

Multiple Linear and Polynomial Regression with Statistical Analysis

Given a set of data of measured (or observed) values of a dependent variable: y_i versus n independent variables $x_{1i}, x_{2i}, \dots, x_{ni}$, multiple linear regression attempts to find the “best” values of the parameters a_0, a_1, \dots, a_n for the equation

$$\hat{y}_i = a_0 + a_1x_{1,i} + a_2x_{2,i} + \dots + a_nx_{n,i}$$

\hat{y}_i is the calculated value of the dependent variable at point i . The “best” parameters have values that minimize the squares of the errors

$$S = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

In polynomial regression there is only one independent variable, thus

$$\hat{y}_i = a_0 + a_1x_i + a_2x_i^2 + \dots + a_nx_i^n$$

Multiple Linear and Polynomial Regression with Statistical Analysis

Typical examples of multiple linear and polynomial regressions include correlation of temperature dependent physical properties, correlation of heat transfer data using dimensionless groups, correlation of non-ideal phase equilibrium data and correlation of reaction rate data.

The software packages enable high precision correlation of the data, however statistical analysis is essential to determine the *quality of the fit* (how well the regression model fits the data) and the *stability of the model* (the level of dependence of the model parameters on the particular set of data).

The most important indicators for such studies are the *residual plot* (quality of the fit) and *95% confidence intervals* (stability of the model)

Regression and Analysis of “Heat of Hardening” Data

No.	x_1	x_2	x_3	x_4	y
1	7	26	6	60	78.7
2	1	29	15	52	74.3
3	11	56	8	20	104.3
4	11	31	8	47	87.6
5	7	52	6	33	95.9
6	11	55	9	22	109.2
7	3	71	17	6	102.7
8	1	31	22	44	72.5
9	2	54	18	22	93.1
10	21	47	4	26	115.9
11	1	40	23	34	83.8
12	11	66	9	12	113.3
13	10	68	8	12	109.4

Woods *et al*(1932) investigated the integral heat of hardening of cement as a function of composition. The independent variables represent *weight percent* of the clinker compounds: x_1 -tricalcium aluminate ($3CaO \cdot Al_2O_3$), x_2 -tricalcium silicate ($3CaO \cdot SiO_2$), x_3 -tetracalcium alumino-ferrite ($4CaO \cdot Al_2O_3 \cdot Fe_2O_3$), and x_4 - β -dicalcium silicate ($3CaO \cdot SiO_2$). The dependent variable, y is the **total heat evolved** (in calories per gram cement) in a 180-day period.

Regression and Analysis of “Heat of Hardening” Data

No.	x_1	x_2	x_3	x_4	y
1	7	26	6	60	78.7
2	1	29	15	52	74.3
3	11	56	8	20	104.3
4	11	31	8	47	87.6
5	7	52	6	33	95.9
6	11	55	9	22	109.2
7	3	71	17	6	102.7
8	1	31	22	44	72.5
9	2	54	18	22	93.1
10	21	47	4	26	115.9
11	1	40	23	34	83.8
12	11	66	9	12	113.3
13	10	68	8	12	109.4

Calculate the coefficients of a linear model representation of y as function of x_1 , x_2 , x_3 , and x_4 , calculate the variance and the correlation coefficient R^2 and the 95% confidence intervals. Prepare a residual plot.

Consider the cases when the model includes and does not include a free parameter.

“Heat of Hardening” Data – Regression and Analysis by Polymath

The screenshot shows the POLYMATH 6.10 Educational Release interface. The main window displays a data table with columns for independent variables (Wpc1, Wpc2, Wpc3, Wpc4) and a dependent variable (hard_heat). The regression analysis panel on the right is configured with the following settings:

- Regression Analysis Graph: Residuals
- Report: Store Model:
- Model Type: Multiple linear (selected)
- Dependent Variable: hard_heat
- Independent Variables: Wpc1, Wpc2, Wpc3, Wpc4
- Through origin:

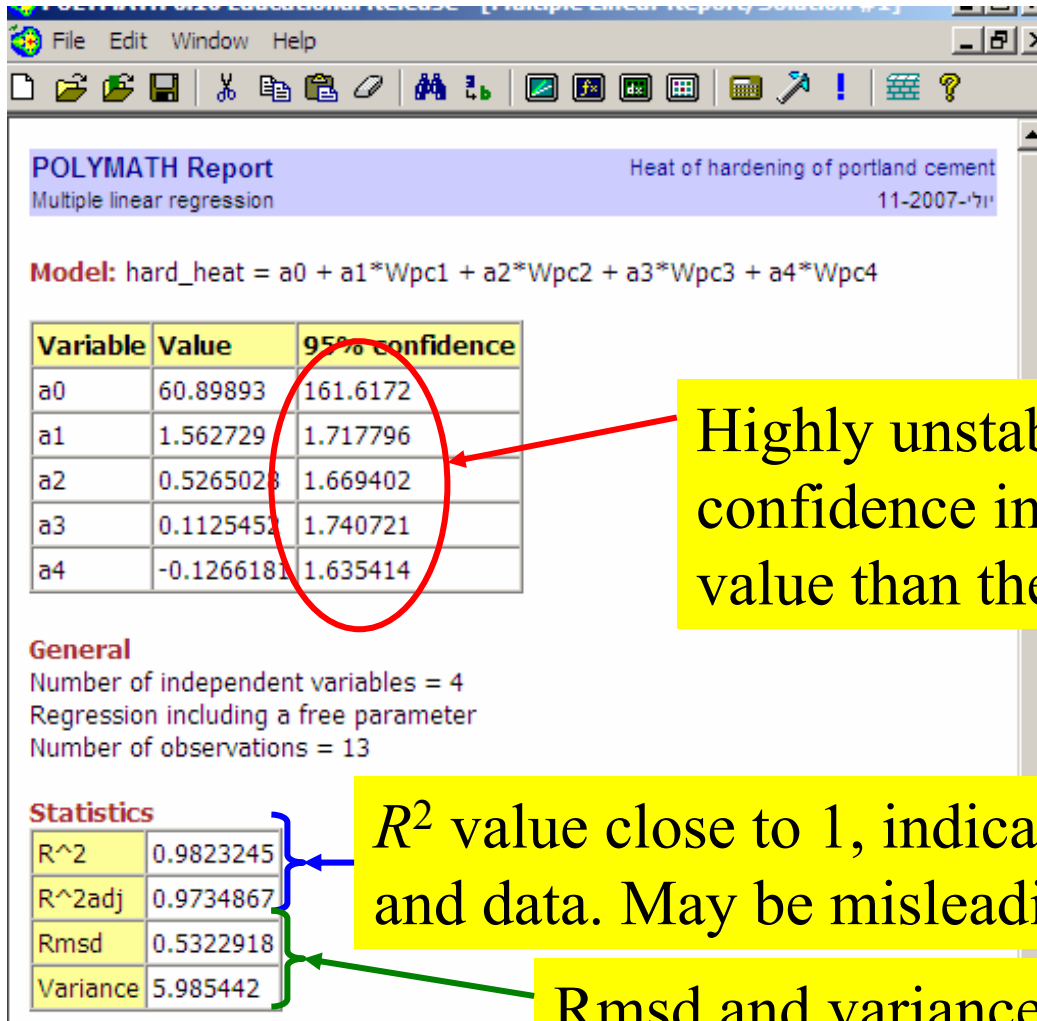
The data table contains the following values:

	Wpc1	Wpc2	Wpc3	Wpc4	hard_heat
01	7	26	6	60	78.7
02	1	29	15	52	74.3
03	11	56	8	20	104.3
04	11	31	8	47	87.6
05	7	52	6	33	95.9
06	11	55	9	22	109.2
07	3	71	17	6	102.7
08	1	31	22	44	72.5
09	2	54	18	22	93.1
10	21	47	4	26	115.9
11	1	40	23	34	83.8
12	11	66	9	12	113.3
13	10	68	8	12	109.4
14					
15					
16					
17					

Non-zero intercept

Model type

“Heat of Hardening” Data – Analysis of the Linear Model that Includes a Free Parameter

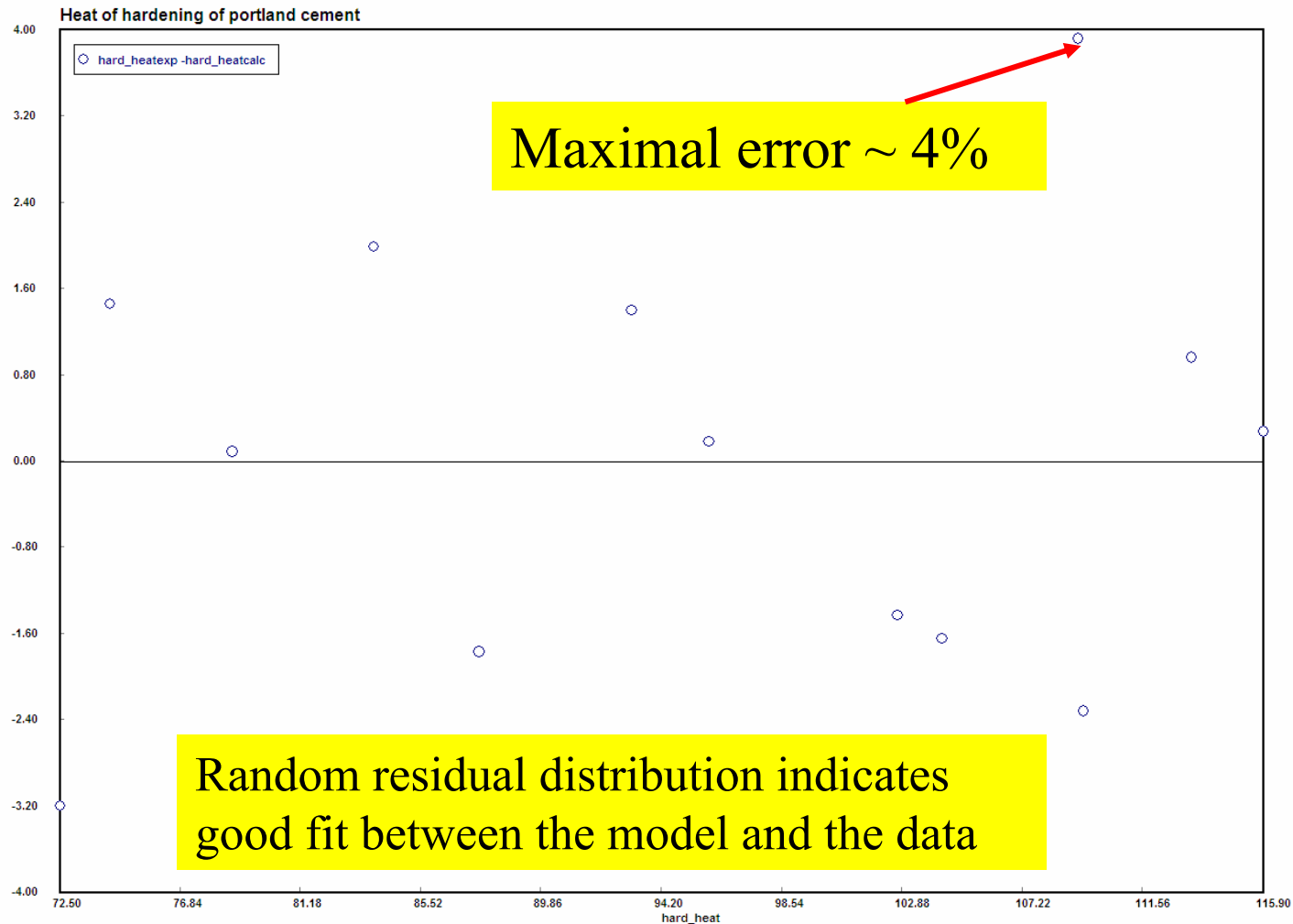


Highly unstable model. All 95% confidence intervals larger in absolute value than the respective parameters.

R^2 value close to 1, indicates good fit between model and data. May be misleading occasionally.

Rmsd and variance values used for comparison between different models

“Heat of Hardening” Data – Residual Plot of the Linear Model that Includes a Free Parameter



“Heat of Hardening” Data – Demonstration of the Harmful Effect of the Instability

Model: $\text{hard_heat} = a_0 + a_1 \cdot \text{Wpc1} + a_2 \cdot \text{Wpc2} + a_3 \cdot \text{Wpc3} + a_4 \cdot \text{Wpc4}$

Variable	Value	95% confidence
a0	60.89893	161.6172
a1	1.562729	1.717796
a2	0.526503	1.669402
a3	0.112545	1.740721
a4	-0.12662	1.635414
Statistics		
R ²	0.982325	
Variance	5.985442	

Last data point removed

Removal of the last data point causes substantial change in all the parameter values.

In case of the parameter a_4 even its sign is changed

Variable	Value	95% confidence
a0	36.20362	170.1033
a1	1.783367	1.784379
a2	0.802752	1.770434
a3	0.328399	1.804248
a4	0.12209	1.720564
Statistics		
R ²	0.983897	
Variance	5.746289	

“Heat of Hardening” Data – A Stable Model is Obtained After Removing the Free Parameter

Model: $\text{hard_heat} = a1 \cdot \text{Wpc1} + a2 \cdot \text{Wpc2} + a3 \cdot \text{Wpc3} + a4 \cdot \text{Wpc4}$

Variable	Value	95% confidence
a1	2.189177	0.4182687
a2	1.154136	0.1082325
a3	0.753295	0.3601112
a4	0.488545	0.093483

Statistics

R ²	0.980656
Variance	5.822523

Last data point removed

Variable	Value	95% confidence
a1	2.151451	0.4078854
a2	1.17869	0.1115862
a3	0.703614	0.3563326
a4	0.487699	0.0901595
Statistics		
R ²	0.983314	
Variance	5.209989	

Correlation of Heat Capacity Data for Ethane

A polynomial has to be fitted to heat capacity data provided by Ingham et al*. This data set includes 41 data points in the temperature range of 100 K – 400 K.

The degree of the polynomial:

$$C_p = a_0 + a_1 T + a_2 T^2 + a_3 T^3 + \dots$$

where C_p is the heat capacity in J/kg-mol·K, T is the temperature in K, and a_0, a_1, \dots are the regression model parameters. The set of parameters which best represents the data, has to be found.

The goodness of fit should be determined based on the variance, the correlation coefficient (R^2), the confidence intervals of the parameters, and the residual plot.

Ingham, H.; Friend, D.G.; Ely, J.F.; "Thermophysical Properties of Ethane"; *J. Phys. Ref. Data* 1991, 20, 275

Heat Capacity Data for Ethane, Fitting a 3rd Degree Polynomial

The screenshot shows the POLYMATH 6.10 Educational Release interface. The main window displays a data table with columns for row number, T_K, and Cp. The regression settings panel on the right is configured for a polynomial fit. The dependent variable is Cp and the independent variable is T_K. The polynomial degree is set to 3. The 'Graph' and 'Residuals' options are checked, and the 'Linear & Polynomial' model type is selected. A red arrow points from the 'Model type' label to the 'Linear & Polynomial' tab, and another red arrow points from the '3' in the polynomial degree list to the 'Model type' label.

	T_K	Cp
01	100	3.5698E+04
02	110	3.6249E+04
03	120	3.6817E+04
04	130	3.7401E+04
05	140	3.8003E+04
06	150	3.8628E+04
07	160	3.9279E+04
08	170	3.9961E+04
09	180	4.0680E+04
10	190	4.1439E+04
11	200	4.2243E+04
12	210	4.3092E+04
13	220	4.3989E+04
14	230	4.4934E+04
15	240	4.5924E+04
16	250	4.6959E+04
17	260	4.8036E+04
18	270	4.9151E+04
19	280	5.0302E+04
20	290	5.1484E+04

Regression Analysis Settings:

- Dependent Variable: Cp
- Independent Variable: T_K
- Polynomial Degree: 3
- Model type: Linear & Polynomial
- Graph:
- Residuals:
- Report: Store Model:
- Through origin:
- Polynomial Integration:
- Polynomial Derivative:

Model type

Heat Capacity Data for Ethane, Fitting a 3rd Degree Polynomial

POLYMATH 6.10 Educational Release - [Polynomial Regression, ...]

File Edit Window Help

POLYMATH Report
Polynomial Regression

Model: $C_p = a_0 + a_1 \cdot T_K + a_2 \cdot T_K^2 + a_3 \cdot T_K^3$

Variable	Value	95% confidence
a0	3.641E+04	549.6705
a1	-44.10631	6.531667
a2	0.4397598	0.0234039
a3	-0.0003702	2.583E-05

General
Degree of polynomial = 3
Regression including a free parameter
Number of observations = 41

Statistics

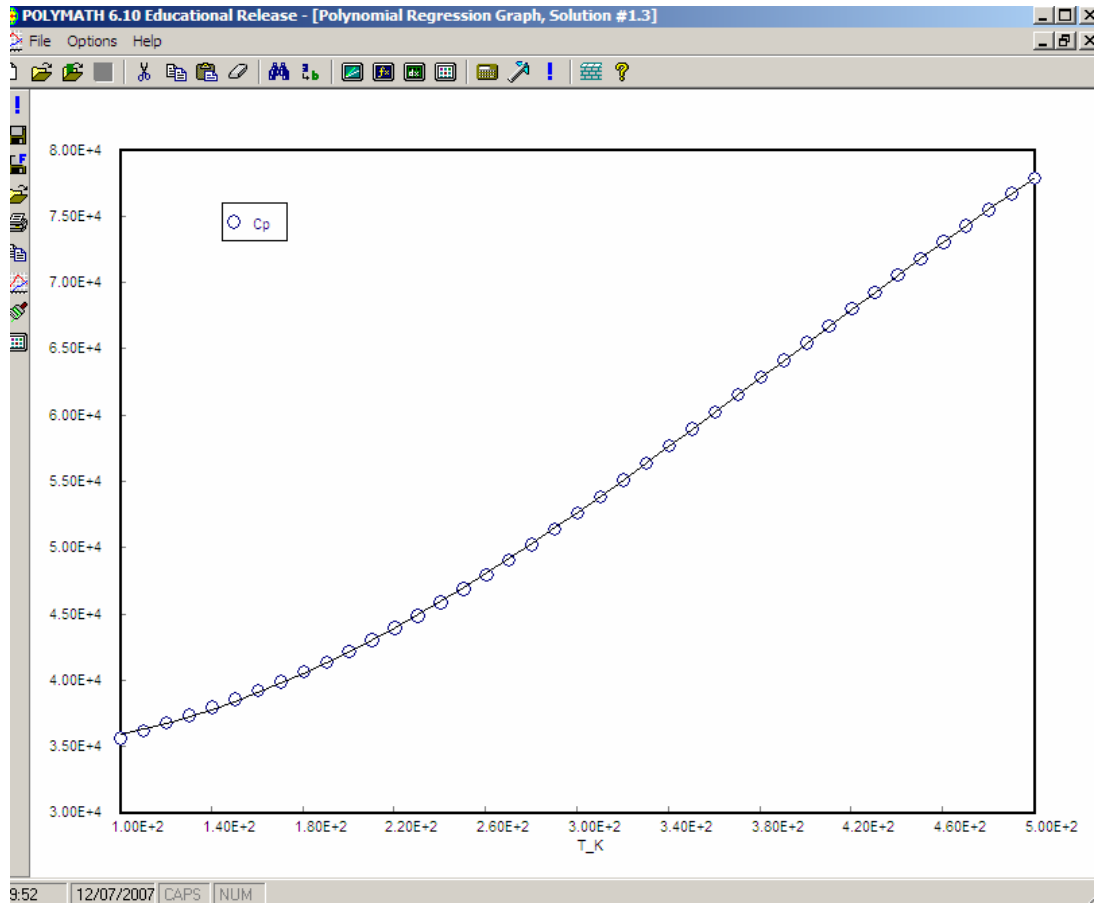
R ²	0.9999407
R ² adj	0.9999359
Rmsd	15.81281
Variance	1.136E+04

All the parameters indicate satisfactory model.

Note, however, the differences of orders of magnitude between the parameter values. This may limit the highest degree of polynomial to be fitted

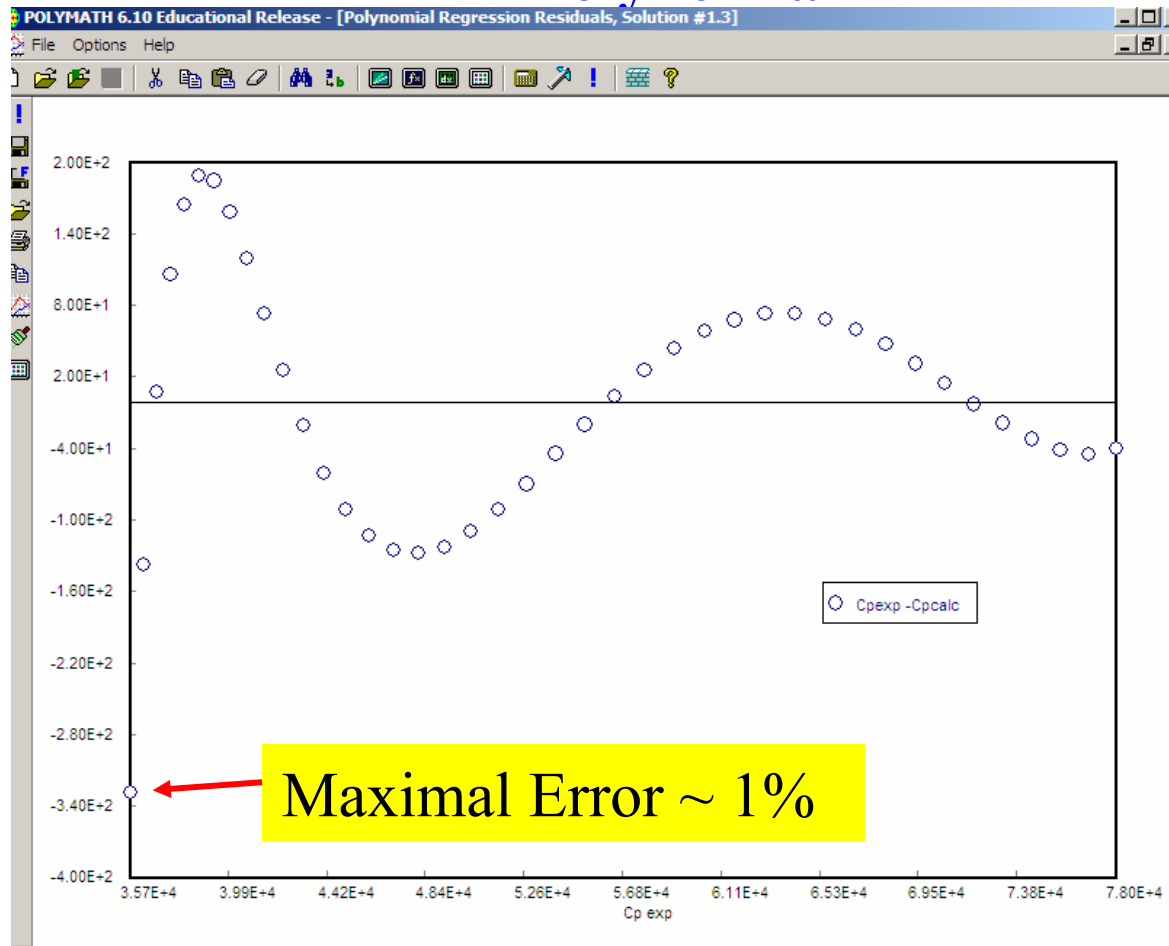
Rmsd and variance values used for comparison between different models

Heat Capacity Data for Ethane, Calculated (3rd Degree Polynomial) and Experimental Values



On the scale of the entire range of the C_p data the fit seems to be excellent

Heat Capacity Data for Ethane, Residual Plot for the 3rd Degree Polynomial



High resolution residual plot shows oscillatory behavior which is not explained by the 3rd degree polynomial

Heat Capacity Data for Ethane, Defining Standardized Temperature Values for High Order Polynomial Fitting

POLYMATH 6.10 Educational Release - [Data Table]

File Program Edit Row Column Format Analysis Examples Win

R001 : C003 Tstd $\times \checkmark$ $= (T_K - 300) / 119.7915$

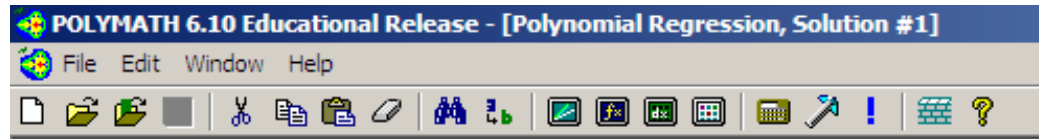
	T_K	Cp	Tstd	C04	C05
01	100	3.5698E+04	-1.669568		
02	110	3.6249E+04	-1.586089		
03	120	3.6817E+04	-1.502611		
04	130	3.7401E+04	-1.419132		
05	140	3.8003E+04	-1.335654		
06	150	3.8628E+04	-1.252176		
07	160	3.9279E+04	-1.168697		
08	170	3.9961E+04	-1.085219		
09	180	4.0680E+04	-1.001741		
10	190	4.1439E+04	-0.9182621		
11	200	4.2243E+04	-0.8347838		
12	210	4.3092E+04	-0.7513054		
13	220	4.3989E+04	-0.667827		
14	230	4.4934E+04	-0.5843486		

Column Statistics

Column name : T_K
Number of rows : 41
Sum : 1.23E+04
Arithmetic mean : 300.
Standard deviation : 119.7915
Variance : 1.435E+04

Column name : Cp
Number of rows : 41
Sum : 2.224E+06
Arithmetic mean : 5.425E+04
Standard deviation : 1.331E+04
Variance : 1.772E+08

Heat Capacity Data for Ethane, Fitting a 5rd Degree Polynomial



POLYMATH Report

Polynomial Regression

Model: $C_p = a_0 + a_1 \cdot T_{std} + a_2 \cdot T_{std}^2 + a_3 \cdot T_{std}^3 + a_4 \cdot T_{std}^4 + a_5 \cdot T_{std}^5$

Variable	Value	95% confidence
a0	5.268E+04	7.53794
a1	1.461E+04	17.87844
a2	1812.29	16.18934
a3	-1041.693	24.12796
a4	-112.7446	6.199
a5	125.1056	7.262645

General

Degree of polynomial = 5

Regression including a free parameter

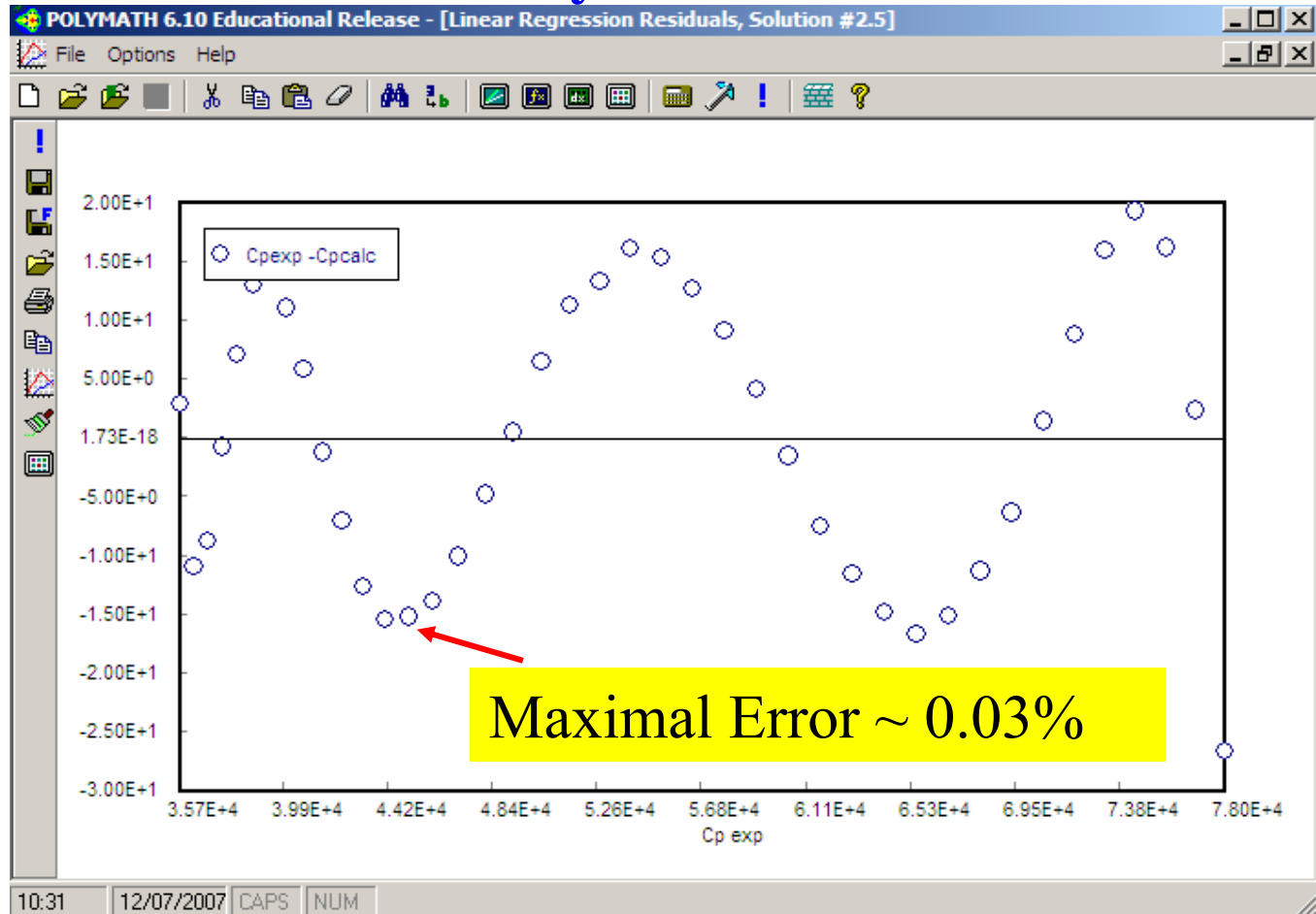
Number of observations = 41

Statistics

R ²	0.9999992
R ² adj	0.9999991
Rmsd	1.834451
Variance	161.6262

Using standardized values yields model parameters of similar magnitude, enables fitting higher order polynomials and improves considerably all the statistical indicators

Heat Capacity Data for Ethane, Residual Plot of a 5th Degree Polynomial



Using standardized independent variable values enables fitting polynomials with *precision higher than justified by the experimental error.*

Modeling Vapor Pressure Data for Ethane

A vapor pressure data set provided by Ingham et al* includes 107 data points in the temperature range of 92 K – 304 K. This temperature range covers almost completely the range between the triple point temperature (= 90.352 K) and the critical temperature ($T_C = 305.32$ K).

The temperature dependence of the vapor pressure should be modeled by the Clapeyron, Antoine and Wagner equations

The Clapeyron equation is a two parameter equation:

$$\ln P = A + \frac{B}{T} \quad \text{where } P \text{ is the vapor pressure (Pa), } T \text{ – temperature (K), } A \text{ and } B \text{ are parameters}$$

*Ingham, H.; Friend, D.G.; Ely, J.F.; "Thermophysical Properties of Ethane"; *J. Phys. Ref. Data* 1991, 20, 275

Modeling Vapor Pressure Data for Ethane by the Clapeyron Equation using Linear Regression

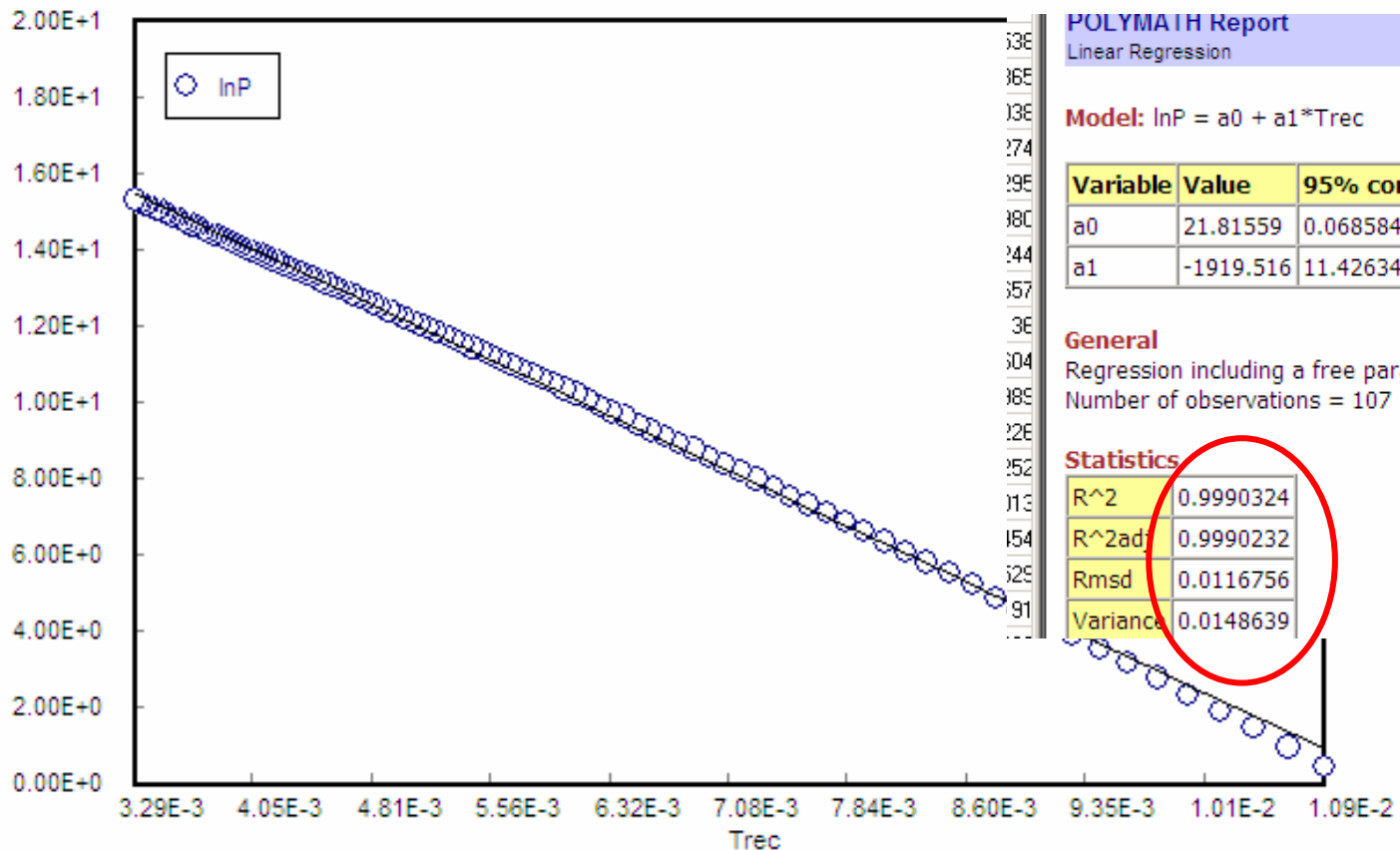
The screenshot displays the POLYMATH 6.10 Educational Release interface. The main window shows a data table with columns for Temperature (T_K), Pressure (P_Pa), and calculated values for T_{rec} and $\ln P$. The regression settings panel on the right is configured for a linear regression model.

Regression settings:
Dependent Variable: $\ln P$
Independent Variable: T_{rec}
Polynomial Degree: 1 Linear
 Through origin
 Polynomial Integration

Annotations on the regression settings:
 $\ln P = \ln(P_Pa)$
 $T_{rec} = 1/T_K$

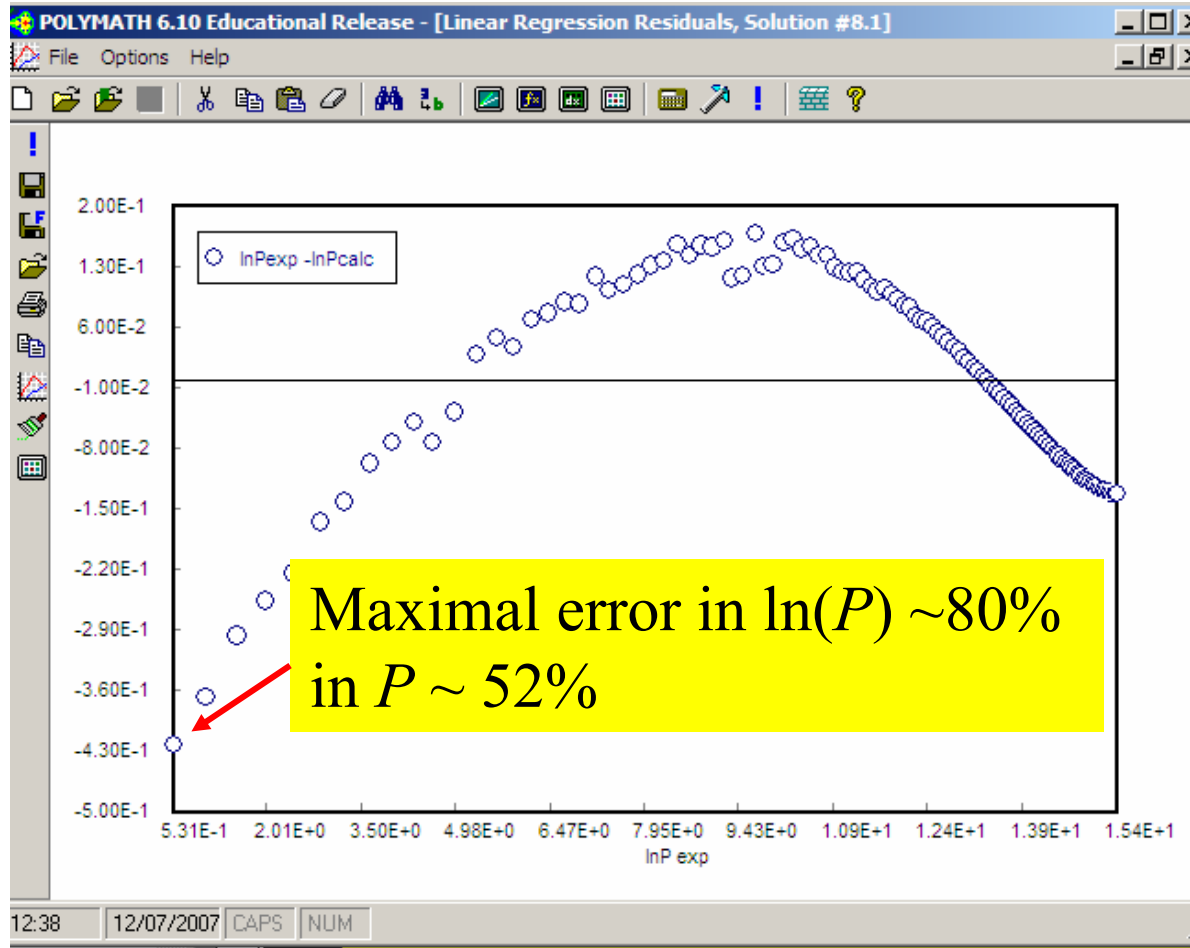
	T_K	P_Pa	Trec	lnP
01	92	1.7	0.0108696	0.5306283
02	94	2.8	0.0106383	1.029619
03	96	4.6	0.0104167	1.526056
04	98	7.2	0.0102041	1.974081
05	100	11	0.01	2.397895
06	102	17	0.0098039	2.833213
07	104	25	0.0096154	3.218876
08	106	37	0.009434	3.610918
09	108	53	0.0092593	3.970292
10	110	75	0.0090909	4.317488
11	112	100	0.0089286	4.60517
12	114	140	0.0087719	4.941642
13	116	200	0.0086207	5.298317
14	118	270	0.0084746	5.598422
15	120	350	0.0083333	5.857933
16	122	470	0.0081967	6.152733

Modeling Vapor Pressure Data for Ethane by the Clapeyron Equation using Linear Regression



All the indicators show good, acceptable fit!

Modeling Vapor Pressure Data for Ethane by the Clapeyron Equation using Linear Regression



The residual plot reveals large unexplained curvature in the data

Modeling Vapor Pressure Data for Ethane by the Antoine Equation using Non-linear Regression

$$\ln P = A + \frac{B}{T + C}$$

The screenshot shows the POLYMATH 6.10 Educational Release interface. The main window displays a data table with columns for T_K, P_Pa, and lnP. The regression settings are configured for a Nonlinear model using the Levenberg-Marquardt (lmqmin) algorithm. The model equation is lnP=A+B/(T_K+C). The initial guess values for parameters A, B, and C are 22, -2000, and 0, respectively. The software also shows a report checkbox and a store model checkbox.

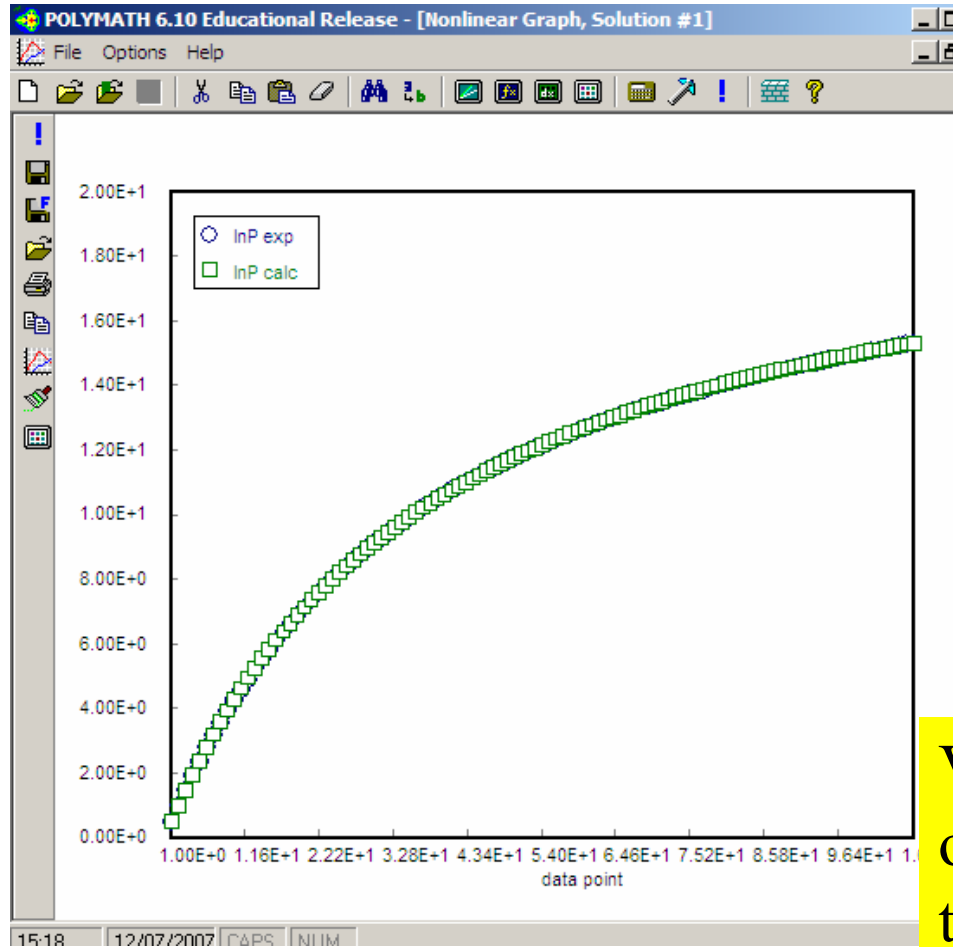
	T_K	P_Pa	lnP	CO
01	92	1.7	0.5306283	
02	94	2.8	1.029619	
03	96	4.6	1.526056	
04	98	7.2	1.974081	
05	100	11	2.397895	
06	102	17	2.833213	
07	104	25	3.218876	
08	106	37	3.610918	
09	108	53	3.970292	
10	110	75	4.317488	
11	112	100	4.60517	
12	114	140	4.941642	
13	116	200	5.298317	
14	118	270	5.598422	
15	120	350	5.857933	
16	122	470	6.152733	
17	124	610	6.413459	
18	126	790	6.670022	

Model parm	Initial guess
A	22
B	-2000
C	0

Model type and solution algorithm

Initial guess from Clapeyron eqn.

Modeling Vapor Pressure Data for Ethane by the Antoine Equation using Non-linear Regression



Model: $\ln P = A + B / (T_K + C)$

Variable	Initial guess	Value	95% confidence
A	22.	20.79896	0.0208029
B	-2000.	-1583.811	6.234624
C	0	-13.88808	0.2627736

Nonlinear regression settings

Max # iterations = 64

Tolerance = .0001

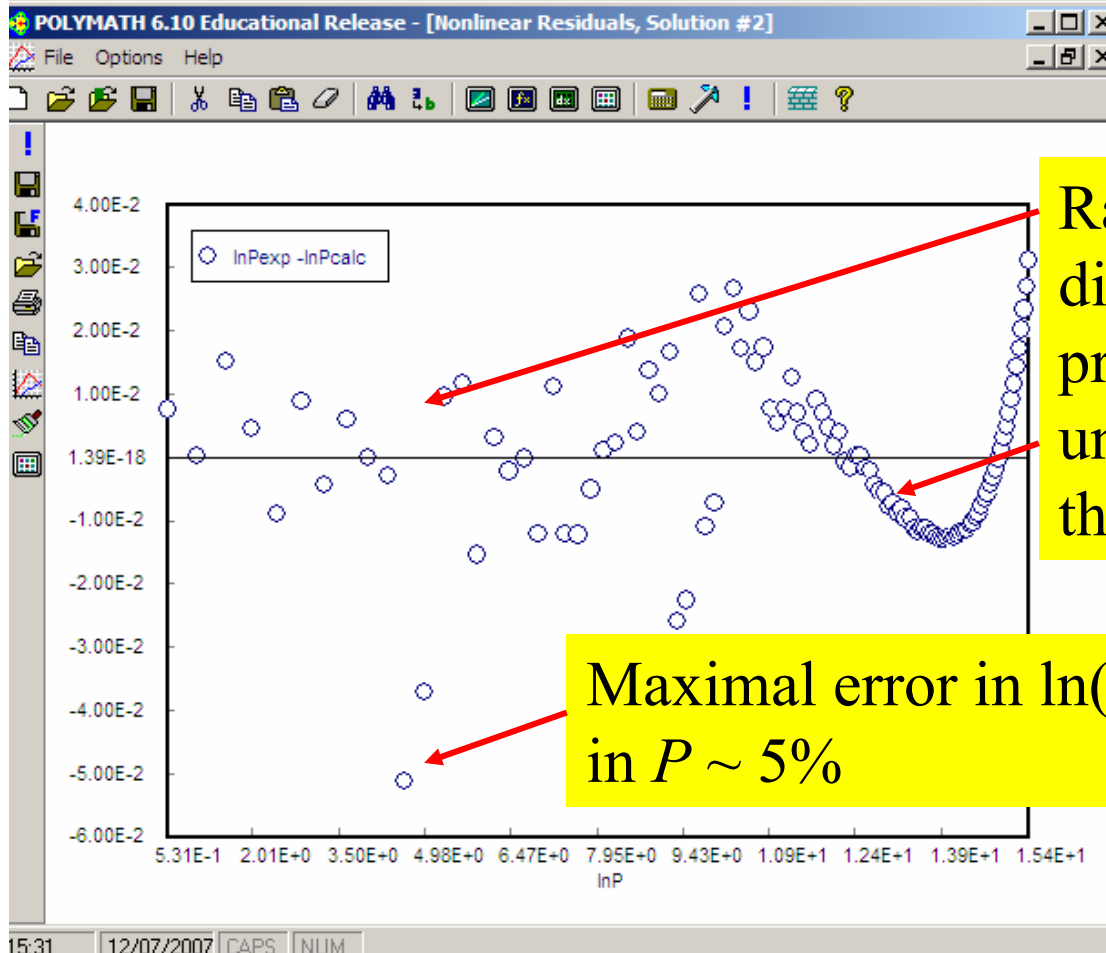
Precision

R ²	0.9999885
R ² adj	0.9999883
Rmsd	0.0012724
Variance	0.0001782
Chi-Sq	1.853625

Variance smaller by 2 orders of magnitude than Clapeyron

Experimental and calculated values cannot be distinguished.

Modeling Vapor Pressure Data for Ethane by the Antoine Equation using Non-linear Regression



Random residual distribution in the low pressure range, unexplained curvature in the high pressure range

Maximal error in $\ln(P) \sim 1\%$
in $P \sim 5\%$

Modeling Vapor Pressure Data for Ethane with the Wagner Equation

$$\ln P_R = \frac{a\tau + b\tau^{1.5} + c\tau^3 + d\tau^6}{T_R}$$

Where $T_R = T/T_C$ is the reduced temperature $P_R = P/P_C$ is the reduced pressure and $\tau = 1 - T_R$.

For ethane $T_C = 305.32$ K, $P_C = 4.8720 \times 10^6$ Pa

In order to obtain the model parameters using linear regression the following variables are defined:

$$\text{Tr} = T_K / 305.32$$

$$\ln\text{Pr} = \ln(P_Pa / 4872000)$$

$$t = (1 - \text{Tr}) / \text{Tr}$$

$$t15 = (1 - \text{Tr})^{1.5} / \text{Tr}$$

$$t3 = (1 - \text{Tr})^3 / \text{Tr}$$

$$t6 = (1 - \text{Tr})^6 / \text{Tr}$$

Modeling Vapor Pressure Data for Ethane with the Wagner Equation using Multiple Linear Regression

POLYMATH 6.10 Educational Release - [Data Table]

File Program Edit Row Column Format Analysis Examples Window Help

R001 : C009 C12

	T_K	P_Pa	Tr	lnPr	t	t15	t3	t6
01	92	1.7	0.3013232	-14.86839	2.318696	1.938126	1.13187	0.3860338
02	94	2.8	0.3078737	-14.3694	2.248085	1.870275	1.07692	0.3570586
03	96	4.6	0.3144242	-13.87296	2.180417	1.805374	1.024827	0.3302303
04	98	7.2	0.3209747	-13.42493	2.11551	1.743244	0.9754096	0.305383
05	100	11	0.3275252	-13.00112	2.0532	1.683718	0.9285029	0.2823653
06	102	17	0.3340757	-12.5658	1.993333	1.626643	0.8839539	0.2610383
07	104	25	0.3406262	-12.18014	1.935769	1.57188	0.8416217	0.2412748
08	106	37	0.3471767	-11.7881	1.880377	1.519298	0.8013759	0.222958
09	108	53	0.3537272	-11.42872	1.827037	1.468775	0.7630958	0.2059807
10	110	75	0.3602777	-11.08153	1.775636	1.420201	0.7266695	0.1902442
11	112	100	0.3668282	-10.79384	1.726071	1.373471	0.6919932	0.1756574
12	114							0.1621366
13	116							0.1496045
14	118							0.1379897
15	120							0.1272261
16	122							0.1172527
17	124							0.1080131
18	126							0.0994548
19	128							0.0915002

WagnerEq.pol No

16:09 12/07/2007

Regression Analysis Graph

Report Store Model

Linear & Polynomial Multiple linear Nonlinear

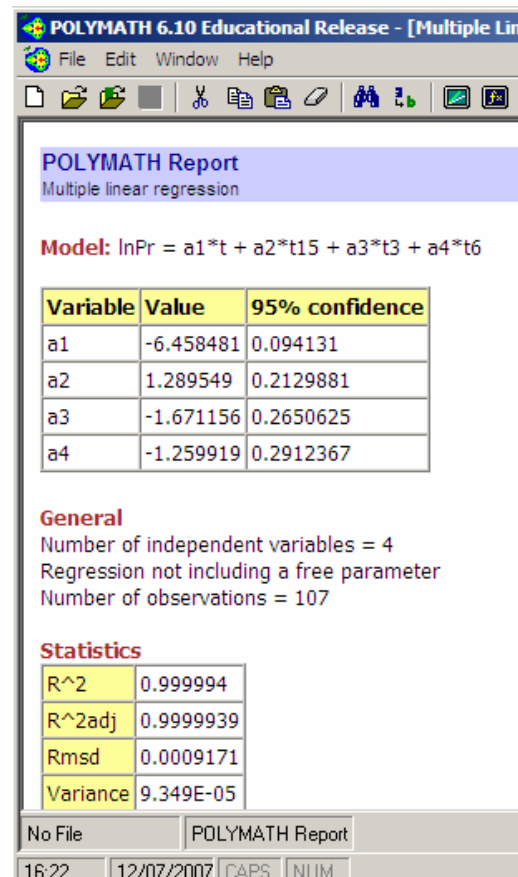
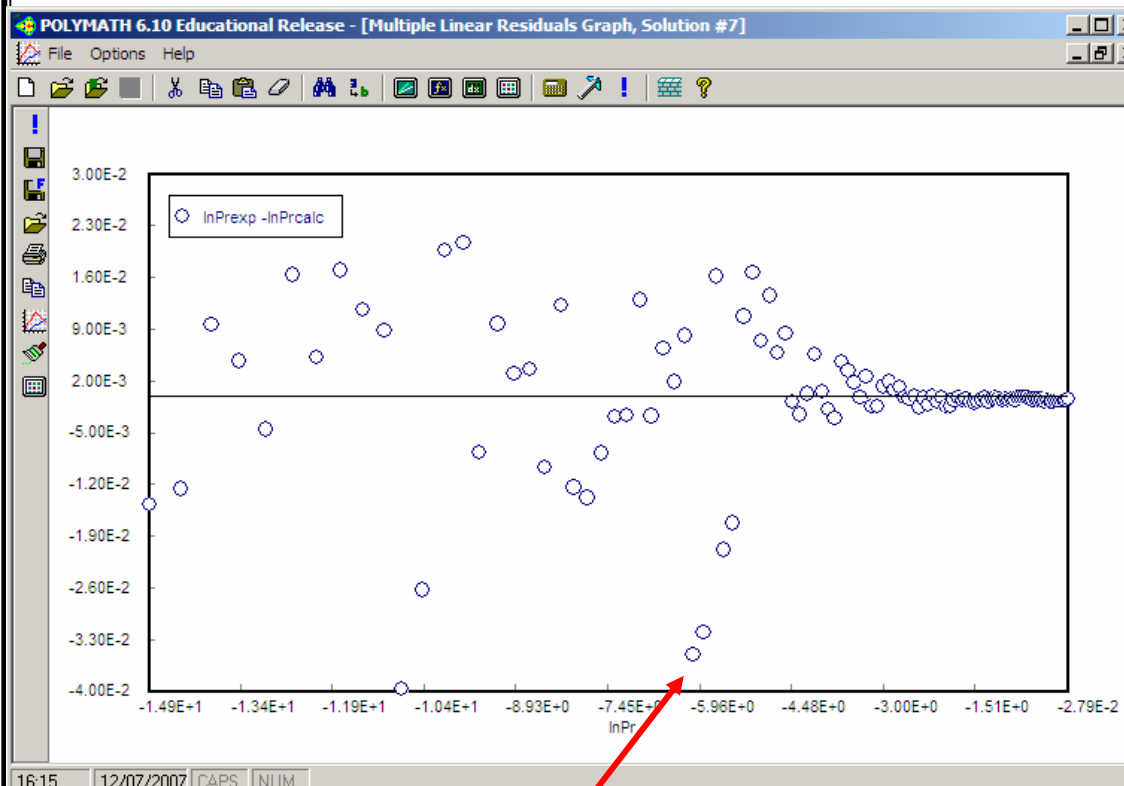
Dependent Variable: lnPr

Independent Variables: T_K, P_Pa, Tr, t, t15, t3, t6

Through origin

$Tr = T_K / 305.32$
 $\ln Pr = \ln(P_Pa / 4872000)$
 $t = (1 - Tr) / Tr$
 $t15 = (1 - Tr)^{1.5} / Tr$
 $t3 = (1 - Tr)^3 / Tr$
 $t6 = (1 - Tr)^6 / Tr$

Modeling Vapor Pressure Data for Ethane with the Wagner Equation using Multiple Linear Regression



Maximal Error in $\ln(Pr) \sim 0.6\%$

Note random residuals distribution in the entire data range