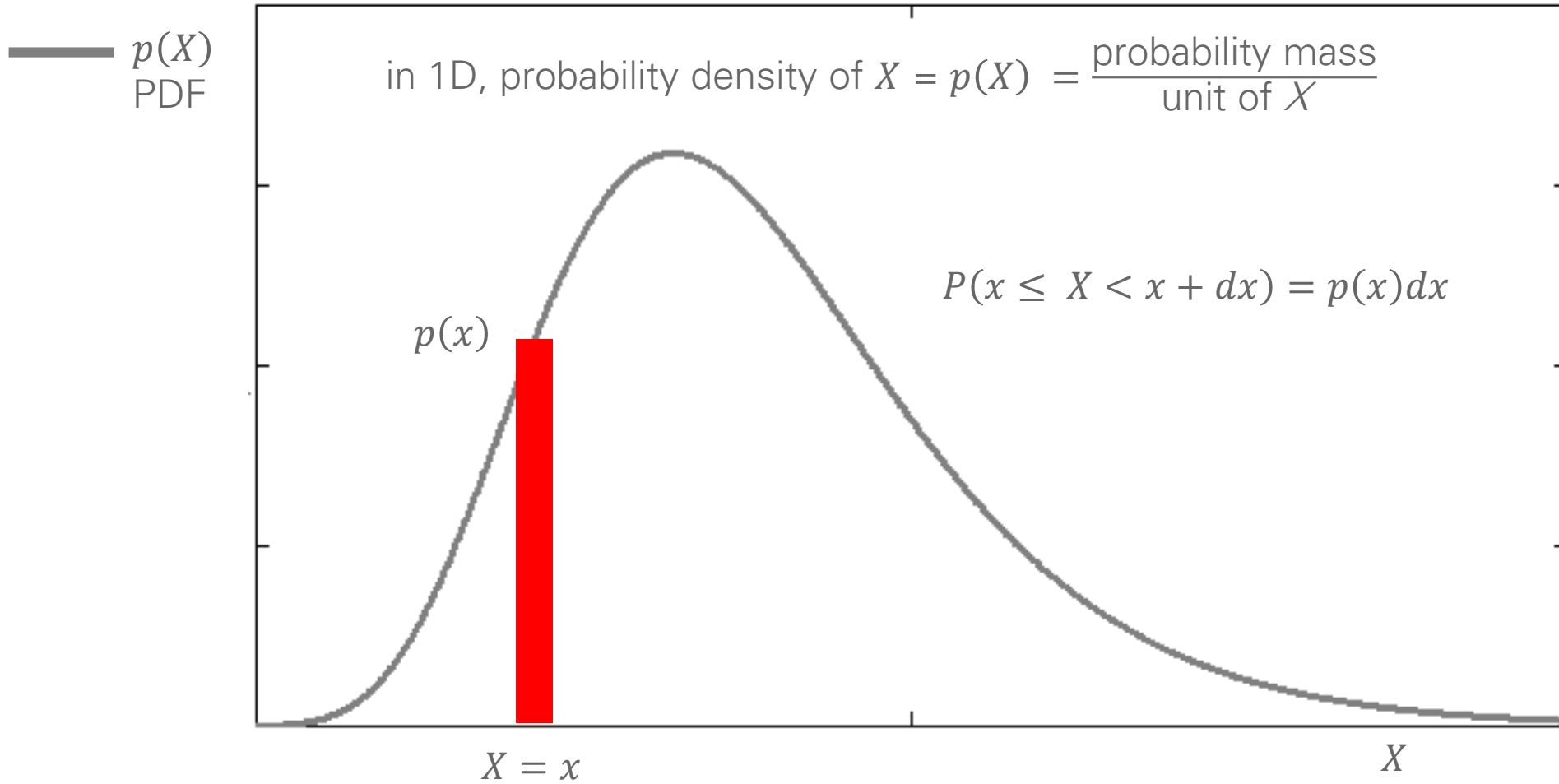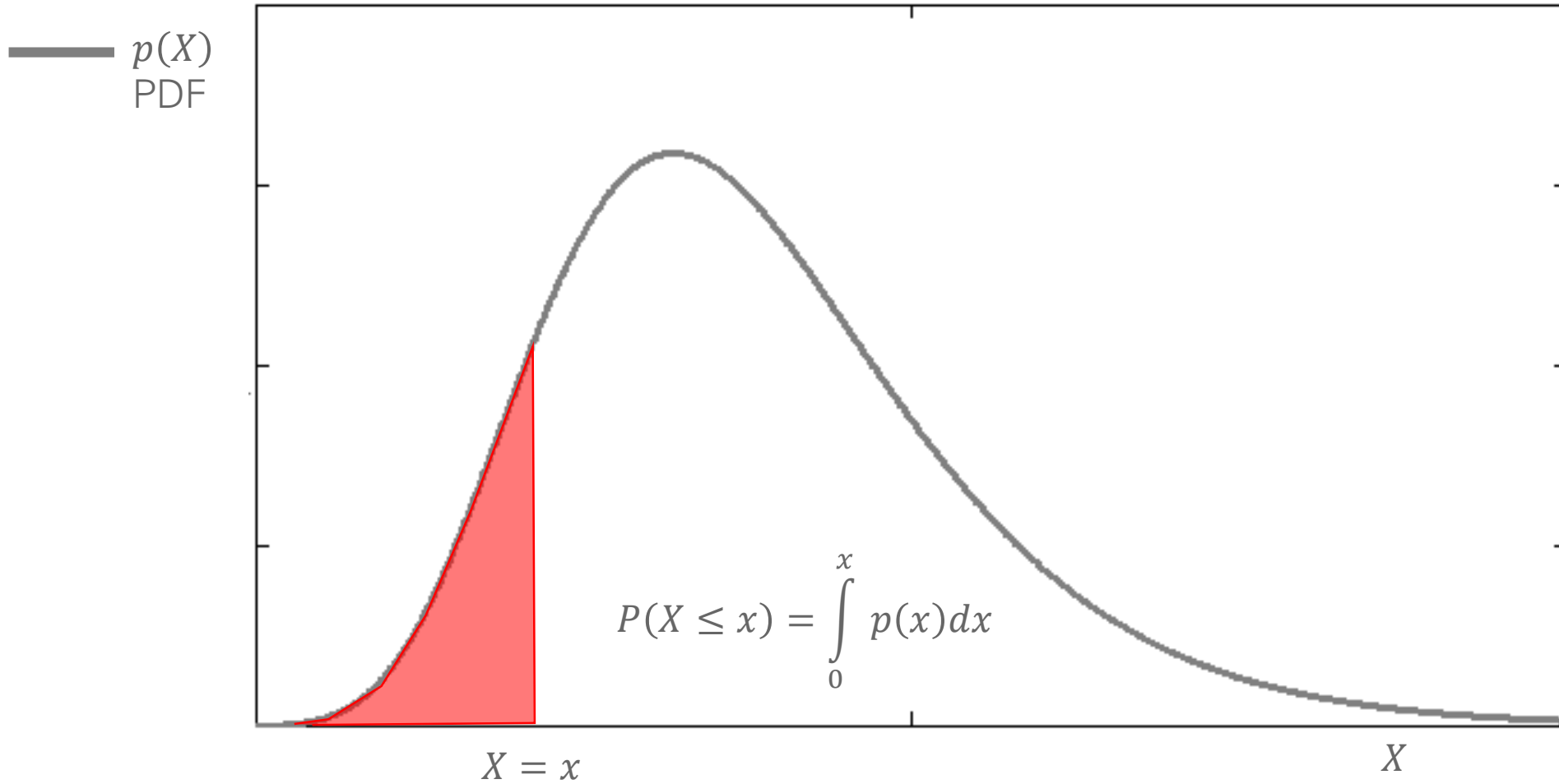for simulating normal distributions

$P(x_i < X < x_{i+1})$

physical quantities that are expected
to be the sum of many independent
processes often have distributions
that are nearly normal
(central limit theorem)

$$p(X)$$
PDF

$$p(x_i) = \frac{n_i}{\sum_{j=1}^{N} n_j}$$

Histogram of sampled values
Exact distribution

$n_i$

$x_1$

$x_i$

$x_N$

# probability density & probability



$p(X)$
PDF

in 1D, probability density of $X = p(X) = \dfrac{\text{probability mass}}{\text{unit of } X}$

$P(x \leq X < x + dx) = p(x)dx$

$p(x)$

$X = x$

$X$

# probability density & probability of non-exceedance



$p(X)$
PDF

$$P(X \leq x) = \int\limits_{0}^{x} p(x)dx$$

$X = x$

$X$

# probability density & probability of non-exceedance



$p(x)$

aka
PDF
$f_X(x)$

$P(X \leq x)$
probability of
non-exceedance

aka
CDF cumulative probability function
$F_X(x)$

$$P(X \leq x) = \int\limits_0^x p(x)dx$$

$X = x$

$X$

# probability density & probability of exceedance

$p(x)$
pdf

$P(X > x)$
probability of
exceedance

aka
$Q_X(x)$
$1 - F_X(x)$
CCDF complementary cumulative
probability function

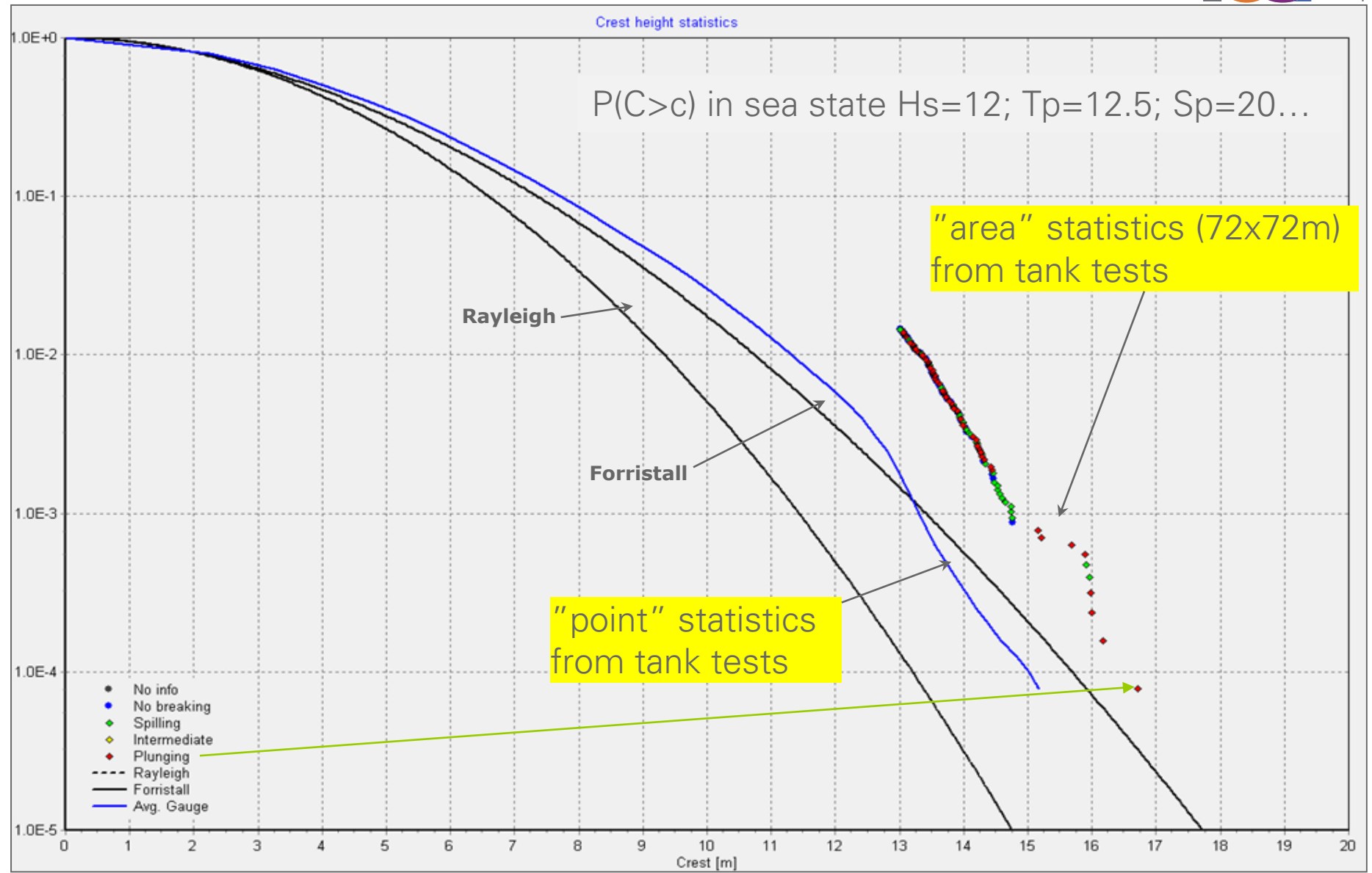$$P(X > x) = \int\limits_{x}^{\infty} p(x)dx$$
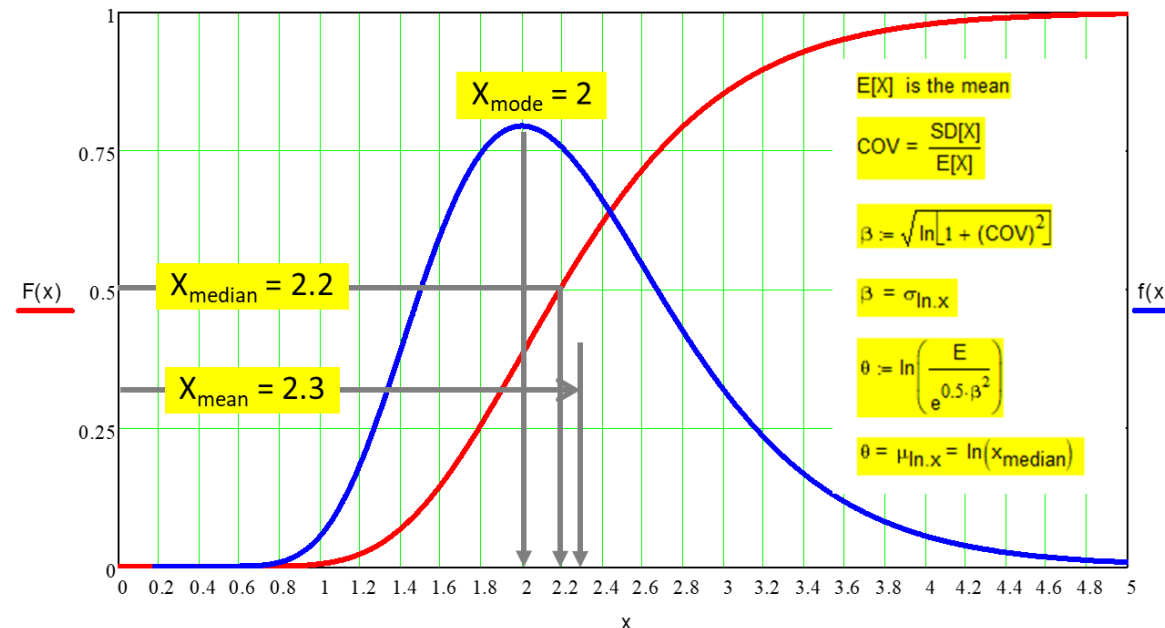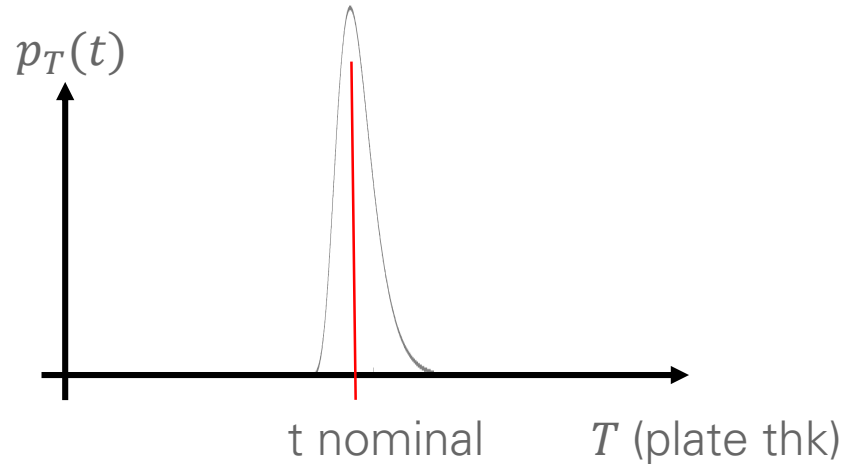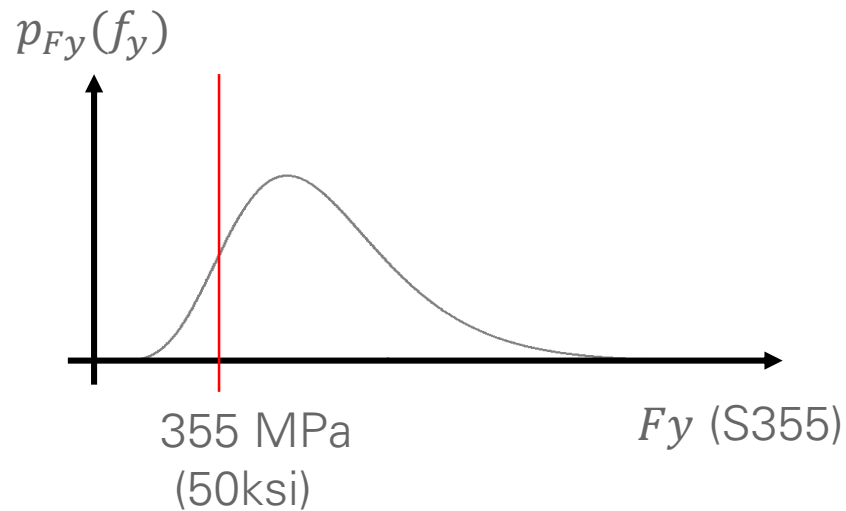
1

0

$X = x$

$X$

# probability of exceedance

$P(C > c)$

probability of exceedance of individual crest ht (C) in a given sea state

log scale rather than linear (0 to 1) shows the tail in more detail at extreme values



Crest height statistics

P(C>c) in sea state Hs=12; Tp=12.5; Sp=20...

"area" statistics (72x72m) from tank tests

Rayleigh

Forristall

"point" statistics from tank tests

- No info
- No breaking
- Spilling
- Intermediate
- Plunging
- Rayleigh
- Forristall
- Avg. Gauge

Crest [m]

International Association of Oil & Gas Producers

# random variables ( v deterministic variables)



$p_{Fy}(f_y)$

355 MPa
(50ksi)

$Fy$ (S355)

$p_T(t)$

t nominal          $T$ (plate thk)

X$_{mode}$ = 2

X$_{median}$ = 2.2

X$_{mean}$ = 2.3

E[X] is the mean

$$COV = \frac{SD[X]}{E[X]}$$

$$\beta := \sqrt{\ln\left[1 + (COV)^2\right]}$$

$$\beta = \sigma_{\ln.x}$$

$$\theta := \ln\left(\frac{E}{e^{0.5 \cdot \beta^2}}\right)$$

$$\theta = \mu_{\ln.x} = \ln(x_{median})$$

F(x)

f(x)

## parameters

location  $\mu$
mean (expectation)
mode (mp)
median (P50)

scale
standard deviation $\sigma$
variance $\sigma^2$
COV  $\sigma/\mu$
dispersion $\beta = \sigma_{logx}$

shape $\xi$
tail properties

# standard probability density functions

**Continuous distributions**
Uniform
Normal
Lognormal
Gamma
Inverse-gamma
Chi-square
Inverse-chi-square
Scaled inverse-chi-square
Exponential
Laplace
Weibull
Wishart
Inverse-Wishart
LKJ correlation
t
Beta
Dirichlet
Logistic
Log-logistic

**Discrete distributions**
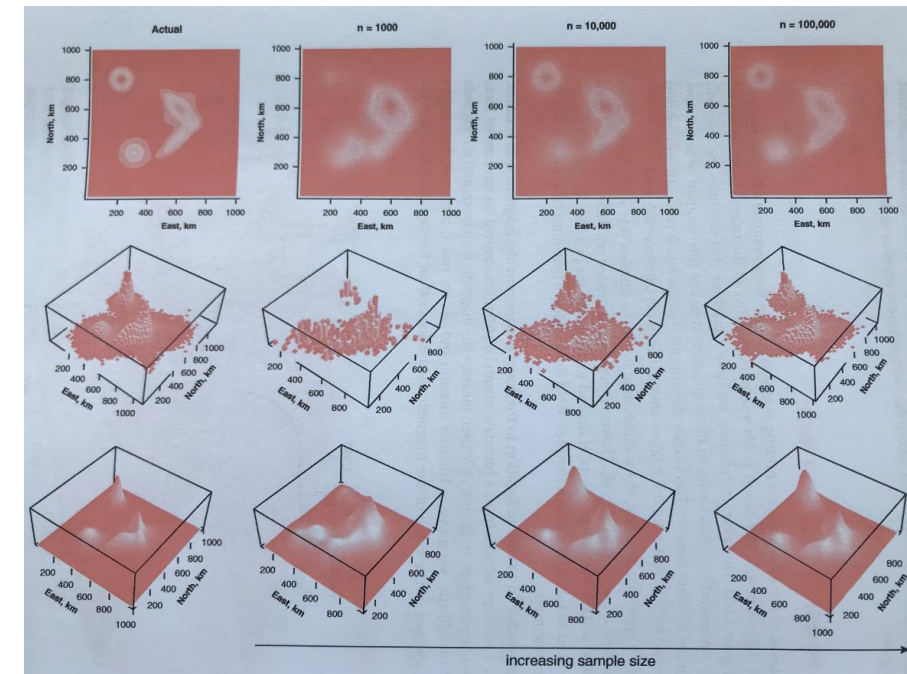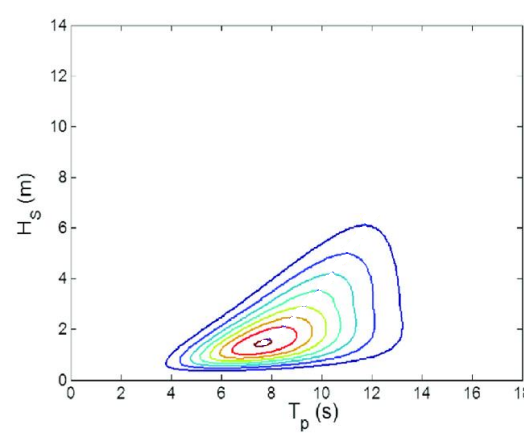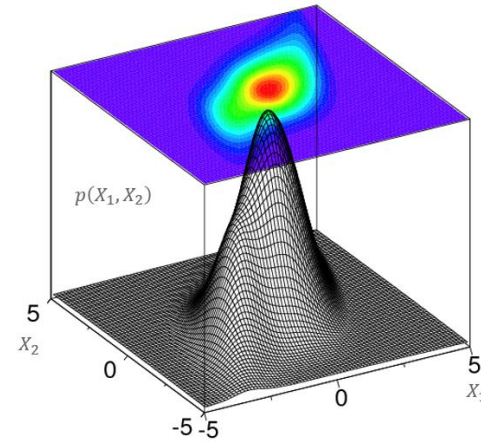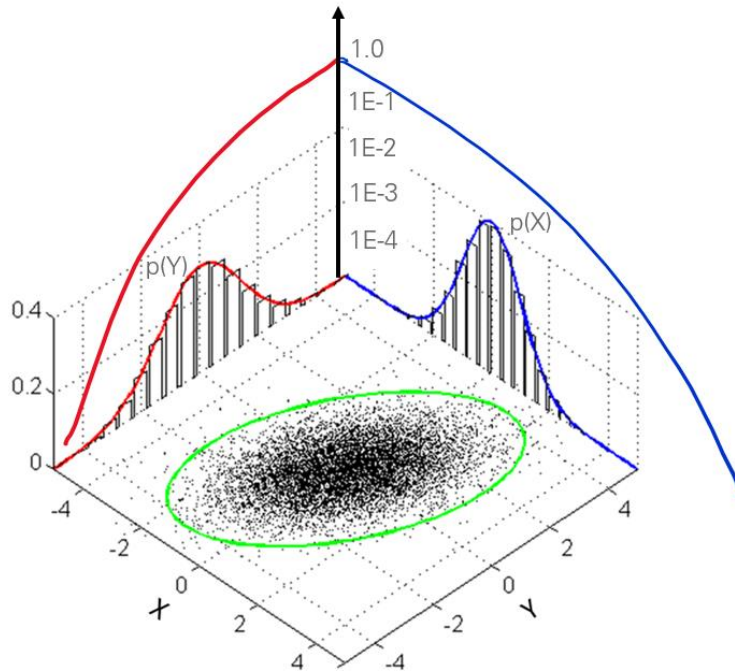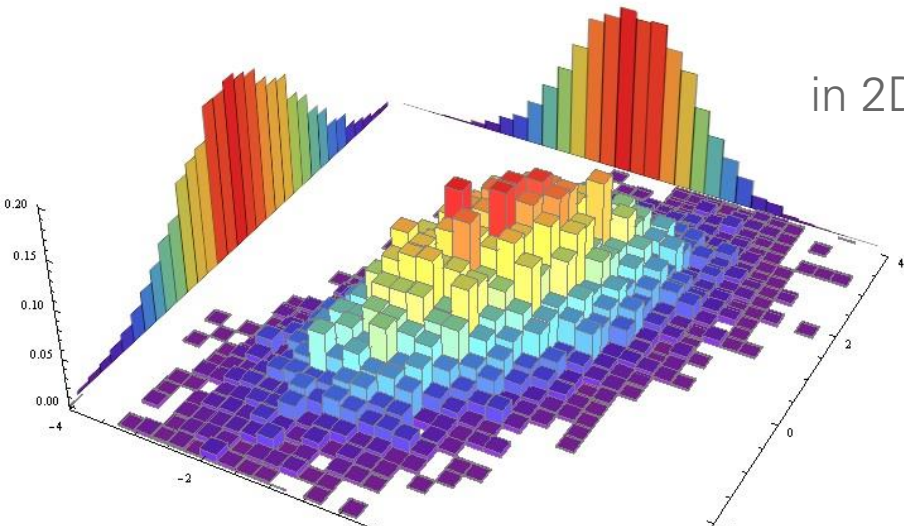Poisson
Binomial
Negative-binomial
Beta-binomial

$$p(x) = \frac{1}{b-a}$$

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

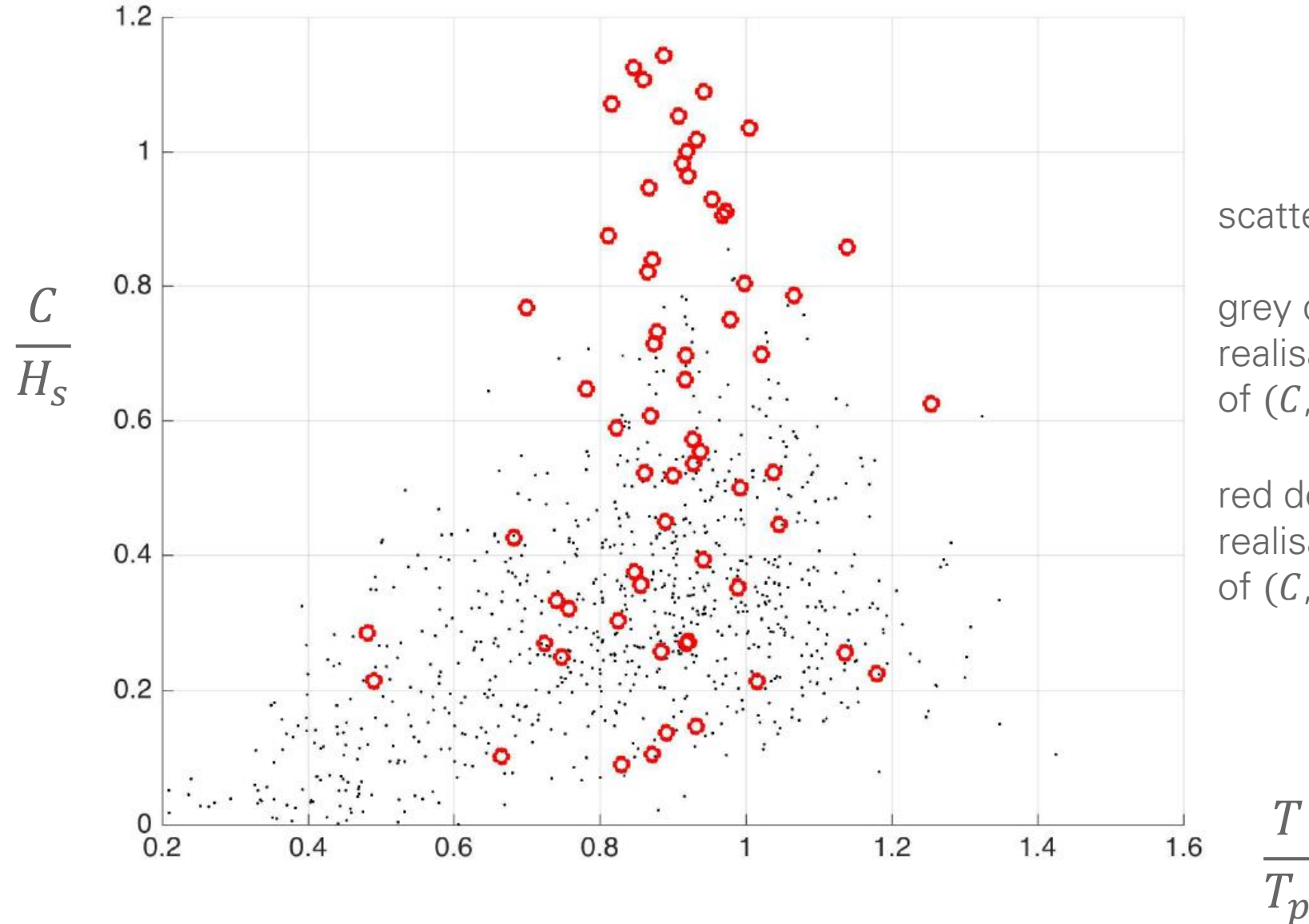$$p(x) = \frac{1}{\sigma x\sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(\log x - \mu)^2\right)$$

# joint probability density function - 2D



in 2D, probability density = $p(X,Y) = \dfrac{\text{probability mass}}{\text{unit of X} \times \text{unit of } Y}$



computing the pdf by sampling using MCMC (HMC)

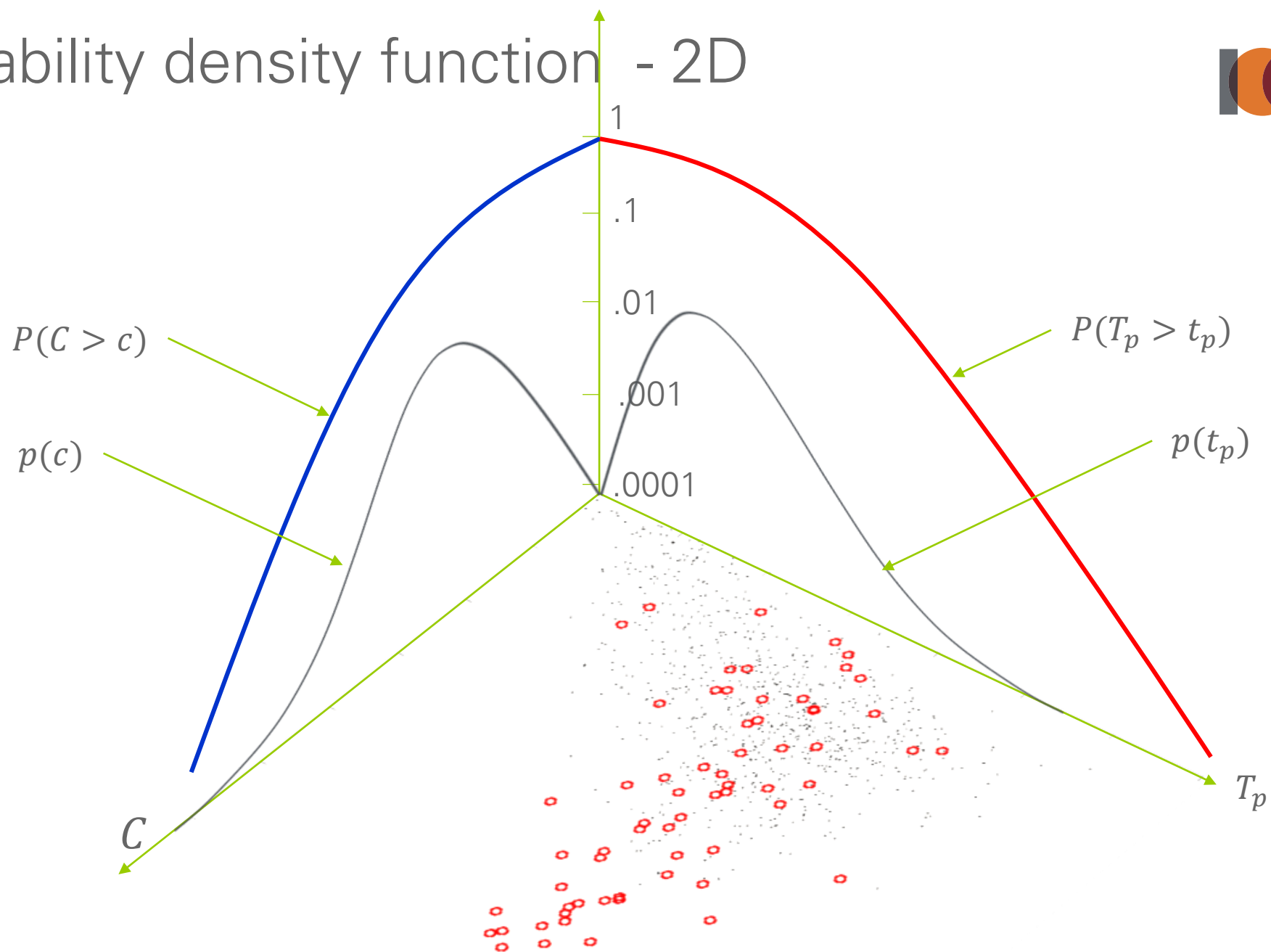# joint probability density function  - 2D



scatter plot

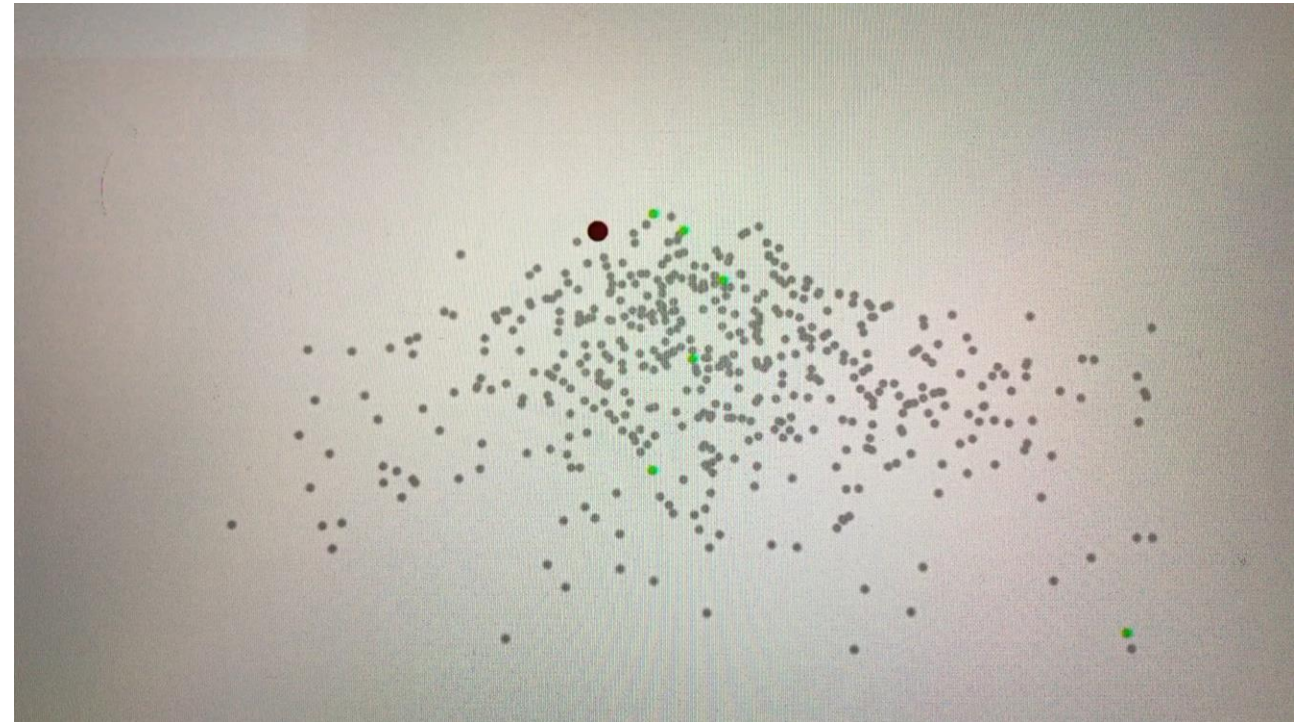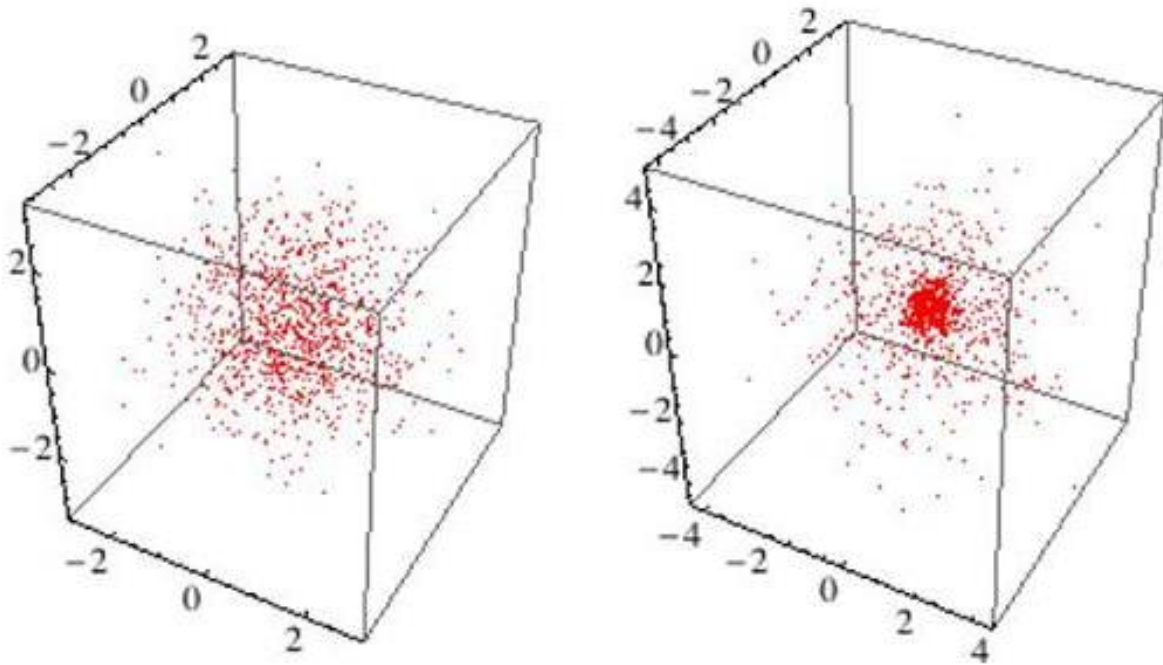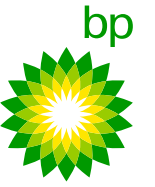grey dots =  Monte Carlo (random) realisations from the joint distribution of $(C, T_p)$

red dots =  stratified Monte Carlo realisations from the joint distribution of $(C, T_p)$

# joint probability density function - 2D

# joint probability density function  - 3D

# joint probability density function - 4D and 12D



$$p(T_p, \sigma_\theta, \gamma, \eta_{swl}, u, \theta_u, w, \theta_w \dots | H_s)$$

# addition rule



$P(A \text{ AND } B)$
$= P(A \cap B)$
$\neq 0$
$= \text{not mutually exclusive}$

$P(A)$

$P(B)$

general addition rule

$P(A \text{ OR } B)$
$= P(A \cup B)$
$= P(A) + P(B) - P(A \cap B)$
$= P(A) + P(B) - P(A \text{ AND } B)$

$P(A \text{ AND } B)$
$= P(A \cap B)$
$= 0$
$= \text{mutually exclusive}$

$P(A)$

$P(B)$

specific addition rule (for mee)

$P(A \text{ OR } B)$
$= P(A \cup B)$
$= P(A) + P(B)$

Venn-pie (pie area=1)

# conditional probability

$P(A|B)$ = probability of event A (eg $2 < X < 4$) occurring given event B (eg $0 < X < 3$) has occurred

$P(E \leq \eta | H_{s_i} = h)$
probability of wave crest elevation $E$ not exceeding $\eta$
given the significant wave height for the $i$ th sea state $H_{s_i}$ equals $h$

$P(L > l | \alpha, \text{storm})$
probability of jacket base shear load $L$ exceeding $l$
given a (random) storm from direction $\alpha$ is occurring

$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)} = \frac{P(A \cap B)}{P(B)} = \frac{\text{intersection}}{\text{normalised}}$$

makes $P(A|B)$ a valid probability

$P(A) = 0.2$

$P(A \cap B) = 0.1$

$P(B) = 0.3$

$P(A \cup B) = 0.4$

Venn-pie (pie area=1)

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.1}{0.3} = 33\%$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.1}{0.2} = 50\%$$

# multiplication rule

conditional probability of event A occurring given event B has occurred

$$P(A|B) = \frac{P(A \text{ AND } B)}{P(B)} = \frac{P(A \cap B)}{P(B)}$$

$$\boxed{P(A \text{ AND } B) = P(A|B) \times P(B)} \quad \text{eqn (1)}$$

$$P(B|A) = \frac{P(B \text{ AND } A)}{P(A)} = \frac{P(B \cap A)}{P(A)}$$

$$\boxed{P(B \text{ AND } A) = P(B|A) \times P(A)} \quad \text{eqn (2)}$$

if event A is independent from event B then

$$P(A|B) = P(A) \qquad \text{eqn (3)}$$
$$P(B|A) = P(B) \qquad \text{eqn (4)}$$

$$\boxed{\begin{array}{l} P(A \text{ AND } B) = P(A) \times P(B) \\ P(\text{B AND } A) = P(B) \times P(A) \end{array}} \quad \begin{array}{l} \text{from (1) \& (3)} \\ \text{from (2) \& (4)} \end{array}$$

# Bayes' rule

from last slide - as event A and event B occur together then

$$P(A \text{ AND } B) = P(B \text{ AND } A)$$

$$P(A|B) \times P(B) = P(B|A) \times P(A)$$

$$\frac{P(A|B) \times P(B)}{P(B)} = \frac{P(B|A) \times P(A)}{P(B)}$$

$$\boxed{P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}}$$

Bayesian inference (probability densities)…

$$p(\boldsymbol{\theta}|y) = \frac{p(y|\boldsymbol{\theta}) \times p(\boldsymbol{\theta})}{p(y)}$$



$P(A) = 0.2$

$P(A \cap B) = 0.1$

$P(B) = 0.3$

$P(A \cup B) = 0.4$

Venn-pie (pie area=1)

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.1}{0.3} = 33\%$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{0.1}{0.2} = 50\%$$

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} = \frac{0.5 \times 0.2}{0.3} = \frac{0.1}{0.3} = 33\%$$

# chain rule

from law of total probability

$$P(A) = \sum_{i=1}^{N} P(A \text{ AND } B) = \sum_{i=1}^{N} P(A \cap B_i)$$

from conditional probability

$$P(A \cap B_i) = P(A|B_i) \times P(B_i)$$

substitute 2nd in 1st

$$P(A) = \sum_{i=1}^{N} P(A|B_i) \times P(B_i)$$

extending from discreet events to continuous random variables

$$P(A) = \int_{B} P(A|X = x) \times p(x) \, dx$$

# marginal probability

aka marginalising or integrating out



marginal probability density of $y$

$$p(y) = \int_0^\infty p(y, \theta) d\theta$$

(area =1)

joint probability density of $y, \theta$

$$p(y, \theta)$$

(vol =1)

$$p(\theta) = \int_0^\infty p(y, \theta) dy$$

(area =1)

marginal probability density of $\theta$

https://www.youtube.com/watch?v=F97Qf5FGt5M

$\mu_y$

$\mu_\theta$

$f$

$b$

$y$

$\theta$

# conditional probability



marginal probability density of $y$

$$p(y) = \int_0^\infty p(y, \theta) d\theta$$

(area =1)

joint probability density of $y, \theta$

$$p(y, \theta)$$

(vol =1)

$$p(\theta) = \int_0^\infty p(y, \theta) dy$$

(area =1)

marginal probability density of $\theta$

probability density of $y$ conditional on $\theta = b$

$$p(y | \theta = b) = \frac{p(y, \theta = b)}{\int_0^\infty p(y, \theta = b) dy}$$

divide by marginal to give a valid probability (area =1)

# $P(H > h)$ in the long-term – by chain rule (aka convolution in ISO 19901-1)

$$P_{annum}(H > h) = \int_0^\infty P(H > h|H_s) \times p(H_s)dH_s$$



$$\frac{H}{H_{mpm}} = \frac{h}{H_{mpm}}$$

probability of exceedance of largest wave height, $P(H > h)$, in a given 3hr sea state $H_s$ (ie in the short term)

annual probability of the given sea state occurring (ie in the long term)

# probability of exceedance for largest in $N$ (random) events

probability of the crest of an individual wave ($E_1$) not exceeding a given value ($\eta$) in a given sea state $H_s = h$ is:

$$P(E_1 \leq \eta | H_s = h)$$

probability of the crest of another individual wave ($E_2$) not exceeding the same value ($\eta$) in $H_s = h$ is:

$$P(E_2 \leq \eta | H_s = h)$$

probability of the larger of crest elevations ($E_1$ and $E_2$) not exceeding a given value ($\eta$) in $H_s = h$ (assuming independence ie far apart in the sea state) is:

$$P\big((E_1 \leq \eta | H_s = h) \text{ and } (E_2 \leq \eta | H_s = h)\big) \;=\; P(E_1 \leq \eta) \times P(E_2 \leq \eta) \;=\; \prod_{i=1}^{2} P(E_i \leq \eta) \;=\; P\left(\max_{i=1,2}[E_i] \leq \eta\right)$$

probability of the largest crest elevation in $N$ waves exceeding a given value ($\eta$) in $H_s = h$ is

$$P\left(\max_{i=1,N}(E_i) > \eta \Big| H_s = h\right) \;=\; 1 - \prod_{i=1}^{N} (E_i \leq \eta | H_s = h)$$

NB probability of the smallest crest elevation in $N$ waves exceeding a given value ($\eta$) in $H_s = h$ is $\prod_{i=1}^{N} P(E_i > \eta)$

# distributions of individual & largest wave ht for a given sea state with duration 3hr (N=1000)



$P(H_i \leq h)$ = parent Rayleigh with $H_s = 10m$

$P(H_i \leq h)$

$P\left(\max_{i=1,1000}[H_i] \leq h\right) = P(H_i \leq h)^N$

$p\left(\max_{i=1,1000}[H] \leq h\right)$

$H_{mpm\_ind} = 0.5H_s$

$H_s = 10m$

$H_{mpm\_3hr} = 1.86H_s$

International Association of Oil & Gas Producers

# distribution of individual crest ht – given a 3hr sea stat



crest heights shown by red dots

3hour simulation of sea state

wave surface elevation ($\eta$)

pdf of individual crest height in a given sea state

crest height ($C$)

distribution of individual crest height in a given sea state

# distribution of largest crest ht – given a 3hr sea state



max crest heights (in 3 hrs) shown by red dots

3hour simulation of sea state

wave surface elevation ($\eta$)

mpm crest height, $C_{mpm}$, in a sea state
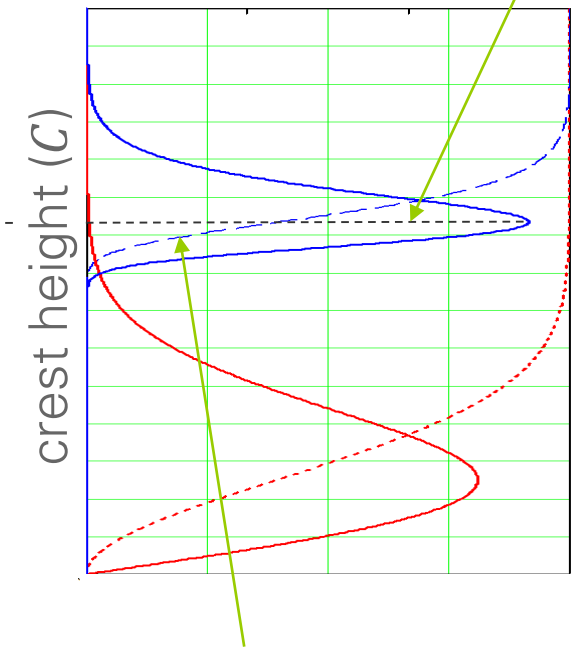
crest height ($C$)

distribution of largest crest height in a given sea state

# statistics of extremes – GEV and GPD

## Generalised Extreme Value $GEV$ distribution

$$P(X \leq x) = exp\left(-\left[1 + \xi\left(\frac{x-\mu}{\sigma}\right)\right]^{-1/\xi}\right) \text{ if } \xi \neq 0$$

$$P(X \leq x) = exp\left(-exp\left(\frac{x-\mu}{\sigma}\right)\right) \text{ if } \xi = 0$$

$\mu$ = location parameter
$\sigma$ = scale parameter
$\xi$ = shape parameter

*$GEV$ is distribution of extreme*



Gumbel $\xi = 0$
parent has exponential tail

Pareto (Fréchet) $\xi > 0$
parent has polynomial tail

Weibull $\xi < 0$
parent has upper end point = $[\sigma + \xi(u - \mu)]/|\xi|$

## Generalised Pareto distribution $GPD$

$$P(X - u > x | X > u) = \left[1 + \frac{\xi(x-u)}{\sigma + \xi(u-\mu)}\right]^{-1/\xi}$$

$\mu$ = location parameter
$\sigma$ = scale parameter
$\xi$ = shape parameter
$u$ = threshold

*$GPD$ is distribution of data points above a threshold*

distribution of extreme = (parent distribution)$^N$
tends to GEV asymptotic distribution as N becomes large

form of extreme depends on the form of the  tail of the parent distribution

# probability of exceedance (extreme values)



annual probability
of exceedance

$$P_{annum}(L > l|\alpha)$$

log scale rather than
linear (0 to 1) shows
the tail in more detail
at extreme values

shifted exponential

generalised Pareto

Wave load ($L$)
normalised by wave load
with RP=100yrs ($L_{100}$)

$L/L_{100}$

International
Association
of Oil & Gas
Producers

# Poisson probability density function

$n$ is the number of storms, the magnitude of which is greater than a given magnitude $m$, over a period of length $t$ is Poisson distributed.

$$p(N|v_m, t) = \frac{(v_m t)^N}{N!} e^{-v_m t} \quad \& \quad P(N \leq n) = \sum_{i=1}^{n} p(i|v_m, t)$$

generally, $t$ is taken equal to 1 year, so that $v_m$ is to be interpreted as the mean annual number of storm occurrences (depends on $m$) - say $v_m = 100$

probability that the time taken for the next storm (with magnitude greater than $m$) to arrive, ie the waiting time $T$, is less than or equal to $t$ is:

$$P(T \leq t) = 1 - \exp(-v_m t) \quad \& \quad p(T) = v_m \exp(-v_m T)$$

if time to next storm is $t$ then number of storms during the waiting time $T$ is zero (ie $N=0$):

$$p(0|v_m, t) = \frac{(v_m t)^0}{0!} e^{-v_m t} = e^{-v_m t} = 1 - P(T \leq t)$$



International Association of Oil & Gas Producers

mean time to next storm = $1/v_m$ years=0.01 years

P50 time to next storm = 0.00691 years

# Poisson spike process

– used to describe time-dependent events (eg wave loading due to discrete but infrequent storms)

– probability that the time to next storm (ie waiting time $T$ is $< t$) has an exponential distribution:

$$P(T \leq t) = 1 - \exp(-\nu t)$$
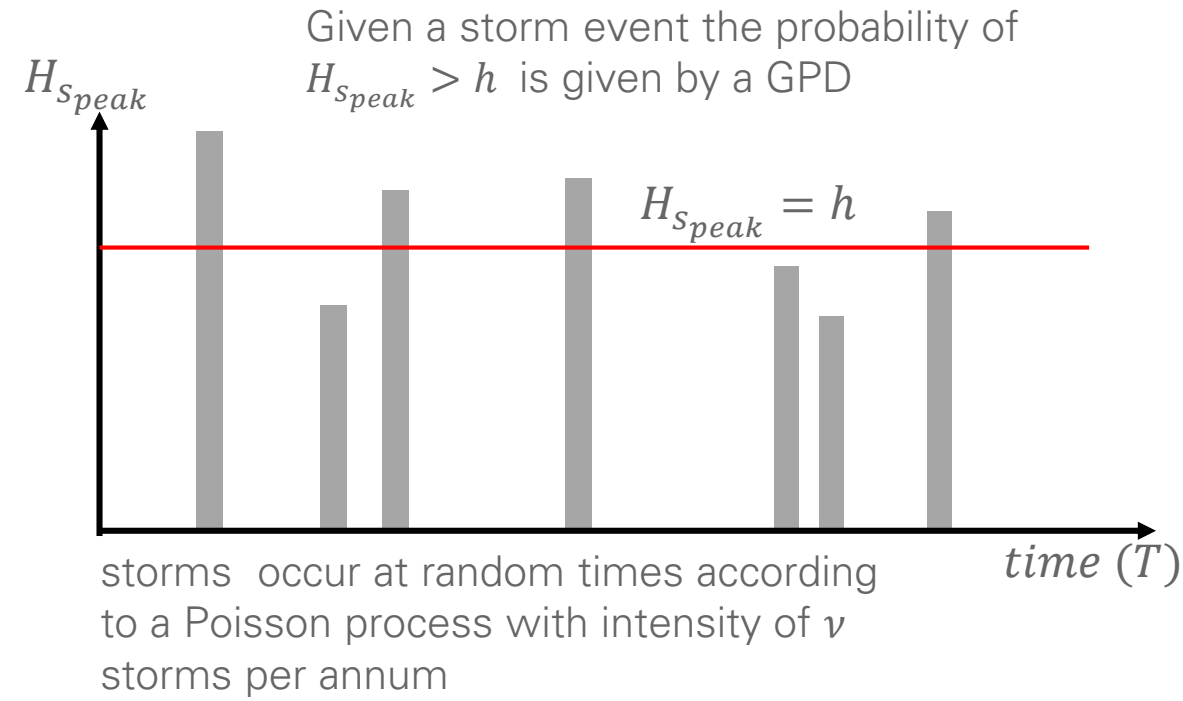
where $\nu$ is the mean annual number of storm occurrences

mean annual number of storms with $H_{s_{peak}} > h$ is

$$\nu_{H_{s_{peak}} > h} = \nu \times P\left(H_{s_{peak}} > h \middle| \mathrm{RS}\right)$$

probability of storms with $H_{s_{peak}} \geq h$ arriving per year is

$$P\left(T \leq 1 \text{ year} \middle| H_{s_{peak}} > h\right) = 1 - \exp\left[-\nu_{H_{s_{peak}} > h} \times 1 \text{ year}\right]$$

$$P_{annual}\left(H_{s_{peak}} > h\right) = 1 - \exp\left[\nu P\left(H_{s_{peak}} > h \middle| \mathrm{RS}\right)\right] \cong \nu P\left(H_{s_{peak}} > h \middle| \mathrm{RS}\right) \text{ for small } \nu P\left(H_{s_{peak}} > h \middle| \mathrm{RS}\right)$$

Given a storm event the probability of $H_{s_{peak}} > h$ is given by a GPD

$H_{s_{peak}}$

$H_{s_{peak}} = h$

storms occur at random times according to a Poisson process with intensity of $\nu$ storms per annum

$time\ (T)$

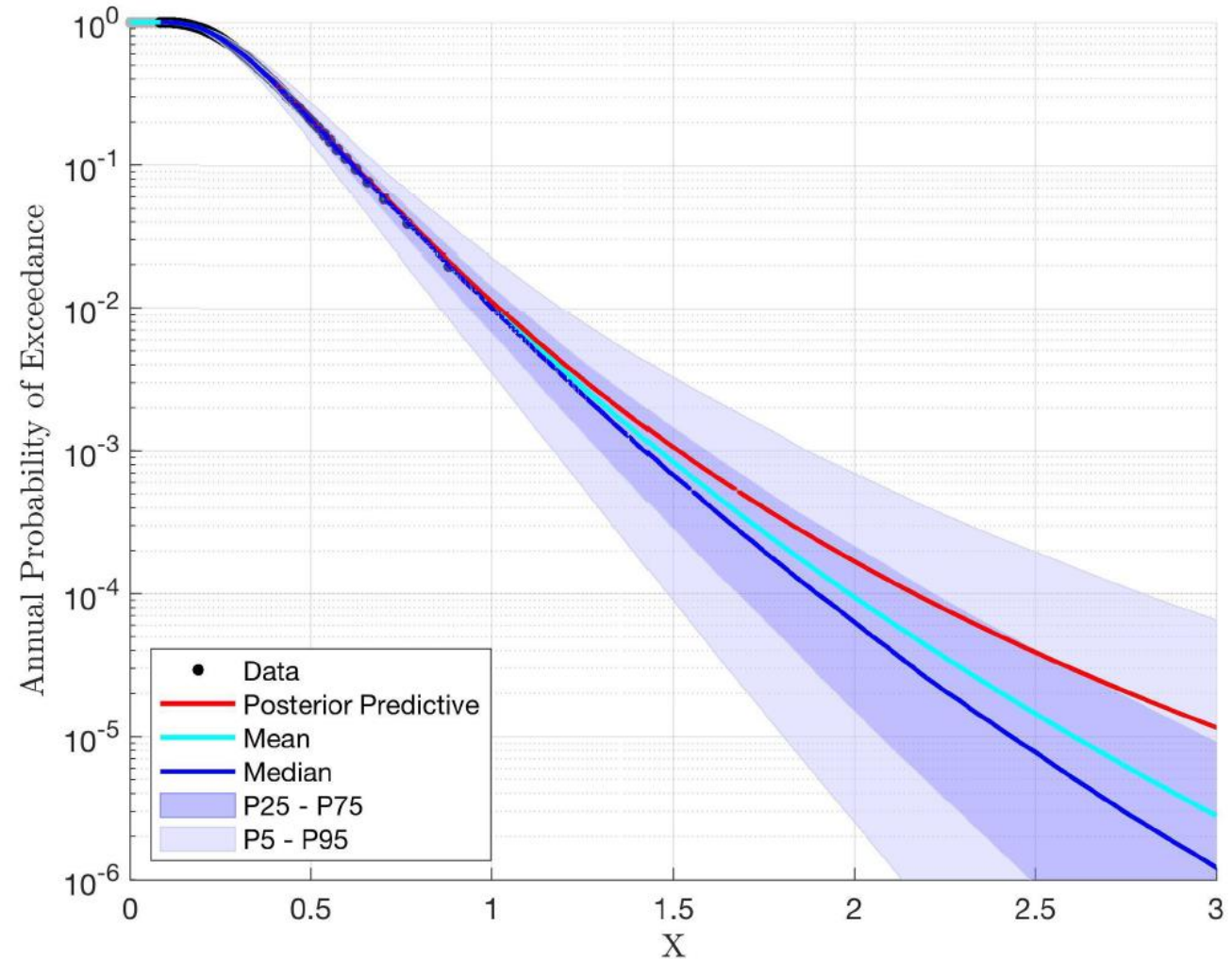# credible interval (for Bayesian inference)

## Quantiles

are points in a distribution that relate to the rank order of values in that distribution.

## Percentiles

are descriptions of quantiles relative to 100; so the 75th percentile (upper quartile) is 75% or three quarters of the way up an ascending list of sorted values of a sample.
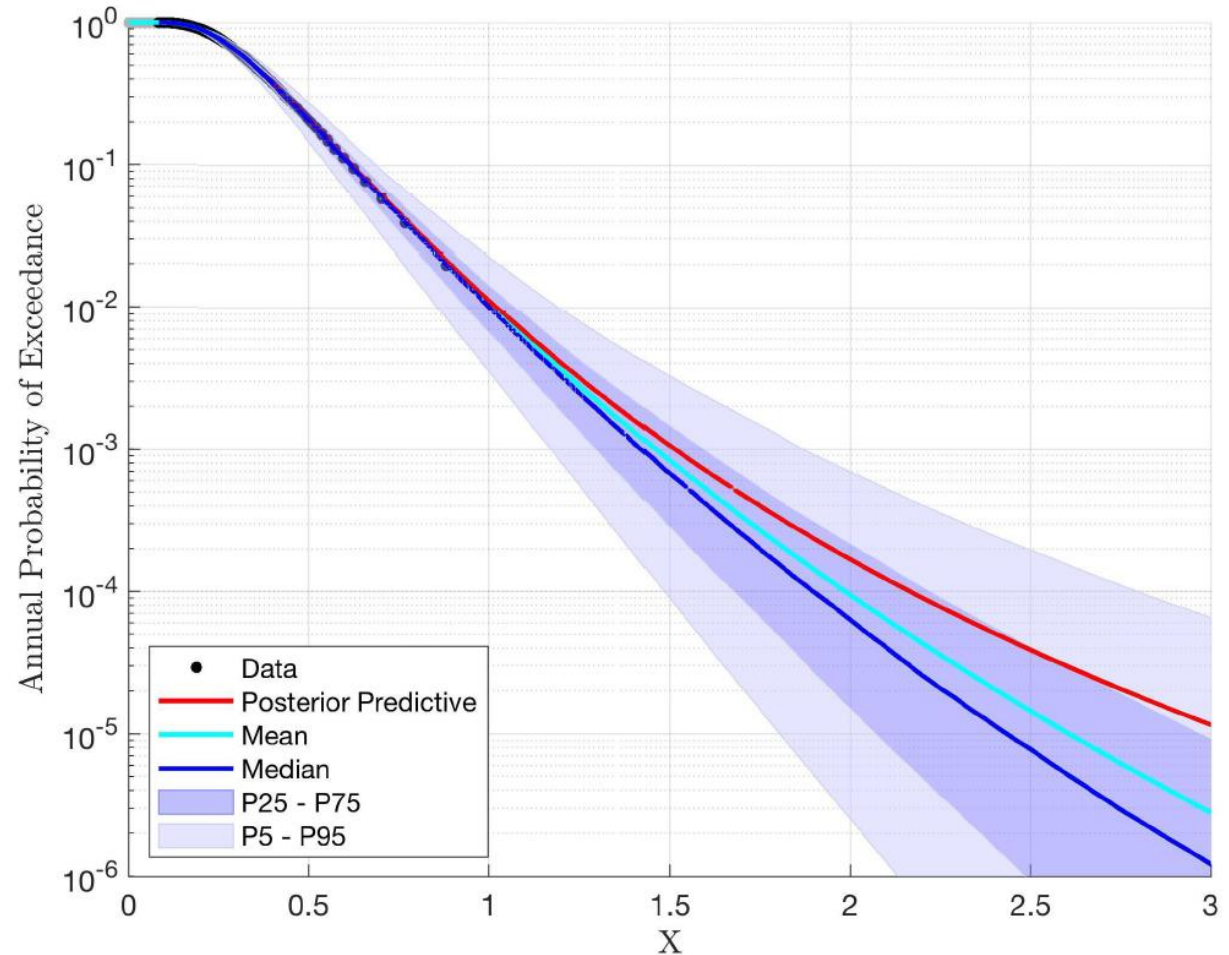
## Credible interval

Credible interval is a "Bayesian confidence interval", but unlike frequentist confidence intervals, credible intervals have a very intuitive interpretation: the 90% credible interval contains the true parameter value ($\theta$) with 90% probability.

# Bayesian inference (1)

Bayesian inference key points…

1) uses prior knowledge of parameter distribution
   ie prior distribution of parameters

2) uses available data together with the prior
   ie posterior distribution of parameters

3) gives the uncertainty explicitly

# Bayesian inference (2)

$$P(\theta|y) = \frac{P(y|\theta) \times P(\theta)}{P(y)}$$

$$posterior = p(\boldsymbol{\theta}|\boldsymbol{h}_{sp\ data}) = \frac{p(\boldsymbol{h}_{sp\ data}|\boldsymbol{\theta}) \times p(\boldsymbol{\theta})}{p(\boldsymbol{h}_{sp\ data})} = \frac{p(\boldsymbol{h}_{sp\ data}|\boldsymbol{\theta}) \times p(\boldsymbol{\theta})}{C} = \frac{likelihood \times prior}{C}$$

where
$\boldsymbol{\theta} = [\mu, \sigma, \xi]$    is the vector of parameters for the GPD
$\boldsymbol{h}_{sp\ data}$        is a vector of values of peak $H_s$ in each storm in the metocean long term simulation

Bayes rule calculates probability densities for $[\mu, \sigma, \xi]$ given the data of peak $H_s$ in each storm
a "continuous family" of GPD fits is obtained, the full posterior distribution is used in the LOADS method

The calculation is performed by sampling using MCMC. MCMC doesn't need to know the denominator as it samples in proportion to the relative magnitude of the posterior rather than the absolute.
The samples are then normalised to give a valid posterior pdf.

# Bayesian inference (3)

$$posterior = p(\mathbf{\theta}|\boldsymbol{h}_{sp\ data}) = \frac{p(\boldsymbol{h}_{sp\ data}|\mathbf{\theta}) \times p(\mathbf{\theta})}{p(\boldsymbol{h}_{sp\ data})} = \frac{p(\boldsymbol{h}_{sp\ data}|\mathbf{\theta}) \times p(\mathbf{\theta})}{C} = \frac{\boxed{likelihood} \times prior}{C}$$

$$p(\mathbf{\theta}|\boldsymbol{h}_{sp\ data}) = \frac{\prod_{i=1}^{N_{data}} p\left(h_{sp\ data_i}|\mathbf{\theta}\right) \times p(\boldsymbol{\theta})}{\int \prod_{i=1}^{N_{data}} p\left(h_{sp\ data_i}|\mathbf{\theta}\right) \times p(\boldsymbol{\theta})d\boldsymbol{\theta}} = \frac{\boxed{\prod_{i=1}^{N_{data}} p\left(h_{sp\ data_i}|\mathbf{\theta}\right)} \times p(\boldsymbol{\theta})}{C}$$

$$p\begin{pmatrix} h_{sp\ data_{i=1}} \\ h_{sp\ data_{i=2}} \end{pmatrix}\mathbf{\theta}\end{pmatrix} = p\left(h_{sp\ data_{i=1}}|\mathbf{\theta}\right) \times p\left(h_{sp\ data_{i=2}}|\mathbf{\theta}\right) = \prod_{i=1}^{2} p\left(h_{sp\ data_i}|\mathbf{\theta}\right)$$

# Bayesian inference (4)

Determine using Bayesian inference with a Generalised Pareto distribution $GPD$

$$P\big(H_{sp} - u > h_{sp}\big|H_{sp} > u\big) = \left[1 + \frac{\xi(h_{sp}-u)}{\sigma+\xi(u-\mu)}\right]^{-1/\xi}$$

$\mu$ = location parameter
$\sigma$ = scale parameter
$\xi$ = shape parameter
$u$ = threshold
$H_{sp}$ = peak significant wave ht in a storm. Say we have 1200 years of data $h_{sp\,data_i}$  $i = 1, N$

$$P\left(H_{sp} - u > h_{sp\,data_i}\Big|H_{sp} > u, \xi, \sigma, \mu\right) = \left[1 + \frac{\xi\big(h_{sp\,data_i}-u\big)}{\sigma+\xi(u-\mu)}\right]^{-1/\xi}$$

Sampling the above for each $h_{sp\,data_i}$  for a range of parameters $\xi_j, \sigma_j, \mu_j$ and then taking the product over $i = 1, N$ gives the likelihood

# $P(H_{sp} > h | \boldsymbol{h}_{sp\ data})$ - posterior predictive prob. exceedance



sampling distribution for future observations of Hs given the GPD parameters

posterior distribution of the parameters given past observations of Hs (ie data)

posterior predictive distribution is a conditional expectation (conditioned on the observed data) weighted by the parameter values from the posterior distribution
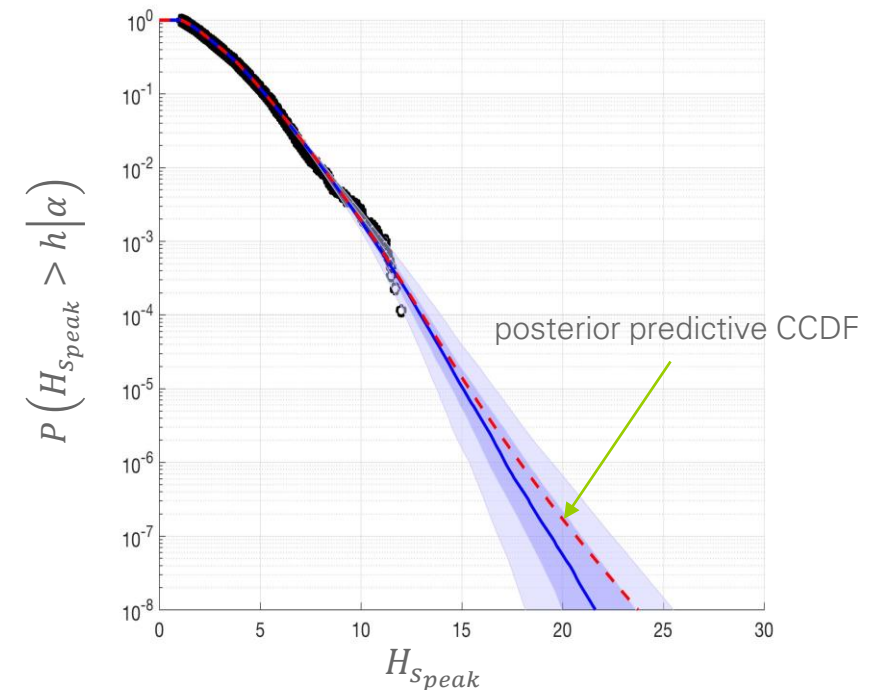
posterior predictive PDF
$$= p(h | \boldsymbol{h}_{sp\ data}) = \int p(h | \boldsymbol{\theta}) \times p(\boldsymbol{\theta} | \boldsymbol{h}_{sp\ data}) d\boldsymbol{\theta} = \mathbb{E}\big(p(h | \boldsymbol{\theta}) | \boldsymbol{h}_{sp\ data}\big)$$
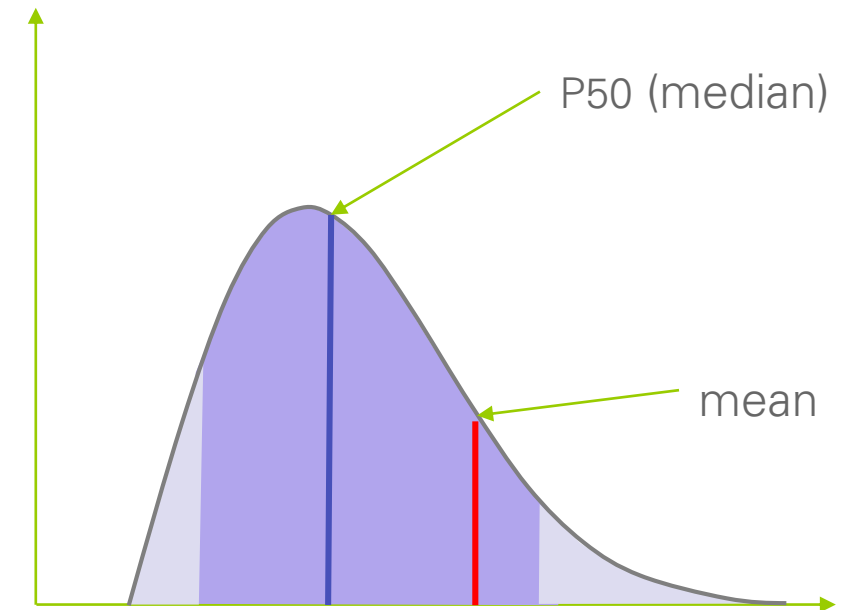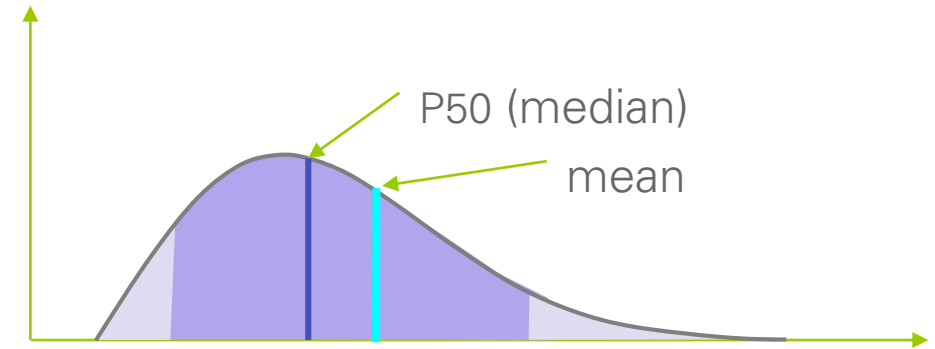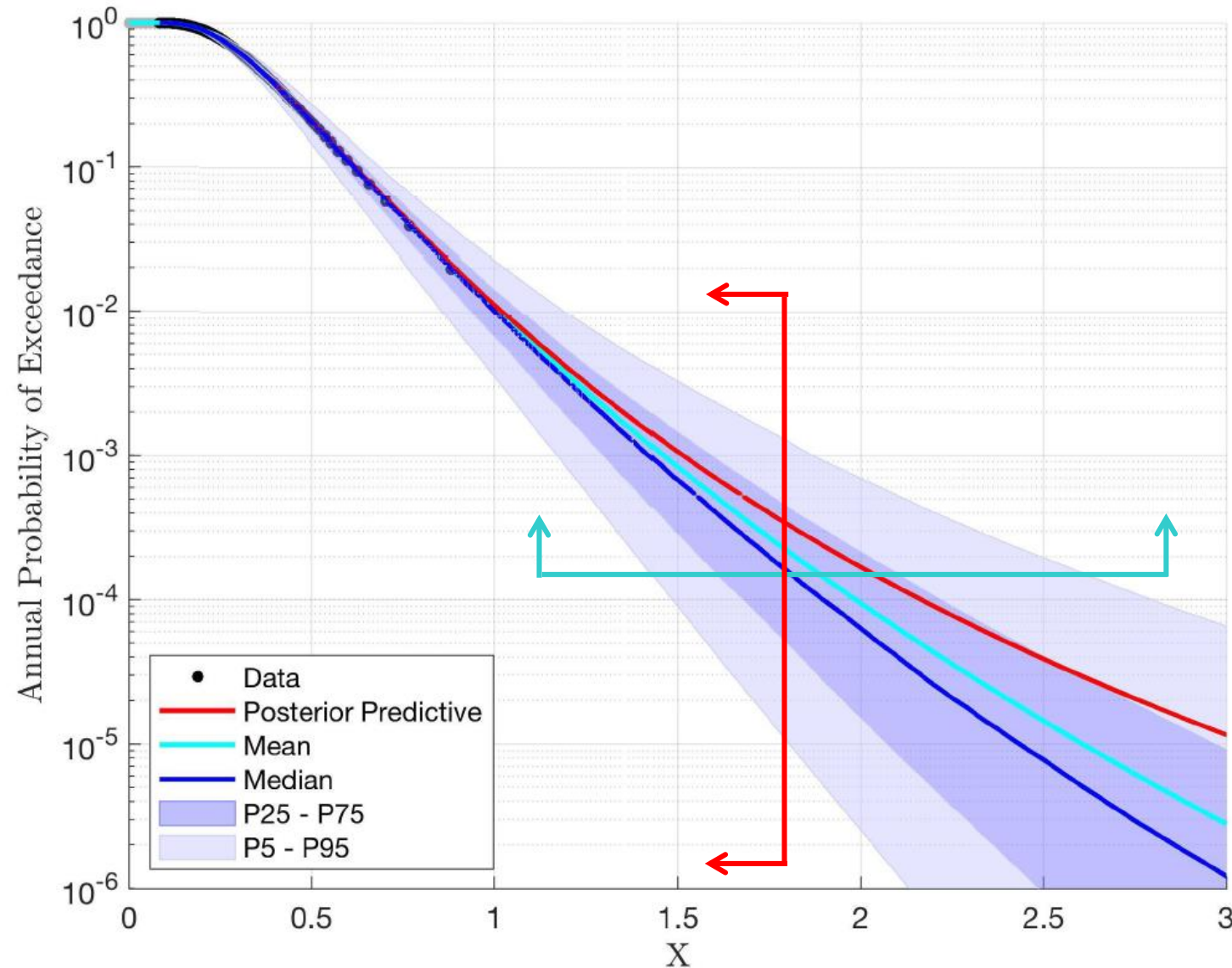
posterior predictive CCDF
$$= P(H_{sp} > h | \boldsymbol{h}_{sp\ data}) = \int_{h=H_{sp}}^{\infty} p(h | \boldsymbol{h}_{sp\ data}) dh$$

the posterior predictive distribution takes into account the uncertainty of the parameter estimates, which is quantified by the posterior distribution.

posterior predictive CCDF

# mean of probabilities  v  mean of values

posterior predictive
prob. of exceedance

# The Case for Using Mean Seismic Hazard

Robin K. McGuire,[a] M.EERI, C. Allin Cornell,[b]
M.EERI, and Gabriel R. Toro,[a] M.EERI

Complete probabilistic seismic hazard analyses incorporate epistemic uncertainties in assumptions, models, and parameters, and lead to a distribution of annual frequency of exceedance versus ground motion amplitude (the "seismic hazard"). For decision making, if a single representation of the seismic hazard is required, it is always preferable to use the mean of this distribution, rather than some other representation, such as a particular fractile. Use of the mean is consistent with modern interpretations of probability and with precedents of safety goals and cost-benefit analysis. [DOI: 10.1193/1.1985447]

## INTRODUCTION

Estimates of earthquake ground motion hazard involve substantial epistemic uncertainty in the mean frequency of exceedance for a given ground motion or, alternatively, in the ground motion for a given mean frequency of exceedance. We believe that the mean estimate of the mean frequency of exceedance should be the standard when a single estimate is necessary. (Please refer to the Addendum for a clarification of the often misused or misunderstood definitions regarding "hazard," "frequency," and "mean." In the context of that addendum we shall adopt the common shorthand convention of "hazard" for the "frequency of exceedance" and "mean hazard" or "mean frequency of exceedance" for the "mean estimate of the frequency of exceedance.")

This epistemic uncertainty has been known and quantified in the United States since the 1970s (see, for example, McGuire 1977). In seismic hazard studies it is preferable to report the complete epistemic distribution of hazard, because this allows any effects of that uncertainty on risk mitigation decisions to be handled in an explicit, quantitative way. This reporting usually takes the form of presenting four or more hazard curves, say, three fractiles (e.g., 0.10, 0.50—or median—and 0.90) plus the mean hazard curve. This reporting position is supported by a finding of the National Research Council Panel on Seismic Hazard Analysis: "Knowledge of earthquake processes and effects in much of the United States is meager, resulting in considerable uncertainty in seismic hazard estimates. No single measure of the seismic hazard (e.g., a mean or median [estimate]) is adequate to represent this basic lack of understanding; therefore, measures of uncertainty must be transmitted as part of a PSHA [probabilistic seismic hazard analysis]." (NRC 1988) (words in brackets added for clarity).

International
Association
of Oil & Gas
Producers