# Regression

Flervariabel & Polynom

# Agenda

- Linjär regression med flera variabler
- Kodstruktur
- Kodexempel
- Polynomregression
- Kodexempel
- Statistik
  - p-value
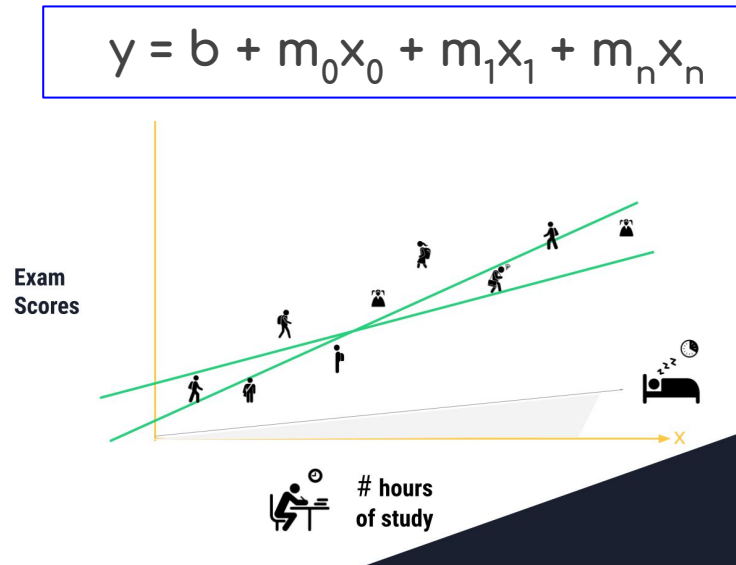  - Multivariable regression - variable selection
  - $R_2$
- Övningar

**RESTOR⟩...**

# Linjär Regression med flera variabler

RESTORY...

# Flervariabel Linjär Regression

$$y = b + m_0x_0 + m_1x_1 + m_nx_n$$

y är nu beroende av två, eller fler, variabler
(years of experience AND years of education)

Varje ny variabel är också oberoende.

Har sin egen riktningskoefficient.



Exam Scores

# hours of study

x

https://images.app.goo.gl/ZuAVS1F2JuiTXN1B8

RESTORY...

# Flervariabel LR



Simple Linear Regression

$$y = b_0 + b_1 * x_1$$

Multiple Linear Regression

Dependent variable (DV)   Independent variables (IVs)

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \ldots + b_n * x_n$$

Constant   Coefficients

**RESTORY...**

# Flervariabler – två oberoende variabler – 3D



bedrooms
bathrooms

price

**RESTORY...**

# The dummy variable trap

- Ordinal to numerical
- Third can be inferred by the absence of the other two (included in $b_0$).
- sklearn model is advanced and will figure that out

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + b_3 * x_3 + b_4 * D_1$$

| | R&D Spend | Administration | Marketing Spend | State | Profit |
|---|---|---|---|---|---|
| 1 | 165349.2 | 136897.8 | 471784.1 | New York | 192261.83 |
| 2 | 162597.7 | 151377.59 | 443898.53 | California | 191792.06 |
| 3 | 153441.51 | 101145.55 | 407934.54 | Florida | 191050.39 |
| 4 | 144372.41 | 118671.85 | 383199.62 | New York | 182901.99 |
| 5 | 142107.34 | 91391.77 | 366168.42 | Florida | 166187.94 |
| 6 | 131876.9 | 99814.71 | 362861.36 | New York | 156991.12 |
| 7 | 134615.46 | 147198.87 | 127716.82 | California | 156122.51 |
| 8 | 130298.13 | 145530.06 | 323876.68 | Florida | 155752.6 |

NY = $D_1$
CA = $D_2$

$D_1 = 1 - D_2$

$b_0$ "includes" $D_2$

RESTORY...

# Kodstruktur

1. Import libraries
2. Import dataset
3. Encode categorical variables
4. Split data from "answer" column
5. Split train/tests
   a. Import "splitter"
6. Training
   a. Import model?
7. Testing

Print options

- `np.set_printoptions(precision=2)`

**RESTORY...**

# Kod

# Övningar

- Ladda ned Iris data
- Omvandla nominell data till numeriska kategorier
- Gör en flervariabelregression för att förutspå Sepal längden.

- Ladda ned Pima indians diabetes data
- Gör en flervariabelregression.

RESTORY...

# Polynom

# Polynom Linjär Regression

**Simple Linear Regression**

$$y = b_0 + b_1 x_1$$

**Multiple Linear Regression**

$$y = b_0 + b_1 x_1 + b_2 x_2 + \ldots + b_n x_n$$

**Polynomial Linear Regression**

$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \ldots + b_n x_1^n$$

RESTORY...

# Polynom LR

$$y = b + m_0x_1 + m_1x_1^2 + ... + m_nx_1^n$$
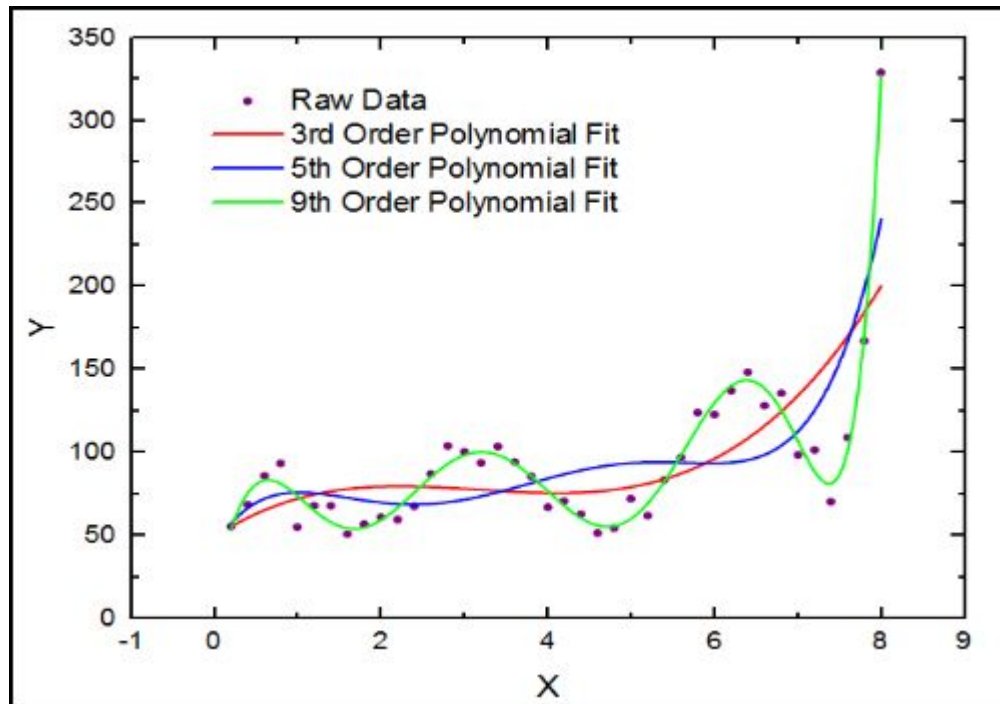
y är beroende av x

x är oberoende

förhållandet mellan dem som en n-graders polynom av x (relationship between them as a n-th degree polynomial of x)

högre grad -> mer komplex modell
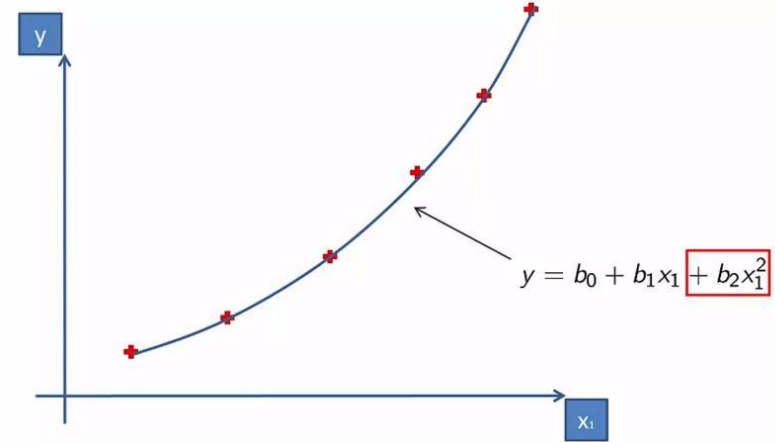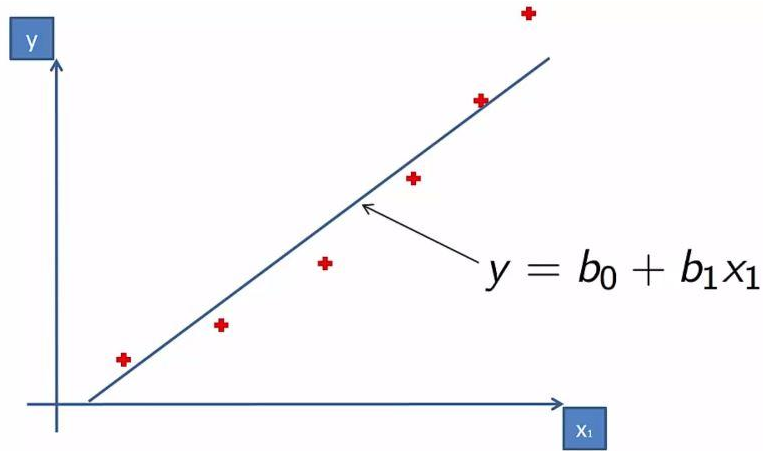
Kan modellera mer komplexa förhållanden.

Linjär - coeffs



https://images.app.goo.gl/qriDKw1nLNMr2fmf9

RESTORY...

# Polynom regression – när behövs det



$$y = b_0 + b_1 x_1$$

$$y = b_0 + b_1 x_1 \boxed{+ b_2 x_1^2}$$

RESTORY...

# Kod

RESTORY...

# Övning

Titta på datan i följande Kaggle exemplet
och försök göra en polynomregression
utan att titta på koden.

Pumpkin data

RESTORY...

# Statistik-ish

RESTORY...

# p-value

*" A p-value is the probability that random chance generated the data, or something else that is equal or rarer"*

- Statistical significance
- Null hypothesis $H_0$ - assume true
- Is it correct? - experiment
  - 0.50
  - halveras …
  - $\alpha$ = 0.05
  - domain dependent

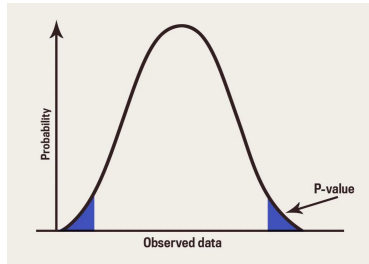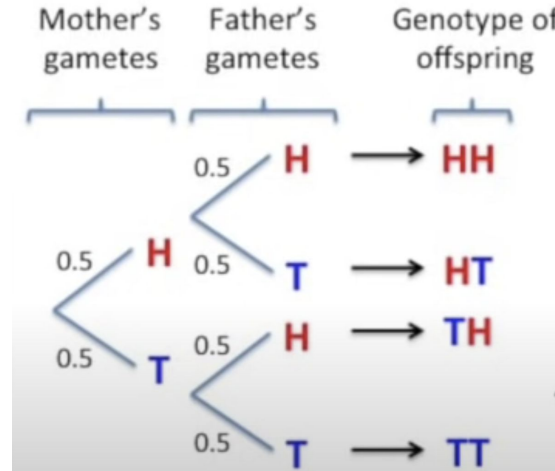**p-value ≠ probability**



Bild ref.



$$\frac{HH}{HH, HT, TH, TT} = \frac{1}{4} = 0.25$$

$$+$$

$$\frac{TT}{HH, HT, TH, TT} = \frac{1}{4} = 0.25$$

StatQuest Youtube vid

RESTORY...

# Multivariable regression - variable selection

Why select?
- GIGO
- Noise
- Explain

1. **All** Variables
   a. Domain knowledge - predictors
   b. Forced - by company
   c. Preparing for no. 2

2. **Backward** Elimination
   a. Set significance level (e.g. sl = 0.05)
   b. Fit full model
   c. Remove predictor with highest **p**-val if **p>sl**.
   d. Refit model
   e. Repeat c & d until no p > sl.

3. **Forward** Selection
   a. Select sl
   b. Fit all simple reg. models for y and X for all x separately.
   c. Select the one with the lowest p-val ($x_n$).
   d. Fit $x_n$ with all other x variables (X-$x_n$) separately for every model. (y and [$x_a$, $x_b$])
   e. If the model with the lowest p-val has a p<sl, add it to the predictor, and repeat d.

RESTORꝨ...

# Multivariable regression – variable selection

Stepwise regression: 2,3,**4**

4. ***Bidirectional*** *Elimination*
    a. Select an entry and stay sl (can be different)
    b.  Add on new variable using Forward Selection (p < sl_enter).
    c. Do all steps of Backward elimination. Keep only variables where p < sl_stay is fulfilled.
    d. Repeat b and c until no new variables can enter or exit.
5. Score Comparison
    a. Set a fit score minimum
    b. Create a model for all possible variable combinations ($2^n-1$)
    c. Select the one with the best score
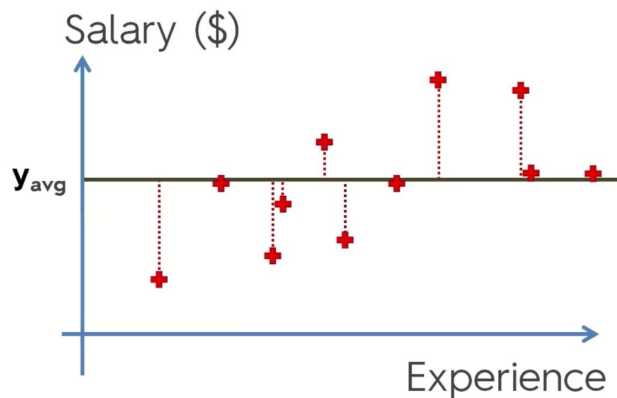
Scikit-learn does this automatically!

RESTORY...

$R^2$

RESTORY...

# $R^2$

How much better than average?

Simple Linear Regression:



$$SS_{res} = SUM\ (y_i - \hat{y_i})^2$$

$$SS_{tot} = SUM\ (y_i - y_{avg})^2$$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

RESTOR〉...

# R² och flera variabler

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$R^2$ – Goodness of fit
(greater is better)

$y = b_0 + b_1 * x_1$

**Problem:**

$y = b_0 + b_1 * x_1 + b_2 * x_2$ ← $+ b_3 * x_3$

$SS_{res} \rightarrow$ Min → $R^2$ will never decrease

Udemy kurs, Lek. 99.

RESTORꙄ...

# Justerad $R^2$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$$Adj\ R^2 = 1 - (1 - R^2)\frac{n - 1}{n - p - 1}$$

p – number of regressors
n – sample size

RESTORY...

# Sammanfattning

- Linjär regression med flera variabler
- Kodstruktur
- Polynomregression
- Statistik
    - p-value
    - Multivariable regression - variable selection
    - $R_2$

RESTORУ...

# Länkar

- StatQuest
- Linear regression

**RESTOR**....