# Skoltech

Skolkovo Institute of Science and Technology

"Experimental Data Processing"

# Topic 2
# "Quasi-optimal approximation under uncertainty"

**Tatiana Podladchikova**
**Term 1B, October 2018**
**t.podladchikova@skoltech.ru**

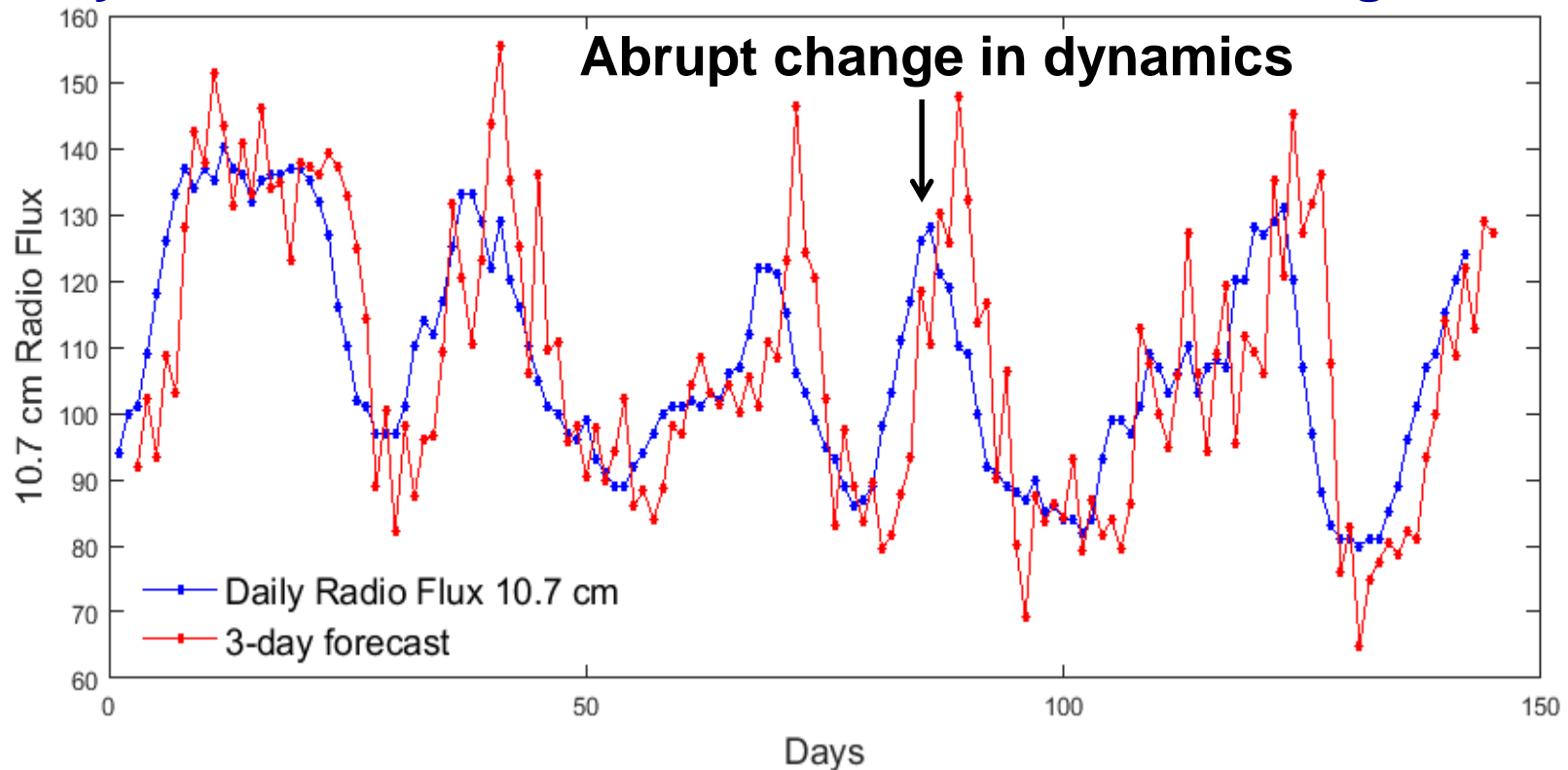# The basis of statistical analysis

**Least-square method and linear regression**

**The LSM method leads to divergence and loses its practical value when a model is inadequate or unknown.**

**Linear regression doesn't provide reliable long-term forecasting**

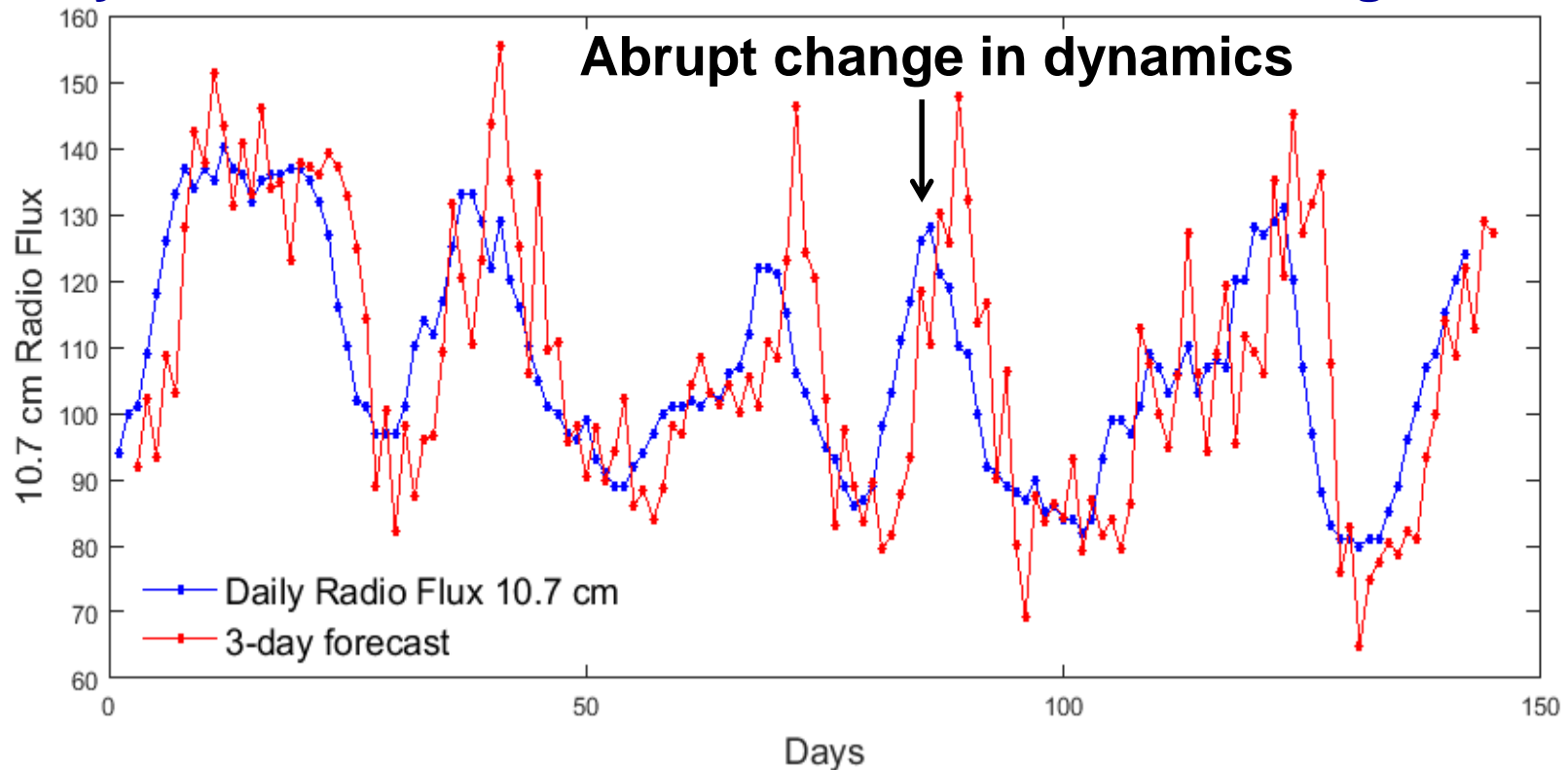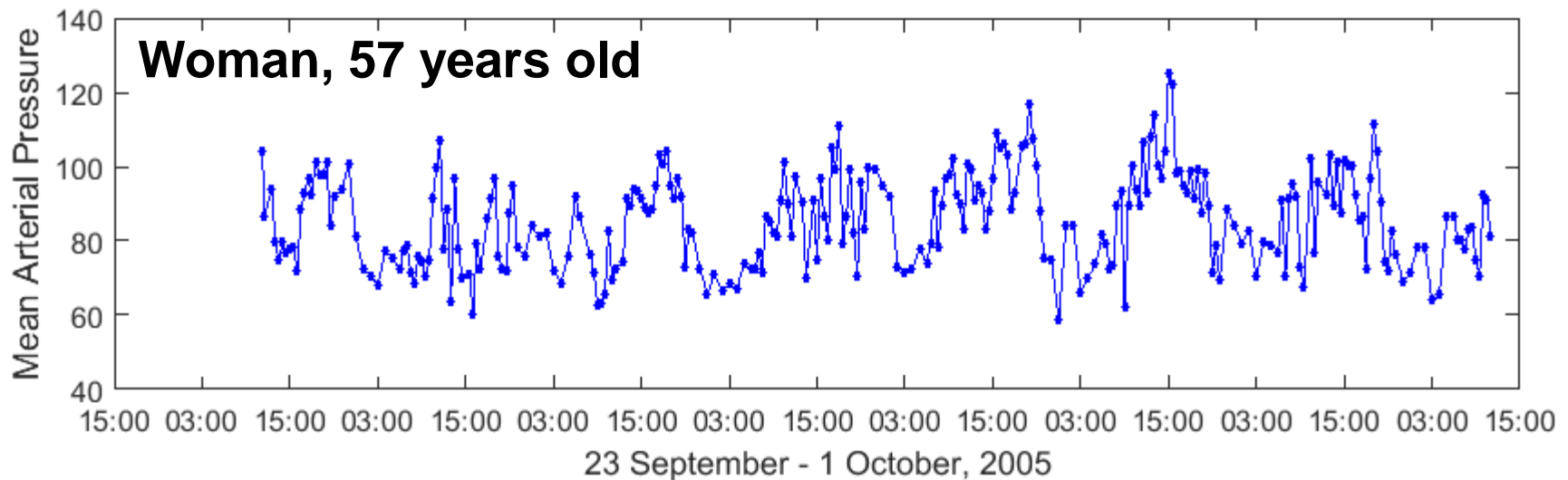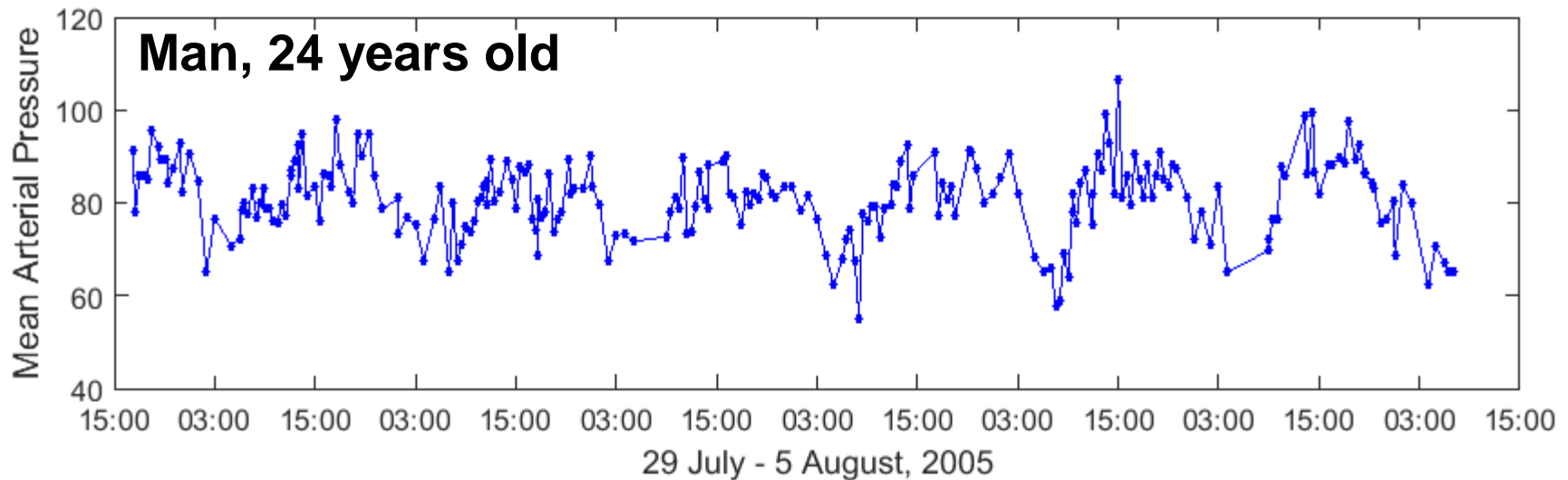# Linear regression doesn't provide long-term forecasting

**3-day 10.7 cm radio flux forecast based on a linear regression**



**Changes in dynamics of a process leads to great increase of forecasting errors**

# Linear regression doesn't provide long-term forecasting

**3-day 10.7 cm radio flux forecast based on a linear regression**
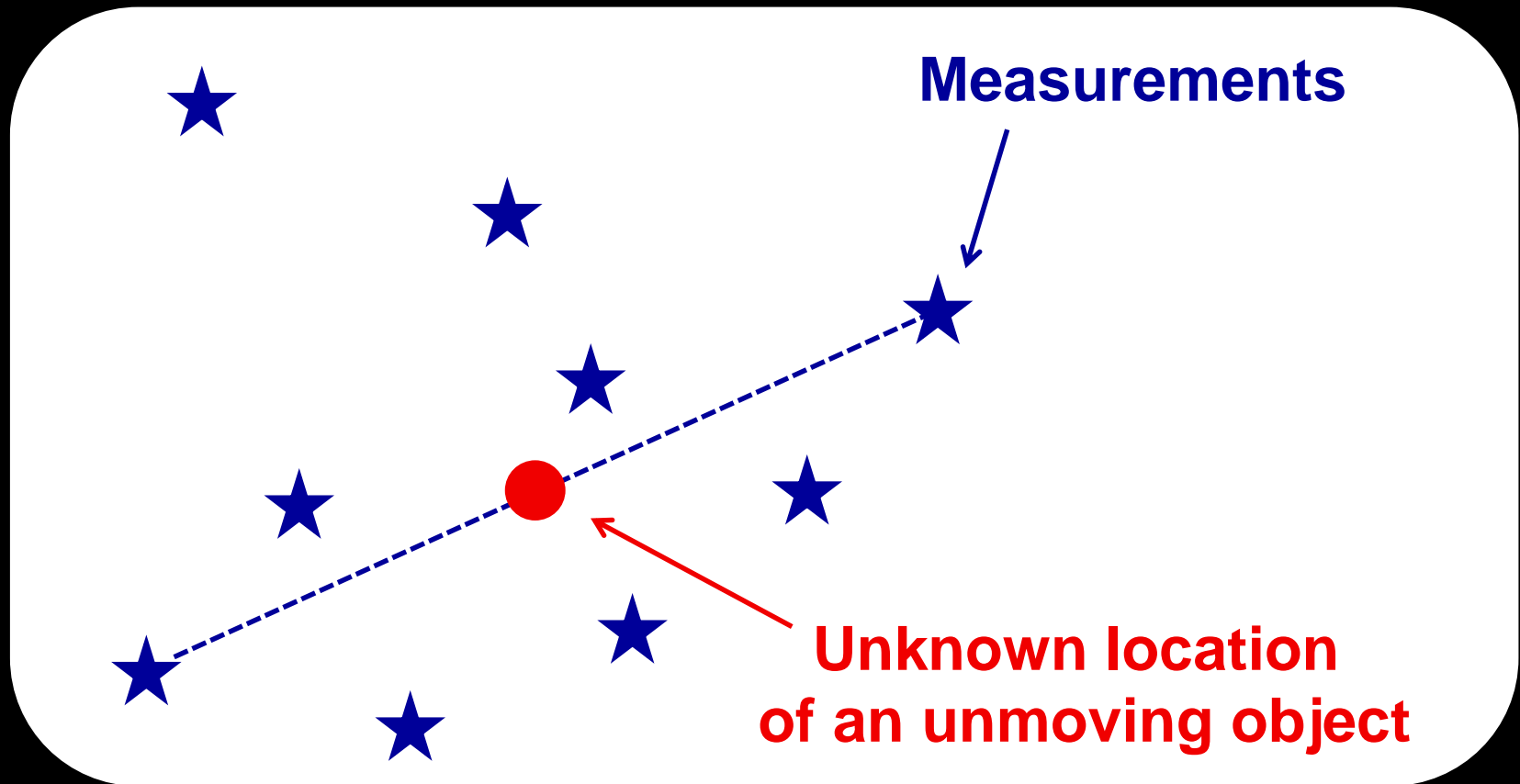


**Abrupt change in dynamics**

**Changes in dynamics of a process leads to great increase of forecasting errors**

**To extract regularities that will allow long-term forecasting we need to smooth data**

# Which regularities can you extract from measurements of mean arterial pressure?

# Estimate the location of an unmoving object



**Measurements**

**Unknown location
of an unmoving object**

Smoothing is weighted averaging of noisy data.
Fluctuation components are self compensated.

# The most popular methods of quasi-optimal estimation
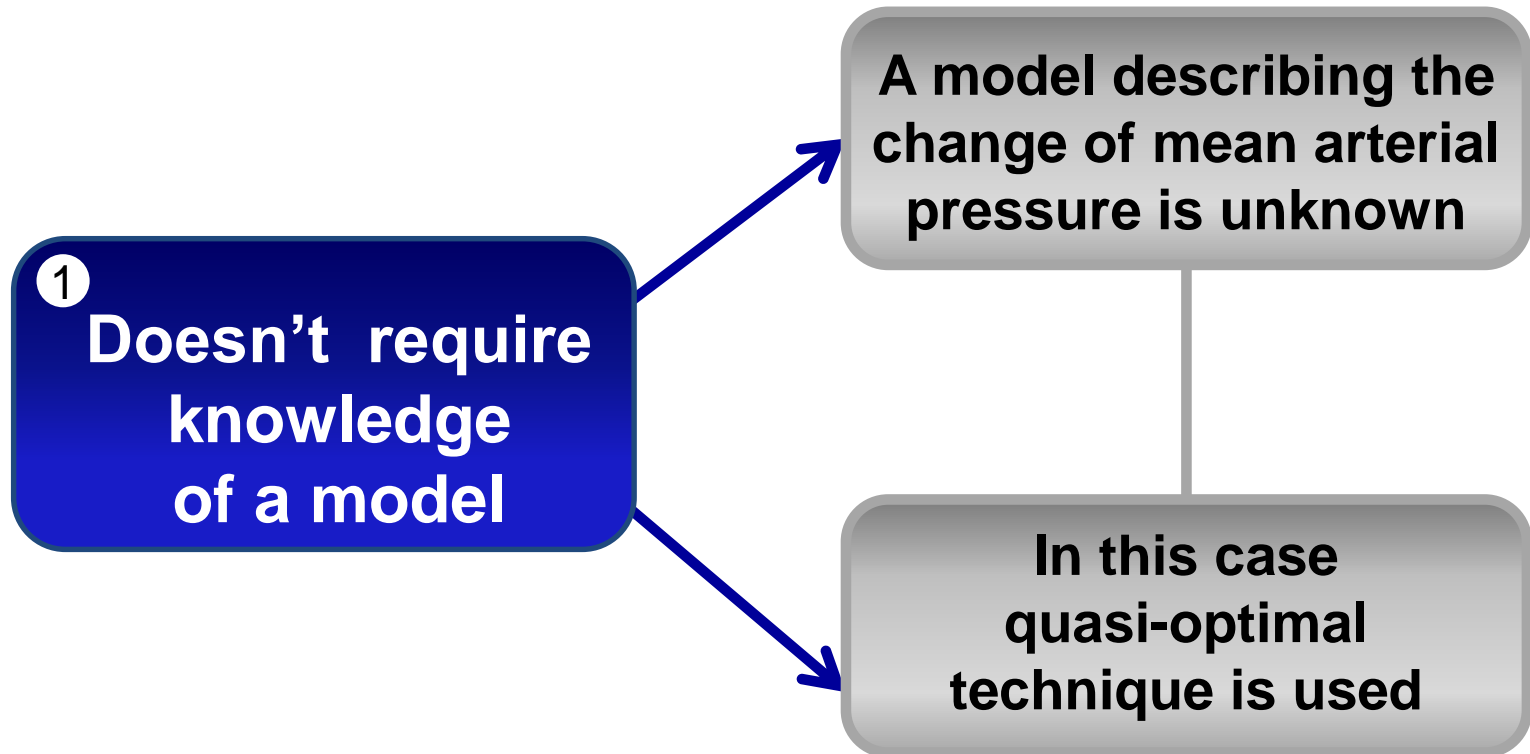
1. **Running mean**

2. **Exponential smoothing**

✓ **Advantages**

✗ **Limitations**

# ✓ Advantages of quasi-optimal estimation methods

**① Doesn't require knowledge of a model**

**A model describing the change of mean arterial pressure is unknown**

**In this case quasi-optimal technique is used**

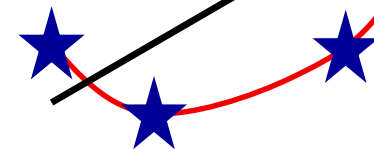# Advantages of quasi-optimal estimation methods

**2** **Robustness**
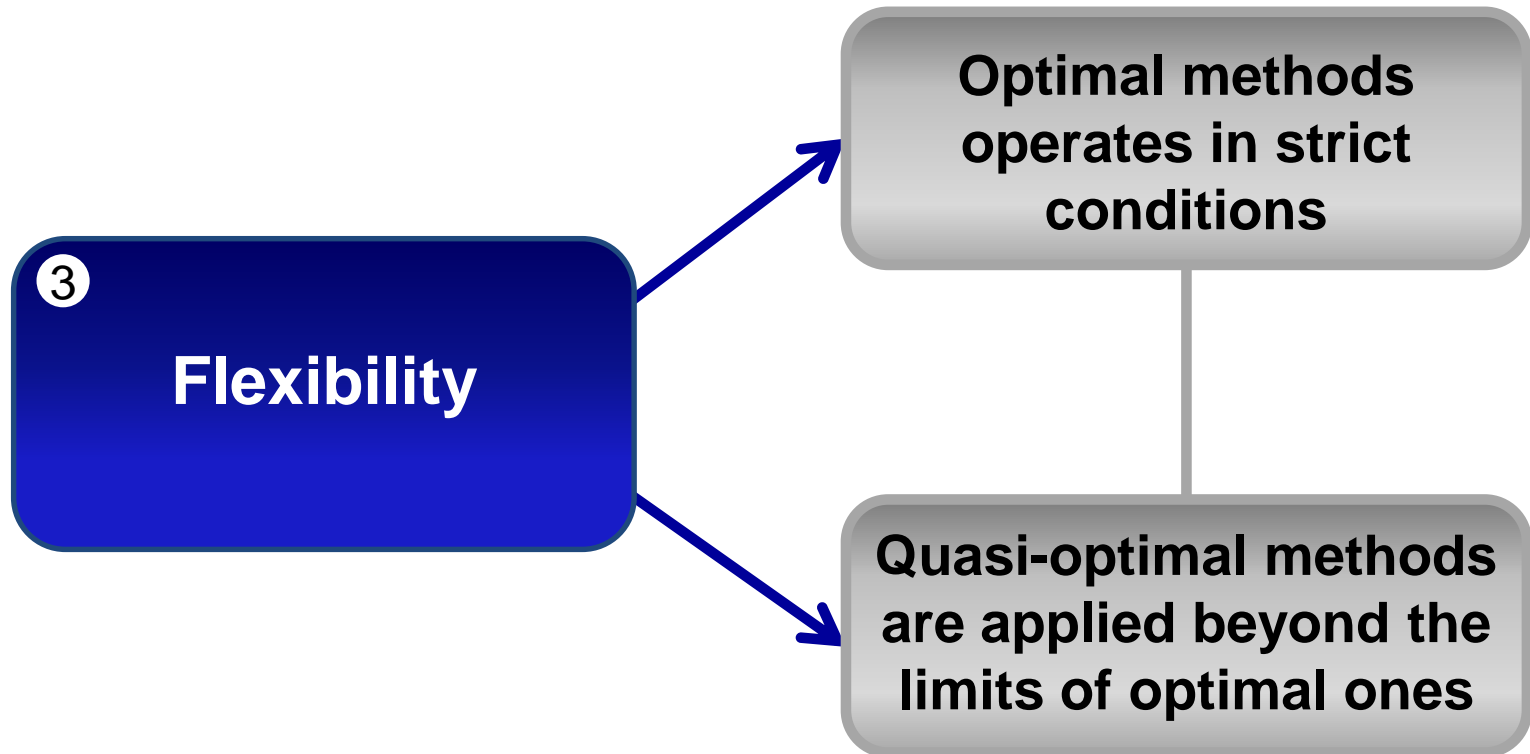
No risk of divergence

**Optimal estimation in conditions of inadequate model**

Divergence. Errors monotonously increase

# Disadvantages of quasi-optimal estimation methods

**1** **Unknown estimation errors**

**Risk of false conclusions**

**It is needed to search for ways of accuracy estimation**

# Which regularities can you extract from smoothed measurements of mean arterial pressure?



**Running mean with window $M = 7$**

# Which regularities can you extract from smoothed measurements of mean arterial pressure?



**Running mean with window $M = 7$**

# Which regularities can you extract from smoothed measurements of mean arterial pressure?



**Running mean with window $M = 7$**

# Which regularities can you extract from smoothed measurements of mean arterial pressure?



**Man, 24 years old**

**③ Maximum at 1-8 p.m.**

**Woman, 57 years old**

**Running mean with window $M = 7$**

# Which regularities can you extract from smoothed measurements of mean arterial pressure?



**Man, 24 years old**

**4 Maximum has two waves**

**Woman, 57 years old**

**Running mean with window $M = 7$**

# Which regularities can you extract from smoothed measurements of mean arterial pressure?



**Man, 24 years old**

**Woman, 57 years old**

**⑤ The dynamics is less regular**

**Running mean with window $M = 7$**

**① Running mean**

**Dilemma of setting goal**

**Greater size of running window**

**GREAT risk of process distortion, but EFFECTIVE filtration of measurement errors**

**Smaller size of running window**

**SMALL risk of process distortion, but LESS EFFECTIVE filtration of measurement errors**

# ② Exponential smoothing

Let's assume that values of process $X$ are characterized by sudden change



Measurements $z_i$

# ② Exponential smoothing

**Let's assume that values of process $X$ are characterized by sudden change**



**Measurements $z_i$**

**Latest measurements have more confidence**

## ② Exponential smoothing

$$\widehat{X}_i = \alpha z_i + (1 - \alpha)\widehat{X}_{i-1}$$

$\widehat{X}_i$
**Smoothed estimate at time $i$**

$\alpha$
**Smoothing constant $\alpha \in (0; 1)$**

$z_i$
**Measurements at time $i$**

$\widehat{X}_{i-1}$
**Smoothed estimate at time $i-1$**

## ② Exponential smoothing

$$\widehat{X}_i = \alpha z_i + (1 - \alpha)\widehat{X}_{i-1}$$

$\widehat{X}_i$
**Smoothed estimate at time $i$**

$\alpha$
**Smoothing constant $\alpha \in (0; 1)$**

$z_i$
**Measurements at time $i$**

$\widehat{X}_{i-1}$
**Smoothed estimate at time $i - 1$**

$$\widehat{X}_i = \alpha z_i + \alpha(1 - \alpha)z_{i-1} + \alpha(1 - \alpha)^2 z_{i-2} + \cdots + \alpha(1 - \alpha)^i z_0$$

**The weight of measurements decreases according to geometric progression or exponential law**

**2** Exponential smoothing: Dilemma of setting goal

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

$\widehat{X}_{i-1}$
**Previous estimate**

$(z_i - \widehat{X}_{i-1})$
**Residual – mismatch between measurement and previous estimate**

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

$\widehat{X}_{i-1}$
**Previous estimate**

$(z_i - \widehat{X}_{i-1})$
**Residual – mismatch between measurement and previous estimate**

**SMALLER $\alpha$, GREATER confidence to the latest estimate, SLOWER reaction to changes**

**Choice of $\alpha$**

**GREATER $\alpha$, GREATER confidence to the latest measurement, FASTER reaction to changes**

**But EFFECTIVE filtration of measurement errors**

**But less EFFECTIVE filtration of measurement errors**

# Comparison of smoothing methods

① **Running mean**

② **Exponential mean**

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

Last $M$ measurements are used

All previous measurements are used

1

# Comparison of smoothing methods

**1** **Running mean**

**2** **Exponential mean**

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

**Equal weights of measurements**

**The weight of measurements decreases according to exponential law**

**2**

# Comparison of smoothing methods

**1** **Running mean**

**2** Exponential mean

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

**Delay of estimation on $\frac{M-1}{2}$ steps**

**Estimation is obtained at last available time moment**

**3**

# Comparison of smoothing methods



**1 Running mean**

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

**2 Exponential mean**

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

$M = 1$

**No filtration of errors**

$\alpha = 1$

**4**

Estimates of both smoothing methods are the same

# Comparison of smoothing methods

**1** **Running mean**

**2** **Exponential mean**

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$\widehat{X}_i = \widehat{X}_{i-1} + \alpha(z_i - \widehat{X}_{i-1})$$

$M \to \infty$

**Effective filtration of errors. But no reaction to changes in dynamics**

$\alpha \to 0$

**5**

# Sources of estimation errors



Measurements

Unknown location
of an unmoving object

| Source 1: Measurement errors | → | Errors of estimation are related with only measurements errors. Model of motion is accurate |

# Sources of estimation errors

**Divergence. Errors monotonously increase**

**Accurate measurements**

**We assume the object to move uniformly**

**But in fact it is uniformly accelerated motion**

**Source 2: Methodical errors** → **Errors of estimation are related with errors of methods. Model of motion is inaccurate.**

# Source 1: Measurement errors

## 1 Running mean

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$\sigma_{\widehat{X}}^2 = \frac{1}{M^2} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} \sigma_\eta^2$$

$$\sigma_{\widehat{X}}^2 = \frac{\sigma_\eta^2}{M}$$

## 2 Exponential mean

$$\widehat{X}_i = \alpha \sum_{k=0}^{i-1} (1-\alpha)^k z_{i-k} + (1-\alpha)^i z_0$$

$$\lim_{i \to \infty} \sigma_{\widehat{X}}^2 = \lim_{i \to \infty} \left( \alpha^2 \sigma_\eta^2 \sum_{k=0}^{i-1} (1-\alpha)^{2k} \right)$$

$$\sigma_{\widehat{X}}^2 = \sigma_\eta^2 \frac{\alpha}{2-\alpha}$$

# Source 1: Measurement errors

**1 Running mean**

**2 Exponential mean**

$$\sigma_{\widehat{X}}^2 = \frac{\sigma_\eta^2}{M}$$

$$\sigma_{\widehat{X}}^2 = \sigma_\eta^2 \frac{\alpha}{2 - \alpha}$$

$M = 1$ ← **No filtration of errors** → $\alpha = 1$

$M \to \infty$ ← **Effective filtration of errors. But no reaction to changes in dynamics** → $\alpha \to 0$

# Source 2: Methodical errors of running mean

**Running mean**

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$z_k = X_k + \eta_k$$

**Estimation error** $\widetilde{X_i}$

$$X_i - \widehat{X}_i = X_i - \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} X_k - \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i-\frac{M-1}{2}} \eta_k$$

# Source 2: Methodical errors of running mean

**Running mean**

$$\widehat{X}_i = \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} z_k$$

$$z_k = X_k + \eta_k$$

**Estimation error** $\widetilde{X_i}$

$$X_i - \widehat{X}_i = X_i - \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} X_k - \frac{1}{M} \sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} \eta_k$$

**Estimation error** $\widetilde{X_i}$

$$\widetilde{X_i} = \Delta_i^X + \Delta_i^\eta$$

Source of $\Delta_i^X$: **methodical errors**

Source of $\Delta_i^\eta$: **Measurement errors**

# Source 2: Methodical errors of running mean

$$\Delta_i^X = X_i - \frac{1}{M}\sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} X_k$$

$$9 = \sum_{k=i-4}^{i+4} 1$$

$$X_i = X_i \frac{1}{M} M = \frac{1}{M}\sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} X_i$$

$$\Delta_i^X = \frac{1}{M}\sum_{k=i-\frac{M-1}{2}}^{i+\frac{M-1}{2}} (X_i - X_k)$$

$$X_i - \widehat{X}_i = -\frac{1}{M}\sum_{k=i-\frac{M-1}{2}}^{i-1} (X_k - X_i) + \frac{1}{M}\sum_{k=i+1}^{i+\frac{M-1}{2}} (X_i - X_k)$$

$$k < i \qquad k > i$$

# Source 2: Methodical errors of running mean

**Left part**

**A process has constant rate**

**Right part**

$X_{i-4}$  $X_{i-3}$  $X_{i-2}$  $X_{i-1}$  $X_i$  $+c$  $X_{i+1}$  $+c$  $X_{i+2}$  $X_{i+4}$  $X_{i+3}$

$$-\frac{1}{M}\sum_{k=i-\frac{M-1}{2}}^{i-1}(X_k-X_i)$$

$$\frac{1}{M}\sum_{k=i+1}^{i+\frac{M-1}{2}}(X_i-X_k)$$

**Methodical error is equal to zero** → **Error of right part compensates that of left part**

# Source 2: Methodical errors of running mean

Left part

$X_i$

$+c$ $\quad$ $-c$

Right part

$X_{i-1}$ $\quad$ $X_{i+1}$

$X_{i-2}$ $\quad$ $X_{i+2}$

$X_{i-3}$ $\quad$ $X_{i+3}$

$X_{i-4}$ $\quad$ $X_{i+4}$

Process rate changes its sign to opposite one

$$-\frac{1}{M}\sum_{k=i-\frac{M-1}{2}}^{i-1}(X_k-X_i)$$

$$\frac{1}{M}\sum_{k=i+1}^{i+\frac{M-1}{2}}(X_i-X_k)$$

Methodical error is doubled $\longrightarrow$ Error of right part doubles that of left part
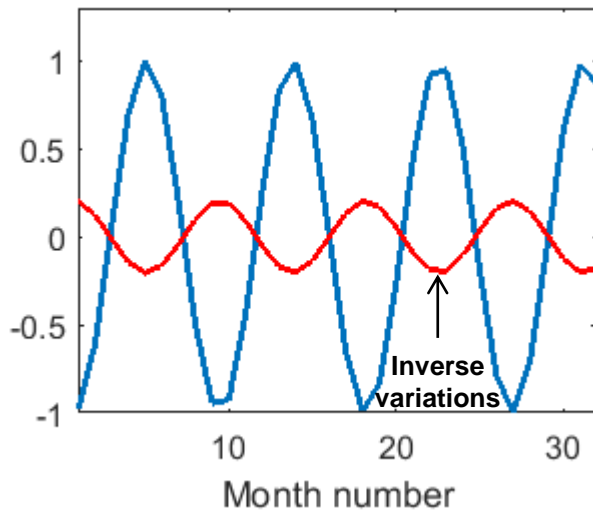
# Analysis of running mean errors

**Running mean may significantly distort the dynamics of the process**

**9-months variation**



Month number

**Inverse variations with periods from 6 to 12 months. Convex curve is replaced by concave curve and vice versa**

# Analysis of running mean errors

**Running mean may significantly distort the dynamics of the process**

**9-months variation**



**13-months variation**
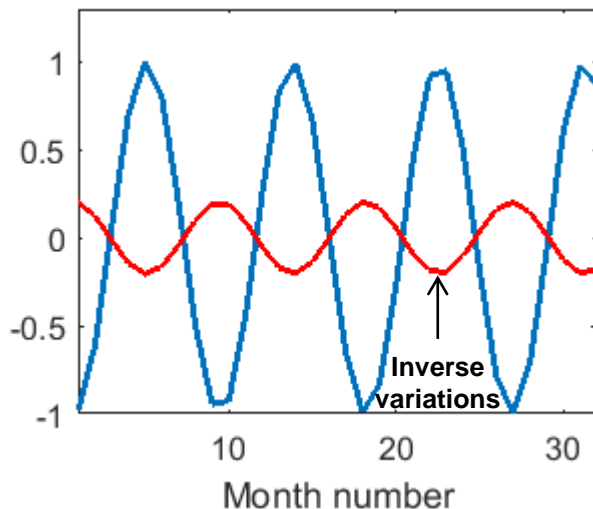


**Inverse variations with periods from 6 to 12 months. Convex curve is replaced by concave curve and vice versa**

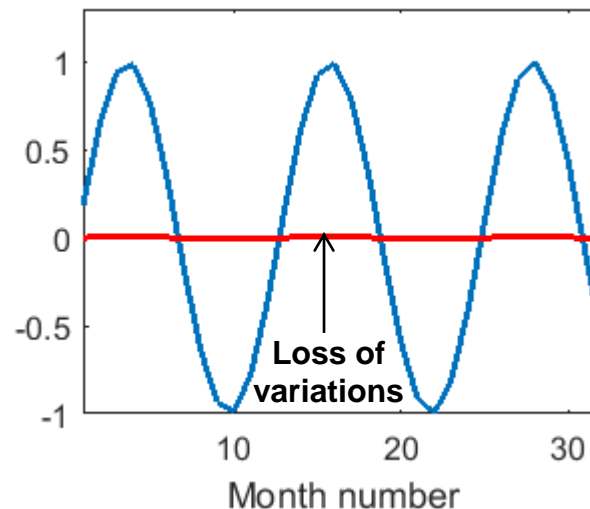**Total loss of 6- and 12-month variations decreasing them to zero**

# Distortion of physics in sunspot cycle 12



**Performed analysis allows us to anticipate the errors of smoothing and getting false conclusions** → **Alternatives in the following topics of course**

# Conclusions

**Don't apply methods in blind to not fall into the trap leading to false conclusions**

**Even if implementation is simple, the method itself requires careful analysis**

# Optimal choice of smoothing constant $\alpha$

Process $X$ is characterized by sudden and unpredictable changes



Measurements $z_i$

# Optimal choice of smoothing constant $\alpha$

Process $X$ is characterized by sudden and unpredictable changes

$$X_i = X_{i-1} + w_i$$

$w_i$ - unbiased random noise with variance $\sigma_w^2$

**Random walk model**

# Optimal choice of smoothing constant $\alpha$

**Optimal $\alpha$ for random walk model**

$$\alpha = \frac{-\chi + \sqrt{\chi^2 + 4\chi}}{2}$$

$$\chi = \frac{\sigma_w^2}{\sigma_\eta^2}$$

$\sigma_\eta^2$ - variance of measurement noise

**Muth J.F. (1960), Optimal properties of exponentially weighted forecasts of time series with permanent and transitory components, J.Amer. Statist. Ass.01960.-Vol.55.-p.299.**

**Full error of smoothing**

$$\sigma_\eta^2 a$$

$$a\sigma_\eta^2 > \sigma_\eta^2 \frac{\alpha}{2 - \alpha}$$

# Optimal choice of smoothing constant $\alpha$

**Optimal $\alpha$ for random walk model**

$$\alpha = \frac{-\chi + \sqrt{\chi^2 + 4\chi}}{2}$$

$$\chi = \frac{\sigma_w^2}{\sigma_\eta^2}$$

$\sigma_\eta^2$ - variance of measurement noise

Muth J.F. (1960), Optimal properties of exponentially weighted forecasts of time series with permanent and transitory components, J.Amer. Statist. Ass.01960.-Vol.55.-p.299.

**Full error of smoothing**

$$\sigma_\eta^2 a$$

$$a\sigma_\eta^2 > \sigma_\eta^2 \frac{\alpha}{2-\alpha}$$

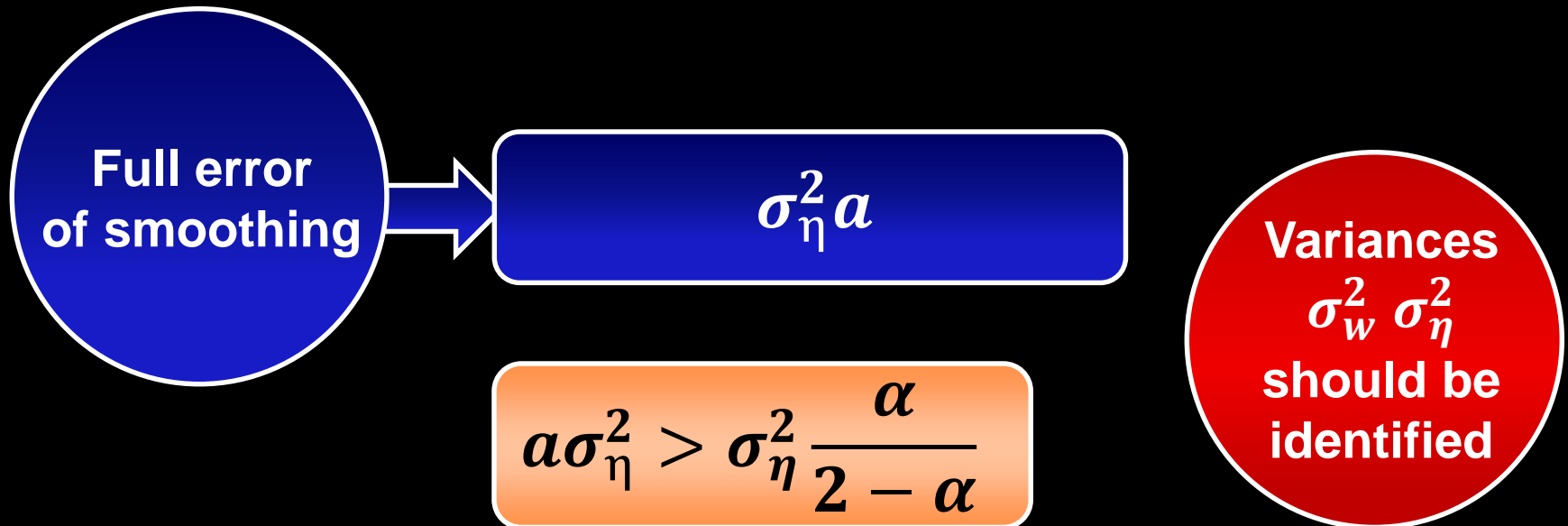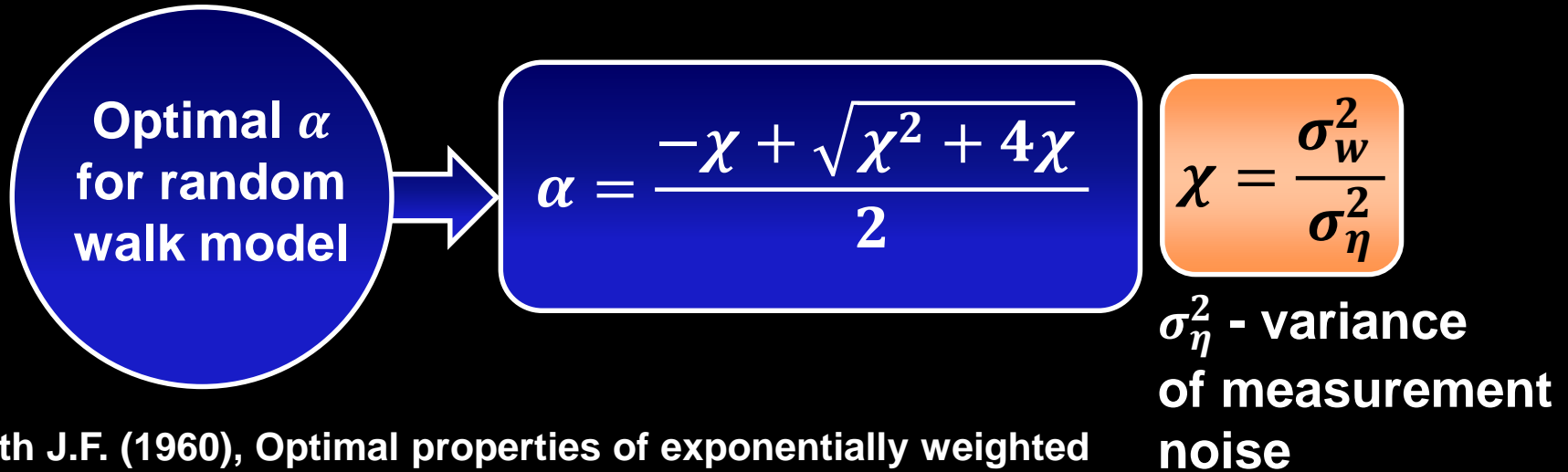**Variances $\sigma_w^2$ $\sigma_\eta^2$ should be identified**

# Identification of noise statistics $\sigma_w^2$ and $\sigma_\eta^2$

| | |
|---|---|
| **Process $X_i$** → | $X_i = X_{i-1} + w_i$ ① |
| **Measurements** → | $z_i = X_i + \eta_i$ ② |
| **Residual $\nu_i$** → | $\nu_i = z_i - z_{i-1}$ ③ |
| **Residual $\rho_i$** → | $\rho_i = z_i - z_{i-2}$ ④ |
| **Residual $\nu_i$** → | $\nu_i = w_i + \eta_i - \eta_{i-1}$ ⑤ |
| **Residual $\rho_i$** → | $\rho_i = w_i + w_{i-1} + \eta_i - \eta_{i-2}$ ⑥ |
| **Math. expectation** → | $E[\nu_i^2] = \sigma_w^2 + 2\sigma_\eta^2$ ⑦ |
| **Math. expectation** → | $E[\rho_i^2] = 2\sigma_w^2 + 2\sigma_\eta^2$ ⑧ |

**Rewrite Eq. 3 using Eq. 1 and 2**

Anderson, W. N., G. B. Kleindorfer, P. R. Kleindorfer, and M. B. Woodroofe (1969), Consistent estimates of the parameters of a linear system, Ann. Math. Stat., 40(3), 2064–2075.

# Identification of noise statistics $\sigma_w^2$ and $\sigma_\eta^2$

**Process $X_i$** → $X_i = X_{i-1} + w_i$ ①

**Measurements** → $z_i = X_i + \eta_i$ ②

**Residual $\nu_i$** → $\nu_i = z_i - z_{i-1}$ ③

**Residual $\rho_i$** → $\rho_i = z_i - z_{i-2}$ ④

**Residual $\nu_i$** → $\nu_i = w_i + \eta_i - \eta_{i-1}$ ⑤

**Residual $\rho_i$** → $\rho_i = w_i + w_{i-1} + \eta_i - \eta_{i-2}$ ⑥

**Math. expectation** → $E[\nu_i^2] = \sigma_w^2 + 2\sigma_\eta^2$ ⑦

**Math. expectation** → $E[\rho_i^2] = 2\sigma_w^2 + 2\sigma_\eta^2$ ⑧

**Rewrite Eq. 3 using Eq. 1 and 2**

$$E[\nu_i^2] \approx \frac{1}{N-1}\sum_{k=2}^{N} \nu_k^2$$

$$E[\rho_i^2] \approx \frac{1}{N-2}\sum_{k=3}^{N} \rho_k^2$$

**Consistent estimates $\sigma_w^2$ and $\sigma_\eta^2$ are obtained by solving system of equations (7,8)**