



---

2. MARTS 2017

# STATISTICS

---

SØREN E. LARSEN  
SENIOR RESEARCHER

U N E R S I T E T



# NORMALIZATION: WHY DO THIS.

---

To reduce variation in yearly loads from climate: runoff.

A large part of the variation due to climate (hydrology) is removed.

To facilitate trend analysis and target testing with reduced year-to-year variation.

Trend tests with increased power.



# FORMULAR FOR NORMALIZATION:

---

Based on the log-log relationship between yearly load and yearly runoff:

$$L_{iN} = \exp\left(\log L_i \cdot \frac{\hat{\alpha} + \hat{\beta} \cdot \log \bar{q}}{\hat{\alpha} + \hat{\beta} \cdot \log q_i}\right) \cdot \exp(0.5 \cdot \text{MSE}). \quad (2.6)$$



# FORMULAR FOR MSE:

---

$$\text{MSE} = \frac{1}{n-2} \sum_{i=1}^n (x_i - \hat{x}_i)^2,$$



# TREND ANALYSIS:

---

Statistical methods used in a total of 270 time series of normalized . Most of the analyses done automatically. Actually 5\*270 time series for trend analysis. Testing of ceiling only for 270 time series.

Mann-Kendall non-parametric trend method for testing for a significant monotone trend in the normalized time series.

This is a fairly robust method although autocorrelation can deflate the power of the test. Calculate p-value based on ranking data. Delivers test-statistic, p-value, Sen's slope estimator, estimate of intercept, and a confidence interval for the slope.

This non-parametric method can be used on both "raw" nutrient time series, normalized time series and runoff (climate) time series.

Software used: SAS, SPLUS (EnvStats), (R-Library: EnvStats).



# EXAMPLE

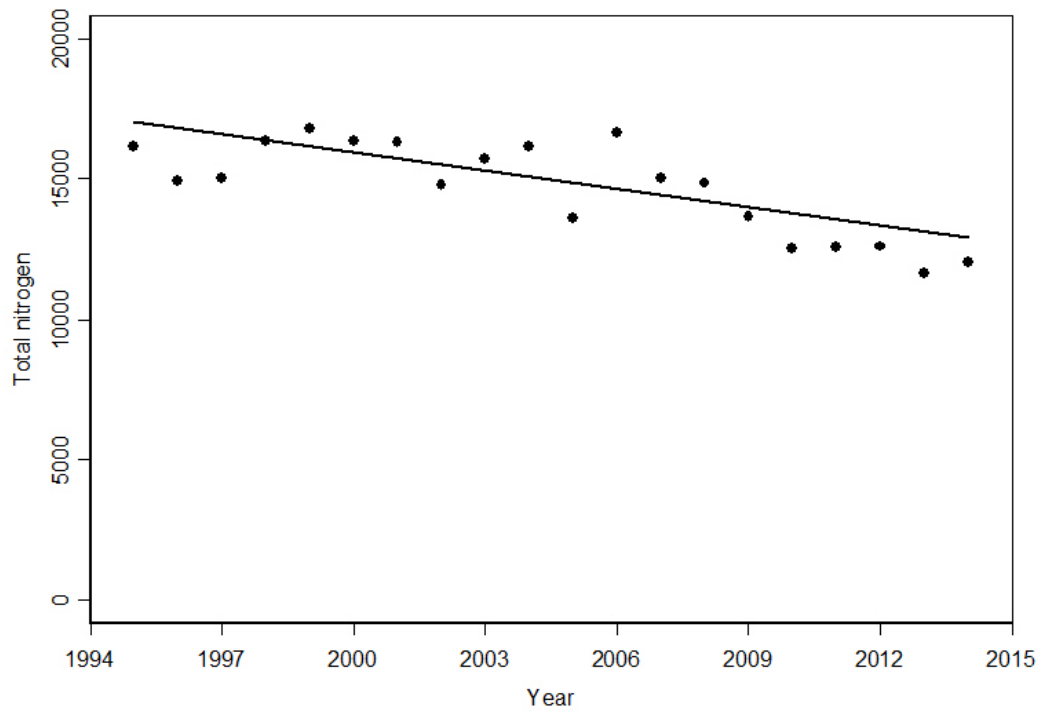
---

```
> Results of Hypothesis Test
> -----
> Null Hypothesis:          tau = 0
> Alternative Hypothesis:  True tau is not equal to 0
> Test Name:               Kendall's Test for Trend
>                           (with continuity correction)
> Estimated Parameter(s):  tau    = -0.5578947
>                           slope  = -217.594
>                           intercept = 1.726966e4
> Estimation Method:      slope:  Theil/Sen Estimator
>                           intercept: Conover's Estimator
> Data:                   y = normTN
> Parent of Data:         --
> Sample Size:            20
> Test Statistic:         z = -3.40665
> P-value:                 0.0006576544
> Confidence Interval for: slope
> Confidence Interval Method: Gilbert's Modification
>                           of Theil/Sen Method
> Confidence Interval Type: two-sided
> Confidence Level:        95%
> Confidence Interval:     LCL= -326.1601
>                           UCL= -146.6605
```



# FIGURE OF TREND EXAMPLE

---





## CHANGE IN %:

---

Estimating the change in normalized nutrient load can be done by the non-parametric Sen's slope estimator.

The method assumes a constant change, i.e. a linear trend. So use the Sen's slope estimator if the trend can be assumed to be fairly linear.

If the trend is not linear, fit a non-linear model or use start-end difference, or the loess model.





# FORMULAE:

---

$$100 \cdot \frac{(n - 1) \cdot \hat{\beta}}{\hat{\alpha}}$$

$$100 \cdot (\text{end} - \text{start}) / \text{start}$$

For some times series, the start value, the end value or both can deviate too much from the general trend; if so, an approach using the average value of, for instance, the first 3 years and the last 3 years would reduce the influence of single years.

Or use model estimates for 1995 and 2014.



## EXAMPLE:

---

› Change (1) = -21,77%

› Change (2) = -19,24%



# T-TESTS:

---

- › Using a standard Students T-test for testing the difference in load between two periods.
- › Either 1997-2003 and 2012-2014
- › Or 1997-2003 and 2010-2014.
- › Assume a Gaussian distribution. Difficult to test formally because of size of samples.
- › Variance inhomogeneity is accounted for by using the Welch-Satterthwaite equation for adjusting degrees of freedom.



## EXAMPLE:

---

- › T-test 1997-2003; 2012-2014
  - › Test value = 7.35
  - › P-value =  $<0.0001$
  - › Difference = 126574
- 
- › T-test 1997-2003; 2010-2014
  - › Test value = 6.34
  - › P-value =  $<0.0001$
  - › Difference = 109269



# TARGET TESTING

---

For normalized time series with a non-significant trend the equation in formula 6.1 (next slide) can be used to calculate the adjusted mean nutrient load and evaluate this value against the target value.

For normalized time series with a significant linear trend the equation in 6.4 (next slide again) should be used. Use this one also for time series with step trends or non-linear models/trends.



# TARGET TESTING: FORMULAE

---

> Calculate the statistic

$$\text{> } \bar{x}_{AD} = \bar{x} + 1.645 \cdot SE, \quad (6.1)$$

where  $\bar{x}$  is the mean of all values in the time series and SE is the standard error (SE = standard deviation divided by square root of N = number of observations in the time series) , and finally 1.645 is the 95% percentile in a Gaussian distribution with mean 0 and variance equal 1.



# SE FOR MEANS:

---

> Mean calculated on the basis of n years:

>

$$> \hat{x}_{2014}^n = \frac{1}{n} \sum_{i=1}^n x_i$$

$$> SE(\hat{x}_{2014}^n) = \sqrt{\frac{1}{(n \cdot (n-1))} \sum_{i=1}^n (x_i - \hat{x}_{2014}^n)^2}$$



# TARGET TESTING: FORMULAE

---

Calculate the statistic

$$\bar{x}_{AD} = \widehat{L}_{nN} + t_{n-2,0.05} \cdot \text{SEr}, \quad (6.4)$$

where  $t_{n-2,0.05}$  is the 95% percentile in a  $t$ -distribution with  $n-2$  degrees of freedom.

$$\text{SEr} = \sqrt{MSE} \cdot \sqrt{1/n + \text{year}_n^2 / \sum_{i=1}^n \text{year}^2} \quad (6.3)$$





# EXAMPLE 1:

---

Using whole time series:

- › Estimated value = 31583
- › SE = 685
- › Test value = 32770

Using last 5 years:

- › Estimated value = 31737
- › SE = 818
- › Test value = 33083

Using last 3 years:

- › Estimated value = 30949
- › SE = 1191
- › Test value = 32908



# STEP TREND:

---

If time series show two distinct trends (trend reversal) use two or more linear regressions to model the time series. The change-point can either be determined by visual inspection of the time series plot or by knowledge of changes in the catchment or by a statistical method.

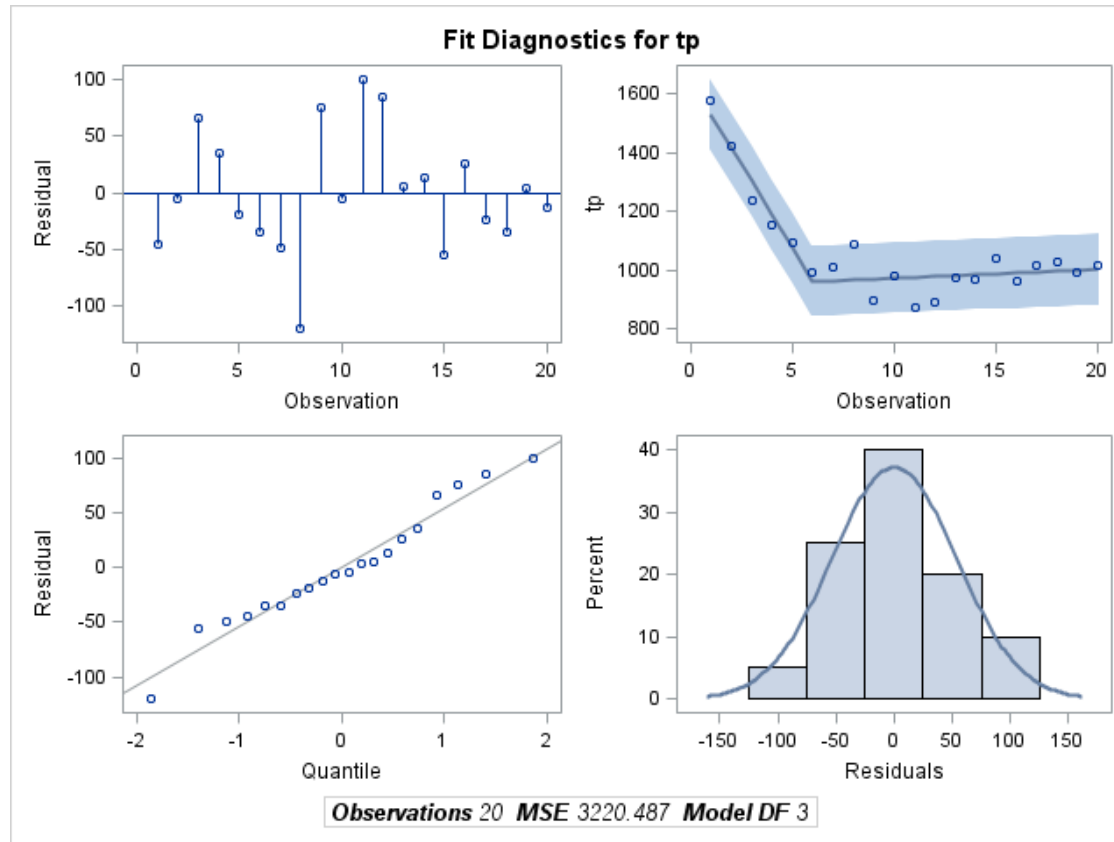
Model examples:

$$L_{Ni} = \begin{cases} a + b \cdot Y, & \text{if } Y < T \\ a + b \cdot Y + d \cdot (Y - T), & \text{if } Y \geq T \end{cases}$$

$$L_{Ni} = \begin{cases} a_1 + b_1 \cdot Y, & \text{if } Y < T \\ a_2 + b_2 \cdot Y, & \text{if } Y \geq T \end{cases}$$



# EXAMPLE OF A STEP TREND.





# EXAMPLE 2:

---

Using whole time series:

- › Estimated value = 885
- › SE = 31
- › Test value = 939

Using last 5 years:

- › Estimated value = 1004
- › SE = 11
- › Test value = 1023

Using last 3 years:

- › Estimated value = 1013
- › SE = 10
- › Test value = 1030

Step trend:

- › Estimated value = 1002
- › SE = 21
- › Test value = 1038



# TRAFFIC LIGHT EVALUATION

---

## Red:

- > If  $\bar{x} > T$ , i.e. the average nutrient load over the considered period is above target value.

## Yellow:

- > If  $\bar{x} < T$  and if  $\bar{x}_{AD} > T$ , i.e. the null hypothesis of target test is accepted but the average load is less than the target.

## Green:

- > If  $\bar{x}_{AD} < T$ , i.e. the null hypothesis of the target test is rejected.



# UNCERTAINTY: DANISH CALCULATIONS (TN)

---

- > Variance in time series:  $\sigma^2 = \sigma_Y^2 + \sigma_{Between Y}^2$
- > It's the total variance that is used for statistical analyses.
  
- > Calculation of an estimate of the uncertainty for the yearly load.
- > Measured area: 169 stations (55% of total danish area).

- > Uncertainty (%) = 
$$\frac{100}{\sum_{i=1}^{169} X_i} \sqrt{\sum_{i=1}^{169} (bias_i \cdot X_i)^2 + (precision_i \cdot X_i)^2}.$$

- > Precision is defined as the relative stand deviation.

- |                            |                      |                 |
|----------------------------|----------------------|-----------------|
| > 0-50 km <sup>2</sup> :   | Bias: -1% til -3%;   | Precision: 1-3% |
| > 50-200 km <sup>2</sup> : | Bias: -0.7% til -3%; | Precision: 1-3% |
| > >200 km <sup>2</sup> :   | Bias: -1% til -4%;   | Precision: 2-5% |



# RESULT:

---

Bias: -1 to -4%

Precision: 0.7 to 1.2 %

Uncertainty: 0.7 to 1.3%.

For one station:

Bias: -1 to -3%

Precision: 3 to 5 %

Uncertainty: 3.2 to 5.8%.



# UNMEASURED AREA:

---

- › Modelling concentrations and runoff in 1286 small catchments.
- › L= F(Ndiffus, Retention lakes, Retention stream, Retention groundwater, Point sources)

	Bias (%)	Præcision (%)
Model	-15 to 25	12 to 15
Retention lake	5	40
Retention stream	-10	40
Retention total	-5	40
Point sources: industry	-1 to -3	5 to 10
Point sources: wastewater	-1 to -3	5 to 10
Point sources: fish farms	-1 to -3	15 to 20
Point sources: rain water	-5	40





# RESULT:

---

Bias: 20 to 28%

Precision: 0.8 to 2.0 %

Uncertainty: 1.2 to 2.2%

For one small catchment:

Bias: 27%

Precision: 15 to 20%

Uncertainty: 31 to 34%.



# RESULT: TOTAL DANISH AREA

---

- › Bias: 7.4 to 12.8 %
- › Precision: 0.5 to 1.1 %
- › Uncertainty: 7.4 to 12.8 %



# TOTAL UNCERTAINTY FORMULAS

---

In DUET-H/WQ the uncertainty for individual measurements is estimated by the formula

$$EP = \sqrt{E_Q^2 + E_C^2 + E_{PS}^2 + E_A^2 + E_{DPM}^2}, \quad (4.1)$$

where according to Harmel et al. (2009)

$E_Q^2$ =Uncertainty for the discharge measurement ( $\pm\%$ )

$E_C^2$ =Uncertainty for sample collection ( $\pm\%$ )

$E_{PS}^2$ =Uncertainty of sample preservation/storage ( $\pm\%$ )

$E_A^2$ =Uncertainty from laboratory analysis ( $\pm\%$ )

$E_{DPM}^2$ =Uncertainty from data processing and data management ( $\pm\%$ ).



# TOTAL UNCERTAINTY FORMULAS

---

Then the total uncertainty for aggregated data can be estimated by

$$EP_{total} = \frac{100}{\sum_{i=1}^n x_i} \sqrt{\sum_{i=1}^n \left( x_i \cdot \frac{EP_i}{100} \right)^2} \quad (4.2)$$

and  $EP_{total}$  is given as  $\pm\%$ .



# TESTING REDUCTION TARGETS, EXAMPLE

---

Actual situation	Inference from data	Diagnosis	Truth	Probability
$\mu \geq T$	Less than T	"Positive"	False	$\alpha$ (type I error)
$\mu < T$	Higher than T	"Negative"	False	$\beta$ (type II error)
$\mu \geq T$	Higher than T	"Negative"	True	$1 - \alpha$
$\mu < T$	Less than T	"Positive"	True	$1 - \beta$