# An approach to the extreme value distribution of non-stationary process

## Hang CHOI & Jun KANDA

*Institute of Environmental Studies*

*Graduate School of Frontier Sciences*

*The University of Tokyo*

# 01 Theoretical Frameworks of EVA

1) Stationary Random Process (Sequence)
- Distribution Ergodicity
- (In)dependent and Identically Distributed random variables (I.I.D. assumption)
→ strict stationarity
* Conventional approach

2) Non-stationary Random Process (Sequence)
- (In)dependent but non-identically Distributed random variables (non-I.I.D. assumption)
→ weak stationarity, non-stationarity

3) Ultimate (Asymptotic) and penultimate forms

# 1) I.I.D. random variable Approach

$$X_1, \mathrm{K}, X_n \in F(x), \; Z_n = \max\{X_1, ..., X_n\}$$

$$\lim_{n \to +\infty} P\left(\frac{Z_n - a_n}{b_n} \leq x\right) = \lim_{n \to +\infty} F^n(a_n + b_n x) = G(x) \in \mathrm{F}$$

$$\mathrm{F} = \{ \text{Gumbel, Fréchet, Weibull} \}$$

GEVD, GPD, POT-GPD + MLM, PWM, MOM etc.

R.A. Fisher & L.H.C. Tippett (1928), Limiting forms of the frequency distribution of the largest or smallest member of a sample, *Proc. Cambridge Philosophical Society*, Vol 24, p180~190

Fréchet

Gumbel

Weibull

x

reduced variates y=-log(-log G(x))

## 2) non-I.I.D. random variable Approach

$$X_1 \in F_1(x), \mathrm{K}, X_n \in F_n(x), Z_n = \max\left\{X_1, \mathrm{K}, X_n\right\}$$

$$\lim_{n \to +\infty} P\left(\frac{Z_n - c_n}{d_n} \leq x\right) = \lim_{n \to +\infty} \prod_{j=1}^{n} F_j(d_j + c_j x) = Q(x) \in \mathrm{S}$$

$$F \subset S$$
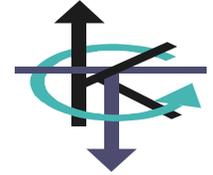
$$\mathrm{S} = \{ \text{Gumbel, Fréchet, Weibull,.....}\}$$

The class of EVD for non-i.i.d. case is much larger.

\* Falk *et al.*, *Laws of Small Numbers: Extremes and Rare Events*, Birkhäuser, 1994

# 3) Ultimate / penultimate form and finite epoch T in engineering practice

In engineering practice, the epoch of interest, $T$ is finite.
*e.g. annual maximum value, monthly maximum value and maximum/minimum pressure coefficients in 10min etc.*

As such, the number of independent random variables $m$ in the epoch $T$ becomes a finite integer, *i.e. $m < \infty$*, and consequently, the theoretical framework for ultimate form is no longer available regardless *i.i.d.* or *non-i.i.d.* case.

Following discussions are restricted on the *penultimate d.f.* for the extremes of non-stationary random process *in a finite epoch T*.
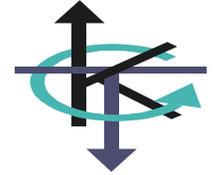
Let $x(t),\ t \in \mathbb{R}^{+} := [0, \infty)$ be the continuous observation record of a non-stationary continuous stochastic process $X(t)$ and assume that $X(t)$ is a mean square differentiable process and hold the following condition.

$$r(t, \tau) := E\big[X(t)X(t + \tau)\big] \Rightarrow r(t, \tau)\log \tau \to 0 \text{ as } \tau \to T$$

As such, how to estimate the extreme value distribution of $X(t)$ in an epoch $T < \infty$ from the continuous observation record $x(t)$?

# 03 Assumptions and Formulation

According to the conventional approach in wind engineering, let assume the non-stationary process $X(t)$ can be partitioned with a finite epoch $T$, in which the partitioned process $X_i(t), (i\text{-}1)T$ $t < iT$, can be assumed as an independent stationary random process, and define the d.f. $F_Z(x)$ as follows.

$$F_Z(x) = P(Z \leq x; \; Z := \sup X(t), \; t \in [0, T < \infty))$$

Then, by the Glivenko-Cantelli theorem and the block maxima approach

$$F_Z(x) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} I_{(-\infty, x]} \left( Z_i := \sup X_i(t) \right) = \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} P(Z_i \leq x)$$

where $I_c(x) = 1$ if $x \in c$ else 0

# 04 i.i.d. random sequence (EQRS) approach of $P(Z_i \leq x)$

By partitioning the interval $[(i-1)T, iT)$ into finite subpartitions $[(j-1)h, jh), 1 \leq j \leq [T/h] = m_i$ in the manner of that

$$P(Z_i \leq x) = \prod_{j=1}^{m_i < \infty} P\left(Z_j^* \leq x;\ Z_j^* := \sup\left(X_i(t), t \in [(j-1)h,\ jh)\right)\right) = F_{Z_i}^{m_i}(x)$$

the required d.f. $F_Z(x)$ can be defined as follows.

$$F_Z(x) := \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} F_{Z_i}^{m_i}(x)$$

If it is possible to assume that all $m_i \approx m$, then

$$F_Z(x) := \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} F_{Z_i}^{m}(x)$$

# 05 lower bound of $F_Z(x)$

From the inequality (geometric mean)<(arithmetic mean), a lower bound of $F_Z(x)$ can be defined as follows:

$$\overset{o}{F_Z}(x) := \lim_{n \to \infty} \prod_{i=1}^{n} F_{Z_i}^{m/n}(x) < \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} F_{Z_i}^{m}(x) = F_Z(x)$$

# 06 alternative definition of $F_Z(x)$

$$\hat{F}_Z(x) := \left( \frac{1}{n} \sum_{i=1}^{n} F_{Z_i}(x) \right)^{m} = \bar{F}_Z^{m}(x)$$

This definition can be found easily in engineering applications and may be reasonable for the case of $F_{Z_1}$ ; L ; $F_{Z_n}$.

Let define quantile functions of each definition as follows.

$$Q(\alpha) := F_Z^{-1}(\alpha), \ \tilde{Q}(\alpha) := \tilde{F}_Z^{-1}(\alpha), \ \hat{Q}(\alpha) := \hat{F}_Z^{-1}(\alpha)$$

Then, by the inequality for the means and the comparison of the distribution of order statistics, i.e. $Z_{n:n}$ and $Z_{1:n}$,

$$Q(\alpha) \le \tilde{Q}(\alpha), \ \hat{Q}(\alpha) < \tilde{Q}(\alpha) \ \ for \ all \ \alpha \in (0,1)$$

$$\begin{cases} \hat{Q}(\alpha) < Q(\alpha) \to \tilde{Q}(\alpha) \ \text{ for large } \alpha \\ Q(\alpha) < \hat{Q}(\alpha) < \tilde{Q}(\alpha) \ \text{ for small } \alpha \end{cases}$$

Therefore, the alternative definition results in smaller variance of extremes.

# Complement for the inequalities of quantile functions

① $Z_{n:n} \sim \prod_{i=1}^{n} F_i^m ,\ \overset{\circ}{Z}_{n:n} \sim \left( \prod_{i=1}^{n} F_i^{m/n} \right)^n = \prod_{i=1}^{n} F_i^m \Rightarrow Z_{n:n}(\alpha) = \overset{\circ}{Z}_{n:n}(\alpha)$

② $\overset{\circ}{Z}_{n:n} \sim \prod_{i=1}^{n} F_i^m ,\ \bar{Z}_{n:n} \sim \left( \frac{1}{n}\sum_{i=1}^{n} F_i \right)^{mn} ,\ \left( \prod_{i=1}^{n} F_i^{m/n} \right)^n < \left( \frac{1}{n}\sum_{i=1}^{n} F_i \right)^{mn} \Rightarrow \overset{\circ}{Z}_{n:n}(\alpha) > \bar{Z}_{n:n}(\alpha)$

③ $Z_{1:n} \sim 1 - \prod_{i=1}^{n}\left( 1 - F_i^m \right);\ \sum_{i=1}^{n} F_i^m - o(F^{2m}),\ \bar{Z}_{1:n} \sim 1 - \left( 1 - \bar{F}^m \right)^n;\ \frac{1}{n^{m-1}}\sum_{i=1}^{n} F_i^m - o(F^{2m})$

$\Rightarrow Z_{1:n}(\alpha) < \bar{Z}_{1:n}(\alpha)$

④ $Z_{1:n} \sim 1 - \prod_{i=1}^{n}(1 - F_i^m),\ \overset{\circ}{Z}_{1:n} \sim 1 - \left( 1 - \prod_{i=1}^{n} F_i^{m/n} \right)^n ,$

$\prod_{i=1}^{n}(1 - F_i^m) < \left( 1 - \bar{F}^m \right)^n < \left( 1 - \prod_{i=1}^{n} F_i^{m/n} \right)^n \Rightarrow Z_{1:n}(\alpha) < \overset{\circ}{Z}_{1:n}(\alpha)$

# 08 Numerical example $\overline{F}(x) = N(0,1)$

( m=10,000, n=5,000, iteration=100 )

$X_i \sim N(0,1)$

$X_i \sim N(0,\sigma_i),\ \sigma_i \sim N(1,0.05)$

$X_i \sim N(0,\sigma_i),\ \sigma_i \sim N(1,0.1)$

$X_i \sim N(0,\sigma_i),\ \sigma_i \sim N(1,0.2)$

09 practical application:
   Annual maximum wind speed in Japan

1) Observation records and Historical annual maximum
   wind speeds at 155 sites in Japan

   Observation Records: JMA records (CSV format)
      1961~1990 : 10 min average wind speed per 3 hours
      1991~2002 : 10 min average wind speed per hour

   Historical annual maximum wind speeds record:
   1929~1999 : A historical annual max. wind speed data set
                    compiled by Ishihara *et al*.(2002)
   2000~2002 : extracted from JMA records (CSV format)

T. Ishihara *et al*. (2002), A database of annual maximum wind speed and corrections for anemometers in Japan, *Wind Engineers, JAWE*, No.92, p5~54 (in Japanese)

# 09 practical application:
## Annual maximum wind speed in Japan

2) The effect of different observation recording format on the basic statistics

Base on the recently opened continuous 1 min average wind speed records (1997. 3~2002. 2), calculating every 10 min average wind speeds, 10 min average per hour and 3 hours, and comparing the basic statistics for each recording format, then the effect of different recording format becomes to be clear.

# 3) Examples of non-stationarity : the basic statistics (1961~2002)

| Site | Coefficient of Variations (C.O.V) | | | |
|---|---|---|---|---|
| | $\mu$ | $\sigma$ | $\gamma_1$ | $\gamma_2$ |
| Abashiri | 5.1% | 5.0% | 11.7% | 14.4% |
| Katsuura | 5.81% | 8.87% | 24.06% | 26.71% |
| Kobe | 5.71% | 8.25% | 12.55% | 16.31% |
| Kumamoto | 7.62% | 5.93% | 14.19% | 32.36% |
| Makurazaki | 4.60% | 6.07% | 25.97% | 47.98% |
| Morioka | 6.99% | 4.43% | 12.09% | 11.55% |
| Oita | 3.42% | 8.41% | 16.49% | 27.41% |
| Shionomisaki | 4.34% | 4.07% | 19.36% | 30.10% |
| Tokyo | 5.90% | 8.06% | 18.90% | 17.88% |
| **Minimum** | **3.42%** | **4.07%** | **11.72%** | **11.55%** |
| **Maximum** | **7.62%** | **8.87%** | **25.97%** | **47.98%** |

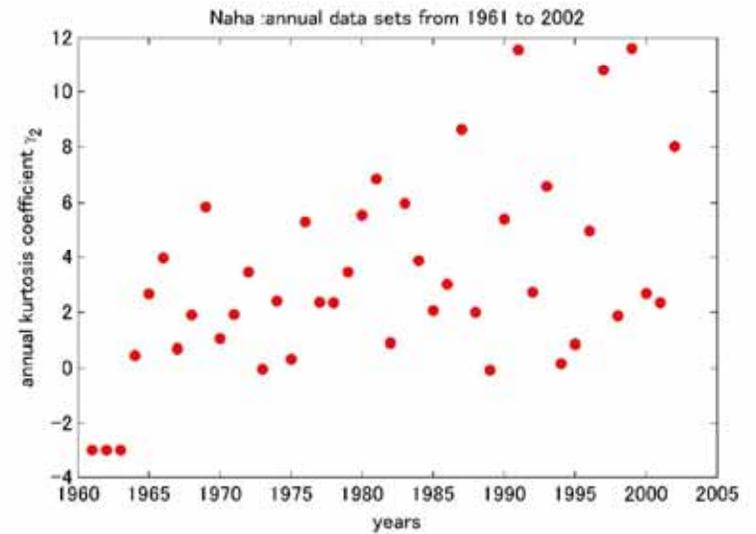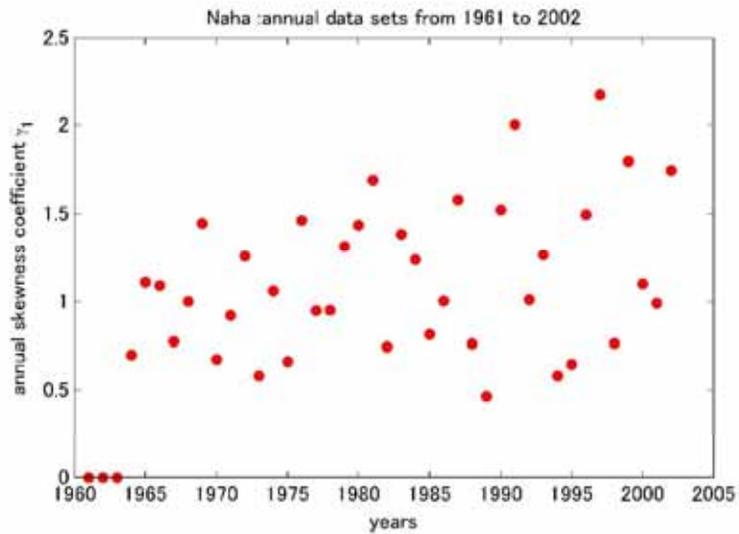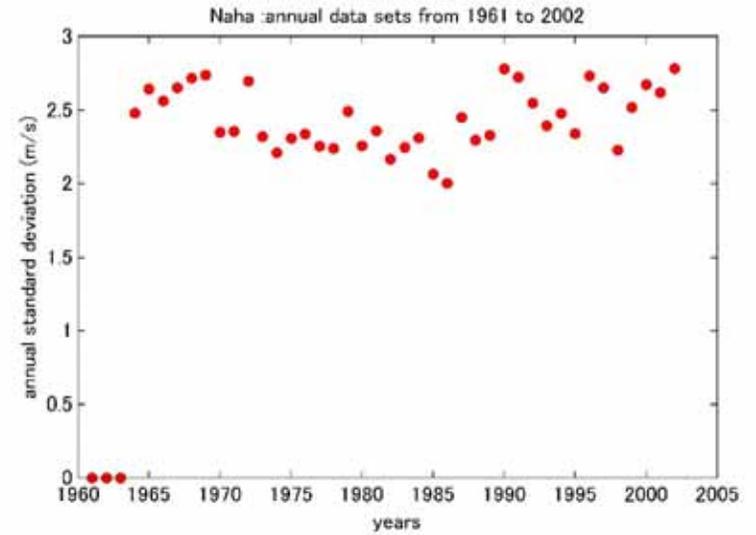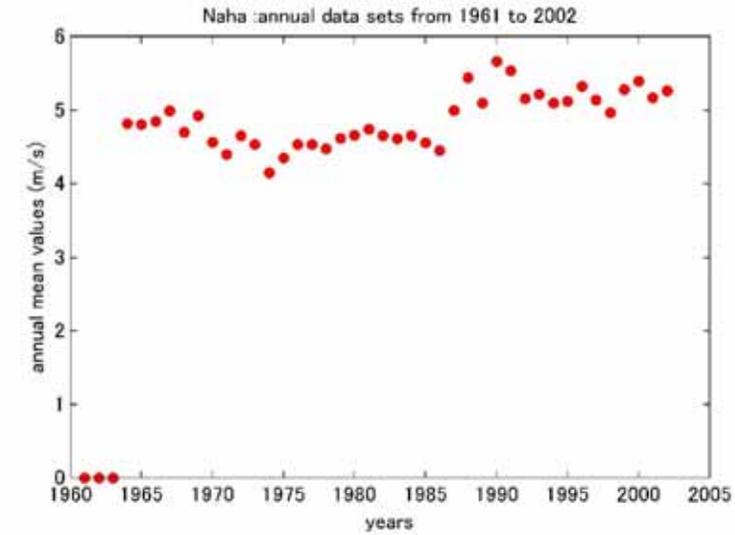Morioka (1961~2002)

# Kobe (1961~2002)

# Shionomisaki (1961~2002)

# Makurazaki (1961~2002)

# Naha (1964~2002)

# 09 practical application: Annual maximum wind speed in Japan

## 4) Approximation of the annual wind speed distribution

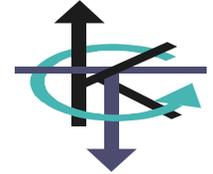Based on the probability integral transformation,

$$\Phi(z_\alpha) = F_{X_i}(x_\alpha) = F_{X_i}\big(g(z_\alpha)\big) = \alpha$$

$$x_\alpha = g(z_\alpha) = a + bz_\alpha + cz_\alpha^2 + dz_\alpha^3$$

$$\Rightarrow \quad F_{X_i}(x_\alpha) = \Phi\big(g^{-1}(x_\alpha)\big)$$

The coefficients $a,b,c$ and $d$ can be estimated from the given basic statistics of annual wind speed, i.e. mean, standard deviation, skewness and kurtosis (Choi & Kanda 2003).

H. Choi and J. Kanda (2003), Translation Method: a historical review and its application to simulation on non-Gaussian stationary processes, *Wind and Struct.*, 6(5), p357~386

# 5) Estimation by Monte Carlo Simulation (MCS)

## 5.1) Simulation methods

① Based on the Spectral representation theorem for stationary stochastic process
→Using a given spectral density function, discrete stationary stochastic process is simulated.
→time consuming method

② Based on equivalent i.i.d. random sequence (EQRS)
→A stochastic process, which can be approximated by Poisson process, is modeled as an i.i.d. random sequence having same quantile function.
→time effective method

## 5.2) Required information for MCS based on EQRS

① $m$ : the number of Independent rv

→ approximated by mean zero crossing rate
(Normal process)

From <span style="color:red">Rice theorem and Poisson approximation</span>, normal quantile function is given as follows:

$$z_\alpha = \sqrt{2\left\{\log \mu_0 T + y_T\right\}}$$

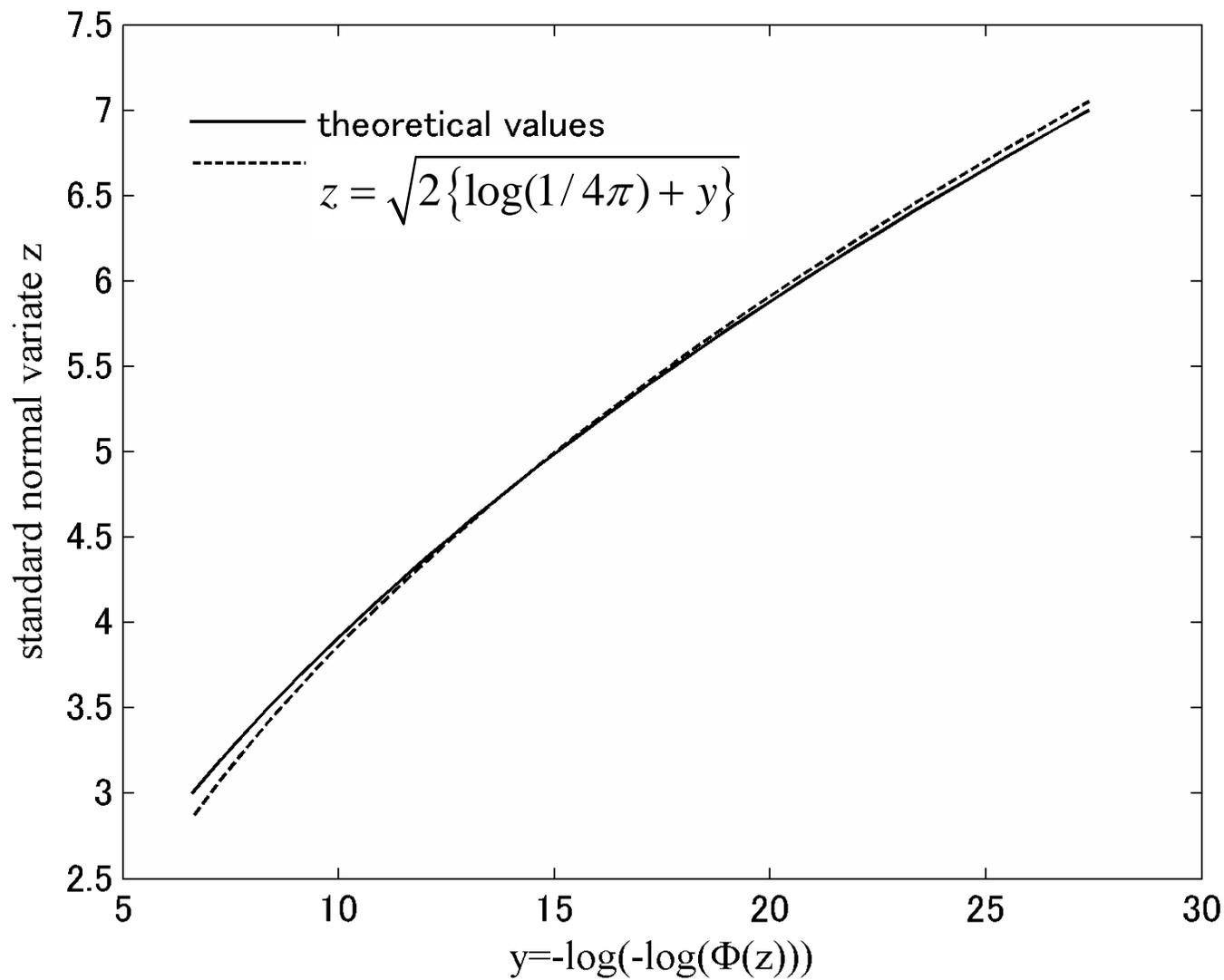in whcih $\mu_0$ : mean zero crossing rate, $y_T = -\log(-\log \alpha)$
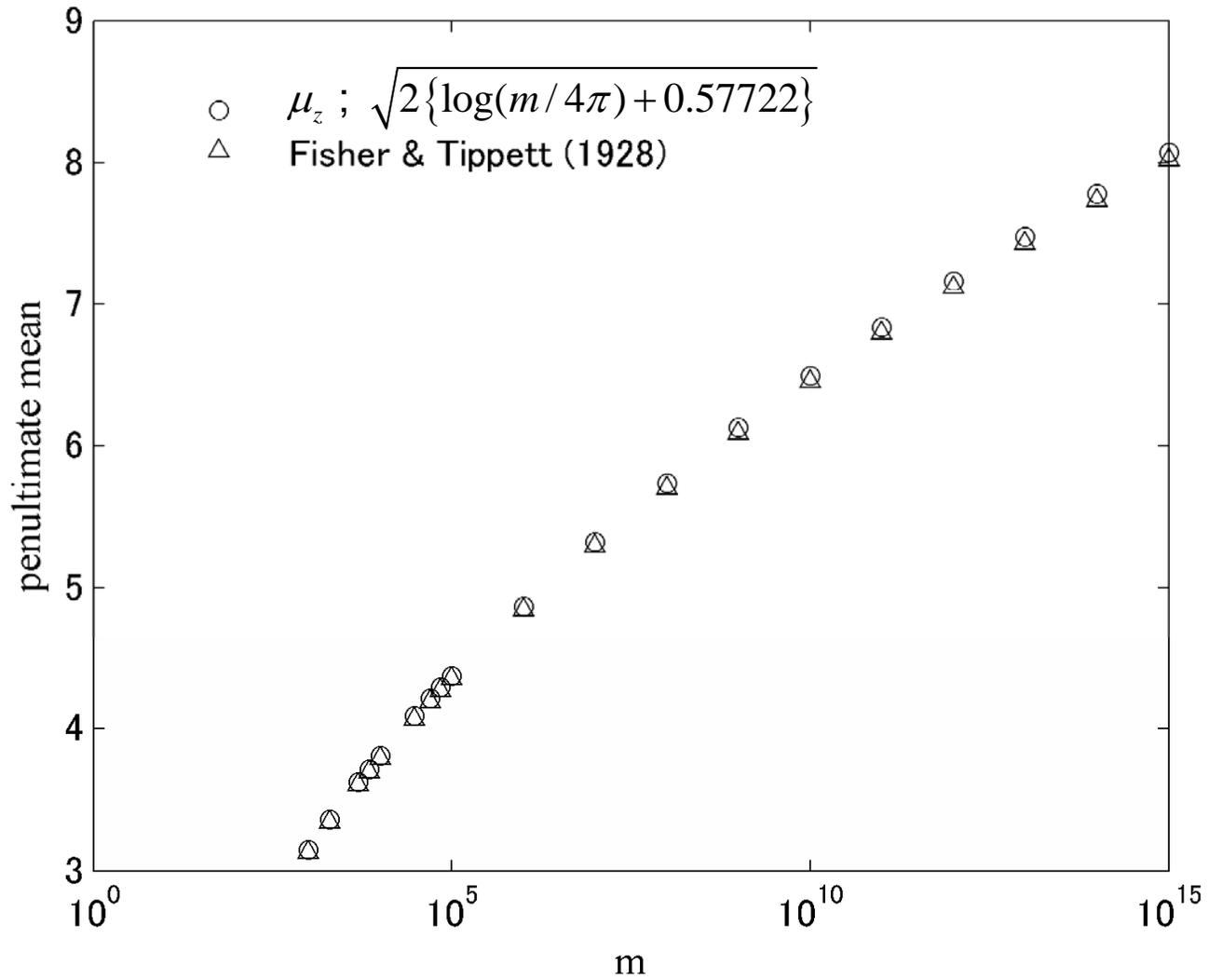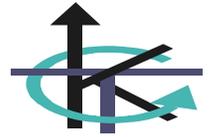
From $\Phi^m(z_\alpha)$,

$$z_\alpha \ ; \ \sqrt{2\left\{\log\left(m/4\pi\right)+y_T\right\}} \quad *$$

By comparison,

$$m \ ; \ 4\pi\mu_0 T$$

* Choi & Kanda (2004), A new method of the extreme value distributions based on the translation method, *Summaries of Tech. Papers of Ann. Meet. of AIJ*, Vol. B1, p23~24 (in Japanese)

The chart shows "standard normal variate z" on the y-axis (ranging from 2.5 to 7.5) versus $y=-\log(-\log(\Phi(z)))$ on the x-axis (ranging from 5 to 30). The legend indicates:
- theoretical values (solid line)
- $z = \sqrt{2\{\log(1/4\pi) + y\}}$ (dashed line)

Figure: Plot of penultimate mean versus $m$ (logarithmic horizontal axis from $10^0$ to $10^{15}$, vertical axis from 3 to 9).

Legend:
- $\circ$   $\mu_z$ ; $\sqrt{2\{\log(m/4\pi)+0.57722\}}$
- $\triangle$   Fisher & Tippett (1928)

(non-Normal process)

To Rice formula for the expected number of crossings, i.e.

$$\nu(x) = \int_0^\infty \dot{x}\, f(\dot{x}\,|\,X=x)\,d\dot{x} = f(x)\cdot\int_0^\infty \dot{x}\, f(\dot{x})\,d\dot{x}$$

applying translation function $g(z)$

$$\nu(x) = \frac{1}{|g'(z)|}\phi\big(g^{-1}(x)\big)\cdot\int_0^\infty \frac{\dot{x}\exp\left\{-\left(\dot{x}/\sigma_{\dot{x}}\cdot g'(z)\right)^2/2\right\}}{\sqrt{2\pi}\cdot\sigma_{\dot{x}}\cdot g'(z)}\,d\dot{x}$$

$$= \frac{\sigma_{\dot{x}}}{2\pi}\cdot\exp\left(-\frac{\{g^{-1}(x)\}^2}{2}\right) = \nu_0\cdot\exp\left(-\frac{\{g^{-1}(x)\}^2}{2}\right) \quad \text{\textcolor{red}{*}}$$

From Poisson approximation

$$x_\alpha = g(z_\alpha) = g\left(\sqrt{2\{\log(\nu_0 T)+y_T\}}\right)$$

With the same manner

$$m \; ; \; 4\pi\nu_0 T$$

* The distribution of $dx/dt$ is assumed as normal distribution and the assumption is reasonable.
*e.g.* H. Choi (1988), *Characteristics of natural wind for wind load estimation*, Master Thesis, Univ. of Tokyo (in Japanese)

# Practical example ($T$=1 year)

- $\nu_0 T$ estimated from long term observation records in Tokyo (1985~1987, Choi 1988)

| Height (m) | 45 | 45 | 46 | 48 | 58 |
|---|---|---|---|---|---|
| $\nu_0 T$ | 2883 | 2665 | 2545 | 2739 | 2307 |
| Height (m) | 63 | 79 | 93 | 187 | |
| $\nu_0 T$ | 3124 | 2586 | 3115 | 2384 | |

$$m = 4\pi\nu_0 T = 4\pi \cdot (2300 \sim 3000) \rightarrow 30,000$$

② parent distribution function for each year
  →Generalized bootstrap method＋
    Translation method (Probability Integral Transform)

For the year $i$,

$$F_i(x = g_i(z)) = \Phi(z), \ g_i(z) = a_i + b_i z + c_i z^2 + d_i z^3$$

$$\{a,b,c,d\}_i = \Psi(\mu_i, \sigma_i, \gamma_{1i}, \gamma_{2i})$$

$$\max(x_1, \mathrm{K}, x_m)_{i=1,\mathrm{K},n} = g_i(\max(z_1, \mathrm{K}, z_m)_{i=1,\mathrm{K},n})$$

# Example : Tokyo (1961~2002)



simulated and historical annual mean

$\mu$

simulated and historical annual rms

$\sigma$

simulated and historical annual skewness

$\gamma_1$

simulated and historical annual kurtosis

$\gamma_2$

◯ Estimated from historical records   ▪ Monte Carlo Simulation

# Correlation between the basic statistics (Tokyo)

| | $\mu$ | $\sigma$ | $\gamma_1$ | $\gamma_2$ |
|---|---|---|---|---|
| $\mu$ | 1.000 | 0.242 | -0.372 | -0.315 |
| $\sigma$ | 0.242 | 1.000 | 0.541 | 0.433 |
| $\gamma_1$ | -0.372 | 0.541 | 1.000 | 0.946 |
| $\gamma_2$ | -0.315 | 0.433 | 0.946 | 1.000 |

Such correlation characteristics between the basic statistics are regenerated by Cholesky decomposition of correlation matrix.

# Regeneration of correlation characteristics

correlation coefficients of simulated ones and given values

| | $\mu$ | $\sigma$ | $\gamma_1$ | $\gamma_2$ |
|---|---|---|---|---|
| $\mu$ | 1.000 | 0.242 (0.242) | -0.389 (-0.372) | -0.343 (-0.315) |
| $\sigma$ | 0.242 (0.242) | 1.000 | 0.565 (0.541) | 0.476 (0.433) |
| $\gamma_1$ | -0.389 (-0.372) | 0.565 (0.541) | 1.000 | 0.958 (0.946) |
| $\gamma_2$ | -0.343 (-0.315) | 0.476 (0.433) | 0.958 (0.946) | 1.000 |

(·): given correlation coefficient

# No. of Simulation：n=100 year x 1000 times



correlation between annual mean and rms

○ : annual mean and standard dev. from historical records

# No. of Simulation：n=100 years x 1000 times



correlation between annual skewness and kurtosis

○：skewness and kurtosis from historical records

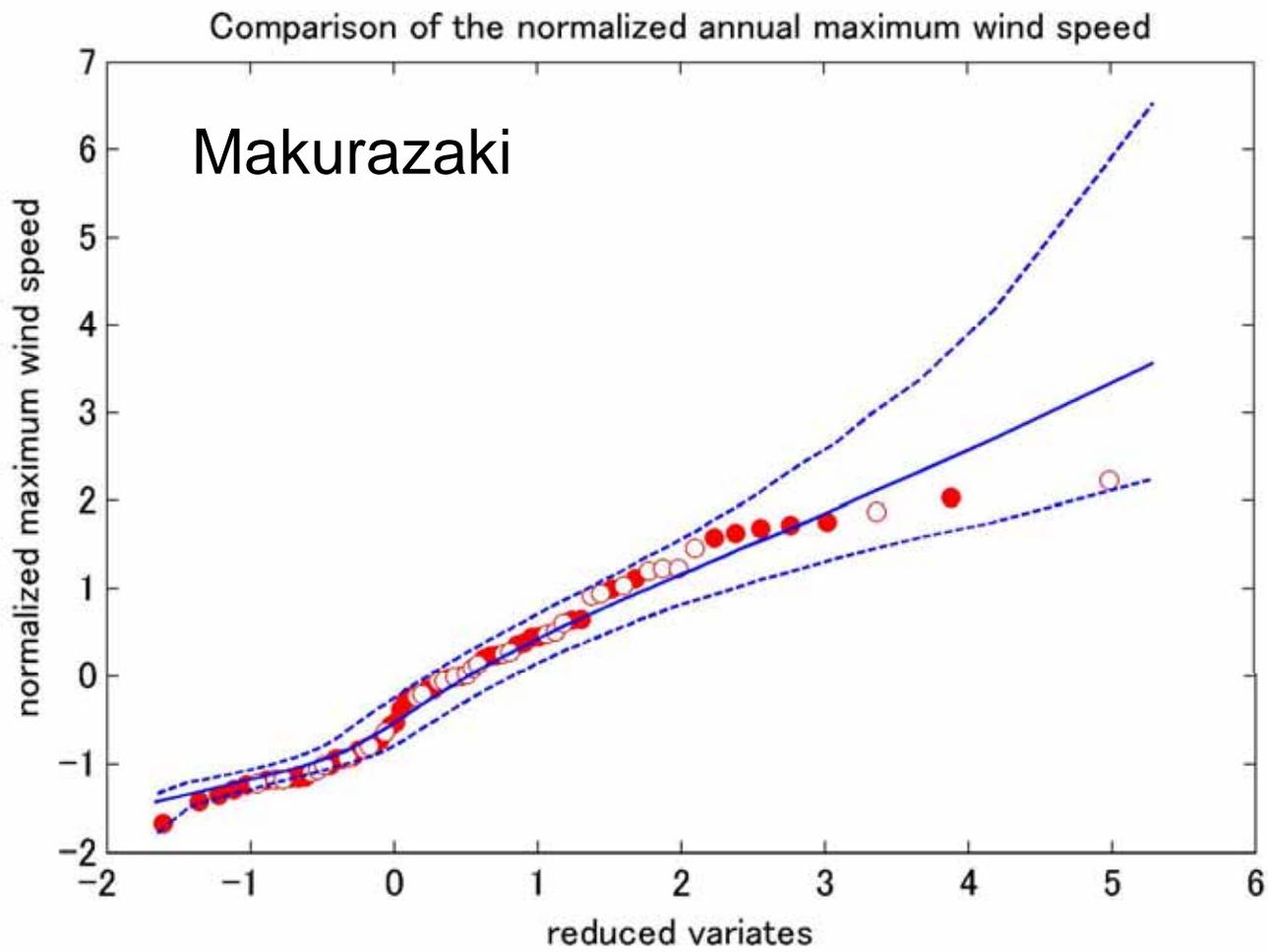# 6) Comparison of the quantile functions from MCS and historical records in normalized form

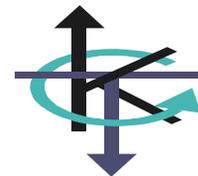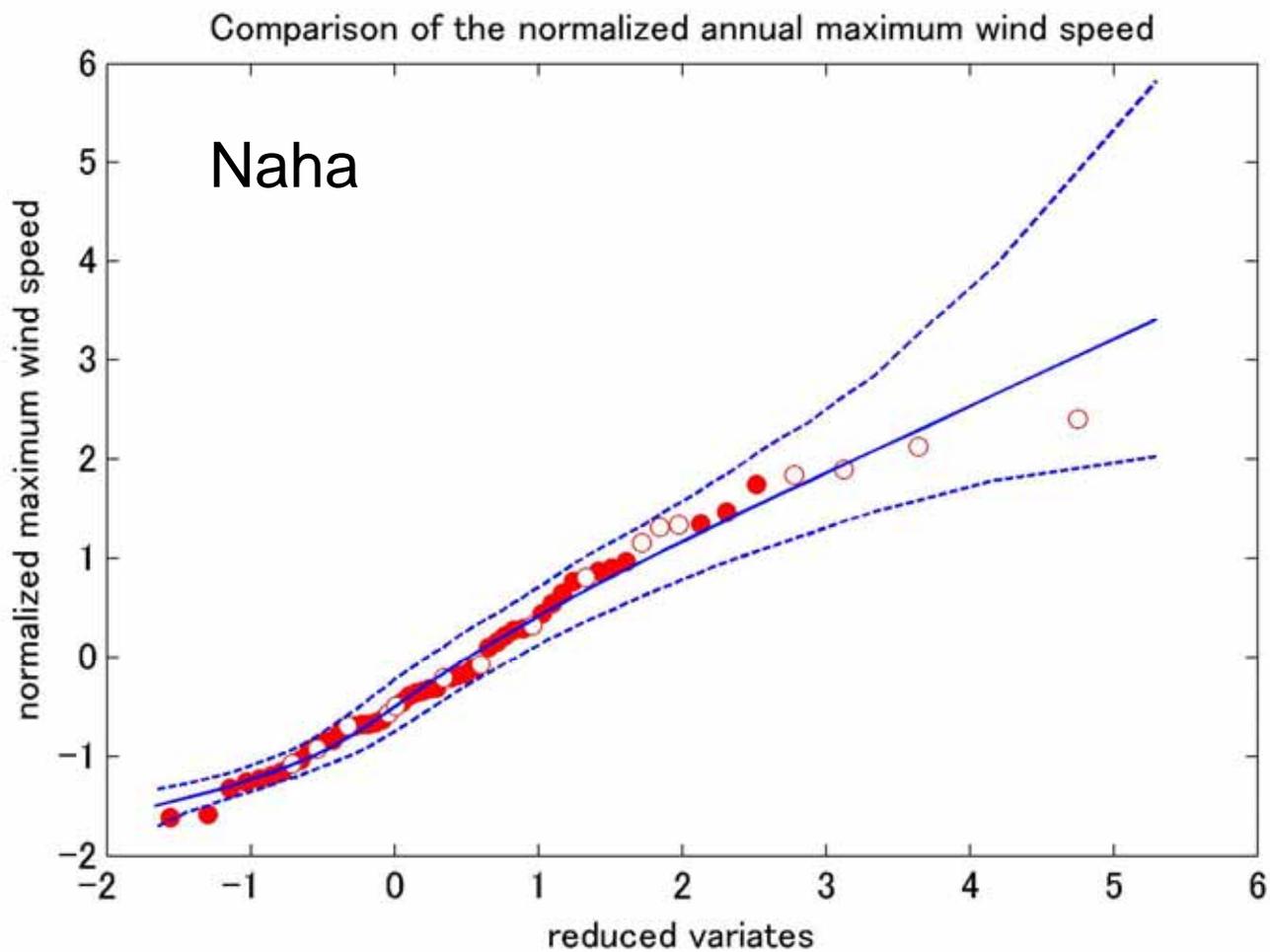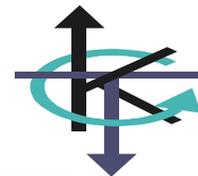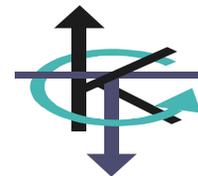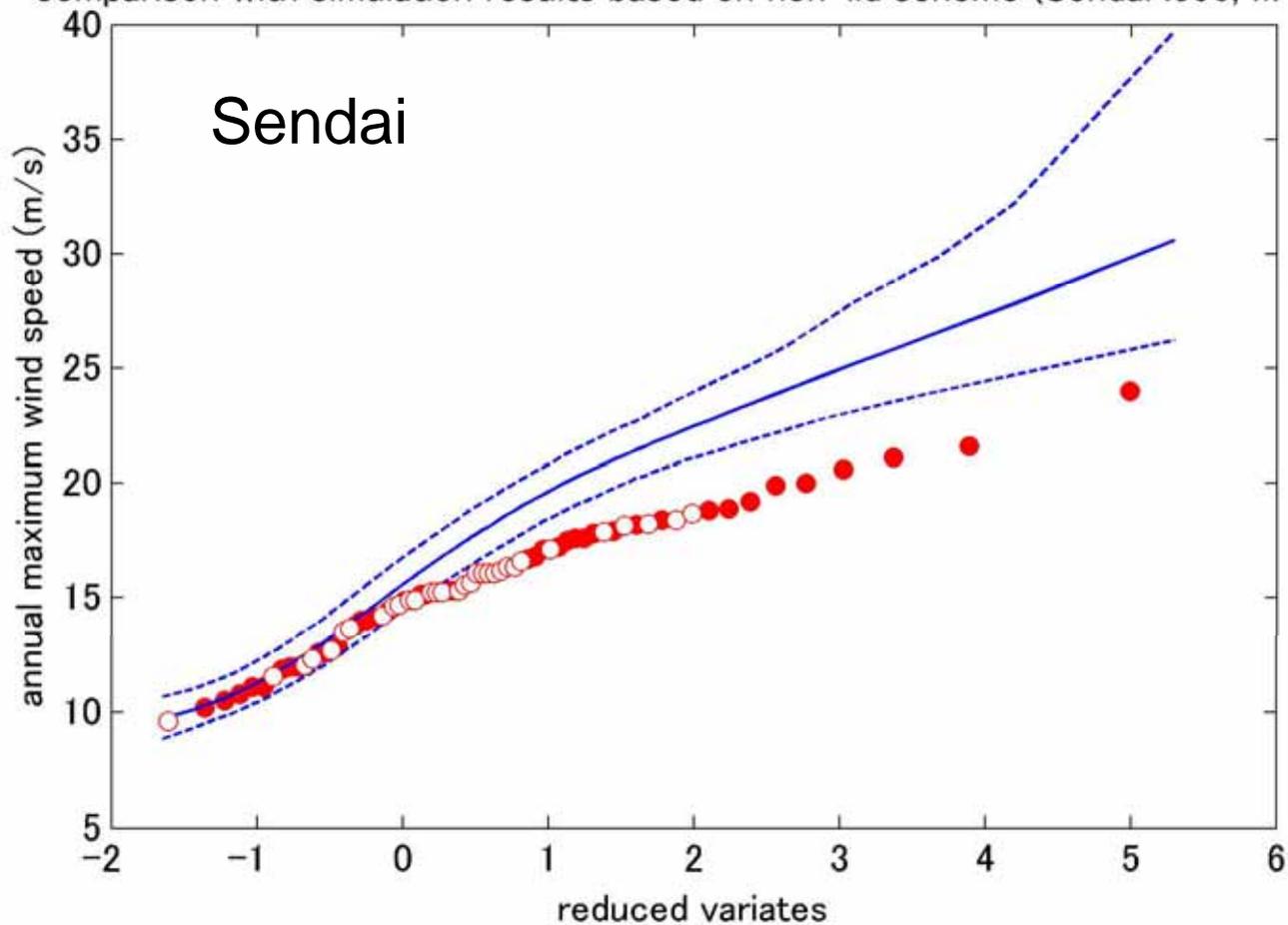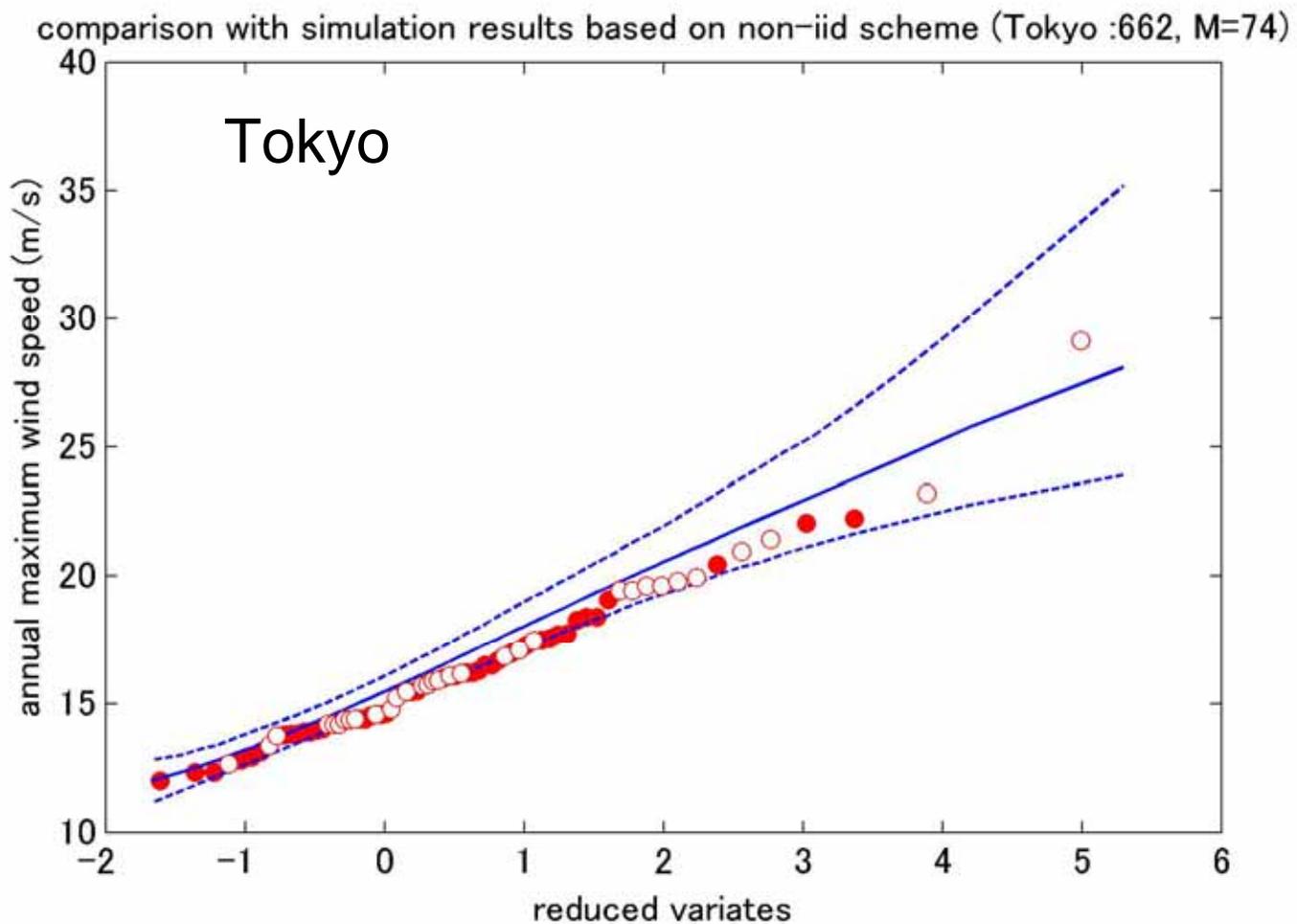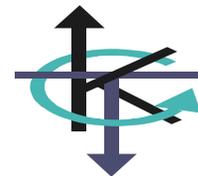Comparison of the normalized annual maximum wind speed

Aomori

- - - - - 95% confidence intervals

——— Mean of 1000 samples

$(Z-a_{m,n})/b_{m,n}$

normalized maximum wind speed

reduced variates

○ : before 1961   ● : after 1962

Comparison of the normalized annual maximum wind speed

Sendai

Comparison of the normalized annual maximum wind speed

Tokyo

Comparison of the normalized annual maximum wind speed

Kobe

Comparison of the normalized annual maximum wind speed

Shionomisaki

Comparison of the normalized annual maximum wind speed

Makurazaki

Comparison of the normalized annual maximum wind speed

Naha

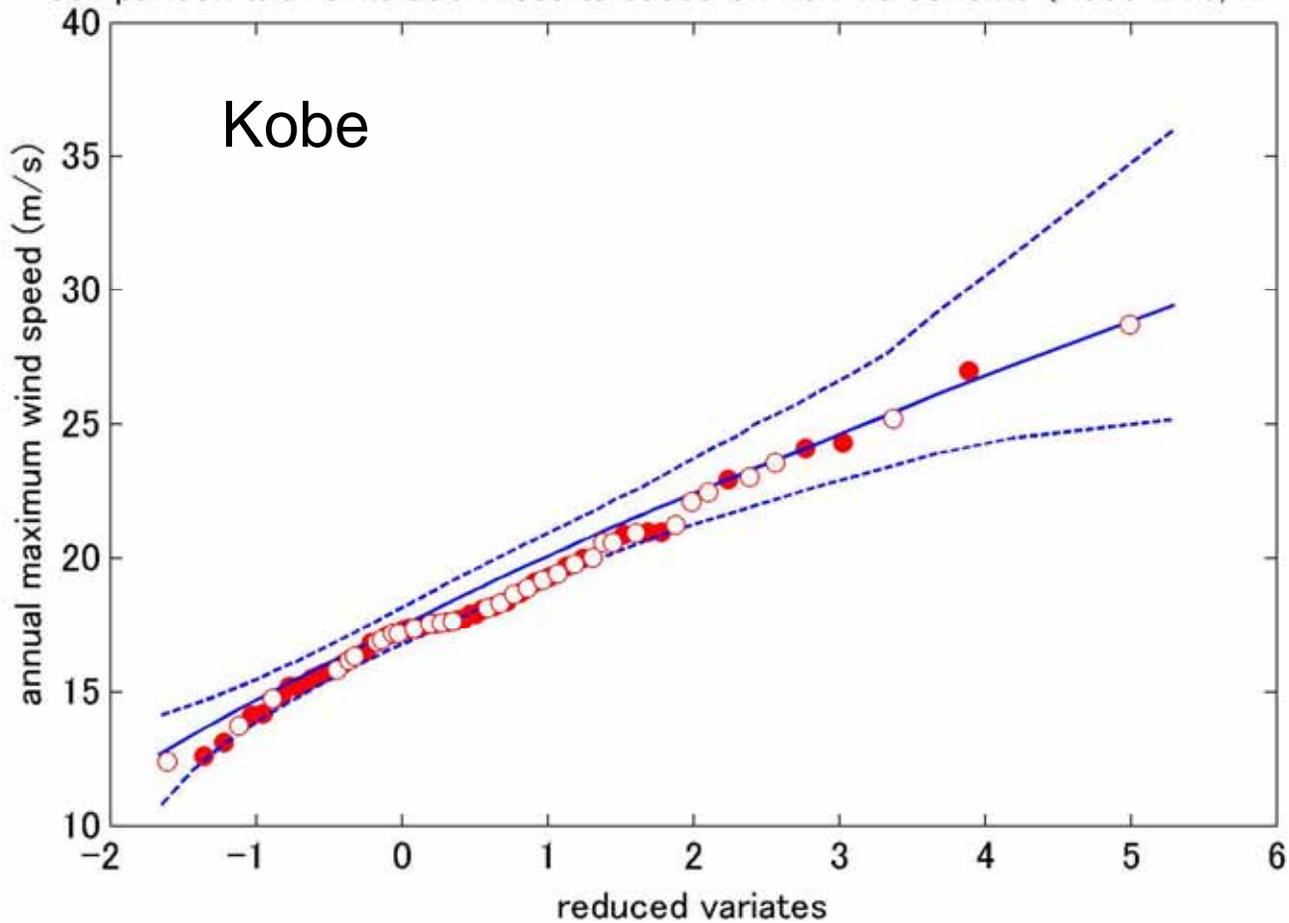# 7) Comparison of the quantile functions from MCS and historical records in full scale



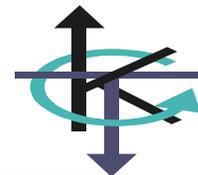comparison with simulation results based on non-iid scheme (Aomori :575, M=74)

Aomori

comparison with simulation results based on non-iid scheme (Sendai :590, M=74)

Sendai

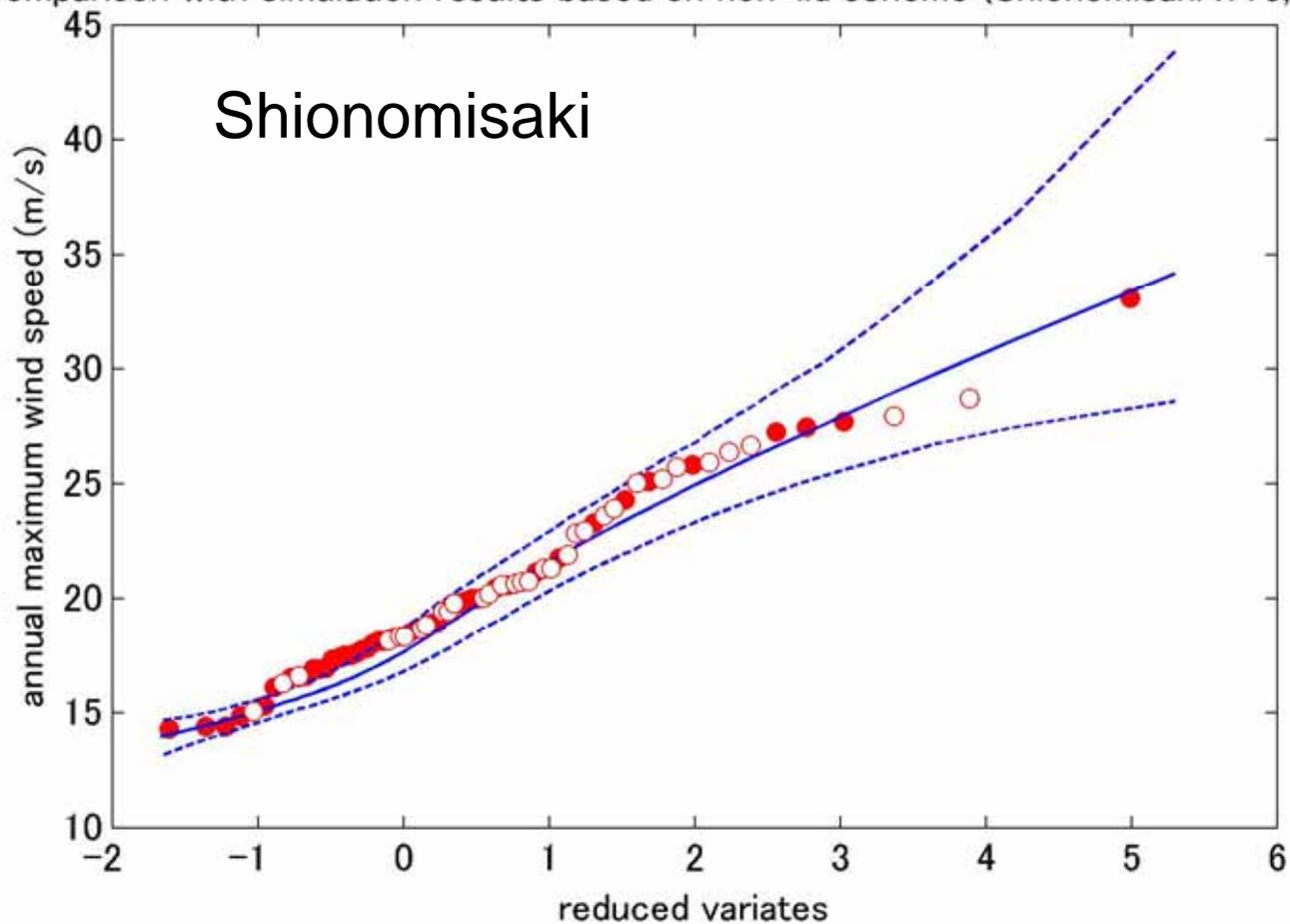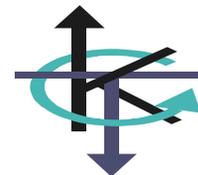comparison with simulation results based on non-iid scheme (Tokyo :662, M=74)

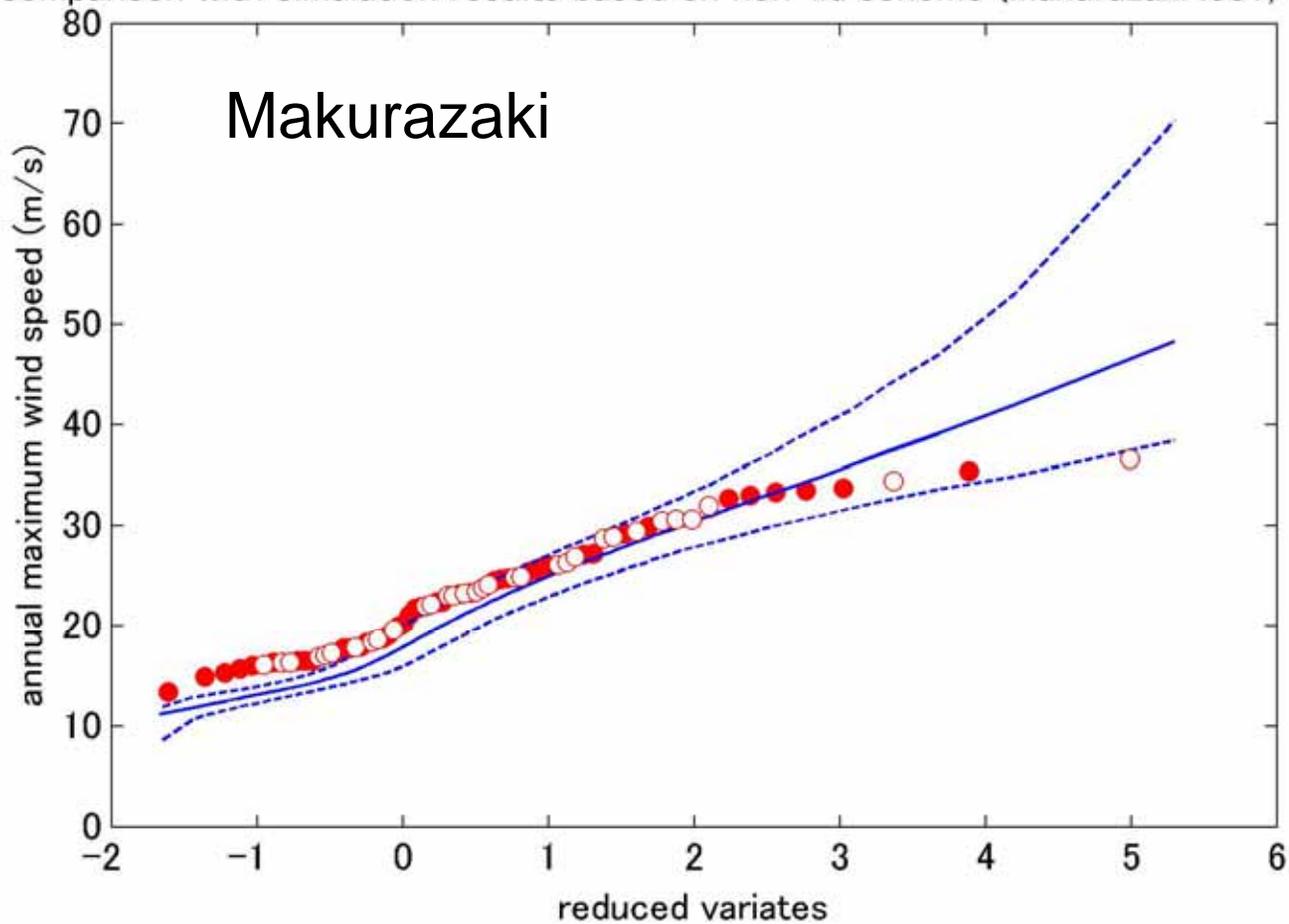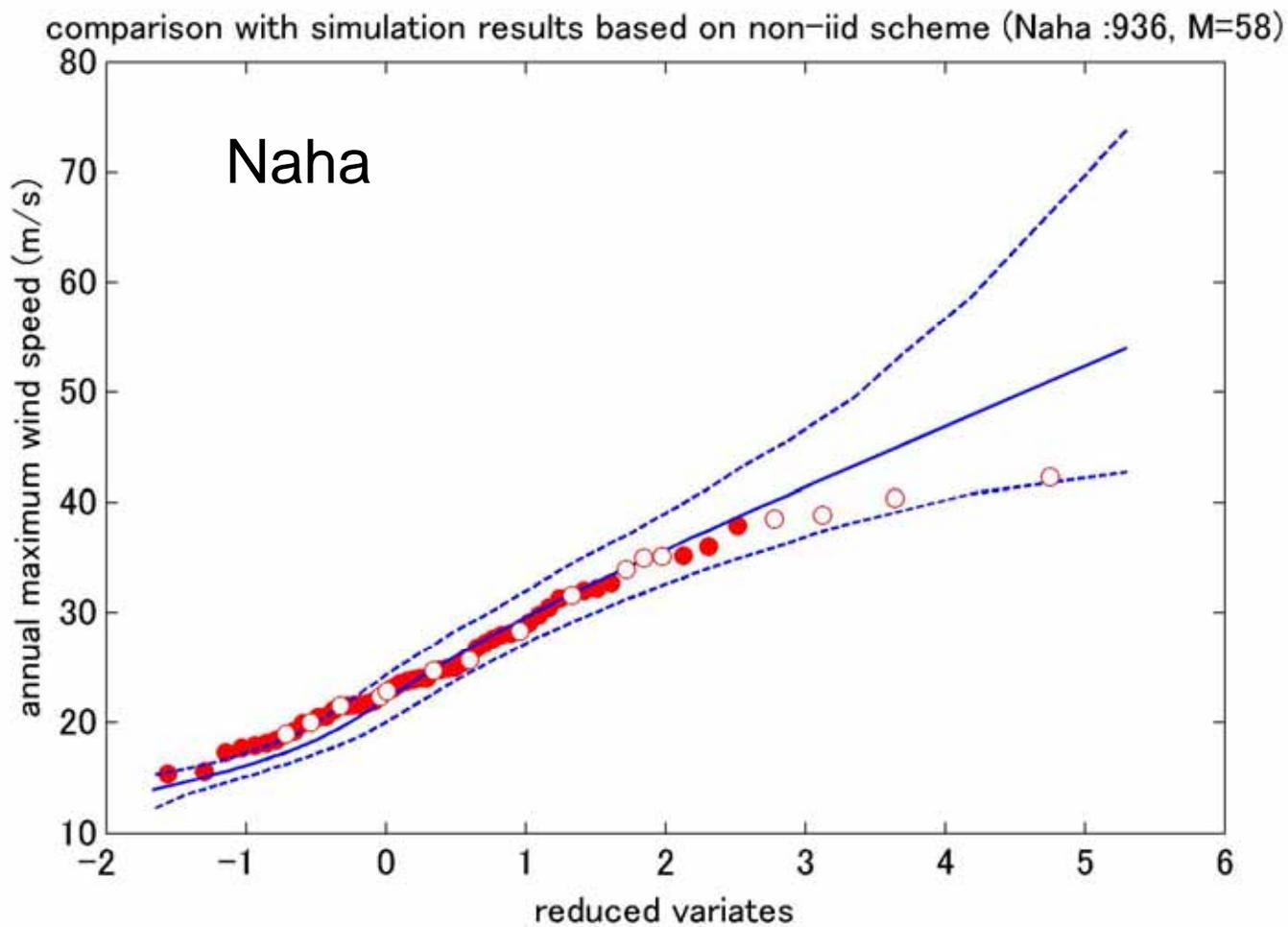comparison with simulation results based on non-iid scheme (Kobe :770, M=74)

Kobe

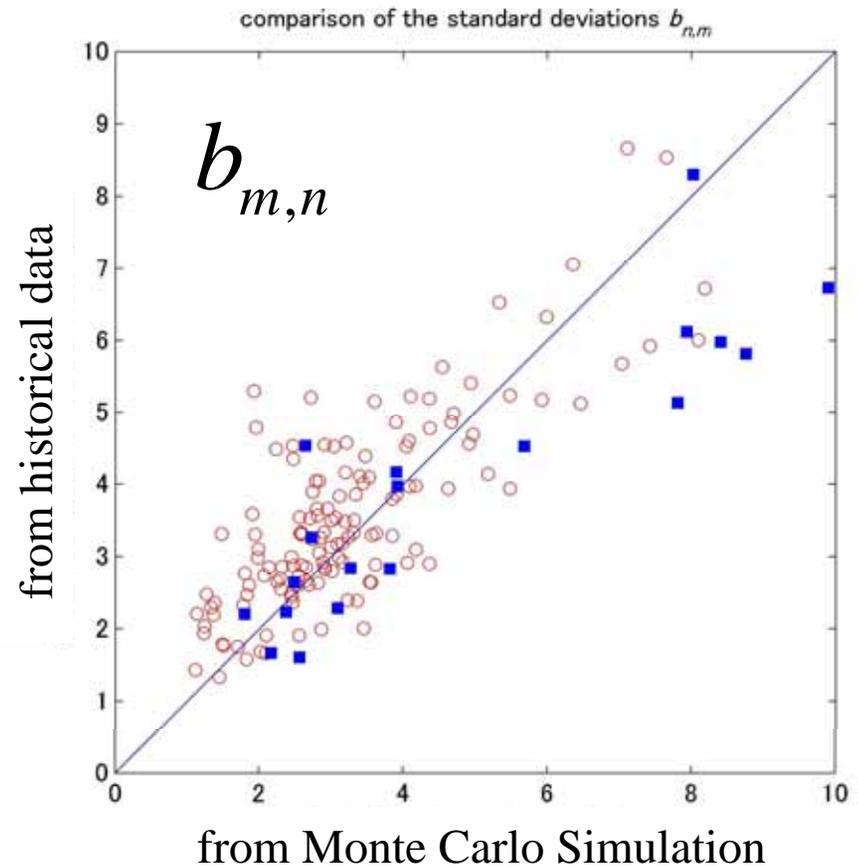comparison with simulation results based on non-iid scheme (Shionomisaki :778, M=74)

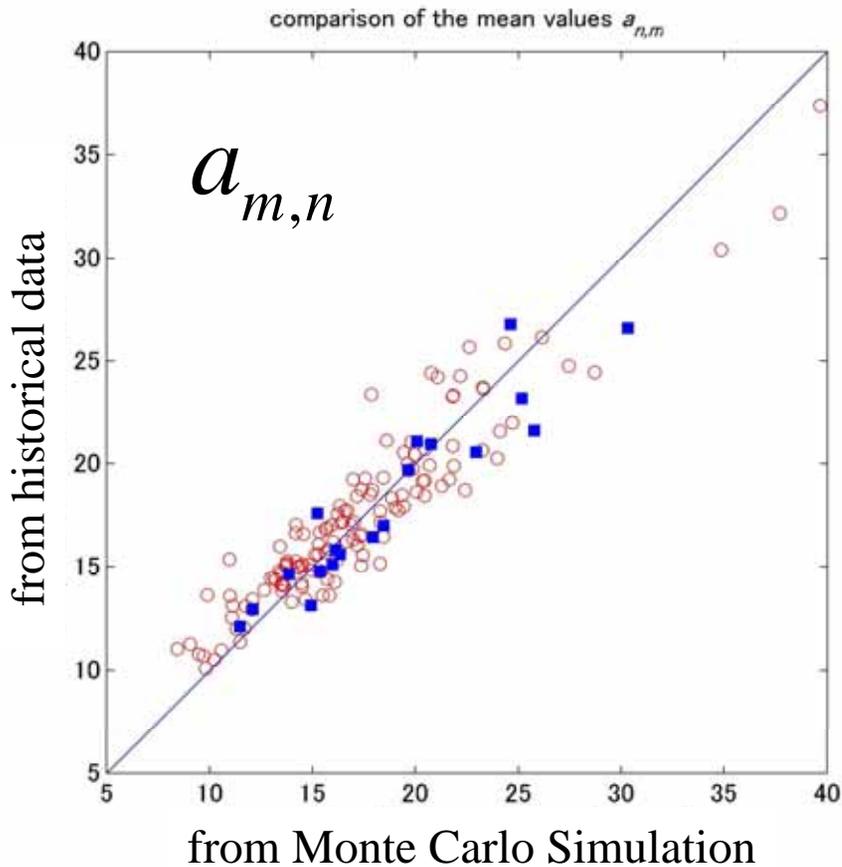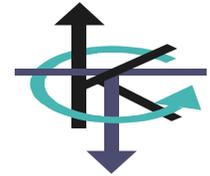Shionomisaki

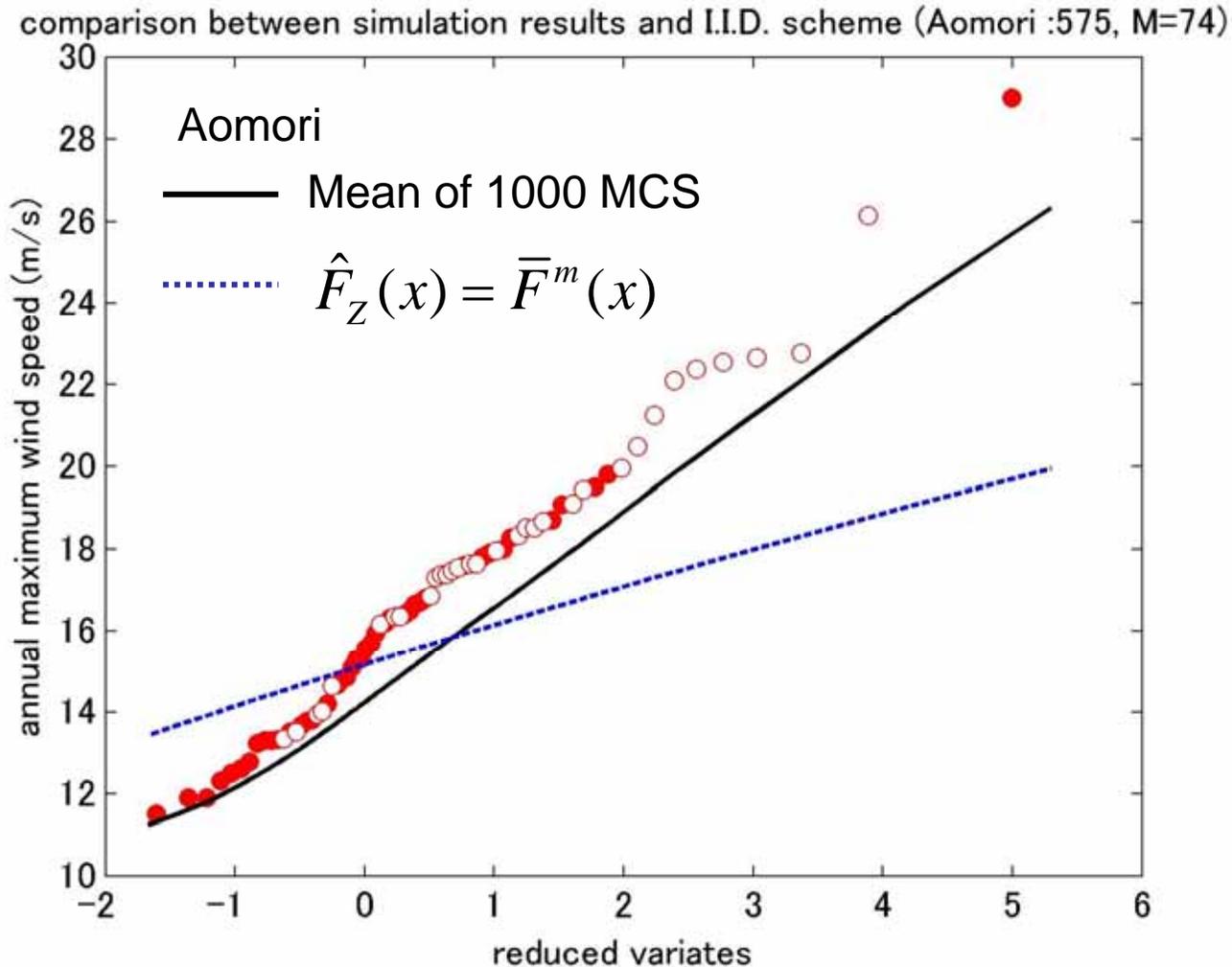comparison with simulation results based on non-iid scheme (Makurazaki :831, M=74)

Makurazaki

comparison with simulation results based on non-iid scheme (Naha :936, M=58)

Naha

# 8) Comparison of the attraction coefficients from MCS and the historical records
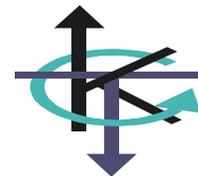


comparison of the mean values $a_{n,m}$

$a_{m,n}$

from historical data

from Monte Carlo Simulation

comparison of the standard deviations $b_{n,m}$

$b_{m,n}$

from historical data

from Monte Carlo Simulation

$\bigcirc$ : n>50 (136 sites), $\blacksquare$ : n ≤ 50 (19 sites)
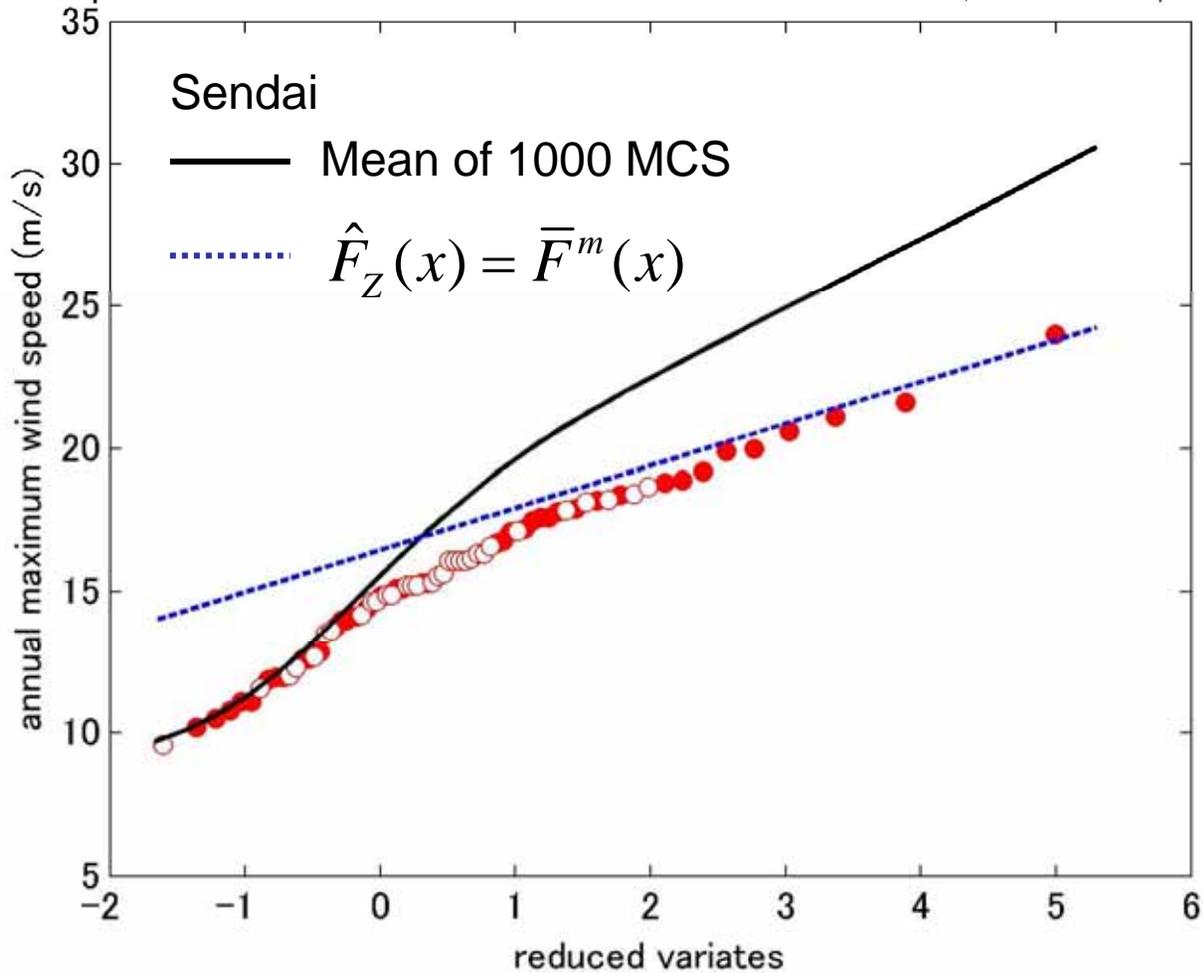
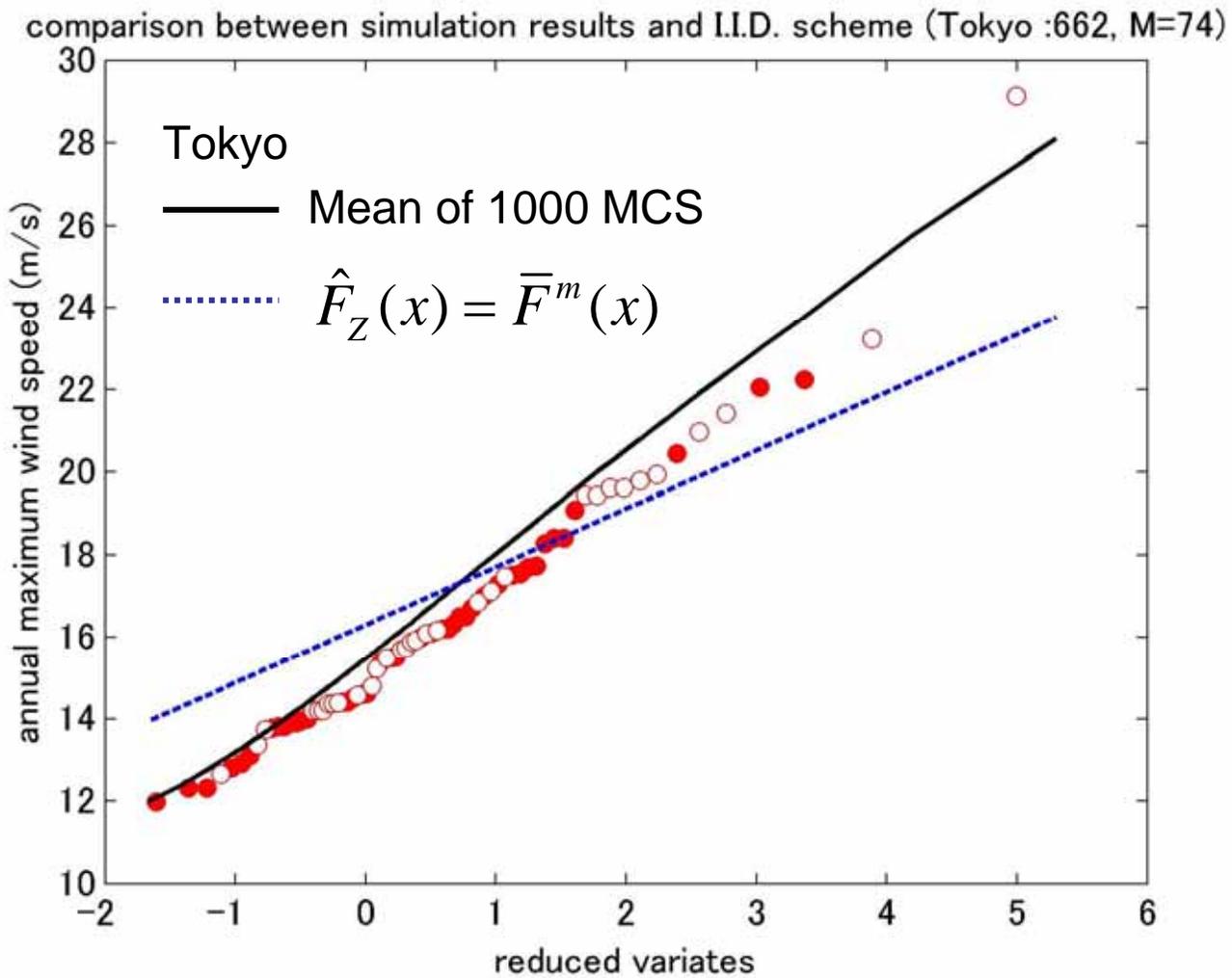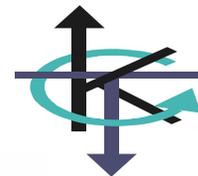$$\lim_{n\to\infty} a_{m,n} = a_m \ , \ \lim_{n\to\infty} b_{m,n} = b_m$$

# 9) Defect of the alternative definition $\hat{F}_Z(x) = \overline{F}^m(x)$



comparison between simulation results and I.I.D. scheme (Aomori :575, M=74)

Aomori

—— Mean of 1000 MCS

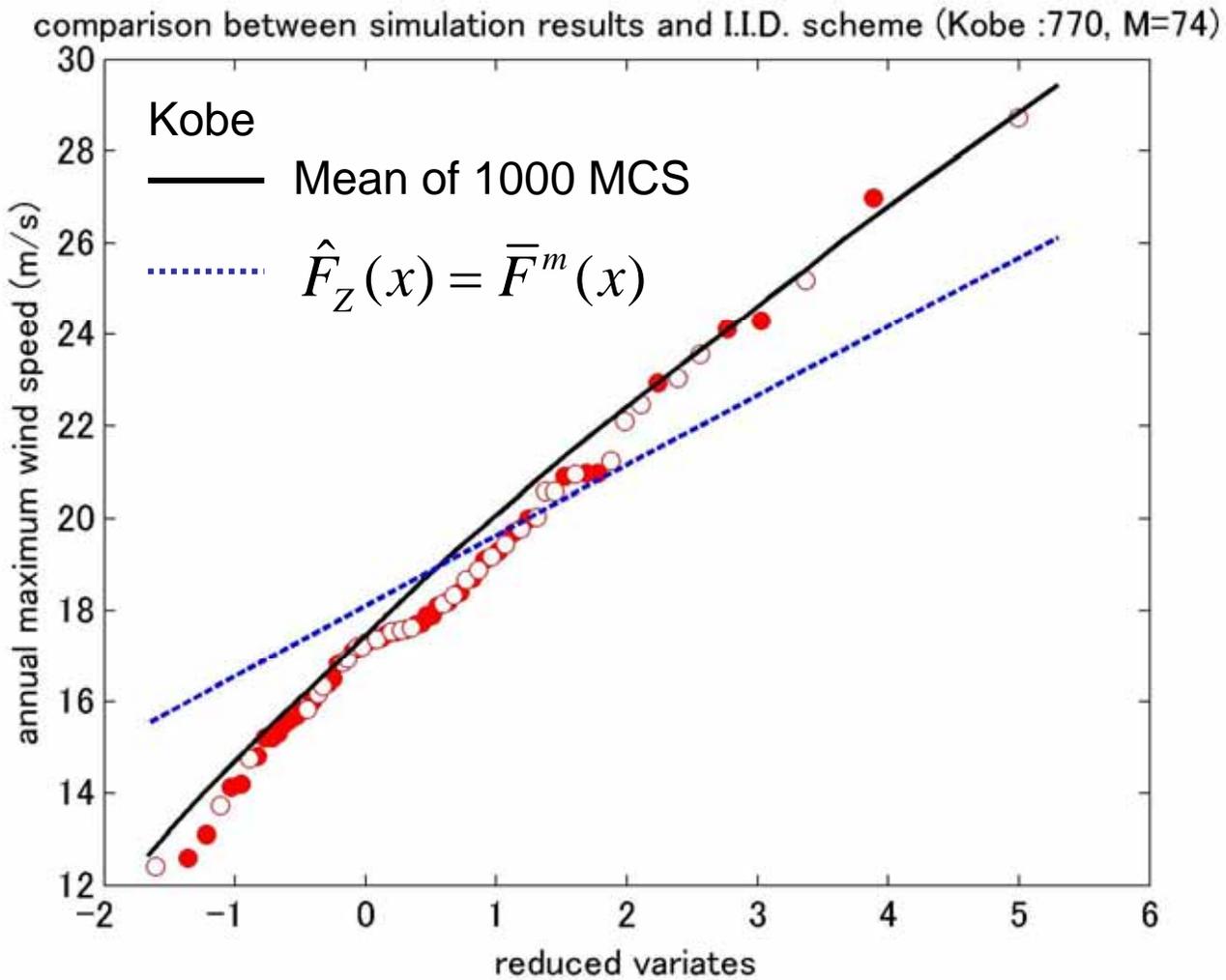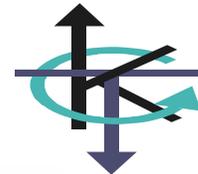········· $\hat{F}_Z(x) = \overline{F}^m(x)$

comparison between simulation results and I.I.D. scheme (Sendai :590, M=74)

Sendai
—— Mean of 1000 MCS
·········· $\hat{F}_Z(x) = \overline{F}^m(x)$

comparison between simulation results and I.I.D. scheme (Tokyo :662, M=74)

Tokyo

——— Mean of 1000 MCS

$\cdots\cdots\cdots\quad \hat{F}_Z(x) = \overline{F}^m(x)$
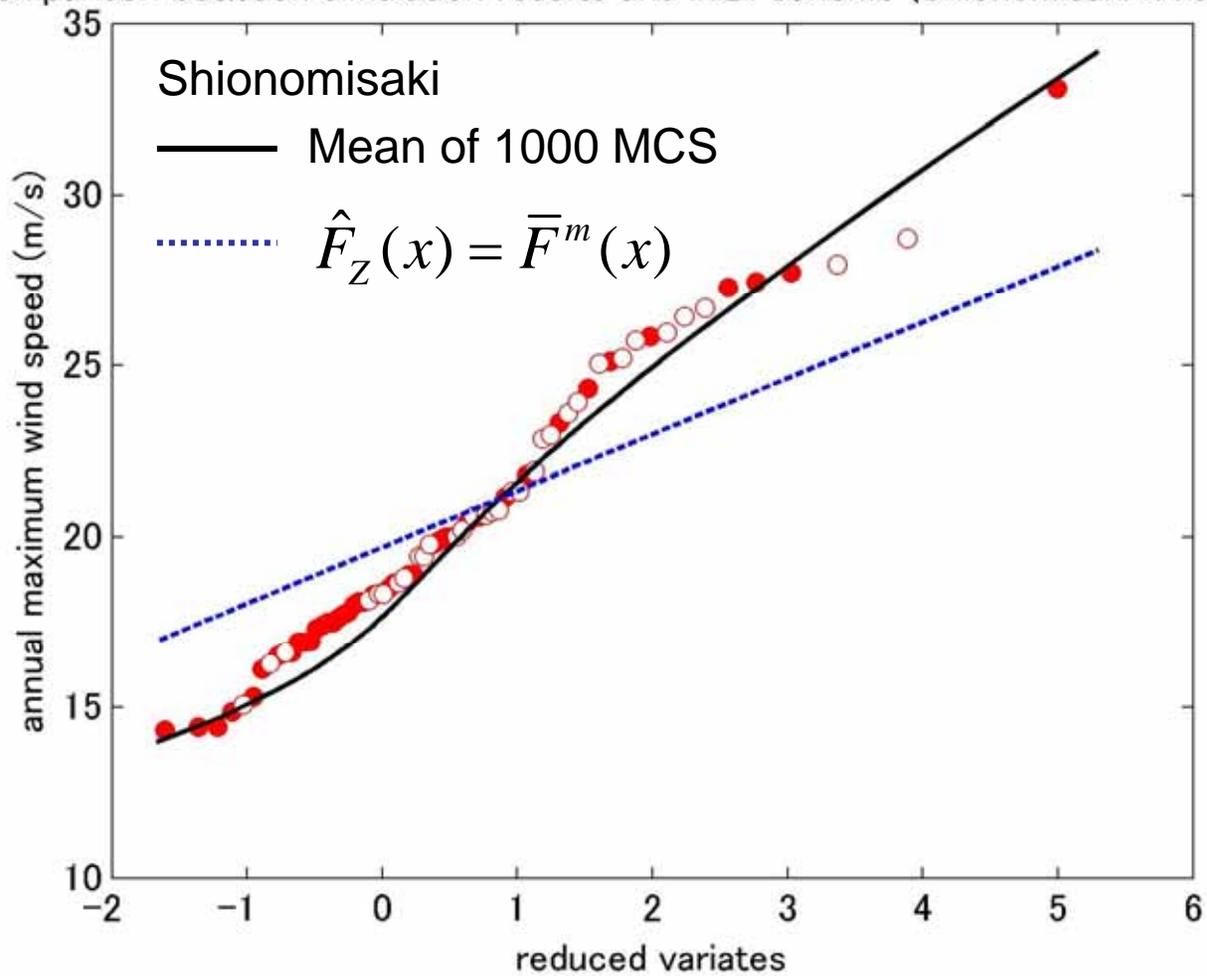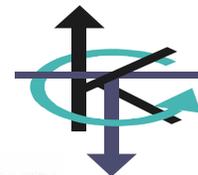
annual maximum wind speed (m/s)

reduced variates

comparison between simulation results and I.I.D. scheme (Kobe :770, M=74)

Kobe

—— Mean of 1000 MCS

$\cdots\cdots \hat{F}_Z(x) = \overline{F}^m(x)$

annual maximum wind speed (m/s)

reduced variates

Shionomisaki

—— Mean of 1000 MCS

$$\cdots\cdots \hat{F}_Z(x) = \overline{F}^m(x)$$

annual maximum wind speed (m/s)

reduced variates

comparison between simulation results and I.I.D. scheme (Makurazaki :831, M=74)

Makurazaki

—— Mean of 1000 MCS

$$\hat{F}_Z(x) = \overline{F}^m(x)$$

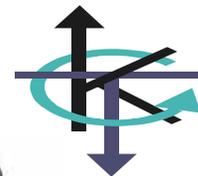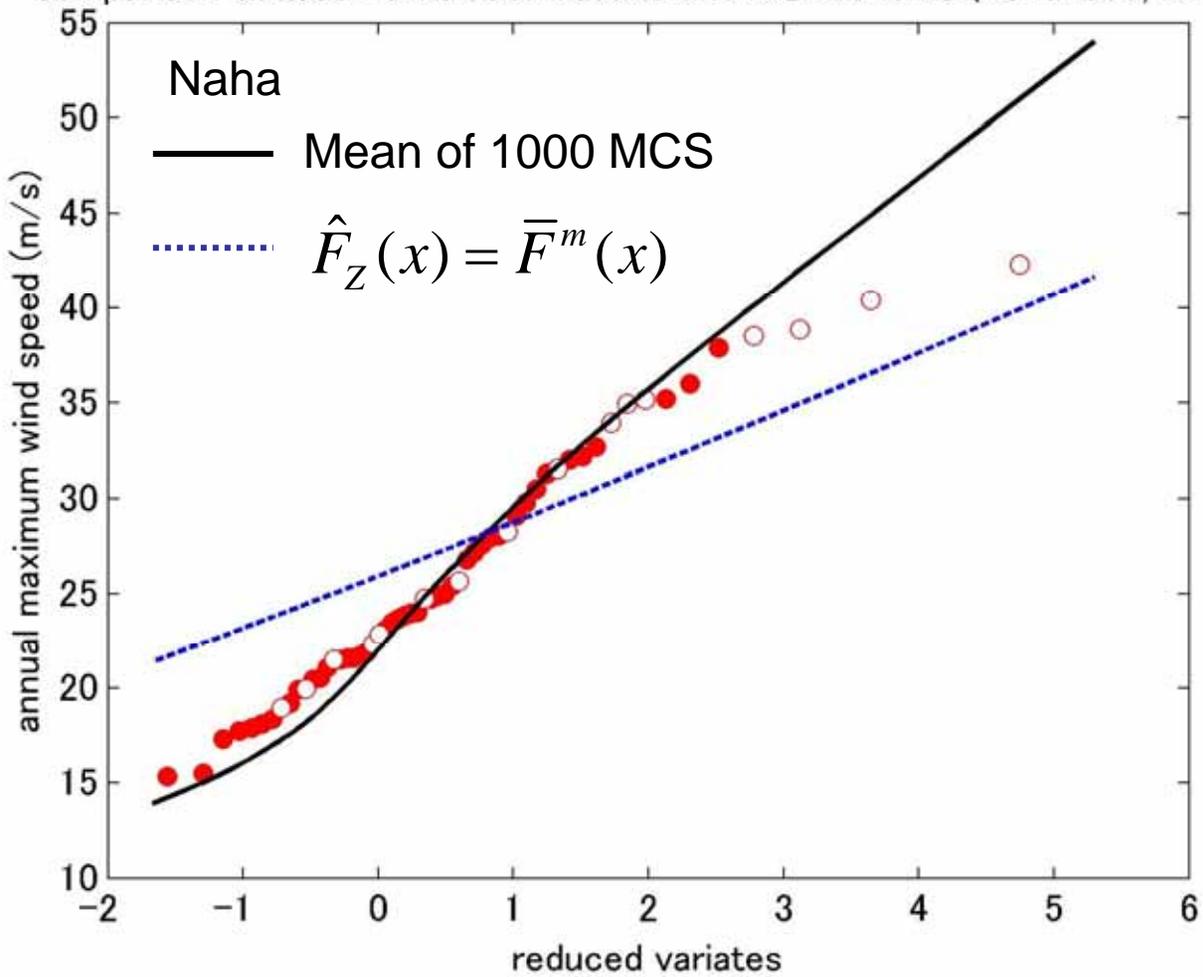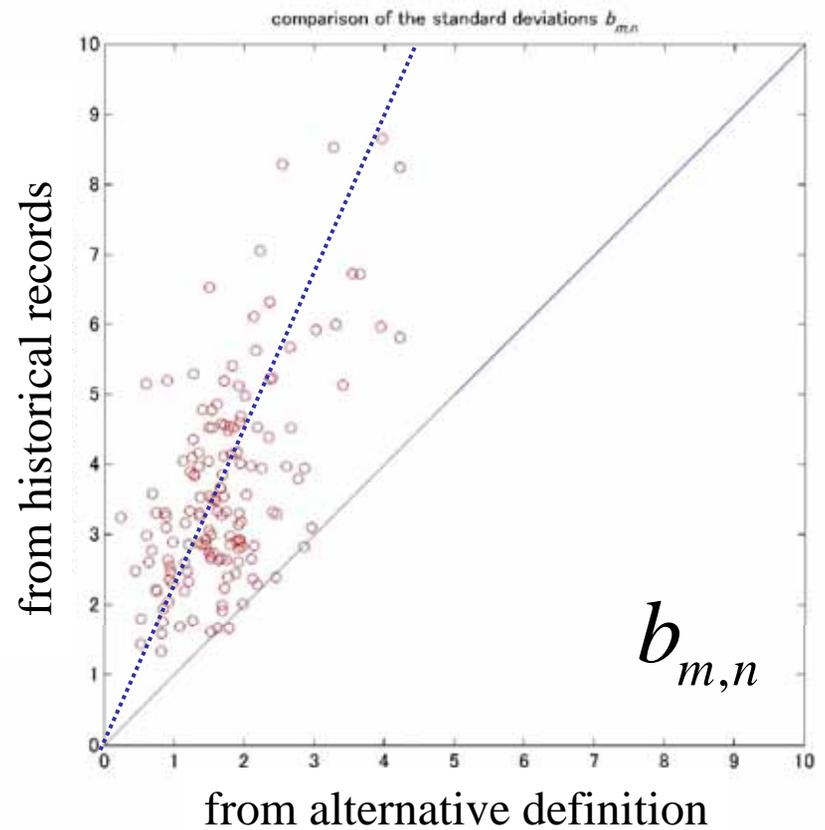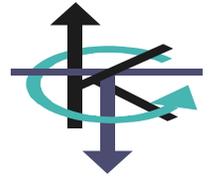comparison between simulation results and I.I.D. scheme (Naha :936, M=58)

Naha

—— Mean of 1000 MCS

$\cdots\cdots \hat{F}_Z(x) = \overline{F}^m(x)$

annual maximum wind speed (m/s)

reduced variates

# 10) Comparison of the attraction coefficients from the alternative definition and the historical records



comparison of the mean values $a_{m,n}$

from historical records

from alternative definition

$a_{m,n}$

comparison of the standard deviations $b_{m,n}$

from historical records

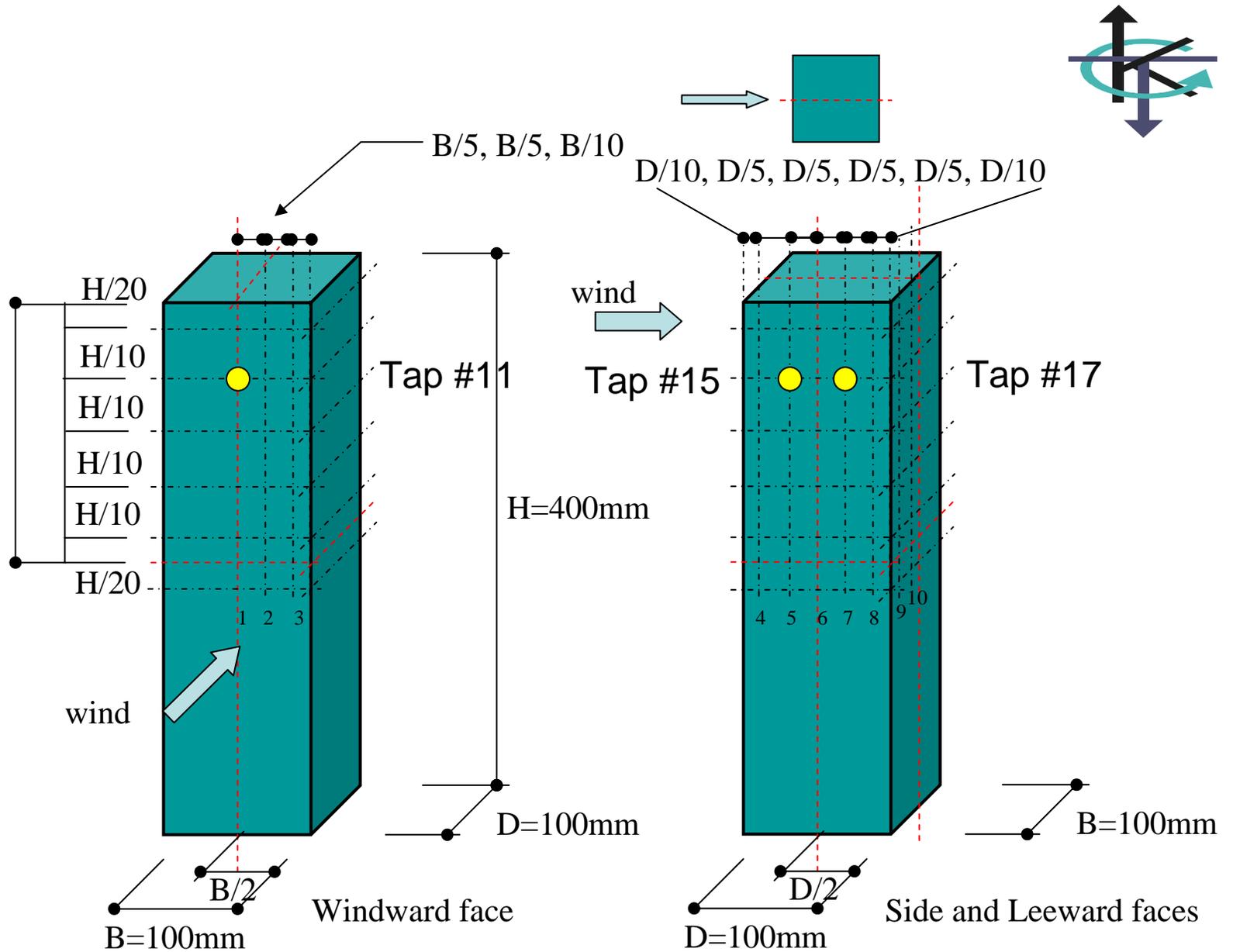from alternative definition

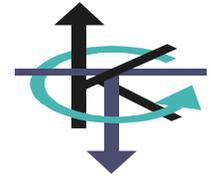$b_{m,n}$

# 10 The cult of isolated statistics and The law of large number

How many extreme values should be used to estimate an extreme value distribution?
( case study for max/min pressure coefficients)*



MATHEMATICA

$\frac{1}{n}(x_1+...+x_n) \rightarrow E(X)$

HELVETIA 80

BURKARD WALTENSPÜL    1994    COURVOISIER

* Choi & Kanda (2004), Stability of extreme quantile function estimation from relatively short records having different parent distributions, *Proc. 18th Natl. Symp. Wind Engr.*, p455~460 (in Japanese)

B/5, B/5, B/10

D/10, D/5, D/5, D/5, D/5, D/10

H/20

H/10

H/10

H/10

H/10

H/20

Tap #11

Tap #15

Tap #17

wind

H=400mm

1  2  3

4  5  6  7  8  9 10

wind

D=100mm

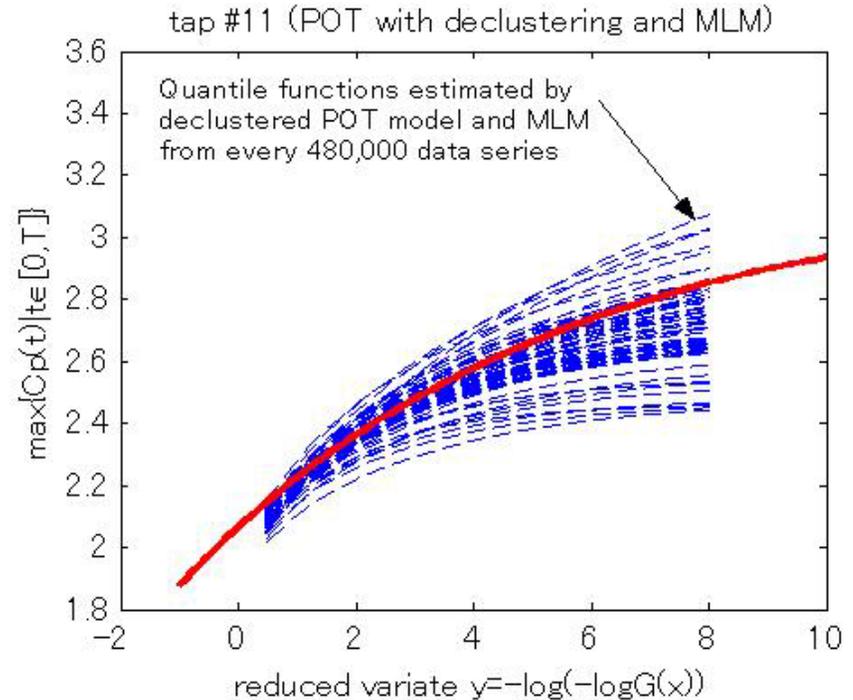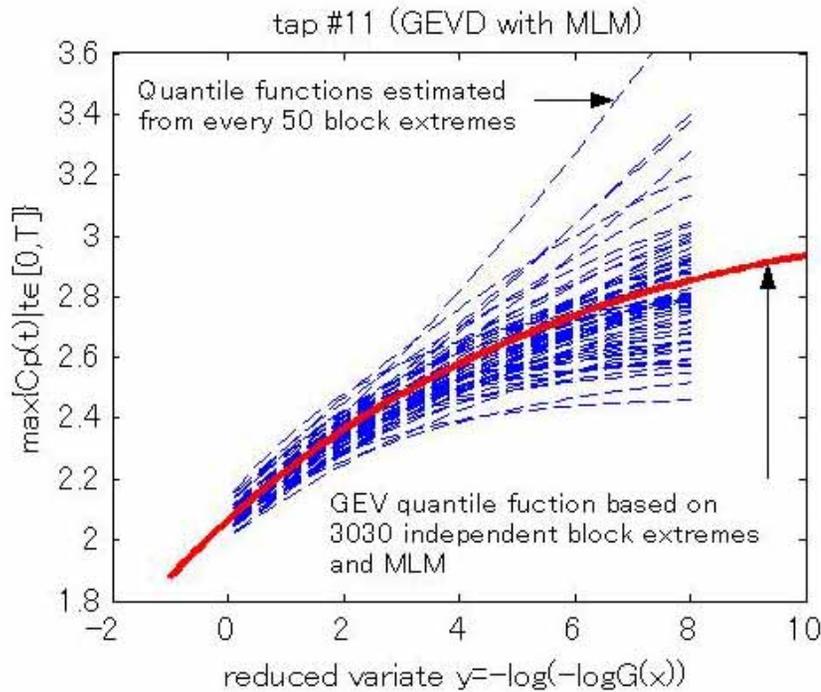B=100mm

B/2

D/2

Windward face

Side and Leeward faces

B=100mm

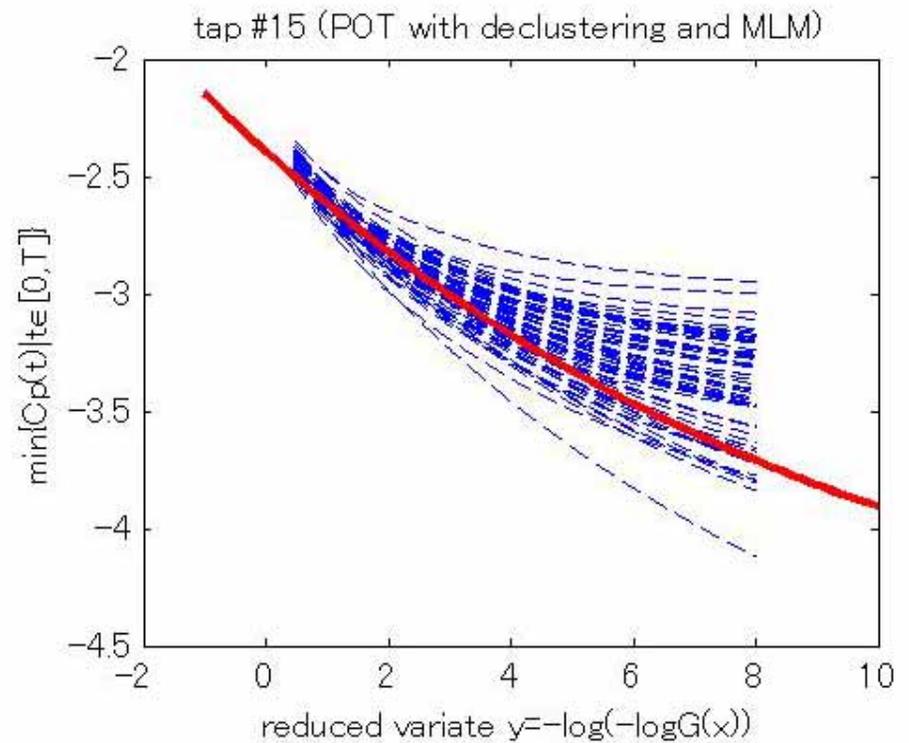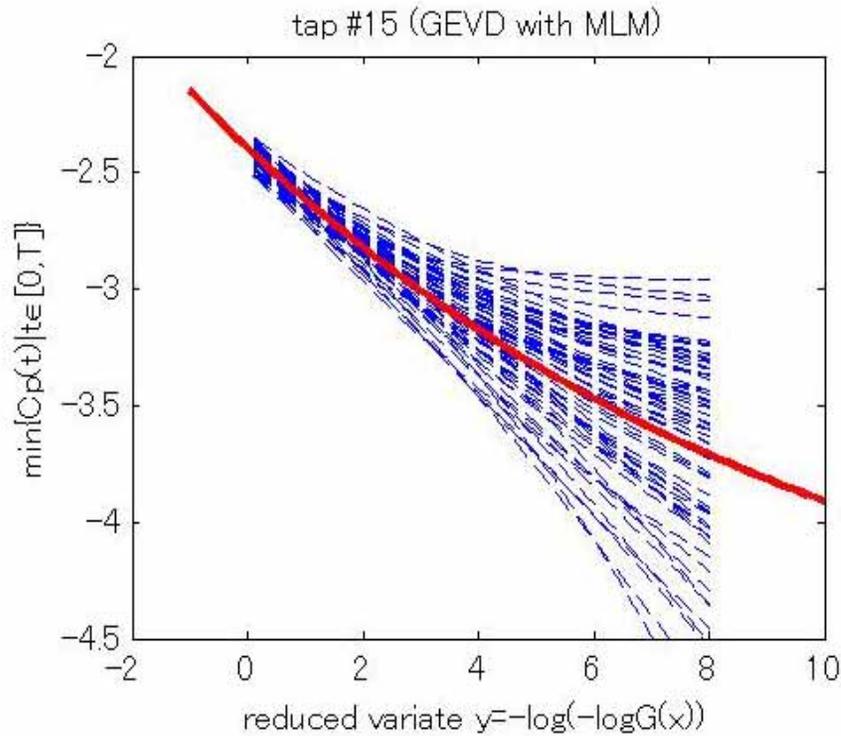D=100mm

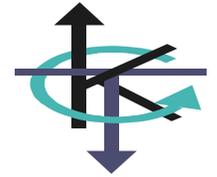# Which one is the best estimation?

$$C_p(t) = \frac{p(t) - p_s}{1/2\rho U_H^2}$$

$p(t)$ : instant total pressure, $p_s$ : static pressure

$\rho$ : air density, $U_H$ : wind velocity at model height



* 50 blocks contain about 480,000 discrete data

# Which one is the best estimation?



tap #15 (GEVD with MLM)

tap #15 (POT with declustering and MLM)

# 10 The cult of isolated statistics and The law of large number

We never be free from the law of large number.

The statistician and the scientists/technologist need to understand that models are necessarily simplifications of the system being modelled; that they are , an absolute sense, wrong; that they are certainly provisional, but nonetheless are useful and necessary for successful quantitative thinking.

from J.A. Nelder (1986), Statistics, Science and Technology – *The address of the president, delivered to the* Royal Statistical Society *on Wednesday, April 16th, 1986, J. R. Statist. Soc. A* 149(2), p109~121