Marcel Dettling

Institute for Data Analysis and Process Design

Zurich University of Applied Sciences

marcel.dettling@zhaw.ch

http://stat.ethz.ch/~dettling

ETH Zürich, October 08, 2012

Example: Airline Passengers

Each month, Zurich Airport publishes the number of air traffic movements and airline passengers. We study their relation.



Example: Airline Passengers

Month	Pax	АТМ	
2010-12	1'730'629	22'666	
2010-11	1'772'821	22'579	
2010-10	2'238'314	24'234	
2010-09	2'139'404	24'172	
2010-08	2'230'150	24'377	

۰. Рах Flugbewegungen

Flughafen Zürich: Pax vs. ATM

Smoothing

We may use an arbitrary smooth function $f(\cdot)$ for capturing the relation between Pax and ATM.

- It should fit well, but not follow the data too closely.
- The question is how the line/function are obtained.





Applied Statistical Regression AS 2012 – Week 03 Linear Modeling

A straight line represents the systematic relation between Pax and ATM.

- Only appropriate if the true relation is indeed a straight line
- The question is how the line/function are obtained.



Flughafen Zürich: Pax vs. ATM

Simple Linear Regression

The more air traffic movements, the more passengers there are. The relation seems to be linear, which is of course also the mathematically most simple way of describing the relation.

$$f(x) = \beta_o + \beta_1 x$$
, resp. $Pax = \beta_0 + \beta_1 \cdot ATM$

Name/meaning of the two $\beta_0 =$ "Intercept"parameters in the equation: $\beta_1 =$ "Slope"

Fitting a straight line into a 2-dimensional scatter plot is known as **simple linear regression**. This is because:

- there is just one single predictor variable ("simple").
- the relation is linear in the parameters ("linear").

Model, Data & Random Errors

No we are bringing the data into play. The regression line will not run through all the data points. Thus, there are random errors:

$$y_i = \beta_0 + \beta_1 x_i + E_i$$
, for all $i = 1, ..., n$

Meaning of variables/parameters:

- y_i is the response variable (Pax) of observation i.
- x_i is the predictor variable (ATM) of observation i.
- β_0, β_1 are the regression coefficients. They are unknown previously, and need to be estimated from the data.
- E_i is the residual or error, i.e. the random difference between observation and regression line.

Least Squares Fitting

→ <u>http://hspm.sph.sc.edu/courses/J716/demos/LeastSquares/LeastSquaresDemo.html</u>

Instructions for this demo are down below the graph.



We need to fit a straight line that fits the data well.

Many possible solutions exist, some are good, some are worse.

Our paradigm is to fit the line such that the squared errors are minimal.

Least Squares: Mathematics

The paradigm in verbatim...

Given a set of data points $(x_i, y_i)_{i=1,...,n}$, the goal is to fit the regression line such that the sum of squared differences between observed value y_i and regression line is minimal. The function

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta x_i))^2 = \min!$$

measures, how well the regression line, defined by β_0, β_1 , fits the data. The goal is to minimize this "quality function".

Solution: → see next slide...

Marcel Dettling, Zurich University of Applied Sciences

Solution Idea: Partial Derivatives

• We are taking partial derivatives on the function $Q(\beta_0, \beta_1)$ with respect to both arguments β_0 and β_1 . As we are after the minimum of the function, we set them to zero:

$$\frac{\partial Q}{\partial \beta_0} = 0 \text{ and } \frac{\partial Q}{\partial \beta_1} = 0$$

- This results in a linear equation system, which (here) has two unknowns β_0, β_1 , but also two equations. These are also known under the name *normal equations*.
- The solution for β_0 , β_1 can be written explicitly as a function of the data pairs $(x_i, y_i)_{i=1,...,n}$, see next slide...

Least Squares: Solution

According to the least squares paradigm, the best fitting regression line is, i.e. the optimal coefficients are:

$$\hat{\beta}_{1} = \frac{\sum_{i=1}^{n} (x_{i} - \overline{x})(y_{i} - \overline{y})}{\sum_{i=1}^{n} (x_{i} - \overline{x})^{2}} \quad \text{und} \quad \hat{\beta}_{0} = \overline{y} - \hat{\beta}_{1}\overline{x}$$

- For a given set of data points $(x_i, y_i)_{i=1,...,n}$, we can determine the solution with a pocket calculator (...or better, with R).
- The solution for our example Pax vs. ATM: $\hat{\beta}_1 = 138.8, \ \hat{\beta}_0 = -1'197'682$

$$\rightarrow$$
 lm(Pax ~ ATM, data=airpax)

Least Squares Regression Line



Pax vs. ATM

Is This a Good Model for Predicting the Pax Number from the ATM?

a) Beyond the range of observed data

Unknown, but most likely not...

b) Within the range of observed data

Yes, under the following conditions:

- the relation is in truth a straight line, i.e. $E[E_i] = 0$
- the scatter of the errors is constant, i.e. $Var(E_i) = \sigma^2$
- the data are uncorrelated (from a representative sample)
- the errors are approximately normally distributed

→ Fodder for thougt: 9/11, SARS, Eyjafjallajökull...?

Model Diagnostics

For assessing the quality of the regression line, we need to (at least roughly) check whether the assumptions are met: $E[E_i] = 0$ and $Var(E_i) = \sigma^2$ can be reviewed by:



Model Diagnostics

For assessing the quality of the regression line, we need to (at least roughly) check whether the assumptions are met: Gaussian distribution can be reviewed by:



Normal Q-Q Plot

We will revisit model diagnostics again later in this course, where it will be discussed more deeply.

"Residuals vs. Fitted" and the "Normal Plot" will always stay at the heart of model diagnostics.

Why Least Squares?

History...

Within a few years (1801, 1805), the method was developed independently by Gauss and Legendre. Both were after solving applied problems in astronomy...

Source: → <u>http://de.wikipedia.org/wiki/Methode_der_kleinsten_Quadrate</u>

	1	Beobschtz	ngen des zu	Palermo #,	1. Jen. 4801	von Prof. Pi	nti neu er	ideckten Gallins.
1\$0	1	Mittlere Sounen- Zeit	Gerado Aufitsig in Zeit	GeradeAuf Reigung in Gradon	Nordi Abweich.	Geocentri- iche Länge	Geocentr. Breite	Ort der Sonne Logar. + 20" d. Diftans Abstration 3
Jan.	1 2 3 4 10 11 13	8 43 17. 8 39 4. 8 34 53. 8 34 53. 8 30 41. 8 6 15. 8 6 15. 8 2 17. 7 54 26.	3 27 11, 25 3 26 53, 85 3 26 38, 45 13 26 23, 16 3 25 32, 15 3 25 32, 15 3 25 32, 15 3 25 30, 75 3 25 30, 75	5 51 47 48,8 51 43 27,8 51 39 36,0 51 35 47,3 51 28 1,5 351 28 26,0 51 22 34,5 51 22 34,5	15 37 45 5 15 41 5 5 15 44 31,6 15 47 57,6 16 10 32,0 16 22 49,5	1 23 22 58,3 1 23 19 44.3 1 23 19 44.3 1 23 14 15,5 1 23 17 59,1 1 23 10 37,6 1 23 10 37,6	3 6 42, 1 3 2 24, 9 2 58 9, 9 1 53 55, 6 2 29 0, 6 1 16 59, 7	Z 9 II I 30, 9 9, 9926156 9 I2 2 38, 6 9, 9926156 9 I3 3 26, 6 9, 9926317 9 I4 4 14, 9 9, 992618 9 20 I0 I7, 5 9, 9927641 9 23 I2 I3, 8 9, 9928490
Febr.	17 18 19 21 23 23 31 2 5 8 11	7 35 11, 7 35 11, 7 31 28, 7 24 2, 7 20 11, 7 16 43, 6 58 51, 6 51 52, 6 48 25, 6 44 59, 6 44 59, 6 41 35, 6 31 31, 6 31 31, 6 31 31, 6 31 58,	3 3: 25 55; 1 5 3 26 8; 12 7 3 26 34; 2 7 3 26 34; 2 7 3 26 34; 2 7 3 26 34; 2 3 28 54; 5 3 29 45; 1 3 3 28 54; 5 3 3 28 54; 5 3 3 28 54; 5 3 3 29 45; 1 3 3 0 47; 2 3 3 4 58; 5 3 3 1 59; 6 5 3 3 3 2; 7 5 3 3 3 2 3 3 4 58; 5 3 3 3 7; 5 4 58; 5 5 3 3 7; 5 5 3 5 5 3 5 5 3 5 5 3 5 5 3 5 5 5 5	151 23 45.0 51 23 45.0 51 23 2,3 51 23 2,3 51 41 21,3 51 41 21,3 51 41 21,3 52 13 38.3 52 41 48.9 53 41 48.9 53 41 48.9 53 45 45.9 53 45 45.9 53 45 37.5 53 45 37.5 53 45 37.5 54 43.7,5 54 43.8 54 37.5 54 38.3 55 40 38.1 55 40	16 40 13.0 16 40 14.1 16 58 35.9 17 3 18.5 17 8 5.5 17 43 11.0 17 43 11.0 17 53 36.3 18 55 17 58 57.5 18 15 J.0 18 31 23.2 18 37 25.8	1 23 25 59.2 1 23 34 21.3 1 23 39 1. 1 23 39 1. 1 23 39 1. 1 24 30 9. 1 24 30 9. 1 24 30 7. 1 24 40 19.3 1 24 46 19.3 1 24 54 57.9 1 25 53 29.5 1 26 56 20.0	1 53 38.2 1 53 38.2 1 46 6,0 1 42 28,1 1 38 52,1 1 13 52,1 1 14 16,0 1 15 54,6 1 15 54,5 1	9 29 19 53. 69, 9928309 9 29 19 53. 69, 9934309 10 1 20 40. 39, 9931836 10 3 12 0, 9, 9931836 10 3 12 12, 79, 993388 10 8 16 20, 19, 9933031 10 11 32 12, 79, 9933031 10 11 38 28. 59, 993703 10 15 39 49, 99, 993843 10 16 31 45, 59, 9940751 10 16 33 33, 31, 994373



Carl Friedrich Gauss



Why Least Squares?

Mathematics...

- Least Squares is simple in the sense that the solution is known in closed form as a function of $(x_i, y_i)_{i=1,...,n}$.
- The line runs through the center of gravity $(\overline{x}, \overline{y})$
- The sum of residuals adds up to zero: $\sum_{i=1}^{n} r_i = 0$
- Some deeper mathematical optimality can be shown when analyzing the large sample properties of the estimates $\hat{\beta}_0, \hat{\beta}_1$. This is especially true under the assumption of normally distributed errors E_i .

Gauss-Markov-Theorem

A mathematical optimality result for the Least Squares line

It only holds if the following conditions are met:

- the relation is in truth a straight line, i.e. $E[E_i] = 0$
- the scatter of the errors is constant, i.e. $Var(E_i) = \sigma^2$
- the errors are uncorrelated, i.e. $Cov(E_i, E_i) = 0$, if $i \neq j$

Not yet required:

- the errors are normally distributed: $E_i \sim N(0, \sigma_F^2)$

Gauss-Markov-Theorem:

- Least Squares yields the best linear unbiased estimates

Properties of the Least Square Estimates

Under the conditions above, the estimates are unbiased:

$$E[\hat{\beta}_0] = \beta_0$$
 and $E[\hat{\beta}_1] = \beta_1$

The variances of the estimates are as follows:

$$Var(\hat{\beta}_0) = \sigma_E^2 \cdot \left(\frac{1}{n} + \frac{\overline{x}}{\sum_{i=1}^n (x_i - \overline{x})^2}\right) \text{ and } Var(\hat{\beta}_1) = \frac{\sigma_E^2}{\sum_{i=1}^n (x_i - \overline{x})^2}$$

Precise estimates are obtained with:

- a large number of observations *n*
- a good scatter in the predictor x_i
- an informative/useful predictor, making σ_{E}^{2} small
- (an error distribution which is approximately Gaussian)

Benefits of Linear Regression

• Inference on the relation between y and x

The goal is to understand if and how strongly the response variable depends on the predictor. There are performance indicators as well as statistical test adressing the issue.

• Prediction of (future) observations

The regression line/equation can be employed to predict the PAX number for any given ATM value.

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

However, this mostly will not work well for extrapolation!

Applied Statistical Regression AS 2012 – Week 03 R^2 : The Coefficient of Determination

The coefficient of determination R^2 is also known as *multiple R-squared*. It tells which portion of the total variation is accounted for by the regression line.



Flughafen Zürich: Pax vs. ATM

Computation of R^2

 R^2 is the portion of the total variation that is explained through regression. It is determined as one minus the quotient of the yellow arrow divided by the blue arrow.

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} (\hat{y}_{i} - y_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \overline{y})^{2}} \in [0, 1]$$

The closer to 1 the value is, the tighter the datapoints are packed around the regression line. However, there are no formal criteria which R^2 value needs to be met such that the regression can be said to be useful/valid.

Confidence Interval for the Slope β_1

A 95%-CI for the slope β_1 tells which values (besides the point estimate $\hat{\beta}_1$) are plausible, too. The uncertainty is due to estimation/sampling effects.

95%-CI for
$$\beta_1 : \hat{\beta}_1 \pm qt_{0.975;n-2} \cdot \hat{\sigma}_{\hat{\beta}_1}$$
, resp.
 $\hat{\beta}_1 \pm qt_{0.975;n-2} \cdot \sqrt{\hat{\sigma}_E^2 / \sum_{i=1}^n (x_i - \overline{x})^2}$

Testing the Slope eta_1

There is a statistical hypothesis test which can be used to check whether the slope is significantly different from zero, or any other arbitrary value b. The null hypothesis is:

$$H_0: \beta_1 = 0$$
, bzw. $H_0: \beta_1 = b$

One usually tests two-sided on the 95%-level. The alternative is:

$$H_A: \beta_1 \neq 0$$
, bzw. $H_A: \beta_1 \neq b$

As a test statistic, we use:

$$T = \frac{\hat{\beta}_1}{\hat{\sigma}_{\hat{\beta}_1}}$$
, resp. $T = \frac{\hat{\beta}_1 - b}{\hat{\sigma}_{\hat{\beta}_1}}$, both follow a t_{n-2} distribution.

Reading R-Output

Residual standard error: 59700 on 22 degrees of freedom Multiple R-squared: 0.9487, Adjusted R-squared: 0.9464 F-statistic: 407.1 on 1 and 22 DF, p-value: 1.110e-15

\rightarrow Will be explained in detail on the blackboard!

Testing the Slope eta_1

Practical Example:

Use the Pax vs. ATM data and perform a statistical test for the null hypothesis H_0 : $\beta_1 = 150$. The information from the summary on slide 25 can be used as a basis. Then, also answer:

- a) Explain in colloquial language what was just tested. What is the benefit of this test? What claims could motivate the test?
- b) How does the testing result relate with the 95%-CI that we computed on slide 23? Would we be able to tell the test results from the CI alone?

\rightarrow See blackboard for the answers

Testing the Intercept β_0

An analogous test can be done for the intercept.

- No matter what the test result will be, the intercept should generally not be omitted from the regression model.
- The presence of the intercept protects against possible non-linearities and calibration errors of measurement devices. If it is kicked out of the model, the results are generally worse.
- If theory dictates that there should not be an intercept but it is still significant, take this as evidence that the linear relation does not hold when extrapolating to x = 0.

Prediction

Using the regression line, we can predict the *y*-value for any desired *x*-value. The result is the expectation for *y* given *x*.

$$E[y | x] = \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$
 a.k.a. "fitted value"

Example: With 24'000 air traffic movements, we expect

 $-1'197'682 + 24'000 \cdot 138.8 = 2'133'518$ Passengers

Be careful:

At best, interpolation within the range of observed x-values is trustworthy. Extrapolation with ATM values such as 50'000, 5'000 or even 0 usually produces completely useless results.

Prediction with R

We can use the regression fit object for prediction. The syntax for obtaining the fitted value(s) is as follows:

- > fit <- lm(Pax ~ ATM, data=unique2010)</pre>
- > dat <- data.frame(ATM=c(24000))</pre>
- > predict(fit, newdata=dat)
- 1 2132598

The x-values need to be provided in a data frame, where the variable/column name is identical to the predictor name.

Then, the predict() procedure is invoked with the regression fit and the new x-values as arguments.

Confidence Interval for E[y | x]

We just computed the fitted value $\hat{\beta}_0 + \hat{\beta}_1 x$, i.e. the expected number of passengers for 24'000 ATMs. This is not a deterministic value, but an estimate that is subject to variability.

A 95%-CI for the fitted value at position x is given by:

$$\hat{\beta}_{0} + \hat{\beta}_{1}x \pm qt_{0.975;n-2} \cdot \hat{\sigma}_{E} \cdot \sqrt{\frac{1}{n} + \frac{(x - \overline{x})^{2}}{\sum_{i=1}^{n} (x_{i} - \overline{x})^{2}}}$$

Prediction Interval for \hat{y}

The confidence interval for E[y | x] tells about the variability of the fitted value. It does not account for the scatter of the data points around the regression line and thus does not define a region where we have to expect the observed value. A 95% prediction interval at position is given by:

$$\hat{\beta}_{0} + \hat{\beta}_{1}x \pm qt_{0.975;n-2} \cdot \hat{\sigma}_{E} \cdot \sqrt{1 + \frac{1}{n} + \frac{(x - \overline{x})^{2}}{\sum_{i=1}^{n} (x_{i} - \overline{x})^{2}}}$$

Confidence and Prediction Interval



Pax vs. ATM with Confidence and Prediction Interval

Confidence and Prediction Interval

Note:

Visualizing the confidence and prediction intervals in R is not straightforward, but requires some tedious handwork.

R-Hints:

dat <- data.frame(xx=seq(..., ..., length=200))
pred <- predict(fit, newdata=dat, interval=...)
plot(..., ..., main="...")
lines(dat\$xx, pred[,2], col=...)
lines(dat\$xx, pred[,3], col=...)</pre>

Example: Braking Distance for Railway

Animation "Bremsweg Eisenbahn"

Zurück zur Übersicht Bremsweg Eisenbahn



Source: http://www.bav.admin.ch/aktuell/bremswege/index.html?lang=de

Braking Distance: Data



Braking Distance: Fitting a Straight Line



Braking Distance vs. Speed

Braking Distance: Residual Diagnostics



Braking Distance: Facts

Conclusion from the residual plots:

• A straight line does not correctly describe the relation braking distance vs. speed relation. Energy is proportional to the square of speed, thus we need a quadratic function. $BrDist_{i} = \beta_{0} + \beta_{1} \cdot Speed_{i}^{2} + E_{i}$

bzw. $Y_i = \beta_0 + \beta_1 \cdot x'_i + E_i$, wobei $x'_i = x_i^2 = Speed_i^2$

• The above model still is a simple linear regression problem. There is only one single predictor, the optimal coefficients $\hat{\beta}_0, \hat{\beta}_1$ can be explicitly determined by using the ordinary LS-algorithm.

Braking Distance vs. Speed^2

> fit02 <- lm(weg ~ I(speed^2))</pre>

Braking Distance vs. Speed^2



Braking Distance: Fitting a Parabola

> fit02 <- lm(weg ~ I(speed^2))</pre>



Braking Distance vs. Speed

Curvilinear Regression

We have seen that simple linear regression offers more than just fitting straight lines. Using the ordinary LS-algorithm, we can fit any curvilinear relation, e.g.:

• $Y_i = \beta_0 + \beta_1 \cdot \ln(x_i) + E_i$

•
$$Y_i = \beta_0 + \beta_1 \cdot \sqrt{x} + E_i$$

•
$$Y_i = \beta_0 + \beta_1 \cdot x^{-1} + E_i$$

Note that the predictor gets transformed, i.e. $x'_i = \ln(x_i)$, $x'_i = \sqrt{x_i}$ or $x'_i = (x_i)^{-1}$. Using this notation, it is obvious that these are simple linear regressions.

→ BUT... see next slide!

Marcel Dettling, Zurich University of Applied Sciences

Braking Distance: Some Thoughts

Curvilinear models rarely yield a good fit:

• For modelling braking distance, reaction time should also be accounted for. This yields a multiple regression:

$$BrDist_{i} = \beta_{0} + \beta_{1} \cdot Speed_{i} + \beta_{2} \cdot Speed_{i}^{2} + E_{i}$$

- The scatter of the errors is not constant: the higher the speed, the higher the variation in braking distance.
- There are many applications where the exponent is unknown and needs to be estimated from data:

$$Y_i = \beta_0 + \beta_1 \cdot x^{\beta_2} + E_i$$

Infant Mortality vs. Per-Capita Income



Infant Mortality vs. Per-Capita Income

Infant Mortality vs. Per-Capita Income



Infant Mortality vs. Per-Capita Income

Residual Analysis of Hyperbolic Fit



Residuals vs. Predictor

Conclusions from Residual Analysis

The problem with the previous fit is that the the **power** x^{-1} is **not correct**. Adjusting it by hand is very laborious. If we try a model such as

 $Y_i = \beta_0 + \beta_1 \cdot x_i^{\beta_2} + E_i$

we do **no longer** have a **linear** problem in the parameters. Thus, least squares fitting cannot be applied for such models.

Yet there is a simple, but very powerful trick that often helps:

$$y'_i = \log(y_i), \ x'_i = \log(x_i)$$

solves the problem, see next slide and the blackboard...

Log-Transformation Helps!



log(infant) vs. log(income)

Model and Coefficients

If a straight line is fitted on the log-log-scale,

$$y'_i = \beta'_0 + \beta_1 \cdot x'_i + E'_i$$
, where $y'_i = \log(y_i)$, $x'_i = \log(x_i)$,

this means fitting the following relation on the original scale:

$$y_i = \beta_0 \cdot x_i^{\beta_1} \cdot E_i$$

The meaning of the parameter β_1 is as follows:

If x, i.e. the income increases by 1%, then y, i.e. the mortality decreases by $\hat{\beta}_1 = 0.56\%$. In other words: β_1 characterizes the relative change in y per unit of relative change in x.

Fitting on the Original Scale

Infant Mortality vs. Per-Capita Income



Fitted Values and Intervals

• To obtain a prediction on the original scale, we can just re-exponentiate to invert the log-transformation.

 $\hat{y} = \exp(\hat{y}')$

• **Be careful**: that result is an estimate for the median, but not for the mean E[y | x]. If we require *unbiased estimation*, we need to use a correction factor:

 $\hat{y} = \exp(\hat{y}' + \hat{\sigma}_E^2 / 2)$

• The confidence and prediction intervals are simply:

 $[l,u] \rightarrow [\exp(l), \exp(u)]$

Mean and Median



Infant Mortality vs. Per-Capita Income

Confidence and Prediction Interval



What to Do if y=0 and/or x=0?

- We can only take logarithms if x, y > 0. In cases where the response and/or predictor takes negative values, we should not log-transform. If zero's occur, they need treatment.
- What do we do with either x = 0 or y = 0?
 - do never exclude such data points!
 - adding a constant value is allowed!
- What about the choice of the constant?
 - standard choice: c = 1
 - scale dependent, thus not recommended!

\rightarrow Set c = smallest value > 0!

Another Example: Daily Cost in Rehab



Logged Response Model

We *transform* the *response* variable and try to explain it using a linear model with our previous predictors:

$$Y' = \log(Y) = \beta_0 + \beta_1 x + E$$

In the *original scale*, we can write the logged response model using the same predictors:

$$Y = \exp(\beta_0 + \beta_1 x) \cdot \exp(E)$$

→ Multiplicative model

→ $E \sim N(0, \sigma_E^2)$, and thus, $\exp(E)$ has a lognormal distribution

Also This Transformation Works!



Fitted Values on the Original Scale



Daily Cost in Rehabilitation vs. ADL-Score

Interpretation of the Coefficients

Important: There is no back transformation for the coefficients to the original scale, but still a good interpretation

$$log(\hat{y}) = \hat{\beta}_0 + \hat{\beta}_1 x_1$$
$$\hat{y} = exp(\hat{\beta}_0) exp(\hat{\beta}_1 x_1)$$

An increase by one unit in x_1 would multiply the fitted value in the original scale with $\exp(\hat{\beta}_1)$.

→ Coefficients are interpreted multiplicatively!