



Visual Analytics for Creating and Validating Regression Models

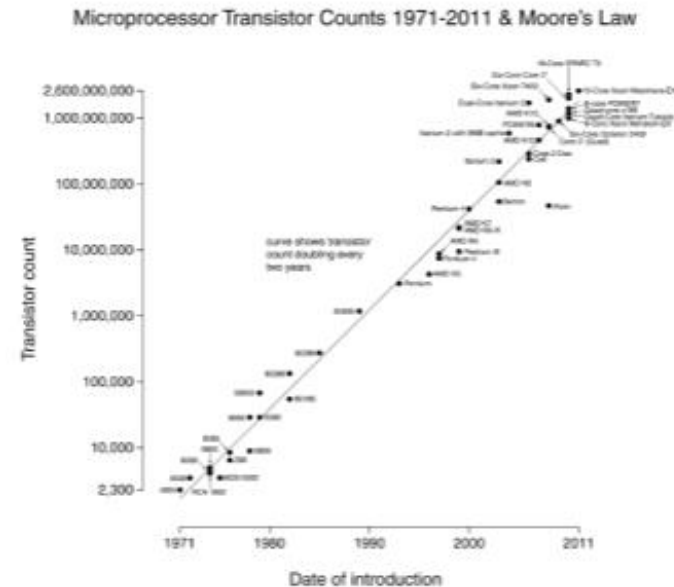
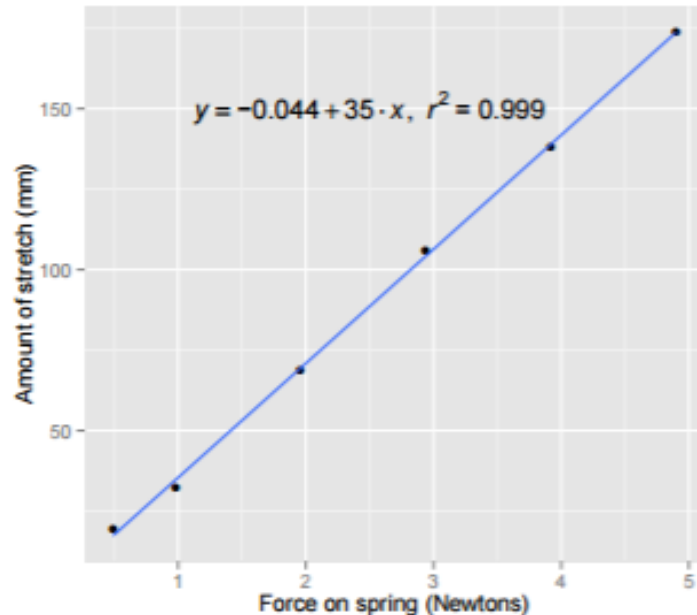
B. Preim | Fakultät für Informatik | Institut für Simulation und Graphik
Otto-von-Guericke-Universität Magdeburg

- Introductory Examples
- Regression Models
- Automatic Building of Regression Models
- Interactive Regression Model Building
 - Visualization Techniques to support Model building
- Partition-Based Regression Analysis
- Validation of Regression Models
- Case Studies
 - Prediction of gas consumption
 - Search for Breast cancer risks

Why we deal with regression in visual analytics?

- Regression is a core part of statistics.
- With nowadays big data, manual evaluation of relations is no longer feasible.
- Automatic regression models exhibit a number of limitations and drawbacks.
- Interactive regression modeling using visualizations of the effects of modeling decisions have a great potential for complex modelling tasks in engineering, climate research, ...
- Successful case studies exist in the Visual Analytics community.

Introductory Examples



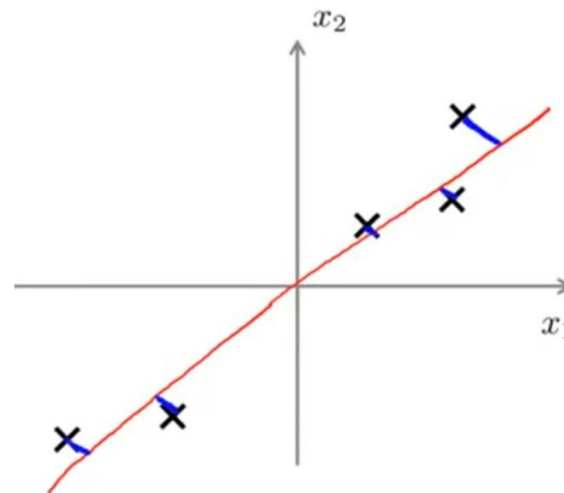
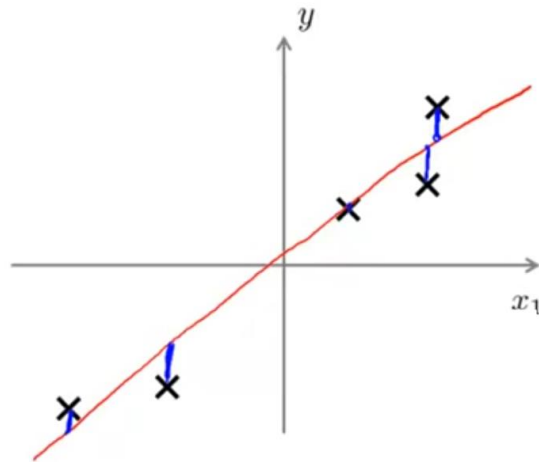
Linear regression is used to empirically verify Hooke's law (relation between the force applied to a spring and the resulting displacement) and to explain how the transistor density developed over 40 years (Moore's law)

From: [Link](#). Right image: is by Wgsimon (own work) [CC-BY-SA-3.0 or GFDL], via Wikimedia Commons.

Introductory Examples

What is the relation between linear regression and PCA?

- In 2D, we search for a 1D representation as a regression line or a 1D vector to represent the 2D data best.
- We minimize distances between the 2D points and the line (in a least square sense).
- The specific error measure is different. In regression modeling, we compute deviations of the true y -value to the model, whereas in PCA we project the data on the 1D vector.



Example 1:

- A gas provider needs accurate predictions of future gas consumption, e.g. to decide on the use of renewable energy
- The gas provider has collected data, relating to temperature, wind speed and direction that may be combined with day of the week, month, season to compute a regression model.

Regression model: $Y = \varepsilon + b_1X_1 + b_2X_2 + \dots$

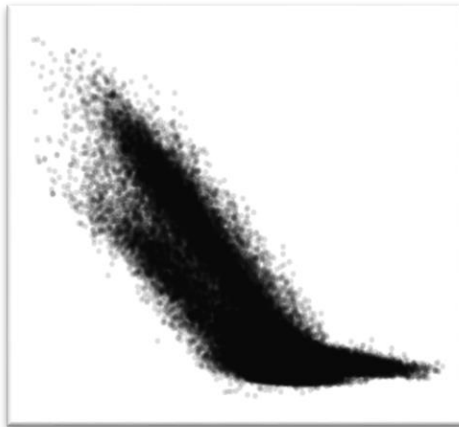
Search for variables that significantly influence our target (gas consumption) and their weight.

Want to know how reliable our model is; i.e. how accurately can we *predict* consumption based on a weather *forecast* (with certain reliability)?

Introductory Examples

Multiple linear regression: gas consumption is predicted based on temperature and time of the day. The resulting curved surface (the model) is compared with the actual data.

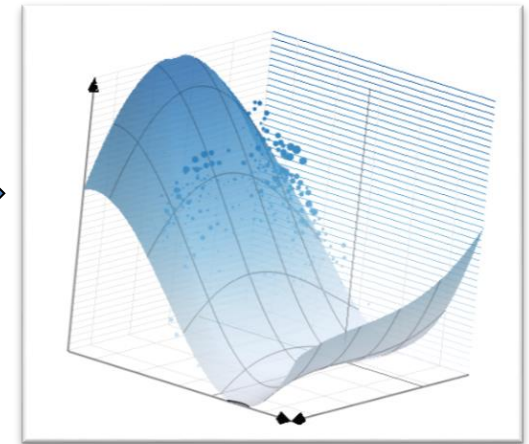
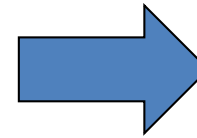
Natural Gas
Consumption



Temperature



Time of Day



Regression Model

(Courtesy of T. Mühlbacher/H. Piringer, VR Vis Vienna)

Introductory Example

Example 2:

Life expectancy in Germany differs up to 6.3 years for males between München and Tübingen compared to Demmin and Bördekreis; 3.5 years for females) (Latzis, 2011, [Link](#))

Data related to income, educational status, unemployment, gross domestic product, air pollution, crime rates, life style (smoking, alcohol consumption, physical activity) and healthcare-related data (density of stroke units, specialized treatment for cardiac diseases), climate-related data is available.

Questions:

- Can we explain the differences in life expectancy?
- Are there factors involved that may be changed?
- Can we predict the amount of improved life expectancy in case certain measures are taken?

Introductory Examples

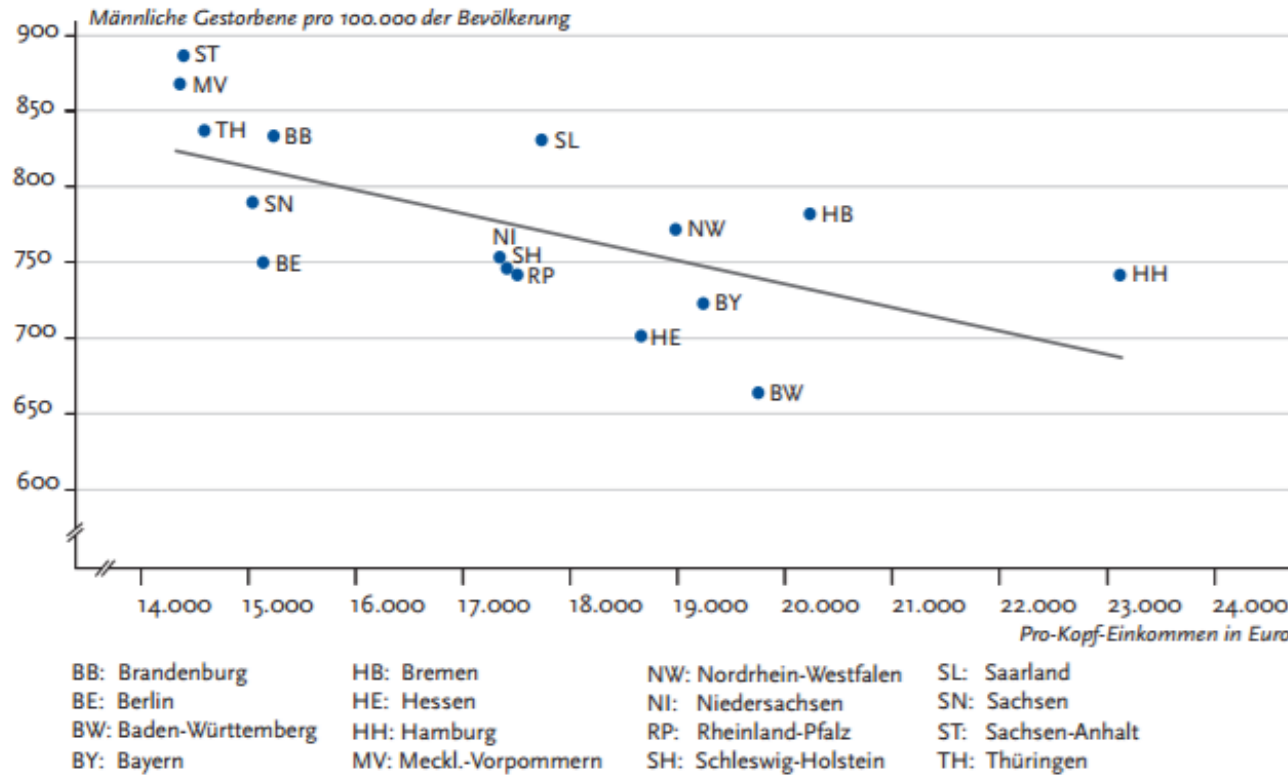
Regression analysis serves to estimate relationships.

It has two major categories of application:

- *Predict* values of the dependent variable given the input parameters (gas consumption, support for decisions in marketing, ...)
- *Explain* how independent variables influence the dependent target variable (health measured as life expectancy). Quantify the strength of the relationship.

For life expectancy example: the strongest positive influence have „Einkünfte je Steuerpflichtigen“, „Verfügbares Einkommen“, „Pflegepersonal im Heim“; the strongest negative influence have „Arbeitslosenquote“, „Keinen Hauptschulabschluss“

Introductory Examples



Mortality (how many persons died within a year) is negatively correlated (-0.65) with income. The relation is highly significant ($p > 0.01$) (From: [Robert-Koch-Institute](https://www.rki.de), see Kuhn, 2006 for more details)

Correlation and causal effects:

- A verified and strong correlation *may be* related to a causal effect. Often it is not.
- Is a low educational level the cause for later health problems or the consequence of already existing health problems?
- Homeless people are in bad health. Did they get homeless because of bad health or is bad health the consequence of being homeless?
- There is a strong significant effect that people with larger shoe sizes die earlier. Is shoe size the true reason?
- Thus, an expert is needed to interpret results of regression modeling. Often further hypothesis-driven research is necessary to investigate the results of exploratory regression modeling.

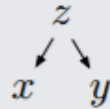
Introductory Examples

Causality and Correlation

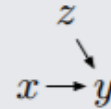
$x \longrightarrow y$

$x \longleftarrow y$

(a) **Causal link:** Even if there is a causal link between x and y , correlation alone cannot tell us whether y causes x or x causes y .



(b) **Hidden Cause:** A hidden variable z causes both x and y , creating the correlation.



(c) **Confounding Factor:** A hidden variable z and x both affect y , so the results also depend on the value of z .

$x \quad y$

(d) **Coincidence:** The correlation just happened by chance

Different reasons for a correlation (From: [Link](#)). The shoe size phenomenon is based on a hidden cause; the gender difference in mortality.

Which types of regression models exist?

- Simple linear models (2 variables)
- Partial linear regression models
- Piece-wise linear regression models
- Non-linear regression models (typically quadratic, cubic)
- Logistic regression models
 - For categorical and binary variables

Important terms:

$$Y = \varepsilon + b_1X_1 + b_2X_2 + \dots$$

Y is called the *regressand*, *target*, response, or dependent variable

The X_i are called *regressor*, explanatory, predictor or independent variable.

The b_i are called *regression coefficients*.

ε is called error term or noise.

Often a vector notation is used:

$$Y = \varepsilon + \mathbf{b}^T \mathbf{X}$$

We restrict to *one* target. Multiple targets – the topic of multivariate regression analysis – is not considered.

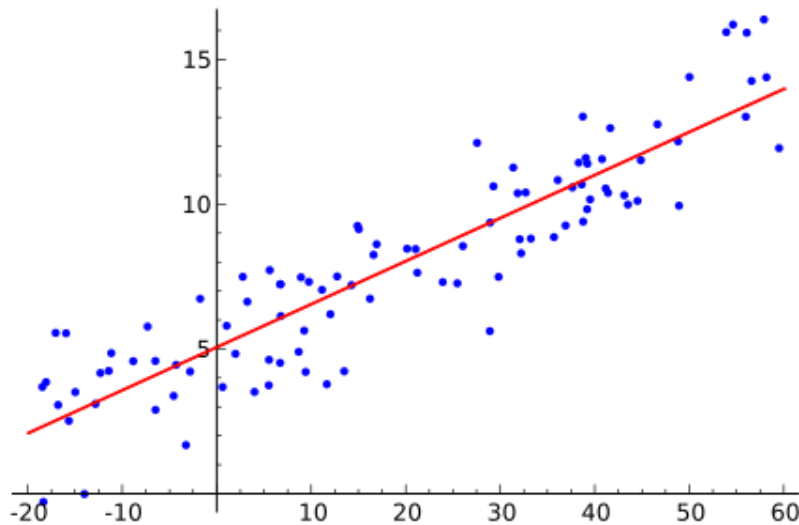
Regression Models

- The determination of specific parameters requires special assumptions, e.g. a linear model.
- The specific parameters are determined based on a certain norm that serves to provide the best fit.
 - Most frequently, the sum of squared differences is minimized („least square“).
 - Least absolute deviation is an alternative.
- Data should be evenly spread on both sides of the regression line/curve.
- The goodness of the fit is determined as correlation coefficient (r ranges between -1 and 1; r^2 is always positive and indicates the variability explained with that model)

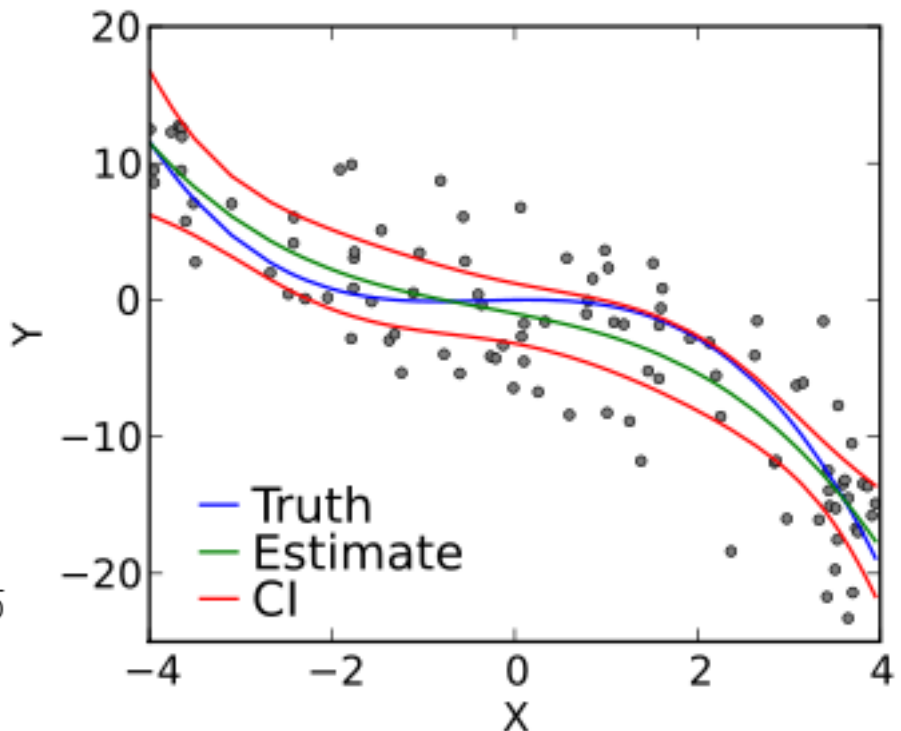
$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right).$$

- Further assumption: explanatory variables are *independent*. Variables that are dependent on each other should be analyzed with care or one of them should be removed. (Think of our examples: Are there suspicious variable pairs?)
- If some variables are perfectly correlated, the regression coefficients cannot be determined (there is no solution).
- R^2 is the most important quality criterion; adjusted R^2 and Aikaike ([Link](#)) are also frequently used.
 - Aikaike coefficient (AIC) incorporates a penalty for complex models. With the same quality of fit, a model with less regressors has a lower (better) AIC.
 - Aikaike prefers less complex models and is less prone to overfitting.
 - A modified version corrects for sample size.

Regression Models



A simple linear regression model. Data was artificially generated with random noise added to a linear term. Least square optimization (From: [Wikipedia](https://en.wikipedia.org/wiki/Least_squares))



A cubic regression model with confidence interval (red). The model (represented by the green line) is applied to test data and does not perfectly fit the truth (blue). (From: [Wikipedia](https://en.wikipedia.org/wiki/Cubic_regression))

Regression Models

- Polynomial regression is a powerful method. A polynomial up to a certain rank is determined to yield an optimal fit.
- Polynomial regression may suffer from overfitting; a risk that is increasing with
 - increasing rank of the polynomial and
 - Small sample sizes. In the extremum: a polynomial of rank n is used to fit a sample of size n . The fit is perfect, but overfitting is extreme.
- Also linear regression may be prone to overfitting when the sample size is small and too many regressors are used ([Link](#)).
 - For a sample size of 1000, no more than 4 regressors should be used

Multicollinearity:

- Means that some explanatory variables are not independent
- When the correlation (R^2) between two variables is > 0.8 possible collinearity exists. If $R^2 > 0.95$ it definitely occurs.
- Another diagnostic test is the Variance Influence Factor (VIF).
- Normal VIF values are < 5 . $VIF > 5$ is suspicious and $VIF > 10$ reliably identifies multicollinearity.

In case of a pair of variables X_1, X_2 being multicollinear, only one of them should be selected for regression modeling.

Scaling variables:

- Variables may have very different units and scales, e.g. studies related to criminality in the US involve the portion of young man (aged 18-24), the portion of non-white persons, average income, population density.
- For expressive regression models, the data needs to be scaled.
- A frequent type of scaling is to center the variables and normalize them w.r.t. the standard deviation σ .

$$\hat{y} = (y - \bar{y}) / \sigma$$

This normalization leads to a new scale with mean 0 and standard deviation 1.

Automatic Building of Regression Models

- Automatic techniques search for a first dependent variable that has a high explanatory power for the target and compute a regression coefficient.
- Iteratively further dependent variables are added until the quality of fit does no longer increase significantly.
- *Stepwise regression* is the most common algorithm (Effroymsen, 1960; Breaux, 1968).
- Statistics packages provide such automatic techniques.

Discussion: In every step a numerical score is evaluated. Several variables may have almost equal scores; thus the decision is somewhat arbitrary. The decision in one step may affect further steps.

Automatic Building of Regression Models

Automatically generated models are objective, may provide reasonable explanations for the dependency of the target variable from the dependent variables.

Drawbacks (compare to decision trees):

- Do not integrate user input. → Results may be partially not interesting or do not consider essential aspects, such as partitions in the data that are less reliable, special effects such as the Arabic oil embargo in the 1970s on economic factors are ignored.
- Since results strongly depend on outliers (that are hard to determine automatically!), the results may be misleading.
- Automatic regression models tend to be unnecessarily complex and may suffer from overfitting.
- User trust is limited.

→ „computer-assisted data analysis“ is more favorable (Henderson, 1981):

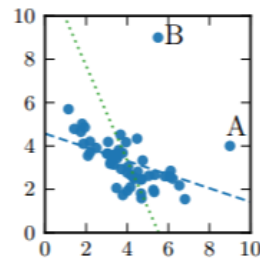
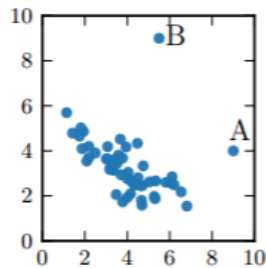
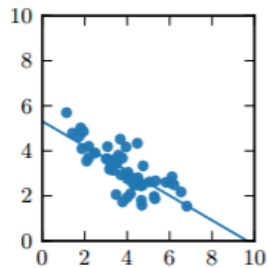
Outliers and influential points:

- Computing the *leverage* of individual points supports the decision to exclude such points from model building.
- The leverage h_i for point i is between 0 and 1. Typical values depend on the sample size n and the number of variables in the regression model p .
- If $h_i > 2p/n$ the corresponding point should be checked.

Interactive Regression Model Building

Outliers and influential points:

- Outliers should be checked w.r.t. influence on the regression slope. They may or may not follow the trend.
- Outliers are often influential if they deviate strongly from the mean of the x-values. The computation of the *leverage effect* considers this.
- Remember: outliers are not necessarily wrong values; they may be caused by (too) low sampling.



A is an influential point whereas B is not (From: [Link](#))

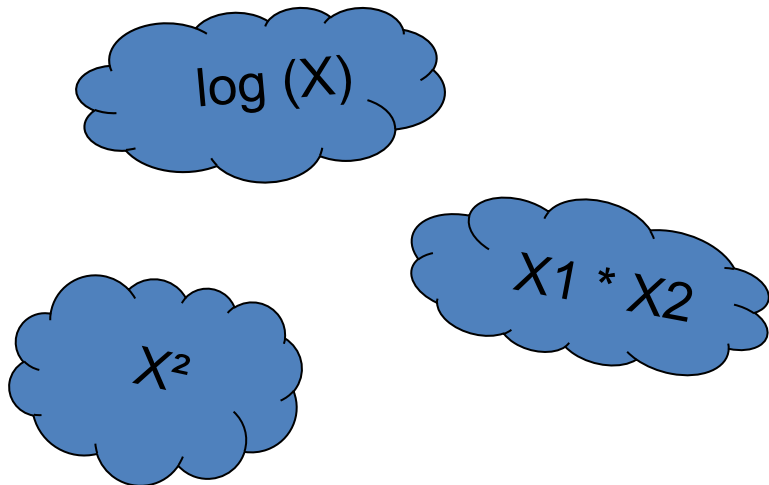
Outliers and influential points:

- An alternative is to compute the prediction with and without a potential outlier and to evaluate the distance (Cook's distance) (Sridharan, 2015)
- Robust regression analysis avoids the need for manual checking of outliers.
 - Basic goal: replace the least square optimization where the penalty for outliers is very large.
 - Most typical approach: RANSAC: Choose a subset of points and compute a model, validate it and record the error. Try other subsets and validate again. Choose the best subset (probably the best subset has no/very few outliers).
 - Drawback: RANSAC has many parameters (which subsets? How they are chosen? When to terminate?) that need to be chosen wisely.

- Various visualizations support user input in model building:
 - Quantile-quantile plots
 - Partial residual plots
 - Partial regression plots
 - Partition-based visualization
- Often, interactive regression model building reveals that derived quantities need to be considered, e.g. horse power/weight in cars to predict oil consumption.

Preparation for Interactive Regression Models:

- Check univariate distributions of explanatory variables. Ideally, they are symmetric, close to a Gaussian.
- If they are not → check whether log-transform, or square-root transform, multiplication of explanatory variables improves the distribution.



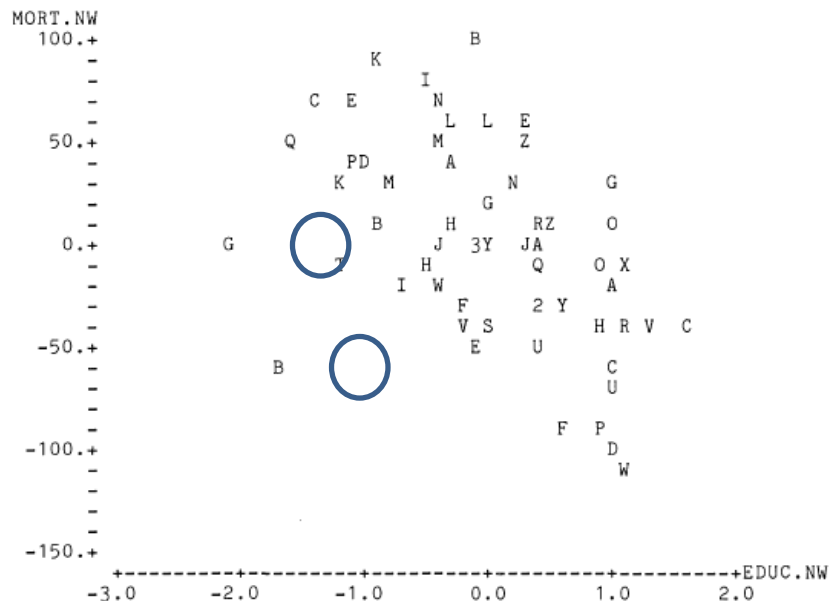
Feature Transformations

Preparation for Interactive Regression Models (cont'd):

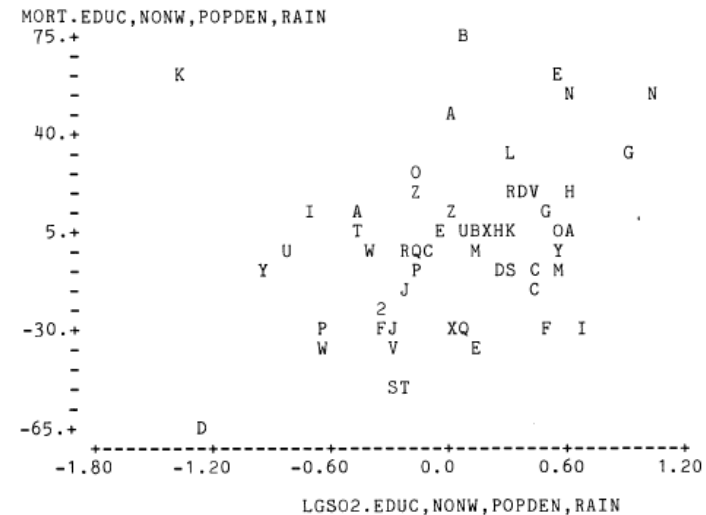
- Eventually replace original data with transformed (re-expressed) data.
- **Example:** Prediction of life expectancy revealed that logarithmized air pollution parameters $\log\text{SO}_2$ and $\log\text{NO}$ are essential parts of a regression model (Henderson, 1981)
- Check for outliers whether they are influential
- Standardize explanatory variables w.r.t. range (0,1)

Summary: Feature selection and transformation as well as handling of outliers are challenging tasks in model building.

Interactive Regression Model Building



Partial regression plot of education level vs. mortality. Two outliers heavily influence the regression analysis (From Henderson, 1981).



Partial regression plot of SO2 vs. Mortality. Data is normalized for education, population density, rain level and proportion of non-white persons (From Henderson, 1981).

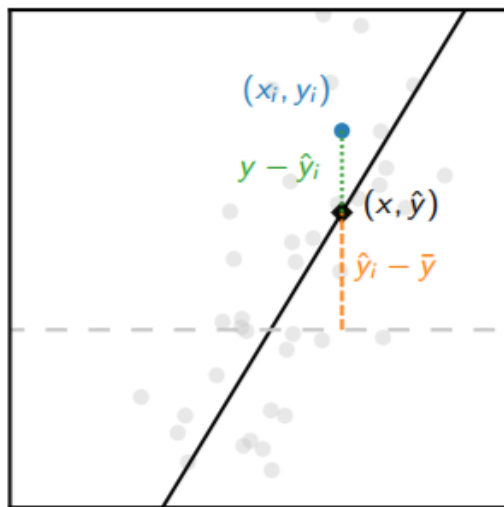
Residual Plots

Residuals: We want to understand how much the prediction of y_i depends on x_i

\hat{y}_i denotes the predicted value according to the model.

Residual is defined as $\hat{y}_i - y_i$

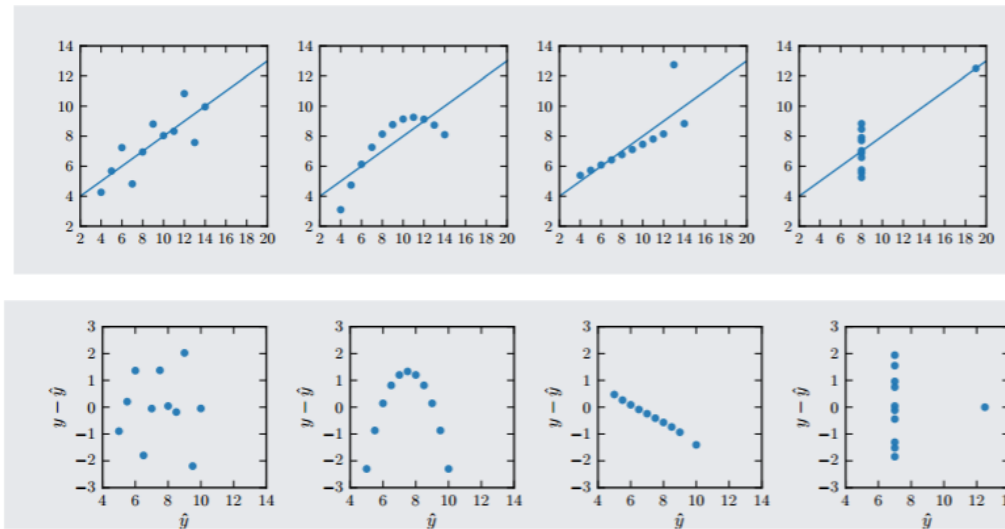
- If the model is good $\hat{y}_i - \bar{y} > \hat{y}_i - y_i$



From [Link](#)

Residual Plots

- Residual plots display the deviation of the data to the regression model – the resulting variation not explained by the model.
- The distribution of the points should look like random noise.
- Otherwise, a systematic error should be checked for, eventually leading to a refined regression model.



Distributions and residual plots (bottom row). The left residual plot shows random noise. The other three systematic errors of the regression model (From [Link](#))

- Residual plots indicate the overall quality of a model and whether there are striking non-linearities.
- Residual plots do not reveal which explanatory variable causes the non-linearity.
- Also outliers are hard to detect → partial residual plots serve these purposes better

Consider the multiple linear regression model

$y = \beta x + \gamma z + u$ and the least square fit

$$\hat{y} = \beta x + \gamma z$$

Regression of y on x and z is carried out in two steps:

- regress y and x on z and compute the residuals \hat{y} and \hat{z} .
- regress \hat{y} on \hat{z} .

The resulting errors correspond to those of the full regression.

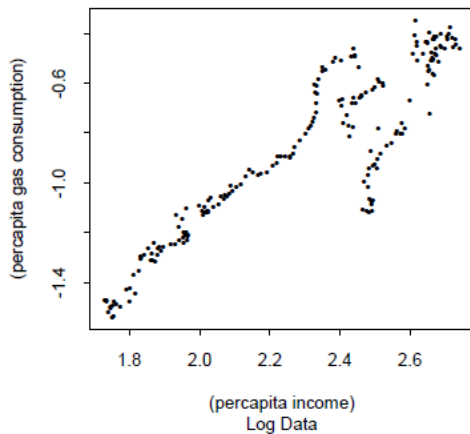
Partial residual plots (next slide) serve to check for problems and outliers, to identify influential points that need investigation.

Thus, the model is verified or refined and can then be used to “answer questions”.

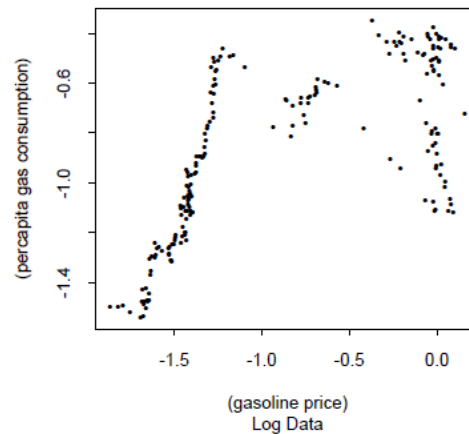
Partial Residual Plots

Consider the following data: gasoline demand, gas consumption and income.

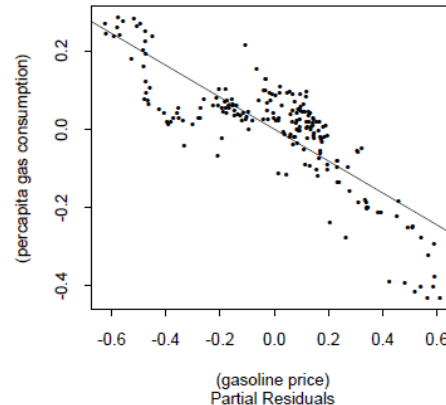
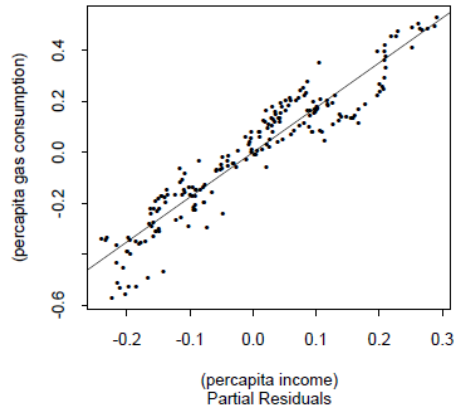
Income vs Gasoline Demand



Price vs Gasoline Demand



The scatterplots do not indicate strong correlations. However, the multiple regression fits quite well.



The partial residual plot of the multiple (linear) regression indicate a good fit (From: [Link](#)).

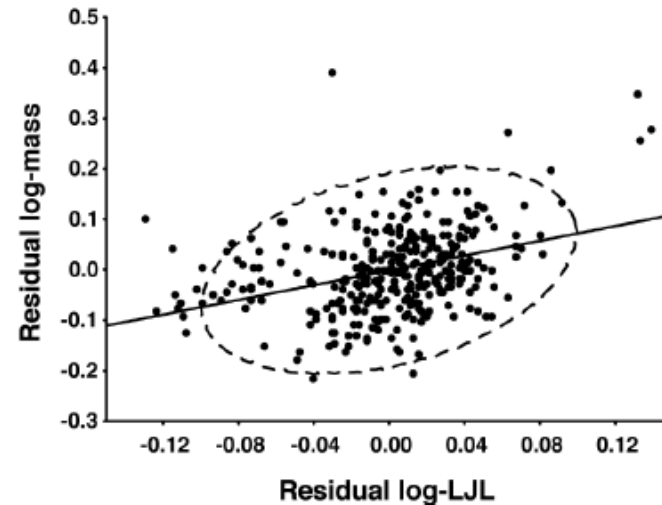
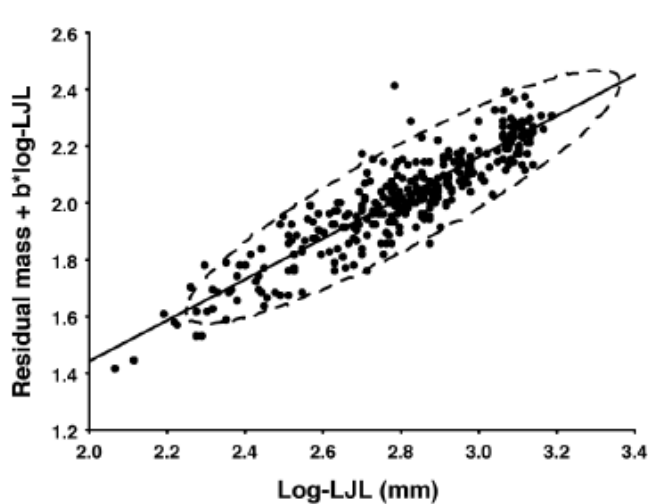
Partial residual plots (Moya, 2008)

- Do not take the variance of the remaining variables fully into account
 - underestimate the scatter of points around the line
 - amount of underestimation depends on the strength of the correlation among the independent variables

Comparison with partial regression plots:

- Non-linear patterns are better detected in partial residual plots
- Residual plots better indicate how the target changes as a consequence of changes in the explanatory variable since the original scale is maintained

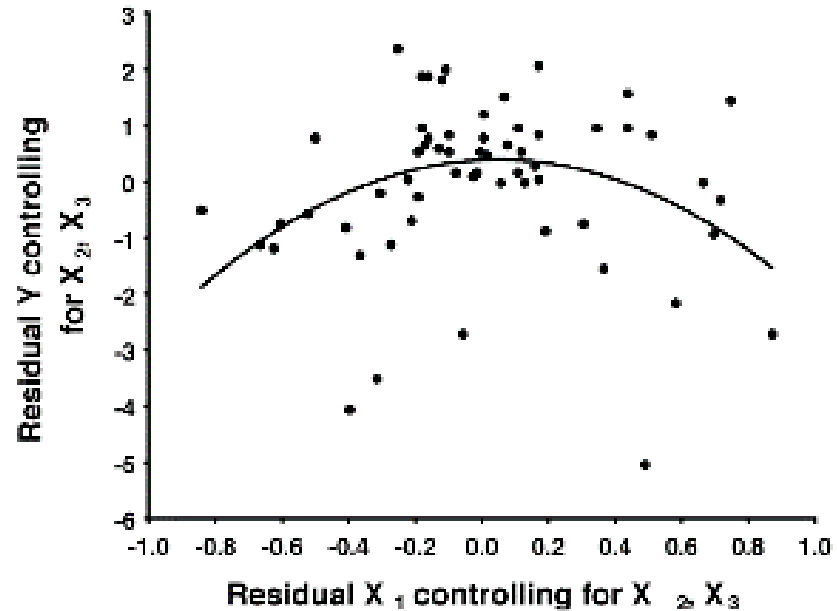
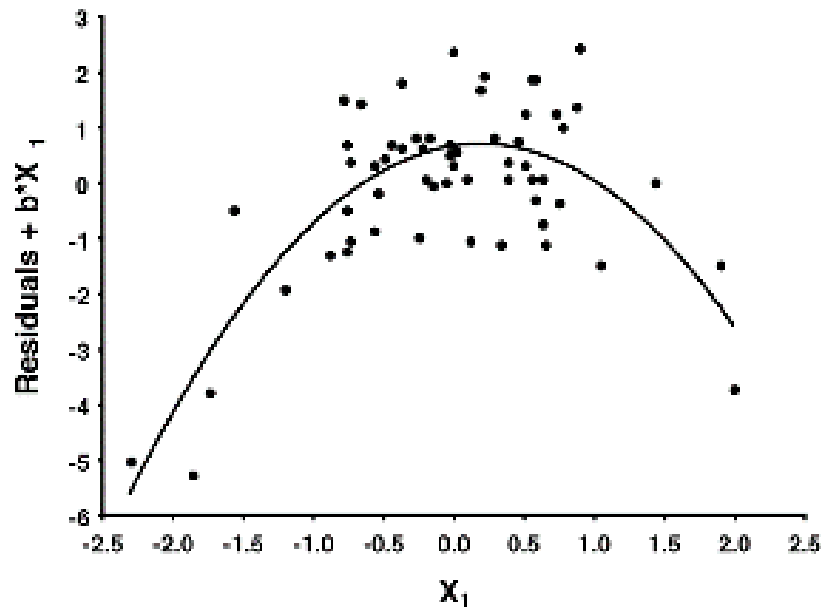
Partial Residual Plots



Comparison between a partial residual plot (left) and a partial regression plot. The 95% confidence ellipse illustrates the different amount of scatter.

Note, that the residual plot maintains the original scale while the points in the regression plot are zero-centered.

Partial Residual Plots



Comparison between a partial residual plot (left) and a partial regression plot for a simulated regression with one quadratic component. The left view indicates the quadratic relation better.

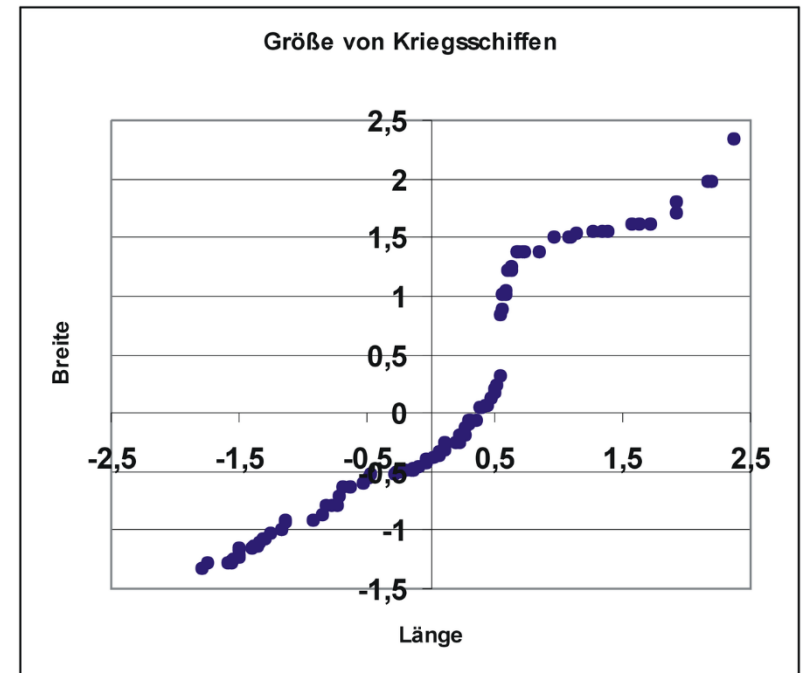
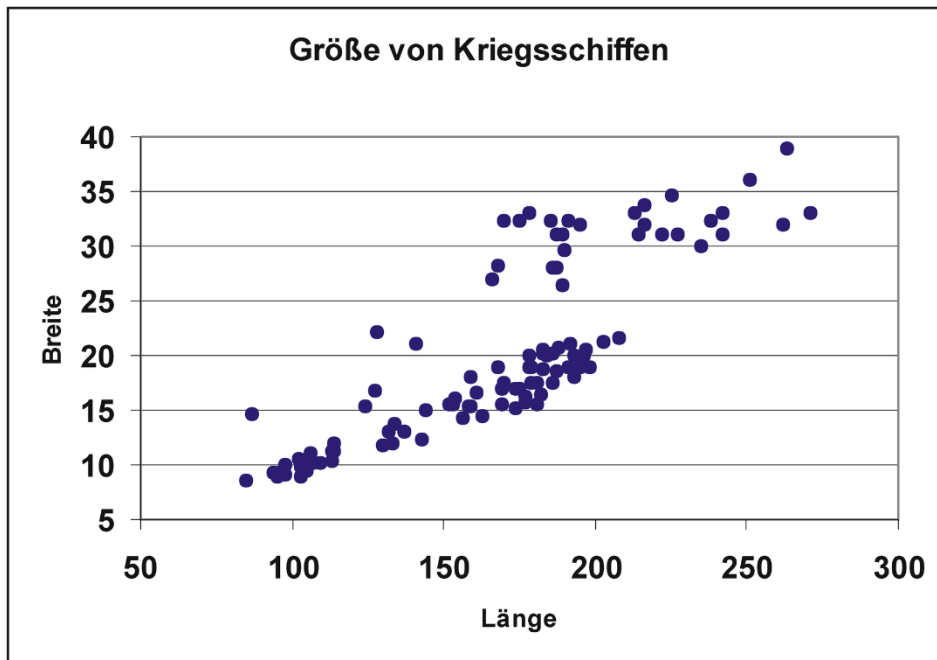
Quantile-Quantile Plots

- Quantile-quantile plots are a graphical method for comparing two distributions by plotting their quantiles.

Interpretation:

- If the two distributions are similar, the points in a Q-Q plot are very close to the diagonal line $y=x$.
- Points are always non-decreasing on x
- Experts can compare the distributions according to skewness, location and scale.
- Q-Q plots serve these purposes better than histograms or raw distributions ([Wikipedia](#))

Quantile-quantile Plots

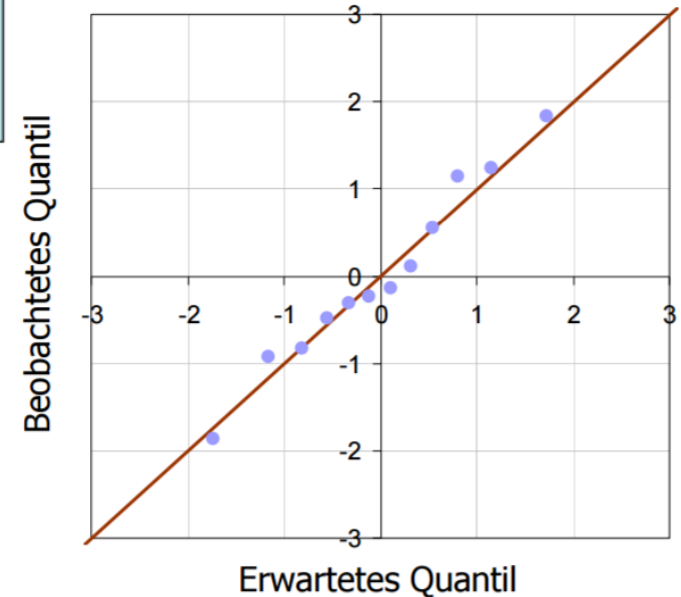


A scatterplot of width and length of 110 ships and a Q-Q plot (From [Wikipedia](#)). A gap between two clusters becomes obvious; also regions where the correlation is low. As a consequence, eventually two subsets should analyzed separately.

Quantile-Quantile Plots

Nr	Datum	Sortiert	z
1	44	25	-1.876
2	61	36	-0.927
3	48	37	-0.841
4	68	41	-0.496
5	25	43	-0.324
6	60	44	-0.237
7	41	45	-0.151
8	37	48	0.108
9	45	53	0.539
10	36	60	1.143
11	43	61	1.229
12	53	68	1.833

$$z = \frac{x - \bar{x}}{s}$$

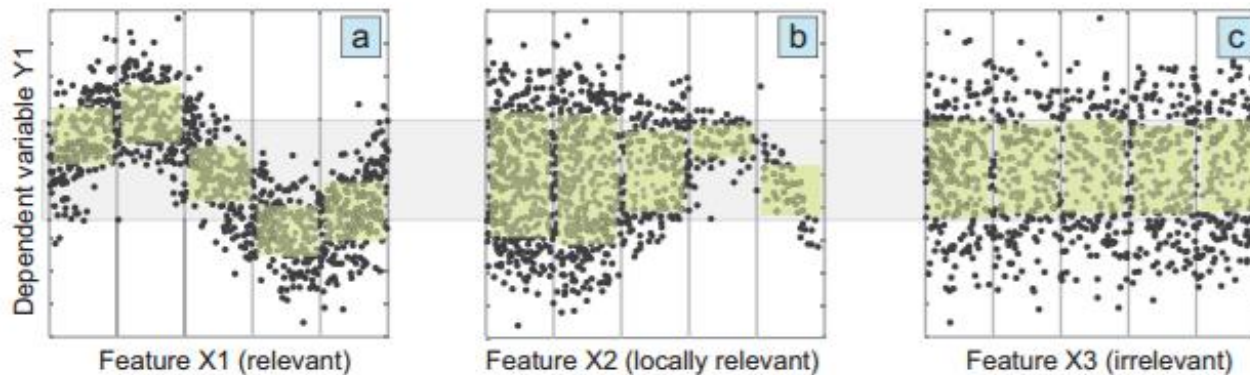


We sort the data, we standardize it (making it a zero means distribution) and compute z-values by dividing through the sample size s . z-values are transformed to quantile numbers (quantile: 0,04 for x means that 4% of the values are below x).

Perfectly correlated values would appear on the main diagonal (From: [Link](#)).

Partition-Based Regression Analysis

- Relations may be locally (strongly) different (remember local regression lines overlaid to scatterplots, Cleveland, 1984).
- Understanding which regions (subsets) should be analyzed separately should be supported.



The left figure represents a global correlation; the middle image a correlation that occurs only locally, whereas the right image shows no relevant correlation at all

(From: Mühlbacher, 2013).

Partition-Based Regression Analysis

$P(T | X_i[X_j])$: A partitioning of X_i and X_j w.r.t. target variable T

A partition is represented as an interval in 1D and an axis-aligned rectangle in 2D.

Partitioning of continuous variables may be performed

- *domain-uniform* (each partition has the same size) or
- *frequency-based* (each partition has the same number of elements)

Frequency-based partitioning has the advantage that all partitions have sufficient elements for identifying statistical significant regression models.

Partitioning may be performed automatically, e.g. with a termination criterion that relates to the minimum significance level.

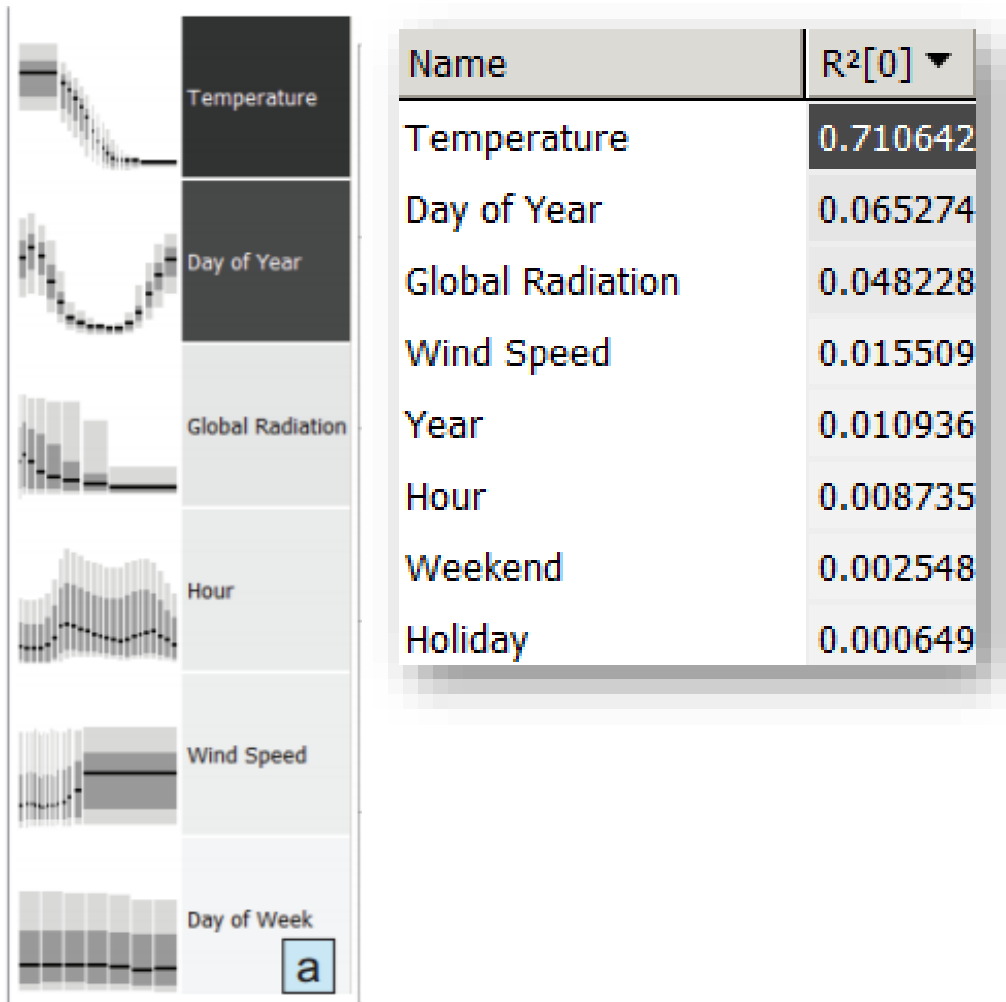
Requirements for visual analytics support (compare Piringer, 2010):

- Relate *known results* to regression models
- Support the *identification of subsets* with a different behavior
- Enable the *separate analysis* of subsets
- Support the analysis of the model behavior at selected points in the HD space
- Support a *comparison of different models* visually and according to regression coefficients, Aikaike, ...
- Support analysis of *large datasets* with many elements and variables and larger regression models
- Support the *validation* of regression models with independent test data
- Perform all actions in *real-time*.

Case Study 1: Analysis of Energy Consumption

- The demand for natural gas is influenced by various factors, including temperature and potentially other weather factors, day of the week and daily hours.
- A precise understanding of the mutual influences of calendric and meteorologic aspects is essential to accurately predict gas consumption and thus to ensure supply and enable cost minimization (Mühlbacher, 2013).
- Partition-based visualization is used, i.e. continuous variables are binned, to identify local correlations (recall variable binned scatterplots)
- The influence of individual features and pairs of features on the target variable (energy consumption) is displayed first (similar to the rank-by-feature framework, which does **not** support regression-based tasks)

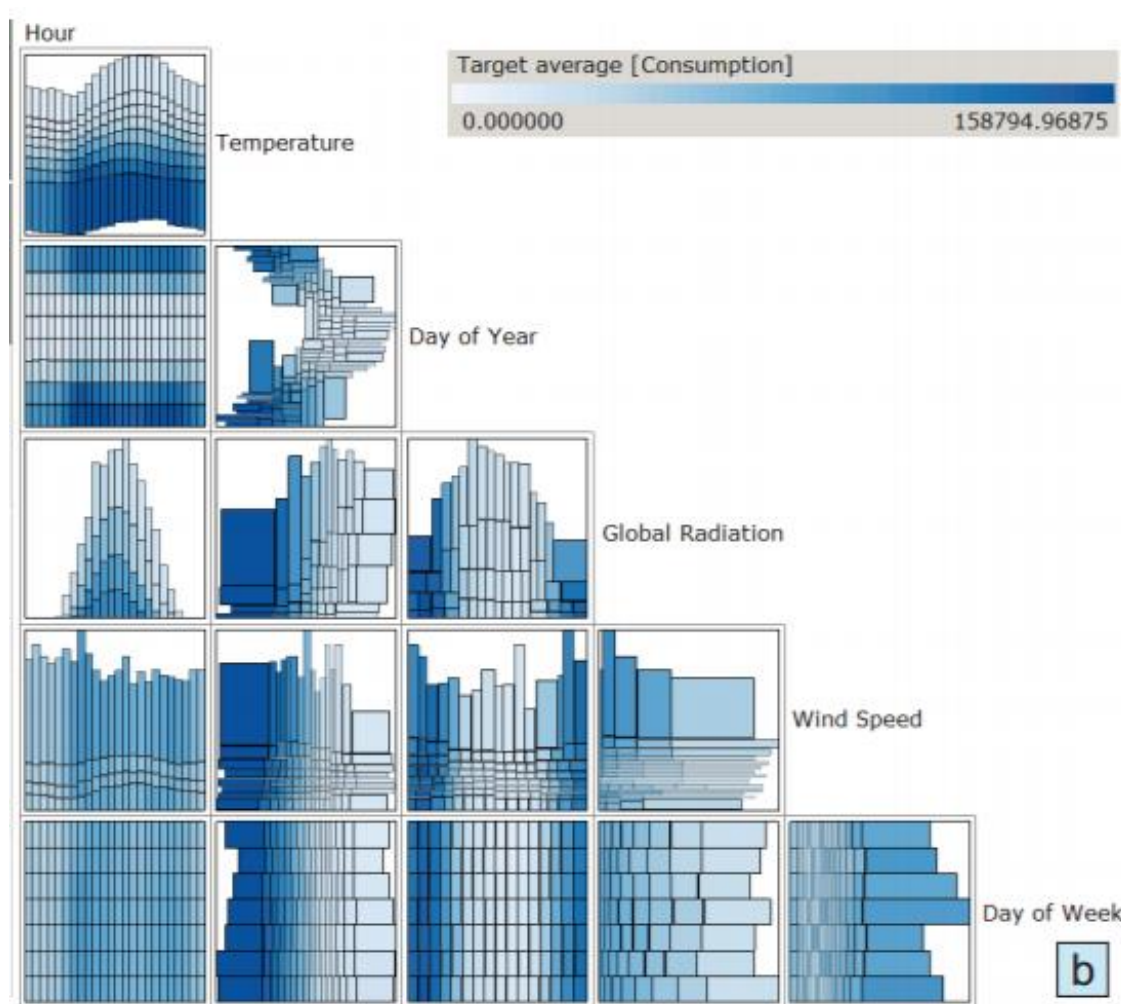
Case Study 1: Analysis of Energy Consumption



Influence of single variables on consumption

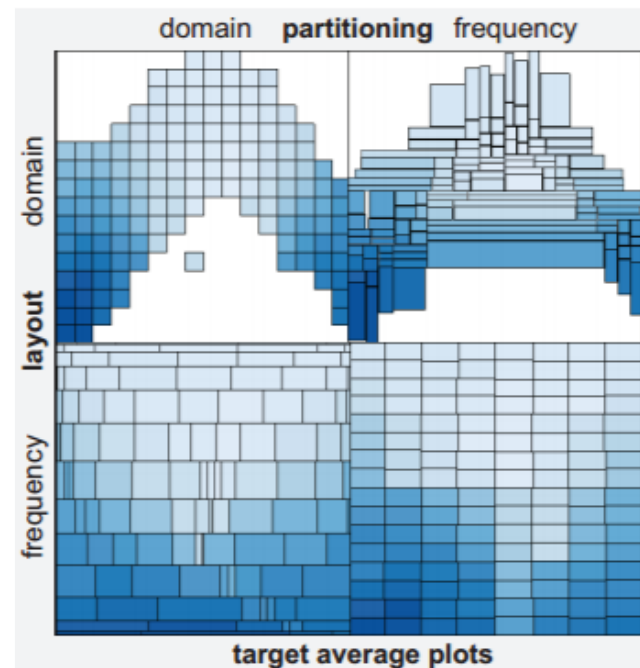
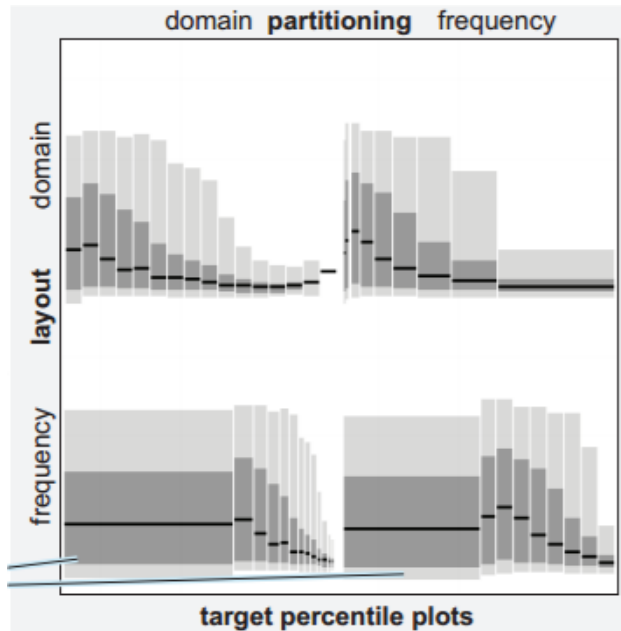
(From: Mühlbacher, 2013).

Case Study 1: Analysis of Energy Consumption



Influence of pairs of variables on consumption (From: Mühlbacher, 2013).

Case Study 1: Analysis of Energy Consumption



(From: Mühlbacher, 2013).

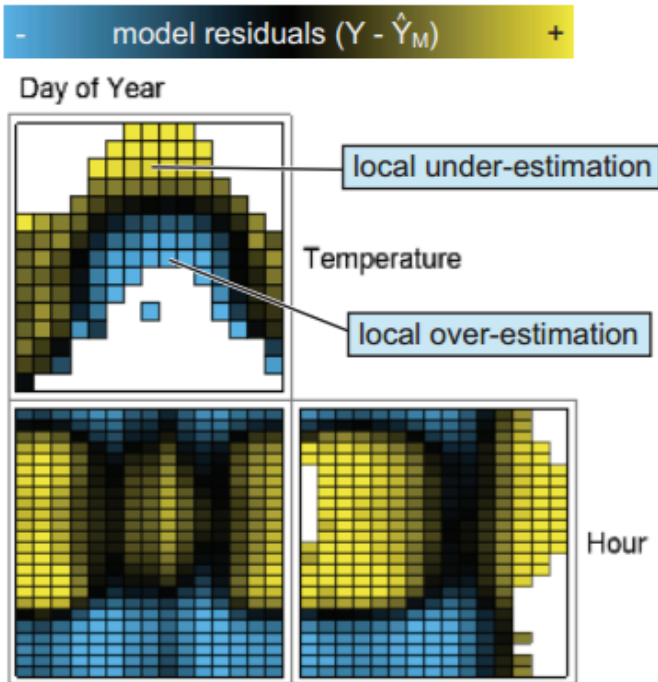
Both the target variable and the independent variable (radiation) may be partitioned based on the *domain* or *frequency*.

When two independent variables are partitioned; one frequency-based and the other domain-based, resulting visualizations are difficult to interpret.

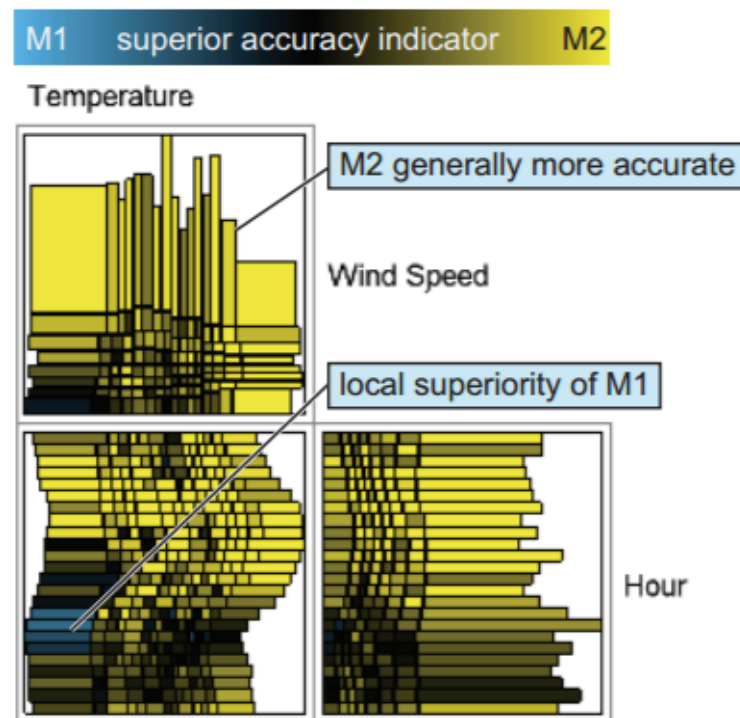
Color encodes the target variable.

Case Study 1: Analysis of Energy Consumption

a) Analyzing prediction bias



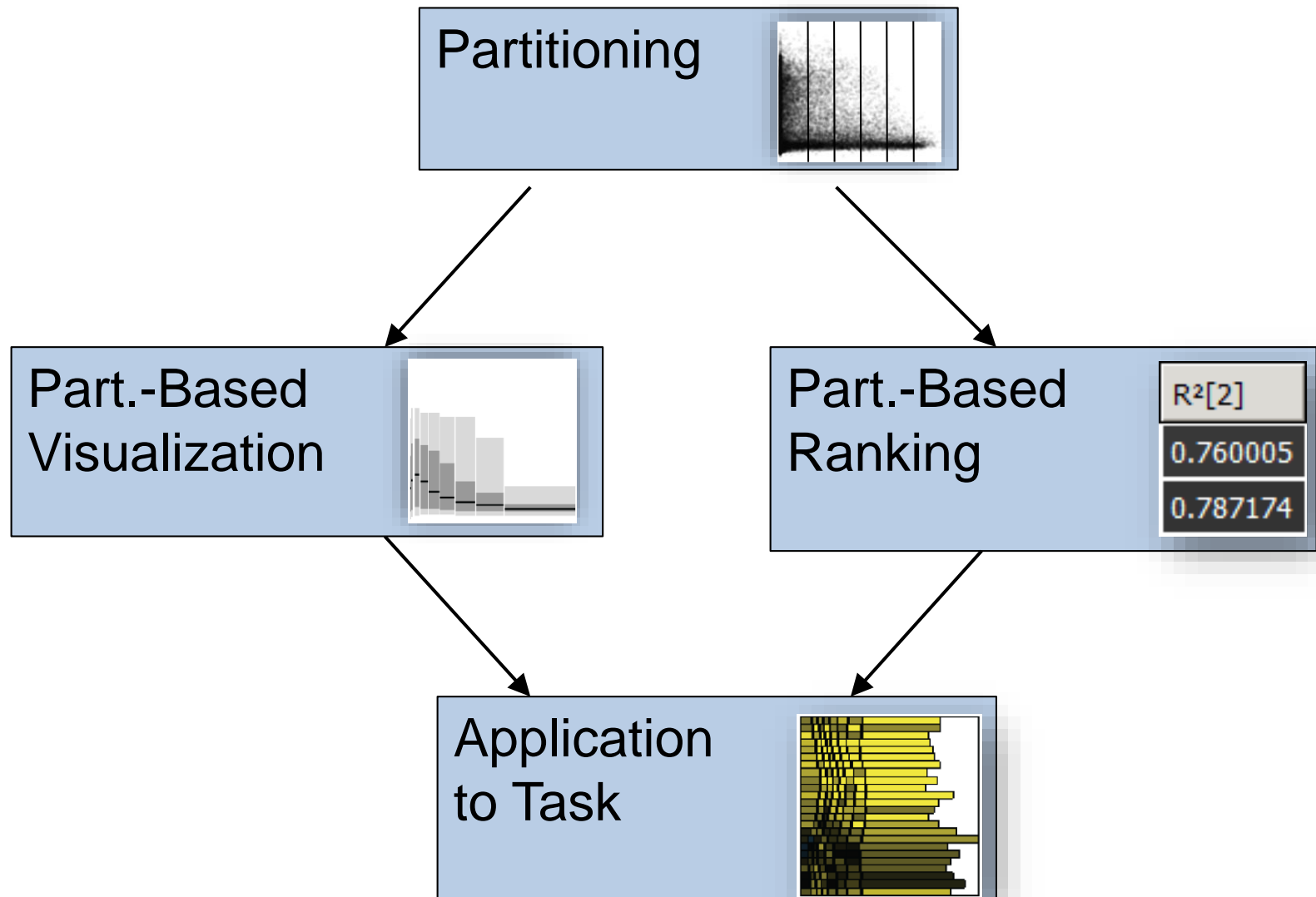
b) Comparing two point-wise predictions



(From: Mühlbacher, 2013).

Regression models may be computed automatically, analyzed and compared. Left: A partition-based partial residual plot reveals local deviations to the model. Right: The deviations of two Modes M_1 and M_2 are compared and displayed.

Case Study 1: Analysis of Energy Consumption



Overall modular architecture (Courtesy of T. Mühlbacher/H. Piringer, VR Vis Vienna)

Case Study 1: Analysis of Energy Consumption

- In addition to analyzing raw features, interaction effects ($X_1 * X_2$) may be analyzed and transformations (log, square, ...) are supported.
- Initial models M_1 may be refined by adding features – guided by the display of their relation to the target.
- User-defined partitioning is supported as well.
- Outliers may be analyzed w.r.t. their influence and removed.
- *Temperature* and have the strongest effects.

Case Study 2: Regression Analysis in Epidemiology

- Consider epidemiology data with rich information potentially related to a target variable that represents a disease
- Enable flexible input of regression formula (including variables) and compute quality of fit for linear regression simultaneously.
 - Epidemiologists have a strong background in statistics and are familiar with regression formula
- Integrate a feature selection step to reduce the dimensions to a manageable subset.
- Consider previous knowledge, e.g. for some diseases, woman before and after menopause (~52 years) should be analyzed separately, age and gender are general confounders for most diseases.

Case Study 2: Regression Analysis in Epidemiology

Regression formula:

- $Z \sim X + Y$, compute all regression coefficients for any possible target Z depending on any pair of X and Y features
- $Z \sim X + Y - \text{age}$, do the same but remove the influence of age, since age is a known confounder
- $\text{Cancer} \sim X + Y$, search for all pairs of features that may explain the occurrence of a cancer disease
- $Z \sim X * Y$, search for all pairs of features that may explain Z considering interactions between the two features

Thus, a text editor is provided, along with a list of possible regression operands and syntax checking.

Case Study 2: Regression Analysis in Epidemiology

Computational complexity: for the breast-related dataset, we have 231 variables. To compute

$$\textit{Breast_fat} \sim X + Y + Z,$$

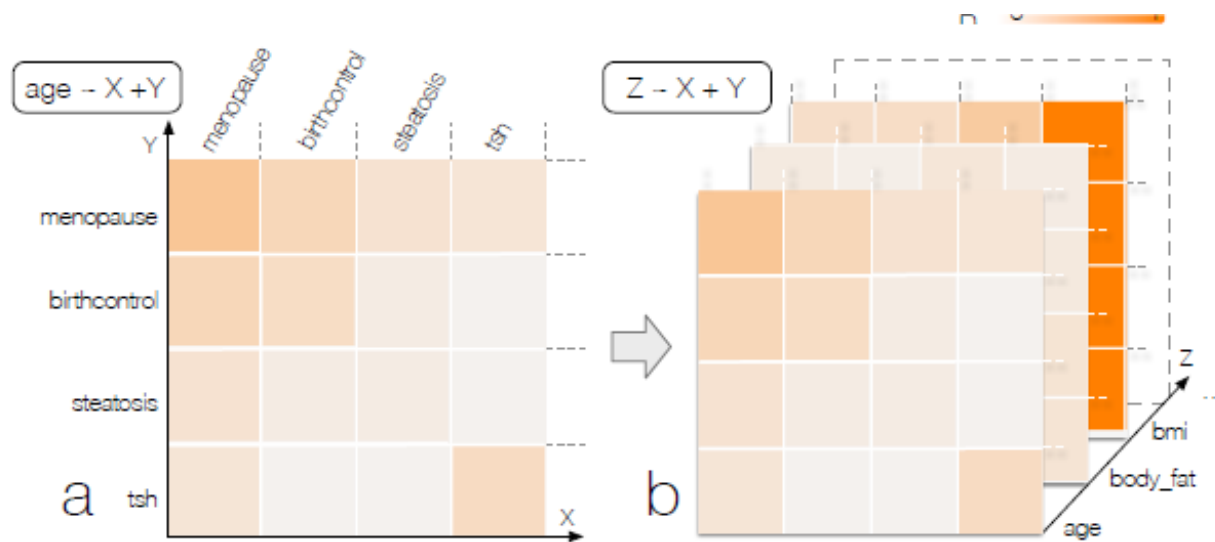
We have to compute 231^3 R^2 values need to be computed involving ~ 1.200 datasets.

Feature selection (CFS):

Based on a merit function and the *correlation-based feature selection* algorithm, the feature set is reduced based on entropy related to the target variable. With an appropriate threshold, feature sets shrink to 10 ... 30 variables.

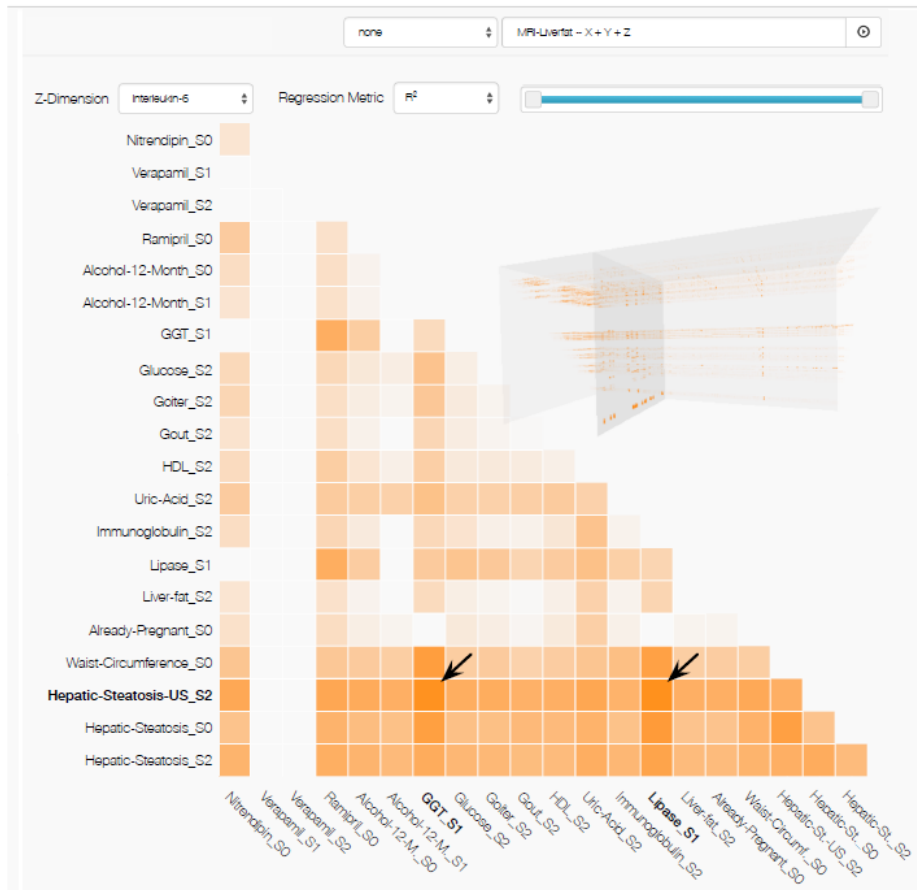
Remember: Such feature selection is risky w.r.t. loosing important information.

Case Study 2: Regression Analysis in Epidemiology



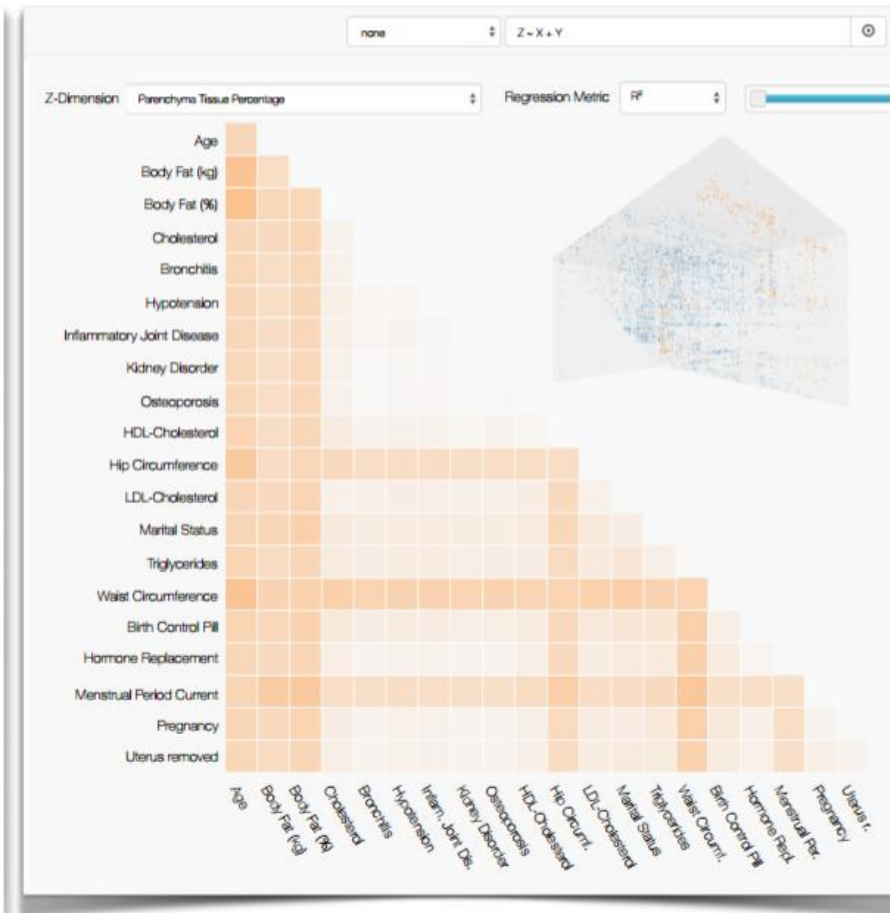
The R^2 value for each regression is mapped to the saturation of color. Left: age as target variable. Right: all variables are considered as target variables, leading to regression cube. Users may slice through the cube, e.g. with the mouse wheel (From: Klemm, 2016)

Case Study 2: Regression Analysis in Epidemiology



A dataset related to hepatic steatosis (fatty liver) is analyzed. It comprises 199 features and is analyzed for man, woman below 52 years and above 52 years. Some enzymes, uric acid and Interleukin-6 levels are found as essential independent variables (From: Klemm, 2016).

Case Study 2: Regression Analysis in Epidemiology



Analysis 2: Breast Fat Data

A dataset with 231 features possibly related to increased breast density was analyzed. Pregnancy status, menstruation period, body fat portion, hip and waist circumference are essential to predict increased breast density.

Surprising fact: strong correlation with kidney disorder. However, only 8 woman are affected. The correlation is not significant (From: Klemm, 2016).

Case Study 2: Regression Analysis in Epidemiology

Discussion so far:

Users welcome fast target selection, simultaneous computation of many regression models, flexibility of input.

However, analysis is restricted to multiple linear regression models.

There are many known, distinctively non-linear relations.

Case Study 2: Regression Analysis in Epidemiology

Workflow for integrating regression (inspired by Seo/Shneiderman):

- Select a target variable *targ*, e.g. a pathology, such as fatty liver
- For all pairs of attributes a_1 , a_2 check how they are related to *targ*. Use a_1 as primary variable, a_2 as *moderator variable* (partition in low, middle, high values).
 - Better splits are available, as used in decision trees, but their computation is time-consuming and groups with few members may arise
- If a_1 has a strong linear correlation to *targ* (use a threshold) emphasize it and overlay fitting lines according to the moderator variable

Case Study 2: Regression Analysis in Epidemiology

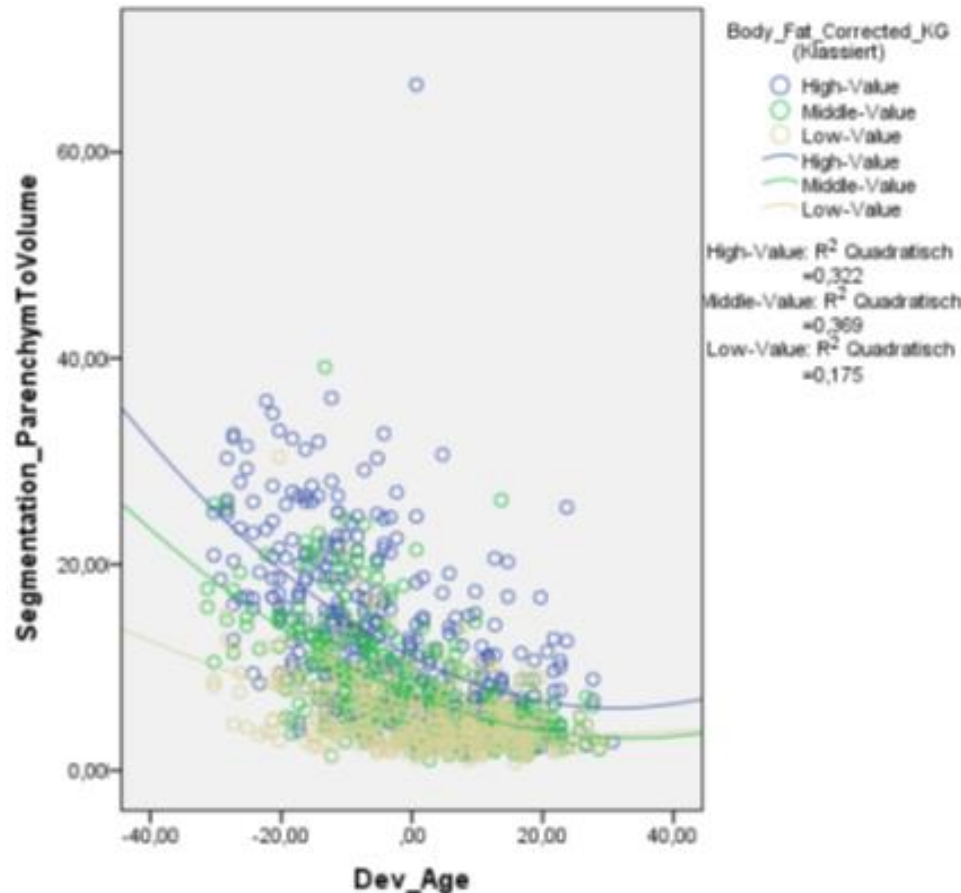
Workflow for integrating regression (continued):

- If a_1 has a low linear correlation to *targ*, check for quadratic correlations to *targ*
- Check for multiple regressions: If a_1 has a strong correlation (linear or quadratic) search for attributes a_2 as moderator variable (only if a_1 and a_2 do not correlate)
 - For multiple regressions the R^2 values have to be corrected
- In case of strong quadratic correlation, fit quadratic curves (capture U-shaped and J-shaped distributions)

Details:

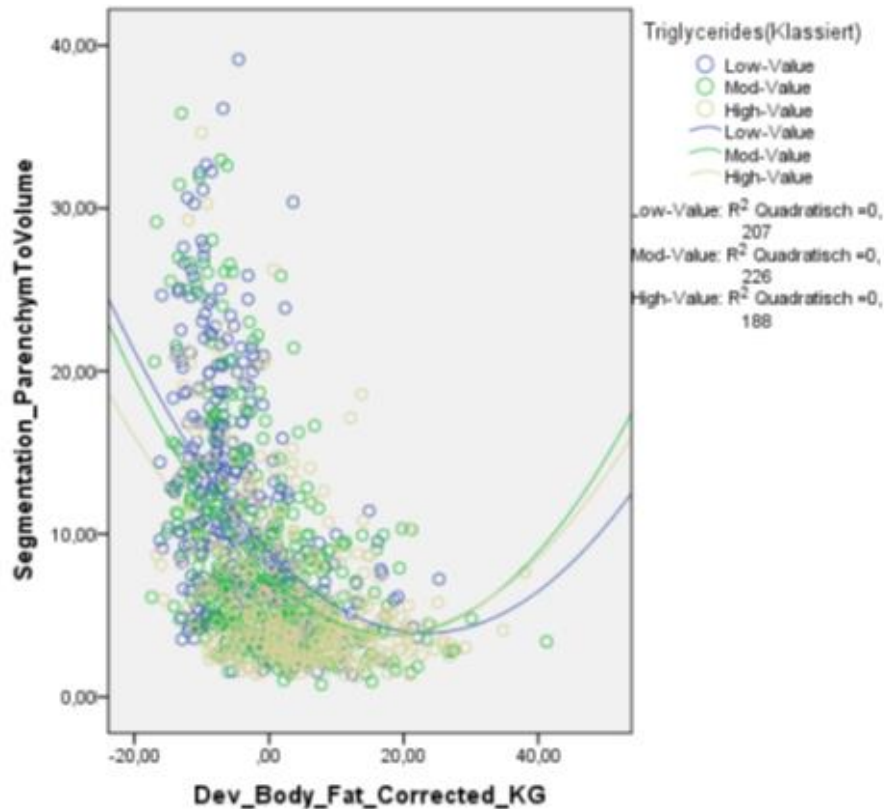
All variables are centralized first and the variance influence factor (threshold 5) is used to check for multicollinearity.

Case Study 2: Regression Analysis in Epidemiology



Age and body fat correlate to increased breast density, a feature with increased risk for breast cancer. For medium and high values of body fat, the correlation is relevant (green and blue lines).
SPSS-Output

Case Study 2: Regression Analysis in Epidemiology



Also the interaction of *triglycerides* and *body fat* with *increased breast density* has a significant quadratic correlation. Influence of body fat is lower, the higher the triglyceride value.

Limitation: Only numeric values considered.

Validating Regression Models

- Validation is essential to transform findings to trusted knowledge.
- A regression model learnt from some training data should be applied to some test data to study its performance.
- The overall data should be split in training data set to learn the model and a test dataset to apply it (like other machine learning techniques, such as classifiers).
- In epidemiology, people start to apply their findings to an independent cohort.

Summary

- Before regression modeling, an understanding of the data and their reliability is important.
- The distribution of each potential explanatory variable and the target should be checked.
- Appropriate features need to be selected, eventually transformed to yield normally distributed data and partitioned to account for local effects.
- Partial plots enable a validation of the influence of variables on the target and the model.
- Visual analytics solutions enable local model assessment, model comparison and analyzing the potential influence of outliers.

Further Information

- Partition-Based Regressing Models [Video](#)
- Application to the Prediction of Gas Consumption [Video](#)
- Predictive Multiple Regression in Excel [Video](#)

References

- Harold J. Breaux. 1968. A modification of Efroymson's technique for stepwise regression analysis. *Commun. ACM* 11, 8 (August 1968), 556-558
- M. A. Efroymson, "Multiple Regression Analysis," In: A. Ralston and H. S. Wilf, Eds., *Mathematical Methods for Digital Computers*, John Wiley, New York, 1960
- H. Henderson and P. Velleman (1981) „Building multiple regression models interactively“. *Biometrics*, 37(2): 391–411.
- Paul Klemm, Kai Lawonn, Sylvia Glaßer, Uli Niemann, Katrin Hegenscheid, Henry Völzke, Bernhard Preim: „3D Regression Heat Map Analysis of Population Study Data“, *IEEE Trans. Vis. Comput. Graph.* 22(1): 81-90 (2016)
- J. Kuhn, A. Zirngibl, M. Wildner, W. H. Caselmann, G. Kerscher, „Regionale Sterblichkeitsunterschiede in Bayern“, *Gesundheitswesen* 2006; 68: 551-556
- N. Latzitis , L. Sundmacher , R. Busse. „Regional Differences in Life Expectancy in Germany at County Levels and their Possible Determinants“, *Gesundheitswesen*, 2011; 73(4): 217-228
- Thomas Mühlbacher, Harald Piringer: "A Partition-Based Framework for Building and Validating Regression Models", *IEEE Trans. Vis. Comput. Graph.* 19(12): 1962-1971 (2013)
- Harald Piringer, Wolfgang Berger, J. Krasser: „HyperMoVal: Interactive Visual Validation of Regression Models for Real-Time Simulation“, *Comput. Graph. Forum* 29(3): 983-992 (2010)
- [Ramesh Sridharan.](http://www.mit.edu/~6.s085/) „Statistics for Research Projects“, Course Notes at MIT, 2015, available at: <http://www.mit.edu/~6.s085/>, accessed: 31/12/2016