# Light years away from a thermodynamic model?

# - Model Discrimination and Validation -

**Dr. Alexander Kud**

**Seminar**
**Thermodynamic data: production, exploitation and impact on process design**
**at „Société Française de Génie des Procédés" and „IFP Energies nouvelles"**
**in Rueil-Malmaison**
**1st April 2016**

BASF Research Department - GME

# Overview

# Sensitivity Analysis

Unscaled dimensionless

$$S_i := \frac{\Theta_i}{f(X,\Theta)}\left(\frac{\partial f(X,\Theta)}{\partial \Theta_i}\right)_{\Theta_{j\neq i}, X}$$

Standardized and unscaled dimensionless

$$\tilde{S}_i := \frac{\Theta_i}{s(f(X,\Theta))}\left(\frac{\partial f(X,\Theta)}{\partial \Theta_i}\right)_{\Theta_{j\neq i}, X}$$

$$X = \left(x_1, x_2 ... x_L\right)^* \qquad x_i \in [x_{i,min}, x_{i,max}] \qquad s = \text{stddev}$$

$$\Theta = \left(\Theta_1, \Theta_2 ... \Theta_{Np}\right)^* \qquad \Theta \in \mathfrak{R}^{Np} \qquad \left(\ ^* \ \text{means the transpose of a vector or a matrix}\right)$$

3

# Sensitivity Analysis

Jacobian matrix  (design matrix)

$$J(X,\Theta) = - \begin{vmatrix} \dfrac{1}{s_1}\left(\dfrac{\partial f(X,\Theta)}{\partial \Theta_1}\right)_{\Theta_{j\neq 1},\,X_1} & \cdots & \dfrac{1}{s_1}\left(\dfrac{\partial f(X,\Theta)}{\partial \Theta_{Np}}\right)_{\Theta_{j\neq Np},\,X_1} \\ \vdots & & \vdots \\ \dfrac{1}{s_M}\left(\dfrac{\partial f(X,\Theta)}{\partial \Theta_1}\right)_{\Theta_{j\neq 1},\,X_M} & \cdots & \dfrac{1}{s_M}\left(\dfrac{\partial f(X,\Theta)}{\partial \Theta_{Np}}\right)_{\Theta_{j\neq Np},\,X_M} \end{vmatrix}$$

$M$ = number of measurements

# Sensitivity Analysis

### Parameter error propagation for a what-if scenario

$$cov(X,\Theta) = \left(J(X,\Theta)^* \cdot J(X,\Theta)\right)^{-1}$$

parameter variance-covariance matrix

$$e_{\%}(\Theta_i) = \frac{\sqrt{cov(X,\Theta)_{ii}}}{\Theta_i} \, 100\,\%$$

expected error of parameter $\Theta_i$

$\left(\,^* \text{ means the transpose of a vector or a matrix}\right)$

# 2. Goodness of Fit

**How to test the power of the most frequently used statistical criteria?**

- **Simulation of measurements with a true and a false model.**

- **Evaluation or analysis of the simulated measurements.**

- **Answer to the question: Which assesment numbers are useful for model discrimination and predictive power ?**

# Choice of Model

Simulation with a heuristic vapor pressure equation. Compound: water (from triple point up to critical point (1))

$$p_{true} := p_c \exp\left\{ \mathring{\Theta}_1 \left( 1 - \frac{T_c}{T} \right) + \mathring{\Theta}_2 \ln\left( \frac{T}{T_c} \right) + \mathring{\Theta}_3 \left[ \left( \frac{T}{T_c} \right)^2 - 1 \right] \right\} = f(T, \mathring{\Theta})$$

$$p_{false} := p_c \exp\left\{ \Theta'_1 \left( 1 - \frac{T_c}{T} \right) + \Theta'_2 \ln\left( \frac{T}{T_c} \right) + \Theta'_3 \left[ \left( \frac{T}{T_c} \right) - 1 \right] \right\} = f(T, \Theta')$$

(1) VDI-Wärmeatlas, 11[th] edition 2013, Springer Vieweg

# Choice of Noise

Realistic simulation assumptions: inhomogeneous variance

$$\varepsilon_i \ \sim \ \mathcal{N}\left(0, \sigma_i^2 \,|\, \sigma_i\right) = v \cdot f\left(T_i, \Theta\right) \qquad v = 0.002 \ \left(rel.\,error\right)$$

$$Y_i \ = \ f\left(T_i, \mathring{\Theta}\right) \ + \ \varepsilon_i \qquad\qquad Y_i \ = \ \text{simulated vapor pressure measurements}$$

$$\mathcal{E}\left(Y\right) \ = \ f\left(T, \mathring{\Theta}\right)$$
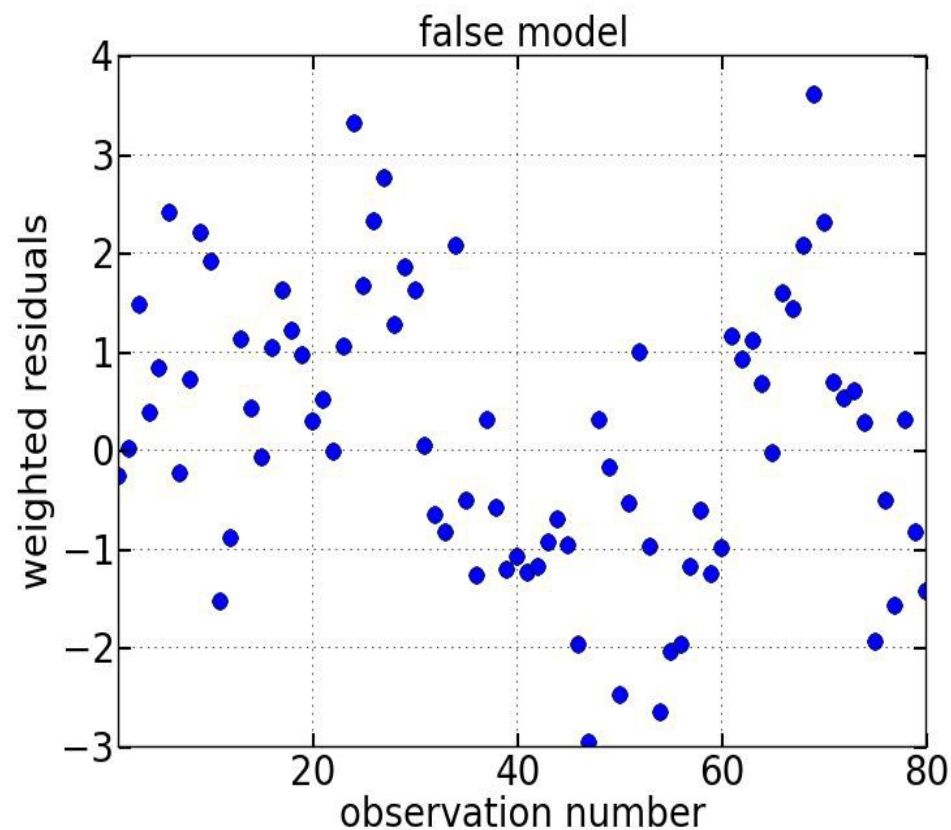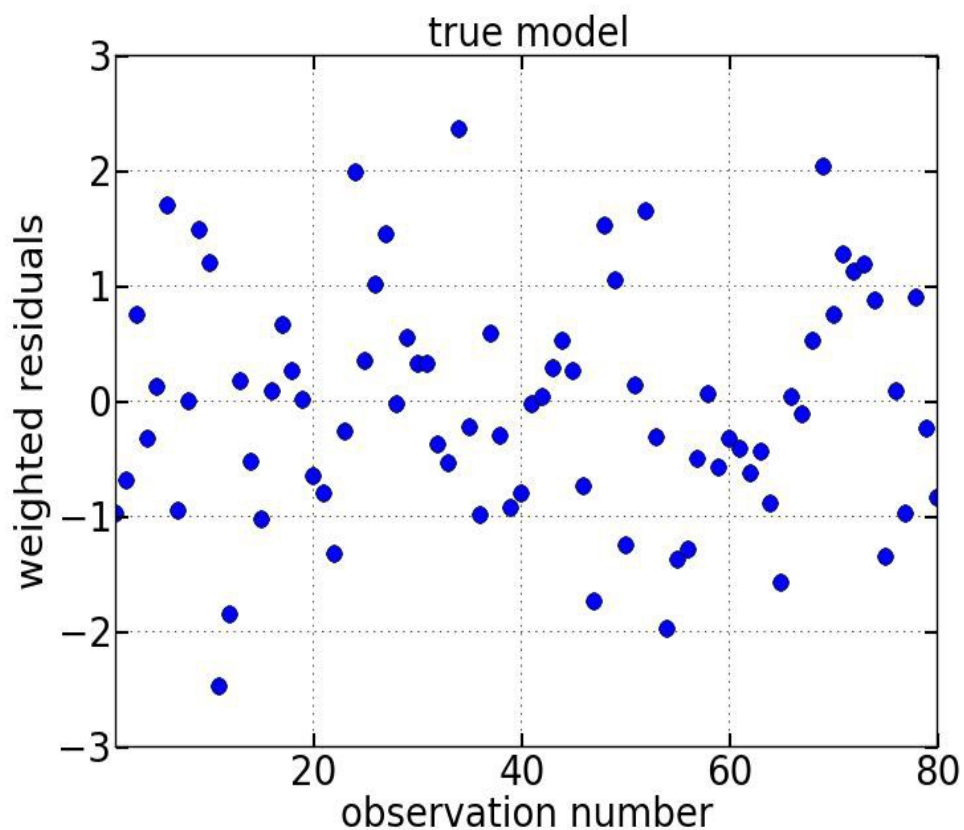
Parameter estimation: Levenberg-Marquardt (Least-Squares)

# Choice of Noise

Plot of Weigted Residuals vs. Measurement Number

for the True and the False Model

# 2. Goodness of Fit Quantities

| # | Test quantity | Equation |
|---|---|---|
| 1 | $SWS$ | $SWS := \sum_i \left( \dfrac{Y_i - \hat{Y}_i}{s_i} \right)^2$ |
| 2 | $\chi_r^2$ | $\chi_r^2 := \dfrac{SWS}{M - Np}$ |
| 3 | $P_\alpha^{(\chi)}$ | $P_\alpha^{(\chi)} = 1 - P(\chi_t^2 > \chi_c^2) = \alpha$ |
| 4 | $AAD\%$ | $AAD\% = 100 \sum_i \left| \dfrac{Y_i - \hat{Y}_i}{Y_i} \right|$ |

# 2. Goodness of Fit Quantities

| # | Test quantity | Equation |
|---|---------------|----------|
| 5 | Model bias 1 | $\text{Bias 1} = 100 \sum_i \dfrac{Y_i - \hat{Y}_i}{Y_i}$ |
| 6 | Model bias 2 | $P_\alpha\left(H_0: (Y = a\hat{Y} + b) \wedge (a=1) \wedge (b=0)\right)$ |
| 7 | *PRE SWS* | $PRE\ SWS := \sum_i \left(\dfrac{Y_i - \hat{Y}_i(\hat{\Theta}_k)}{s_i}\right)^2 \quad (1)$ |
| 8 | $P_\alpha^{(\chi)}(PRE\ SWS)$ | $P_\alpha^{(\chi)} = 1 - P(\chi_t^2 > \chi_c^2) = \alpha$ |

$(1)\quad k = \text{class number}$

# 2. Goodness of Fit Quantities

| # | Test quantity | Equation / Explanation |
|---|---|---|
| 9 | Model Bias 2 for $\hat{Y}_i(\hat{\Theta}_k)$ | $P_\alpha\left(H_0: (Y = a\hat{Y} + b) \wedge (a = 1) \wedge (b = 0)\right)$ |

Remark: The $R^2$ test for non-linear regression is meaningless.

See D.A. Ratkowsky, Handbook of Nonlinear Regression Models. Marcel Dekker, Inc. New York and Basel, 1990. p. 44

# 2. Goodness of Fit Quantities

Bias Test 2

Hypothesis: $$Y = \hat{Y}$$

# 2. Goodness of Fit Quantities

Bias Test 2

Hypothesis:

$$Y \; = \; \hat{Y}$$

$$Y \; = \; a \cdot \hat{Y} \qquad\qquad if \; a \; = \; 1$$

# 2. Goodness of Fit Quantities

Bias Test 2

Hypothesis:
$$Y = \hat{Y}$$

$$Y = a \cdot \hat{Y} \qquad if\ a = 1$$

$$Y = a \cdot \hat{Y} + b \qquad if\ b = 0$$

# 2. Goodness of Fit Quantities

Bias Test 2

Hypothesis:

$$Y = \hat{Y}$$

$$Y = a \cdot \hat{Y} \qquad if \ a = 1$$

$$Y = a \cdot \hat{Y} + b \qquad if \ b = 0$$

$H_0$ Hypothesis

$$P_\alpha\left(H_0 : \left(Y = a\,\hat{Y} + b\right) \wedge (a = 1) \wedge (b = 0)\right)$$

16

# 2. Goodness of Fit – Predictive Power 1

e.g. 3-fold cross validation test

$N^0$ : number of run (randomized)

$\hat{Y}^1$ : predicted response in class 1

$\hat{Y}^0$ : calculated for all data points

$\hat{R}^1_i$ : predicted residual i in class 1

| $N^0$ | $N^1$ | $N^2$ | $N^3$ | | $\hat{Y}^0$ | $\hat{Y}^1$ | $\hat{Y}^2$ | $\hat{Y}^3$ | | remarks |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | | $\hat{R}^0_1$ | $\hat{R}^1_1$ | | | | for prediction |
| 2 | 2 | 2 | 2 | | $\hat{R}^0_2$ | $\hat{R}^1_2$ | | | | for estimation |
| 3 | 3 | 3 | 3 | | $\hat{R}^0_3$ | $\hat{R}^1_3$ | | | | |
| 4 | 4 | 4 | 4 | | $\hat{R}^0_4$ | | $\hat{R}^2_4$ | | | |
| 5 | 5 | 5 | 5 | | $\hat{R}^0_5$ | | $\hat{R}^2_5$ | | | |
| 6 | 6 | 6 | 6 | | ... | | $\hat{R}^2_6$ | | | |
| 7 | 7 | 7 | 7 | | ... | | | $\hat{R}^3_7$ | | |
| 8 | 8 | 8 | 8 | | ... | | | $\hat{R}^3_8$ | | |
| 9 | 9 | 9 | 9 | | $\hat{R}^0_9$ | | | $\hat{R}^3_9$ | | |

17

# 2. Goodness of Fit - Results

| # | Test quantity | True model | False model | Model dis-crimination |
|---|---|---|---|---|
| 1 | $SWS$ (1) | 79 | 162 | O |
| 2 | $\chi^2_r$ | 1.03 | 2.1 | O |
| 3 | $P^{(\chi)}_\alpha$ % | 42 | << $10^{-3}$ | + |
| 4 | $AAD$ % (2) | 0.15 | 0.22 | - |

**+** appriopriate    **O** appropriate if $f$ and database are equal for both models

**-** not appropriate

(1) $\chi^2_c = \chi^2_{1-\frac{\alpha}{2}, f} = 103$      (2) remember: assumed rel. error for measurement 0.2 %

# 2. Goodness of Fit - Results

| # | Test quantity | True model | False model | Model dis-crimination |
|---|---|---|---|---|
| 5 | Model bias 1 | -0.02 | 0.014 | **-** (1) |
| 6 | Model bias 2 | | | |
| | $P_{\alpha/2}(b=0)$ % | 29 | << 10$^{-3}$ | **+** |
| | $P_{\alpha/2}(a=1)$ % | 22 | << 10$^{-3}$ | **+** |

**+** appriopriate    **-** not appropriate

(1)  sensitive for residual structure !  →  indicator for bias  →  residual plot

# 2. Goodness of Fit - Results

| # | Test quantity | True model | False model | Model dis-crimination |
|---|---|---|---|---|
| 7 | $PRE\ SWS$ | 85.8 (1) | 719 (2) | <span style="color:orange">O</span> |
| 8 | $P_\alpha^{(\chi)}(PRE\ SWS)$ | 23 | << 10$^{-3}$ | + |
| 9 | Model Bias 2 for $\hat{Y}_i(\hat{\Theta}_k)$ | | | |
|  | $P_{\alpha/2}(b{=}0)$  % | 29 | << 10$^{-3}$ | + |
|  | $P_{\alpha/2}(a{=}1)$  % | 22 | << 10$^{-3}$ | + |

(1)  $\chi^2_{crit} = \chi^2_{1-\frac{\alpha}{2},f} = 103$     (2)  remember: PRE SWS  for all data points and false model: 162

# 3. Goodness of Parameter

Questions:

- How trustworthy are estimated parameter values?

- What kind of powerful parameter test methods exist?

- Which criteria must be fulfilled for model validation?

- How can model validation be defined?

# 3. Goodness of Parameter

Practical example:

-   Modeling Peng Robinson Equation of State (PR)

    with geometric mean mixing rule for the gas phase.

-   System  $CH_4$ (1) / $H_2O$ (2) with 168 measured data for gas phase.

    L.L. Joffrion and P.T. Eubank; FPE **43** (1988) 263.

-   Independent variables:  $T[K], \rho^*[mol/m^3], \ y_2 : 0.1, \ 0.25, \ 0.5$

    dependent variable:    $p[Pa]$    measured pressure

    overall error  ~ 0.2 %

# 3. Goodness of Parameter

Which parameter are selected in the PR EoS?

$$B(T, y, \Theta) = \sum_i \sum_j y_i y_j b_{ij}(1 - \kappa_{ij}\Theta_1) \longleftarrow$$

$$A(T, y, \Theta) = \sum_i \sum_j y_i y_j \sqrt{a_{ii} a_{jj}}(1 - \kappa_{ij}\Theta_2)$$

$$\kappa_{ij} = \kappa_{ji} := \begin{cases} 1 & if \quad i \neq j \\ 0 & if \quad i = j \end{cases}$$

$\rightarrow$ Analysis of the estimated parameter $\hat{\Theta}_1, \hat{\Theta}_2$

# 3. Goodness of Parameter

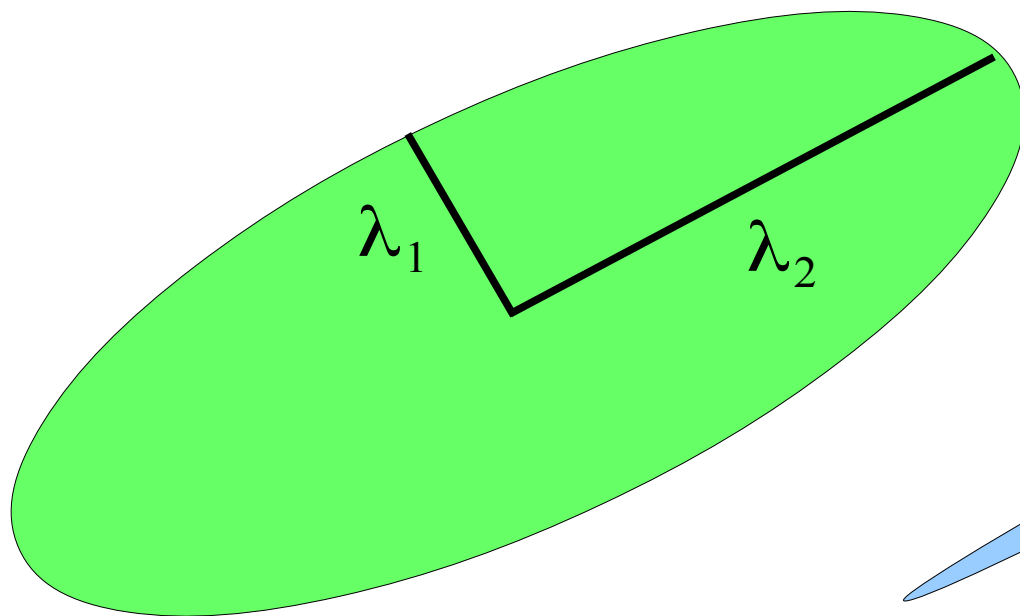| # | Test quantity | Equation / Explanation |
|---|---|---|
| 10 | $\mathcal{R}_{cov}$ | $\mathcal{R}_{cov} = \mathcal{R}ank\left(cov\left(X,\Theta\right)\right)$ |
| 11 | $cond_\lambda$ | $cond_\lambda = \dfrac{\lambda_{max}\left(cov\left(X,\Theta\right)\right)}{\lambda_{min}\left(cov\left(X,\Theta\right)\right)}$ |
| 12 | $e_{\%}\left(\Theta_i\right)$ | $e_{\%}\left(\Theta_i\right) = \dfrac{\sqrt{cov\left(X,\Theta\right)_{ii}}}{\Theta_i}100\,\%$ (1) |
| 13 | $P_\alpha^{(\chi)}\left(var\left(\Theta\right)\right)$ | $P_\alpha^{(\chi)} = 1 - P\left(var\left(\Theta\right) > \chi_c^2\right) = \alpha$ |

(1) better: exact confidence region based on F statistic (appendix)
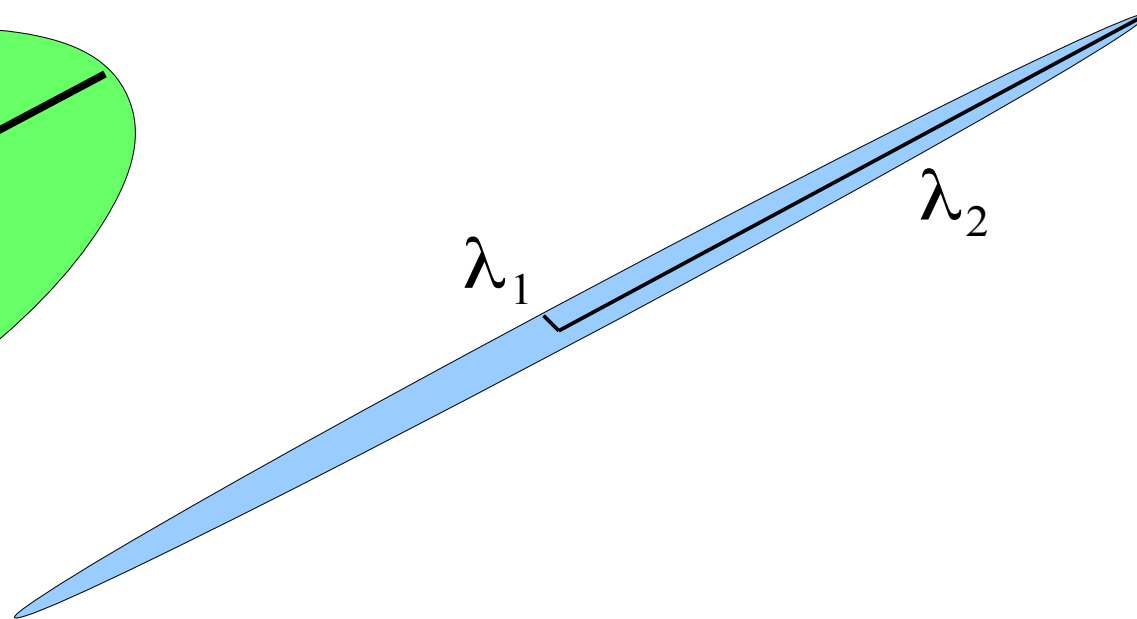
# Explanation of the condition number

$$cond_\lambda = \frac{\lambda_{max}(cov(X,\Theta))}{\lambda_{min}(cov(X,\Theta))}$$

Well conditioned
$\lambda_1$
$\lambda_2$

poorly conditioned
$\lambda_1$
$\lambda_2$

Generally: the eigenvalues $\lambda_i$ are variances of the orthogonal main axis of a hyper ellipsoid

25

# 3. Goodness of Parameter

| # | Test quantity | PR | Model Acceptance |
|---|---|---|---|
| 10 | $\mathcal{R}_{cov}$ | 2 | **+** (1) |
| 11 | $cond_{\lambda}$ | 600 | (2) |
| 12 | $e_{\%}(\Theta_i)$ | 12.2 6.4 | |
| 13 | $P_{\alpha}^{(\chi)}(var(\Theta))/\%$ | 0.2 19.5 | **-** **+** |

The model is not accepted

(1) necessary but not sufficient.    (2)  sensitiv indicator for parameter insufficiency

# 4. Model Validation – Definition

Selection criterion: quantities with statistical constraints

| # | Test quantity | Good-ness of Fit | Model Discrim | Pred. Power 1 | Good-ness of Param. | Model Validation 1 |
|---|---|---|---|---|---|---|
| 3 | $P_\alpha^{(\chi)}$ % | + | + | | | + |
| 6 | Model bias 2 | + | + | | | + |
| 8 | $P_\alpha^{(\chi)}(PRE\ SWS)$ | | + | + | | + |
| 9 | Model Bias 2 for $\hat{Y}_i(\hat{\Theta}_k)$ | | + | + | | + |
| 10 | $\mathcal{R}_{cov}$ | (+) | | | + | + |
| 13 | $P_\alpha^{(\chi)}(var(\Theta))$ | | | | + | + |

# 5. Predictiv Power – Definition

Predictive Power 1

$$Y = f(X, \hat{\Theta})$$

$$X = \left( x_1, x_2 \dots x_L \right)^*$$
$$Y = \left( y_1, y_2 \dots y_Q \right)^*$$
$$\hat{\Theta} = \left( \hat{\Theta}_1, \hat{\Theta}_2 \dots \hat{\Theta}_{Np} \right)^*$$

$X$ independent variables (observation)

$Y$ dependent var.

Predictive Power 2
($\rightarrow$ EoS development)

$$Z = g(X, \hat{\Theta})$$

$$X = \left( x_1, x_2 \dots x_L \right)^*$$
$$Z = \left( z_1, z_2 \dots z_I \right)^*$$

The dependent variable Z is not used for parameterizing $\hat{\Theta}$

# 5. Predictive Power 2

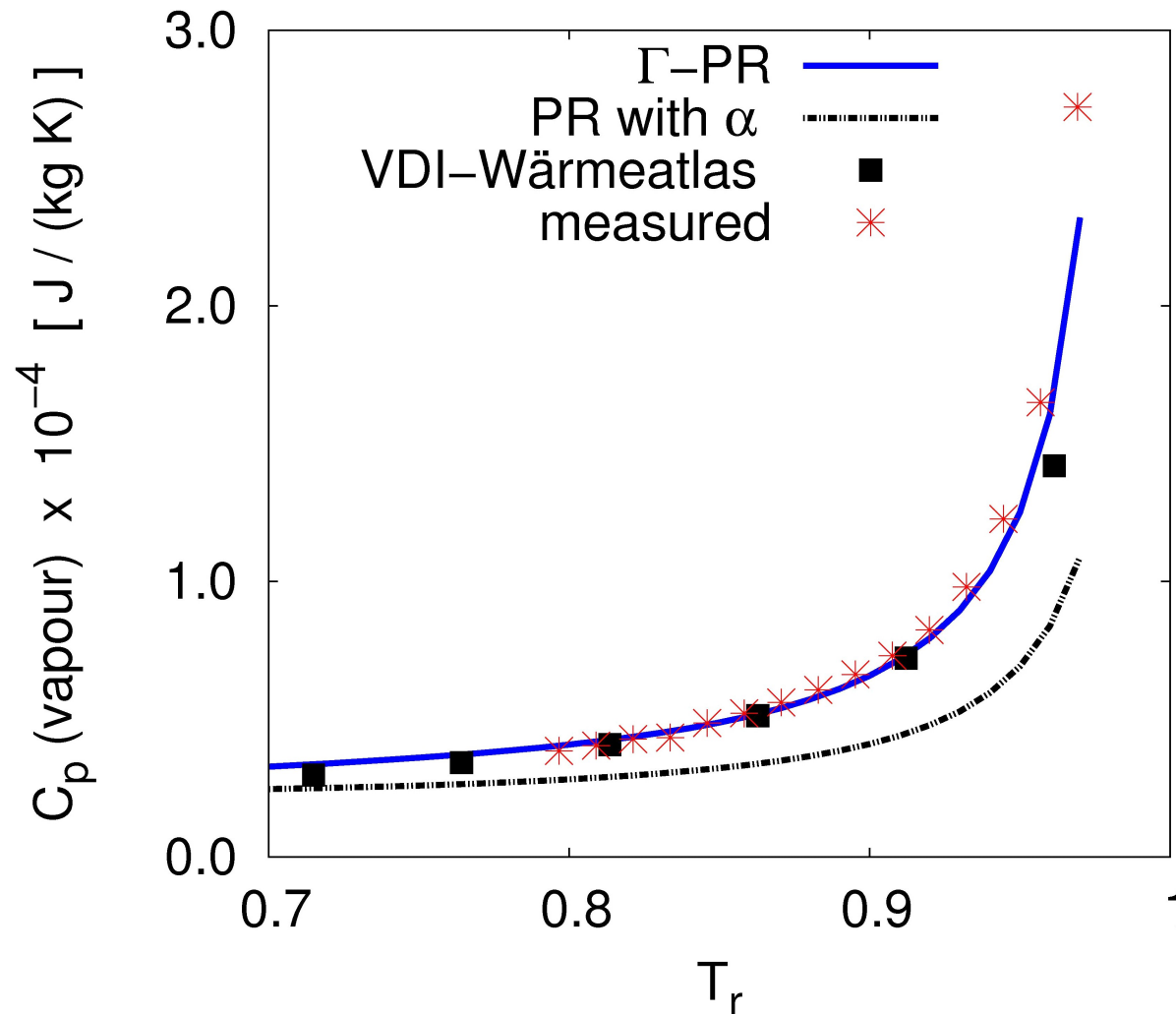Predicted quantity: heat capacity $C_{P,M}$ (1) for vapor and for NH$_3$

$\Gamma - PR$ Calibration of $\hat{\Theta}$ based on pVT data (2) only

$$C_{p,M}(T) = C_{p,id}(T) - T \cdot \int_{\upsilon}^{\infty} \left( \frac{\partial^2 p_M}{\partial T^2} \right)_V dV - R - T \cdot \frac{\left( \frac{\partial p_M}{\partial T} \right)_\upsilon^2}{\left( \frac{\partial p_M}{\partial V} \right)_T}$$

$(1)\ M = \Gamma - PR$ or $PR(\alpha_{Soave})$   (2) VDI-Wärmeatlas, 11$^{th}$ edition 2013, Springer Vieweg

# 4. Predictive Power 2

$\Gamma - PR$ versus $PR$ with Soave correction, $\alpha$

30

# Contributors

Gerhard
Krennrich

Robert
Lee

Michael
Rieger

Simeon
Sauer

# Back-up
# &
# Appendix

# Goodness of Fit - Results

Example: PR Equation of State

| # | Test quantity | Ideal Gas | Virial Equ. | PR |
|---|---|---|---|---|
| 1 | $SWS$ (1) | 5000 | 267 | 119 |
| 2 | $\chi^2_r$ | 30 (2) | 1.6 | 0.72 |
| 3 | $P^{(\chi)}_\alpha$ % | 0 | $< 10^{-4}$ | 99 |
| 4 | $AAD\,\%$ | 0.65 | 0.23 | 0.15 |

(1) $\chi^2_{crit} = \chi^2_{1-\frac{\alpha}{2}, f} \approx 197$

(2) One unit is one „statistical light year"

# Goodness of Fit - Results

Example: PR Equation of State

| # | Test quantity | Ideal Gas | Virial Equ. | PR (1) |
|---|---|---|---|---|
| 5 | Model bias 1 (2) | -0.61 | 0.22 | 0.08 |
| 6 | Model bias 2 | | | |
| | $P_{\alpha/2}(b=0)$ % | 0 | 0 | 29 |
| | $P_{\alpha/2}(a=1)$ % | 0 | 0 | 3 |

(1) Goodness of fit: model accepted.  Goodness of parameter:  model refused

(2)  sensitive for residual structure !  $\rightarrow$  indicator for bias  $\rightarrow$  residual plot

# Goodness of Parameter

$$J(X,\Theta) := -\frac{1}{s_{i=1,2...M}}\left(\frac{\partial f(X,\Theta)}{\partial \Theta}\right)_{X_{i=1,2...M}}$$

$$cov(X,\Theta) = \left(J(X,\Theta)^* \cdot J(X,\Theta)\right)^{-1}$$

$$cov(X,\Theta) = \begin{vmatrix} cov(X,\Theta)_{11} & ... & cov(X,\Theta)_{1Np} \\ & \cdot & \\ & \cdot & \\ & \cdot & \\ cov(X,\Theta)_{Np1} & ... & cov(X,\Theta)_{NpNp} \end{vmatrix}$$

$$s(X,\Theta_i) = \sqrt{cov(X,\Theta)_{ii}}$$

# Goodness of Parameter

The exact conficence region of parameter

$$F(N_p, M-N_p, \alpha) = \frac{[Y-f(X,\hat{\Theta})]^* \mathcal{P}(X,\hat{\Theta})[Y-f(X,\hat{\Theta})]/N_p}{[Y-f(X,\hat{\Theta})]^*[I-\mathcal{P}(X,\hat{\Theta})][Y-f(X,\hat{\Theta})]/(M-N_p)}$$

$$\mathcal{P}(X,\hat{\Theta}) = J(X,\hat{\Theta})\left[J(X,\hat{\Theta})^* J(X,\hat{\Theta})\right]^{-1} J(X,\hat{\Theta})^*$$

Ratkowsky, Handbook of Nonlinear Regression Models. Marcel Dekker, Inc.
New York and Basel, 1990. p. 36

# PR and SRK EoS

$$p_{PR}(T, \upsilon, y, \Theta) := \frac{RT}{\upsilon - B(y, \Theta)} - \frac{A(T, y, \Theta)}{(\upsilon + \delta B(y, \Theta)) \cdot (\upsilon + \varepsilon B(y, \Theta))}$$

# Sensitivity Analysis

## Simulation results for the PR EoS

| Sensitivity index | $\Theta_1$ | $\Theta_2$ |
|---|---|---|
| $S_I$ | 0.037 | 0.022 |
| $\tilde{S}_I$ | 18.3 | 11.1 |

Fundamental theorem in statistic for the probability P

$$P(X \leqslant c) + P(X > c) = 1$$