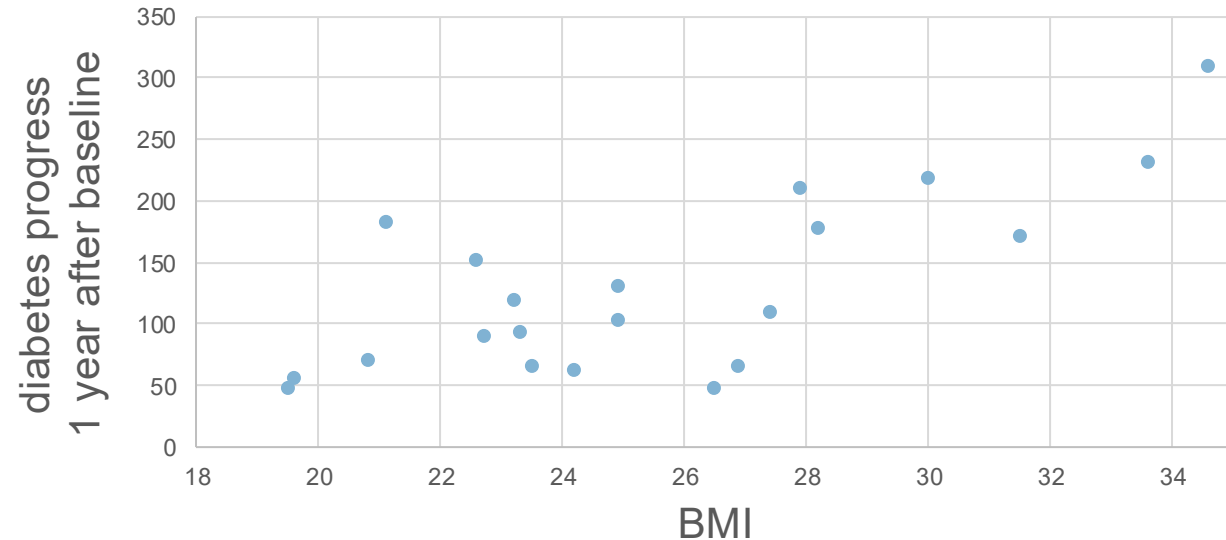


Supervised Machine Learning Regression

Exemplary Regression Problem: Course of diabetes progress depending on patient's BMI

bmi	progress
x	y
27.8	201
22.8	40
35	140
34.6	264
\vdots	\vdots

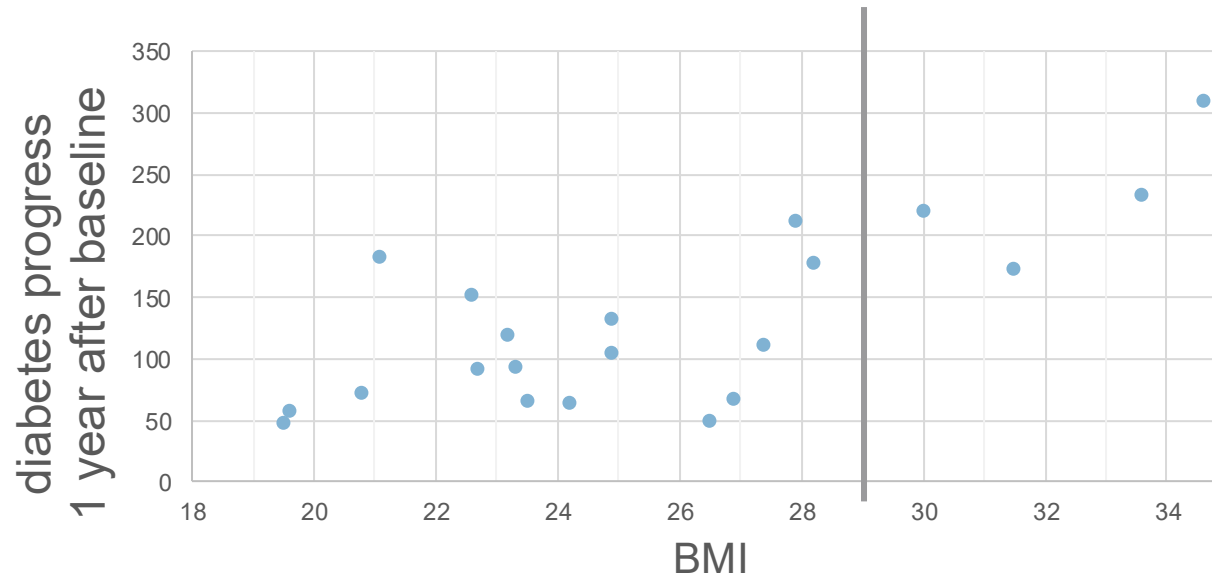


Univariate: 1 independent variable (feature)

Data: <https://www4.stat.ncsu.edu/~boos/var.select/diabetes.html>

Illustration of the Regression Problem

What is the expected course of disease with a BMI of 29?



Instance-Based Learning Applied to a Regression Problem

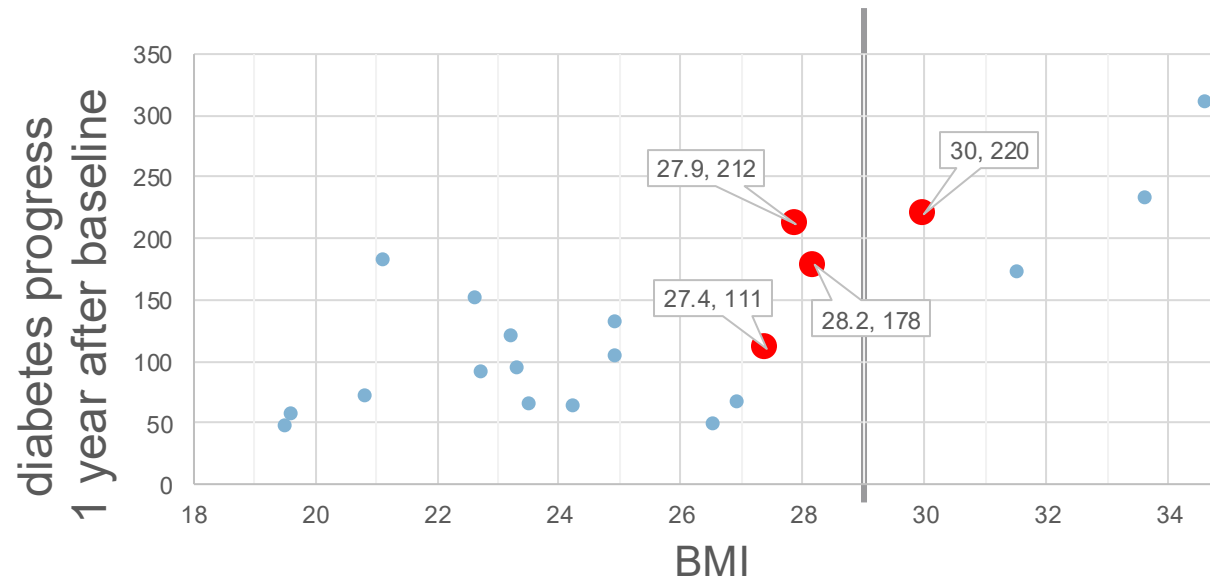
KNN: K-Nearest Neighbour Algorithm

What is the expected course of disease with a BMI of 29?

$k=4$

Distance metric: difference in BMI

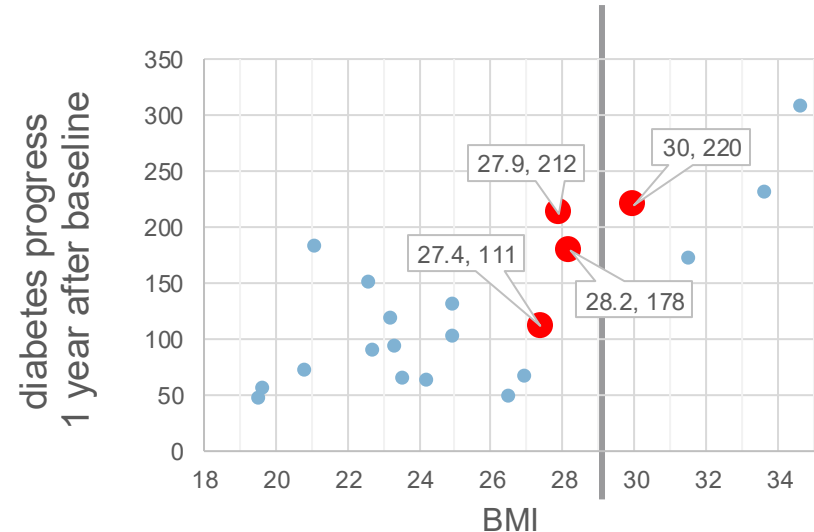
BMI=29 \rightarrow average of the 4 closest points = $(111+212+178+220)/4 = 180.3$



The Nearest Neighbours-Algorithmus Applied to a Regression Problem

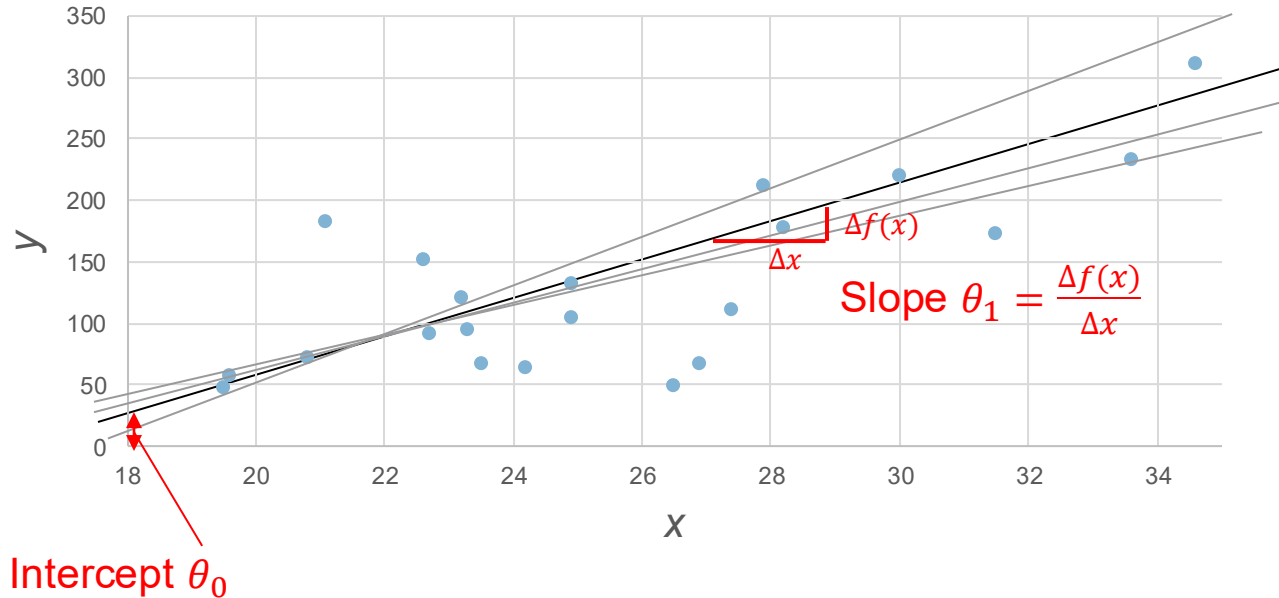
Disadvantages:

- Inference on a new data point **always requires to consider explicitly all training data** (select k closest points and compute the output)
- Is **sensitive to outliers** and noise
- Training samples further away than the k neighbours are ignored for inference → **loss of information** contained in the training data



Univariate linear regression

Mathematical expression which consistently models the relationship between the variables.



Hypothesis: Linear function
with parameters θ_0, θ_1

$$h(x; \theta_0, \theta_1) = \theta_0 + \theta_1 x$$

Goal of the Training/Learning: Parameter values, which describe the relationship best.

«Training» of the univariate linear regression model

Training data: M samples $\{(x_m, y_m)\}$

The outputs y are explained by the hypothesis h plus

the stochastic residuals e : $y_m = h(x_m; \theta_0, \theta_1) + e_m$

The model predicts for the sample x_m : $\hat{y}_m = \theta_0 + \theta_1 x_m$

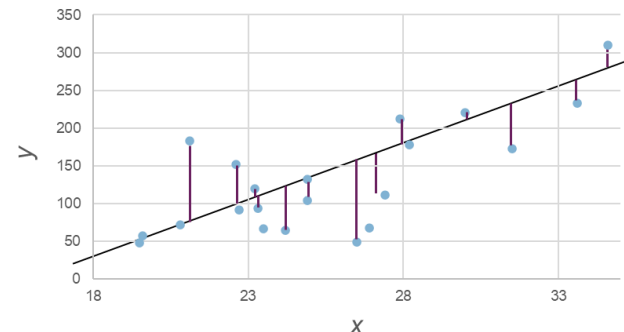
Loss function: Sum of squared **residuals** $\mathcal{L}_{RSS}(\theta_0, \theta_1) = \sum_{m=1}^M e_m^2 = \sum_{m=1}^M (y_m - \hat{y}_m)^2$

Learning: Minimise the cost function

$$J(\theta_0, \theta_1) = \frac{1}{2M} \sum_{m=1}^M (y_m - \hat{y}_m)^2$$

Parameters $\hat{\theta}_0, \hat{\theta}_1$ minimise the cost function

and can then be used for inference on new data samples $\hat{y}_i = h(x_i; \hat{\theta}_0, \hat{\theta}_1)$



Explicit solution of univariate linear regression

Setting $\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = 0$ and $\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = 0$

and solving these for θ_0 and θ_1 yields (without derivation):

$$\hat{\theta}_0 = \mu_y - \theta_1 \mu_x \quad \text{and} \quad \hat{\theta}_1 = \frac{\sum_{m=1}^M (x_m - \mu_x)(y_m - \mu_y)}{\sum_{m=1}^M (x_m - \mu_x)^2}$$

where μ_x is the mean of all $\{x_m\}$ and μ_y is the mean of all $\{y_m\}$ of the training data

Normal Equations

Write the hypothesis in matrix form: $\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{e}$

For the (univariate) diabetes problem:

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{pmatrix} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & \vdots \\ 1 & x_M \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_M \end{pmatrix}$$

The vector of residuals is $\mathbf{e} = \mathbf{y} - \mathbf{X}\boldsymbol{\theta}$

Loss: The sum of squared residuals is

$$\begin{aligned} \mathbf{e}^T \mathbf{e} &= (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \mathbf{y}^T \mathbf{y} - \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\theta} \\ &= \mathbf{y}^T \mathbf{y} - 2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\theta} \end{aligned}$$

Normal Equations

The gradient of the squared residuals (without derivation):

$$\frac{\partial \mathbf{e}^T \mathbf{e}}{\partial \boldsymbol{\theta}} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta}$$

Equating the gradient with zero produces the so-called normal equations:

$$(\mathbf{X}^T \mathbf{X}) \hat{\boldsymbol{\theta}} = \mathbf{X}^T \mathbf{y}$$

which leads to a closed form-expression for the optimal parameters:

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Multivariate linear regression with the normal equations

$$h(\mathbf{x}_m; \theta_0, \theta_1, \dots, \theta_N) = \theta_0 x_{m0} + \theta_1 x_{m1} + \theta_2 x_{m2} + \dots + \theta_N x_{mN} = \boldsymbol{\theta}^T \mathbf{X}_{m,:}$$

with $x_{m0} = 1$ for all $m = 1, \dots, M$

Diabetes example

with $M = 4$, $N = 3$:

	bmi	bp	glu	progress
x_0	x_1	x_2	x_3	y
1	27.8	73	73	201
1	22.8	101	97	40
1	35	79.33	96	140
1	34.6	115	109	264

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1N} \\ 1 & x_{21} & x_{22} & \cdots & x_{2N} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{M1} & x_{M2} & \cdots & x_{MN} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_N \end{pmatrix}$$

Dimensions: $M \times (N + 1)$

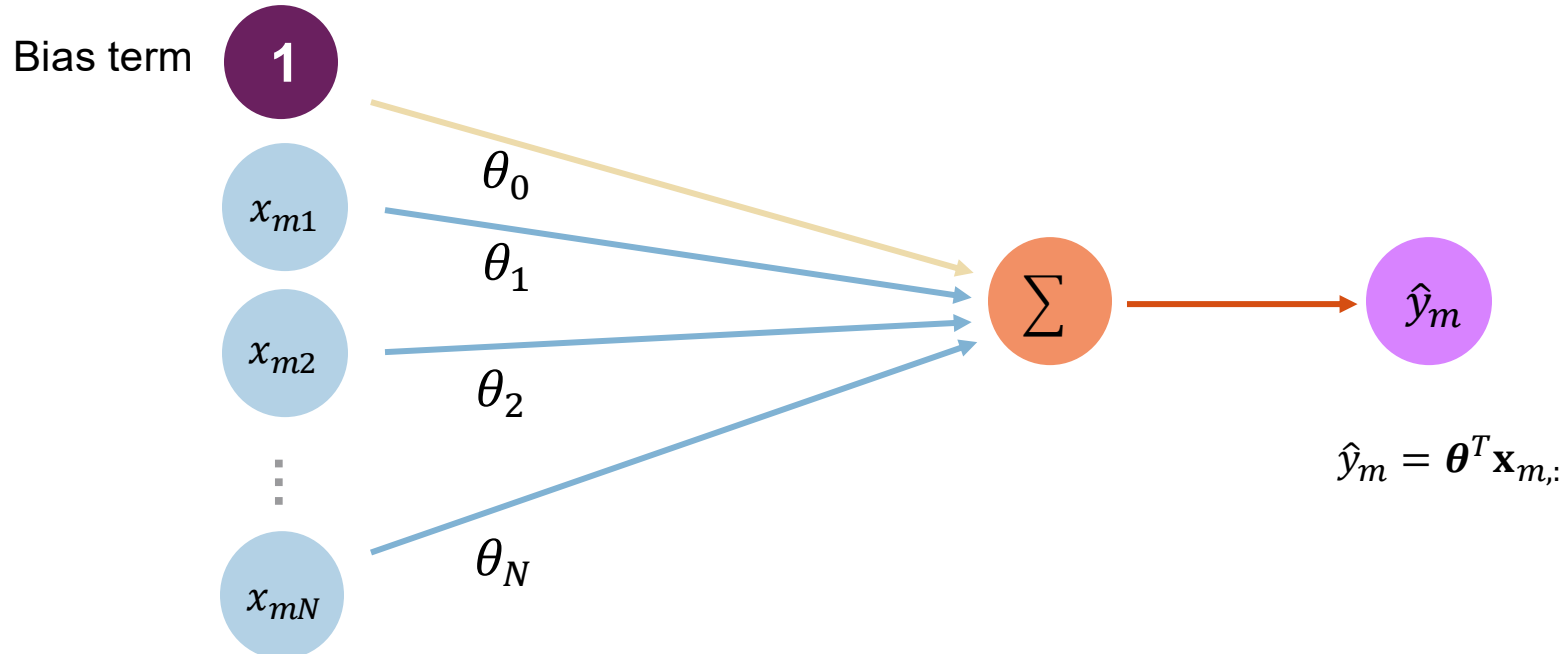
$M \times 1$

$N + 1$

Linear regression for one sample, N features

Forward pass on sample $\mathbf{x}_{m,:}$ (m -th row in X) : \mathbf{x}

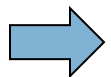
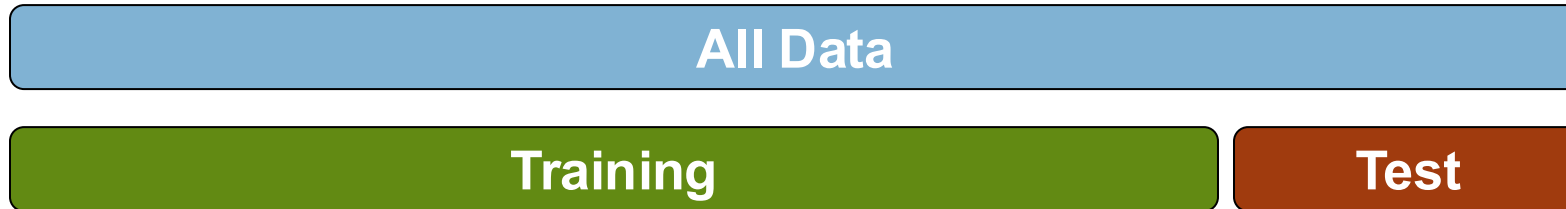
Input $\mathbf{x}_{m,:}$ Parameters/weights linear combination output



Evaluating Regression Models

Can we estimate quantitatively how well the model generalises beyond the training data?

General for supervised machine learning: Split the available data into a training and independent test set



Estimate the generalisation error based on the independent test set

Evaluation Metrics for Regression

Mean Absolute Error: $MAE = \frac{\sum_{i=1}^I |y_i - \hat{y}_i|}{I}$

Mean Squared Error (MSE): $MSE = \frac{\sum_{i=1}^I (y_i - \hat{y}_i)^2}{I}$

Root Mean Squared Deviation: $RMSE = \sqrt{MSE} = \sqrt{\frac{\sum_{i=1}^I (y_i - \hat{y}_i)^2}{I}}$

with I number of samples in the independent test set

Evaluation Metrics for Linear Regression

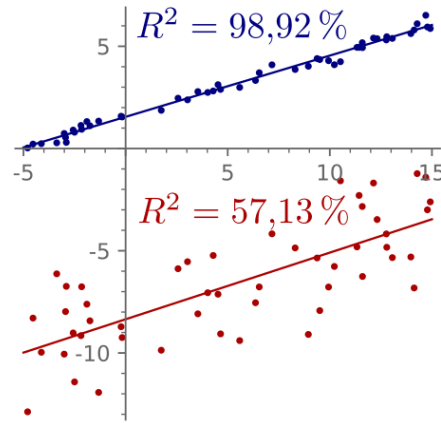
Sum of squares of residuals: unexplained variance (variance of the model's errors)

Total sum of squares (proportional to the variance of the data)

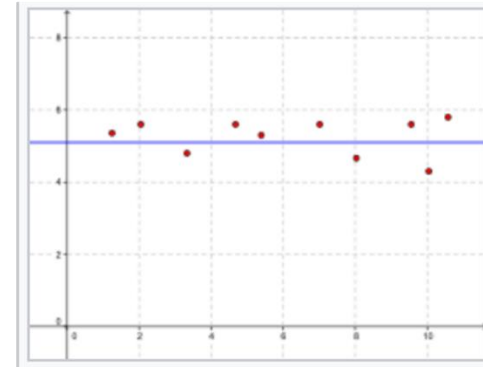
Coefficient of Determination $R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum_{i=1}^I (y_i - \hat{y}_i)^2}{\sum_{i=1}^M (y_i - \mu_y)^2}$

Measures fraction of the variance of the data, which can be explained by the model.

In the best case: $y_i = \hat{y}_i \rightarrow SS_{res} = 0 \rightarrow R^2 = 1$



A baseline model, which always predicts $\mu_y \rightarrow R^2 = 0$



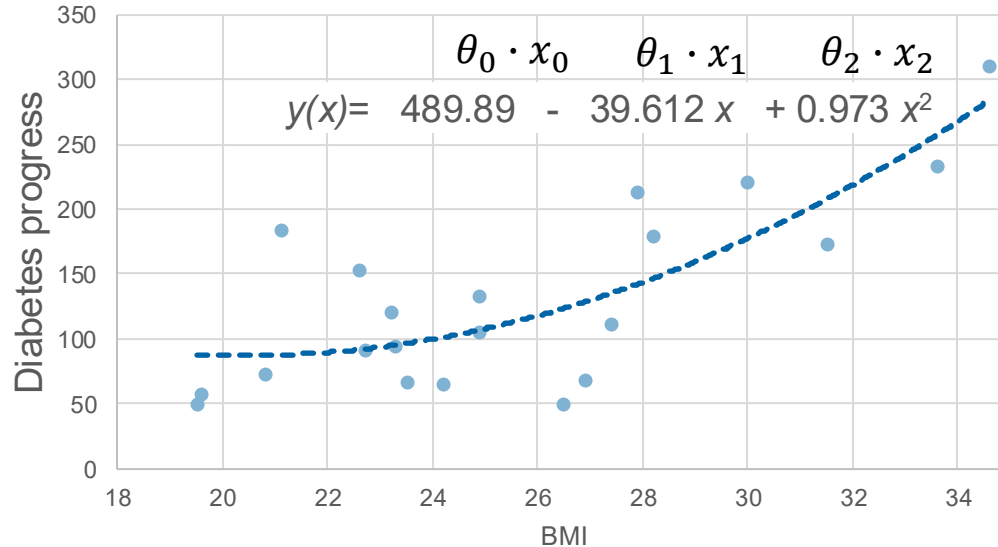
Models which perform worse than the baseline $\rightarrow R^2 < 0$

<https://de.wikipedia.org/wiki/Bestimmtheitsma%C3%9F>

Polynomial Regression

Modelling non-linear relationships using higher order polynomials

e.g. quadratic function:

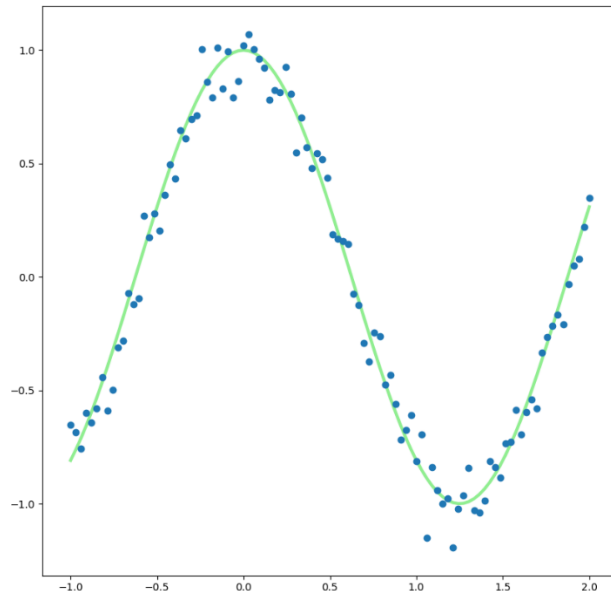


→ linear combination of higher order features.

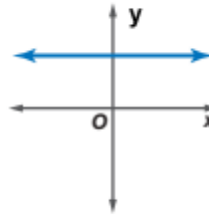
→ treat x, x^2, \dots, x^K as distinct independent features in a multivariate regression model and use the same procedure to solve for the optimal parameters.

Higher Order Polynomials

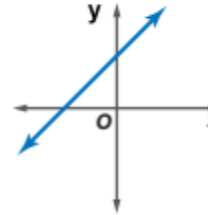
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 (x_2)^3 + \theta_3 x_2^3 x_3 x_4$$



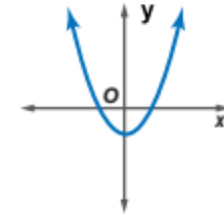
Constant function
Degree 0



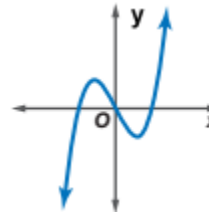
Linear function
Degree 1



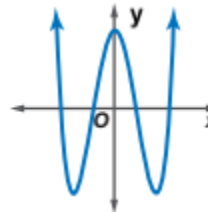
Quadratic function
Degree 2



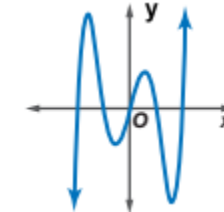
Cubic function
Degree 3



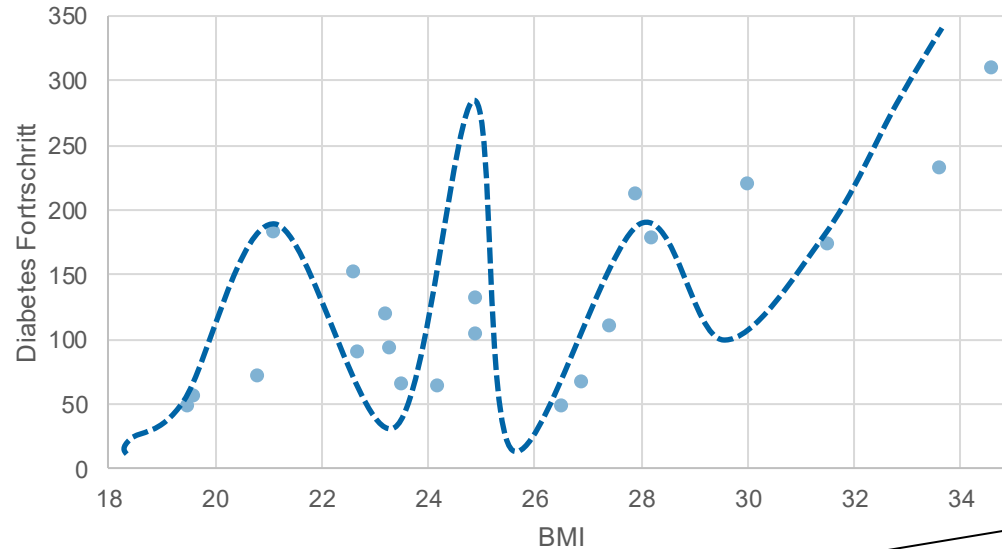
Quartic function
Degree 4



Quintic function
Degree 5



Overfitting



(here the order of the polynome)

Increasing the flexibility of the model (e.g. polynomial order) can, to a certain extent, lead to a better fit with the training data.

Overfitting occurs when the flexibility of the model leads to fitting the noise or outliers rather than representing the generalizing pattern of the data.

Other Aspects of Model Flexibility/Complexity

	$K = 0$	$K = 1$	$K = 3$	$K = 9$
\hat{w}_0	-0.06	1.04	-0.11	0.38
\hat{w}_1		-2.04	11.28	-15.58
\hat{w}_2			-33.18	450.84
\hat{w}_3			22.13	-4228.63
\hat{w}_4				20509.34
\hat{w}_5				-57747.54
\hat{w}_6				97233.05
\hat{w}_7				-96374.54
\hat{w}_8				51856.86
\hat{w}_9				-11684.80

Values of the parameters learned during training of polynomial regression model with varying degree on the synthetic dataset

Overfitting is manifested by large oscillations, i.e. large absolute values of the coefficients

Besides restricting the number of parameters, the model's flexibility/complexity (and therefore overfitting) can be controlled by preventing larger values of the coefficients during model training → **Regularisation**

Regularisation in polynomial regression

Regularization

Hypothesis: $h(x_m; \theta_0, \theta_1, \dots, \theta_N) = \boldsymbol{\theta}^T \mathbf{x} = \theta_0 x_{m0} + \theta_1 x_{m1} + \dots + \theta_N x_{mN}$

Add an additional term to the cost function that penalises large values of the model parameters.

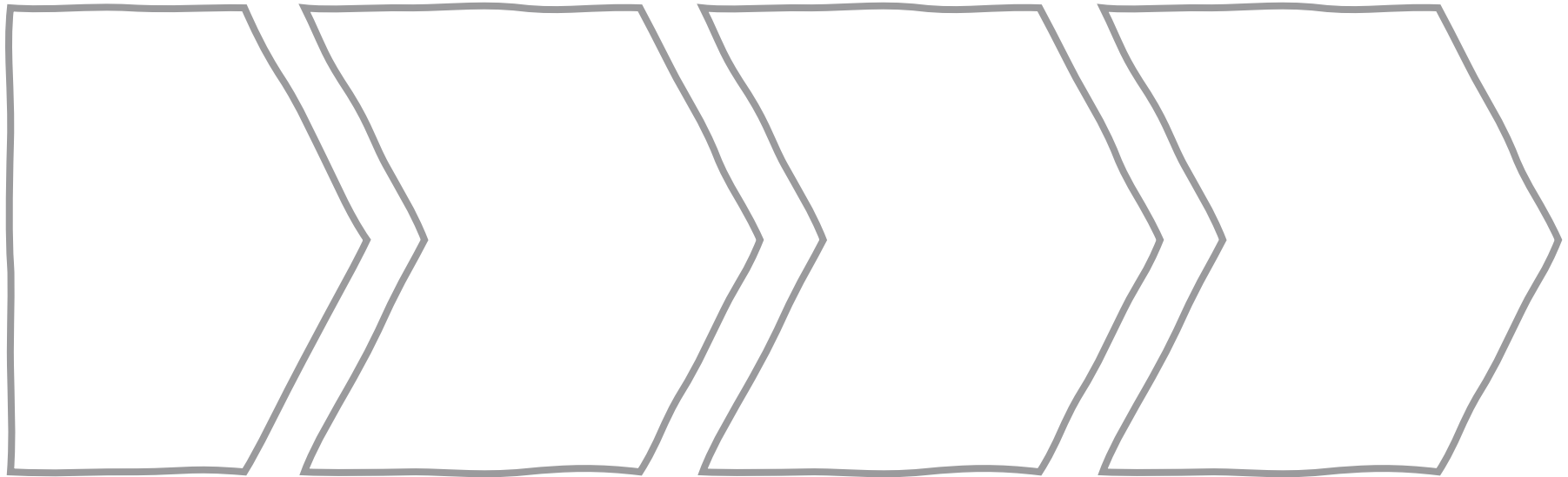
Cost Function of Ridge Regression:

$$J(\{(x_m, y_m)\}, \theta_0, \theta_1, \dots, \theta_N) = \frac{1}{2M} \left[\sum_{m=1}^M (y_m - h(x_m; \theta_0, \theta_1, \dots, \theta_N))^2 + \lambda \sum_{n=1}^N \theta_n^2 \right]$$

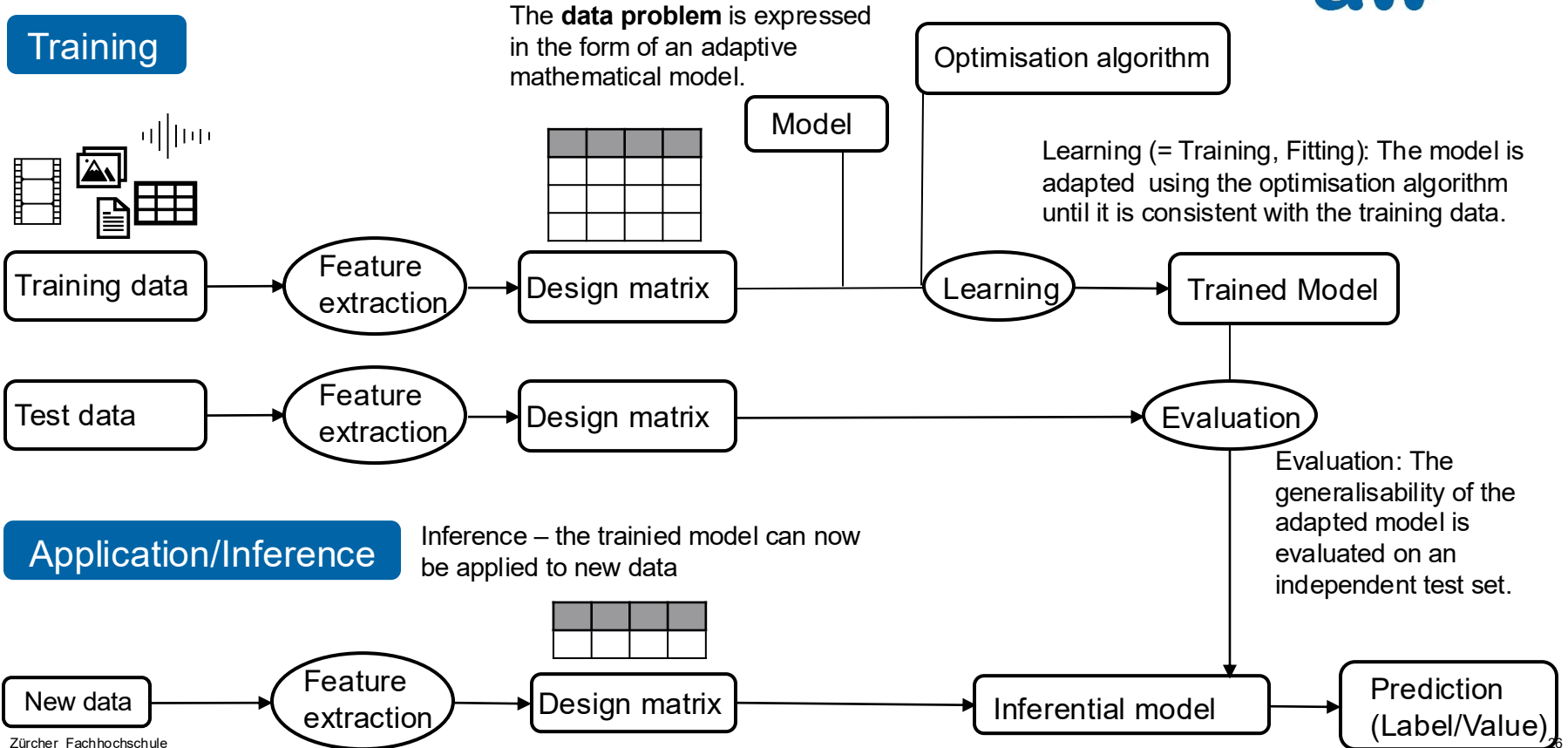
The hyperparameter λ is not optimised during training \rightarrow model selection

Appendix

Structure of the Learning Problem

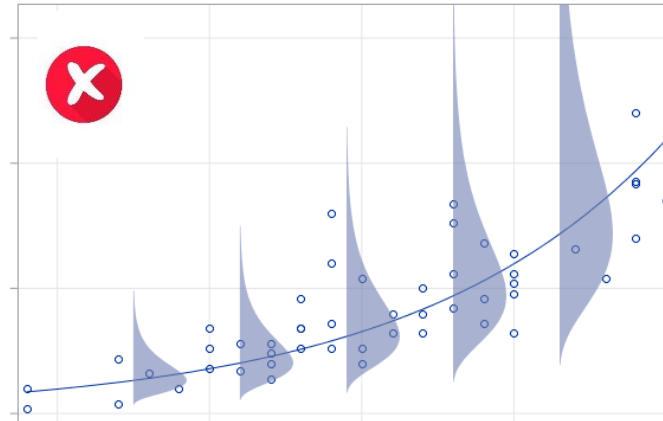
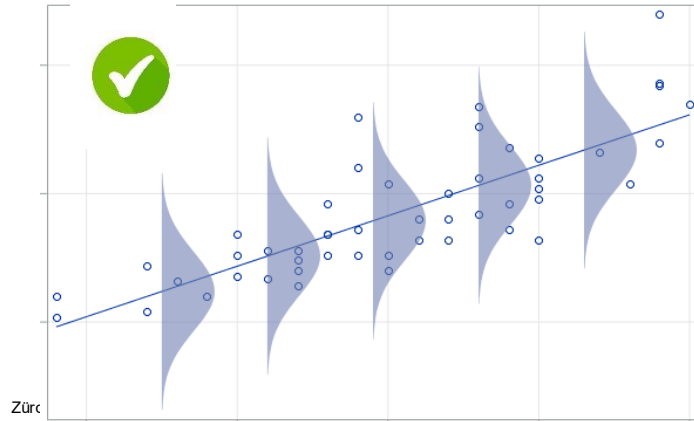
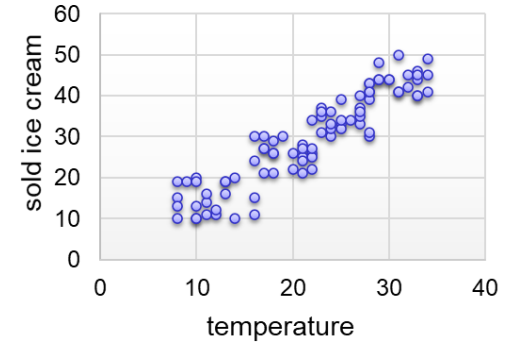


Elements of a «Machine Learning Pipeline»



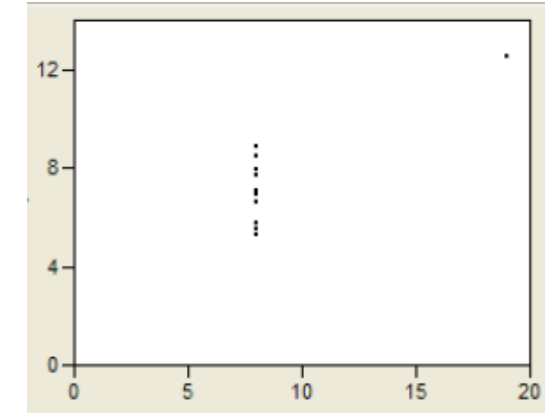
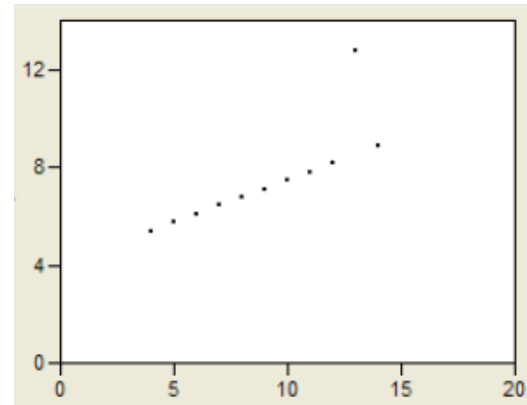
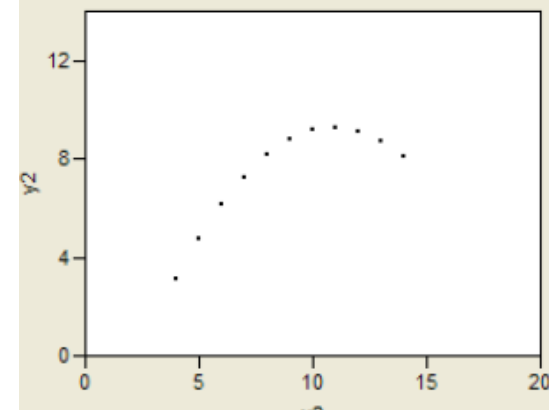
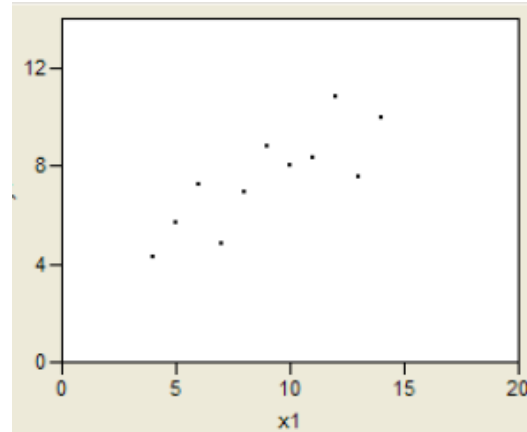
Basic Assumptions of Linear Regression

- **Linearity:** The input and output values have a linear relationship
- **Independence:** The outcome of one sample does not affect the others
- **Normality:** Errors should be normally distributed, i.e. larger deviations from mean should be less likely
- **Equality of Variance ("Homoscedasticity"):** Error distribution should be the same for all input values



Evaluation Through Visual Inspection

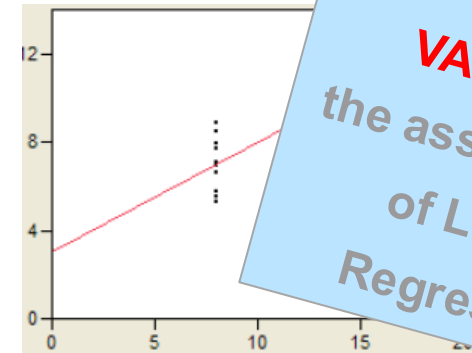
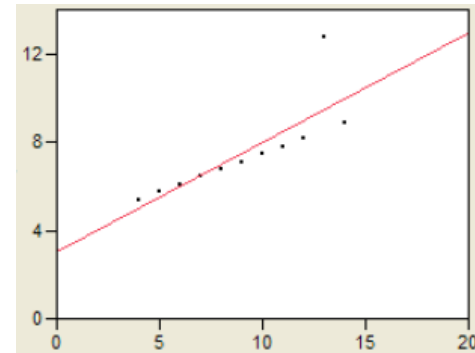
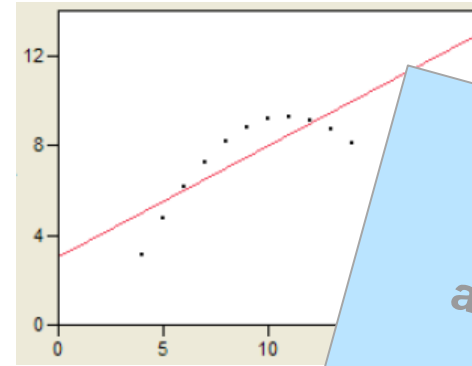
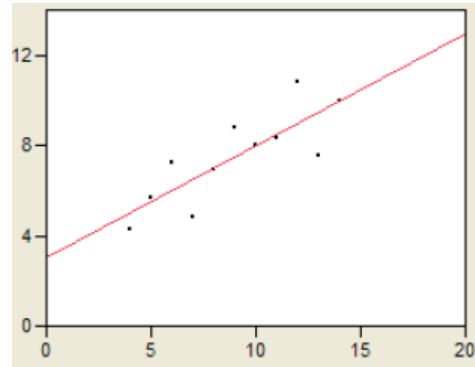
Ancombe's Quartet:



Evaluation Through Visual Inspection

Linear Regression yields same fit line and same mean square error
for all examples

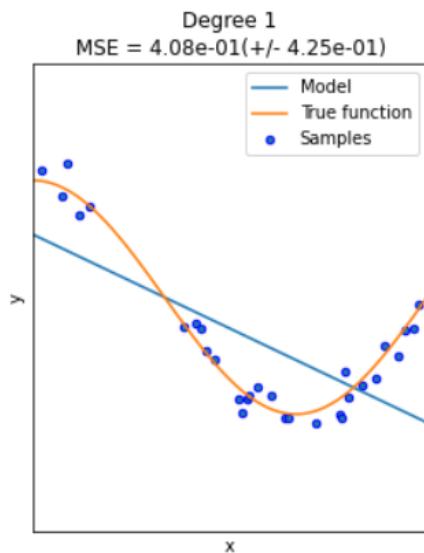
Ancombe's Quartet:



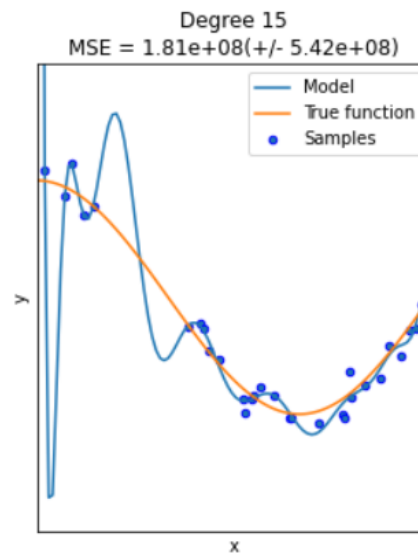
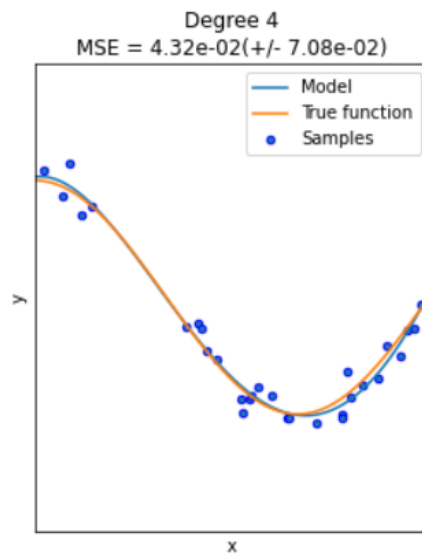
You have to
LOOK
at the data
and
VALIDATE
the assumptions
of Linear
Regression

Higher Order Polynomials might Overfit

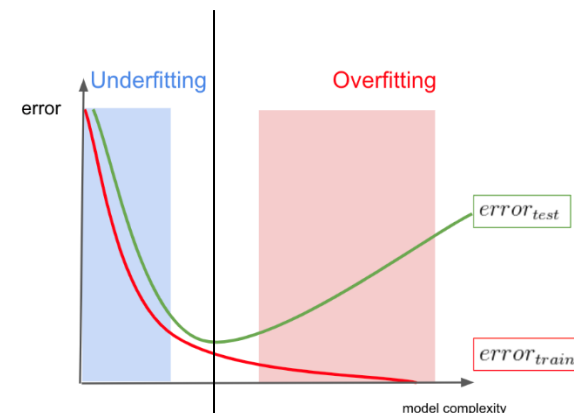
Zürcher Hochschule
für Angewandte Wissenschaften



Underfitting

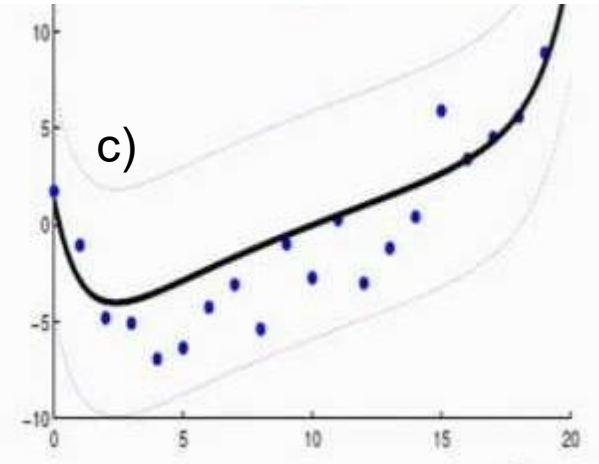
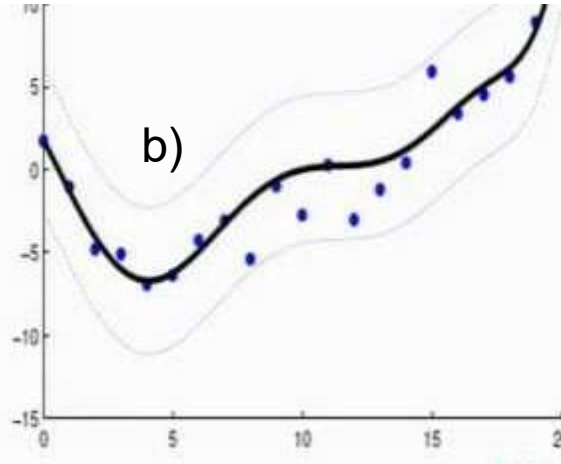
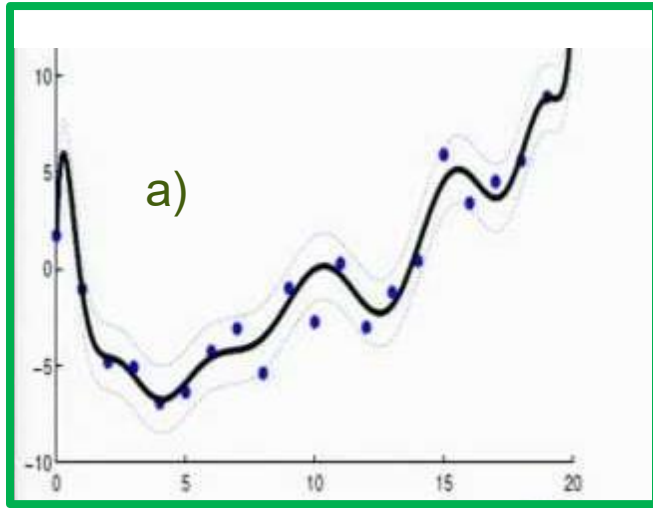


Overfitting



Solution: Regularization Parameter

The images below show regression models with polynomials of degree 14. Which of them has the *lowest* value of regularization parameter λ ?



$$\text{Cost Function: } J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

SOLUTION: Large Regularization Parameter

What happens if we set the regularization parameter λ to a very large value, e.g. $\lambda = 10'000'000$?

$$\text{Cost Function: } J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Then all parameters θ_j will be close to or equal to zero, except for θ_0
Thus, we obtain a horizontal line.

