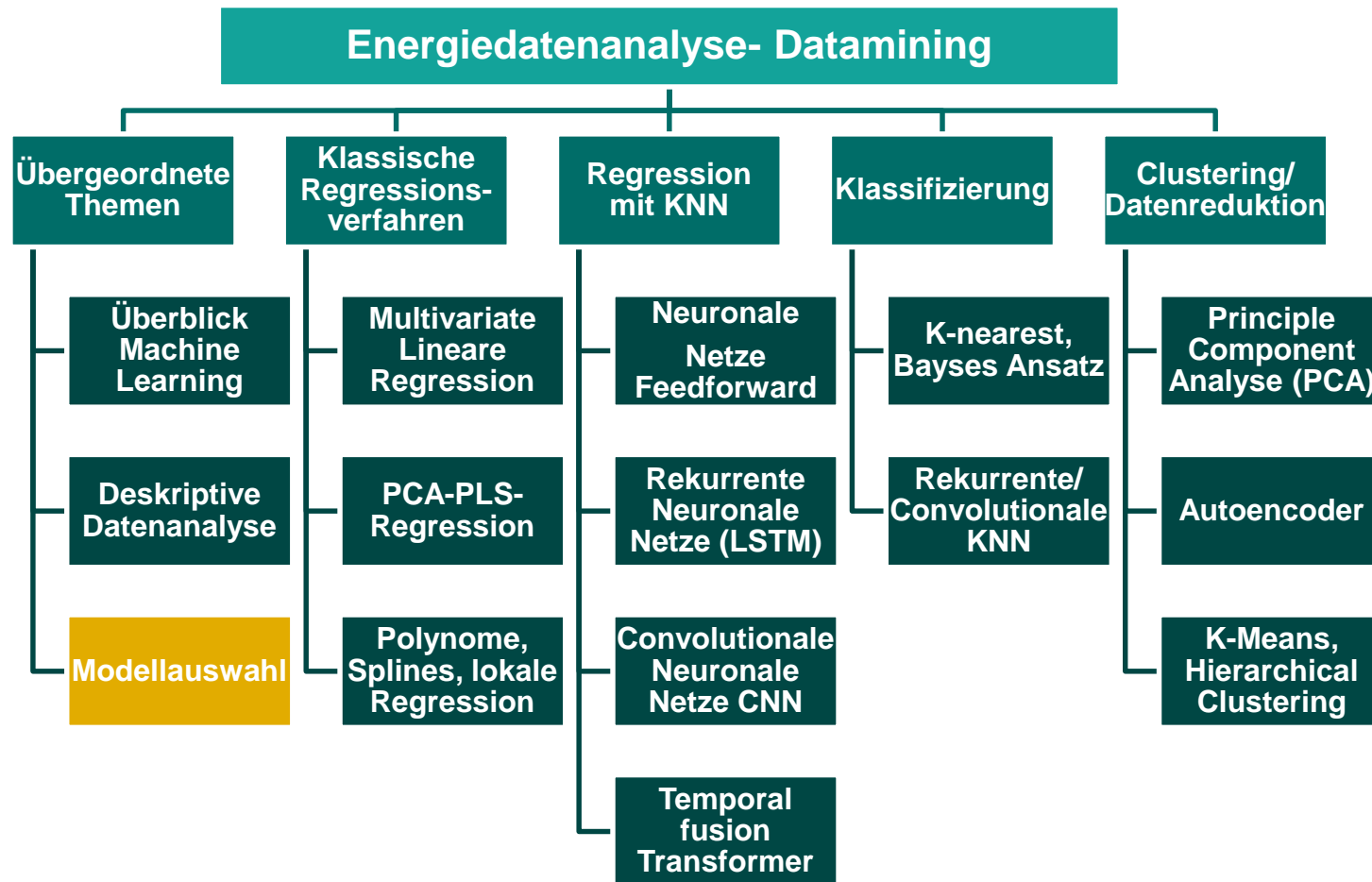


Modellauswahl

Energiedatenanalyse - Datamining

Die Themengebiete der Veranstaltung verknüpfen Modelle des „Machine Learning“ mit energiewirtschaftlichen Fragestellungen



Zielsetzung der heutigen Vorlesung: Einführung in die Verwendung der linearen Regression

Thema	Überblick über Validierungs- und Auswahlverfahren
-------	---

Aufbau der heutigen Vorlesung:

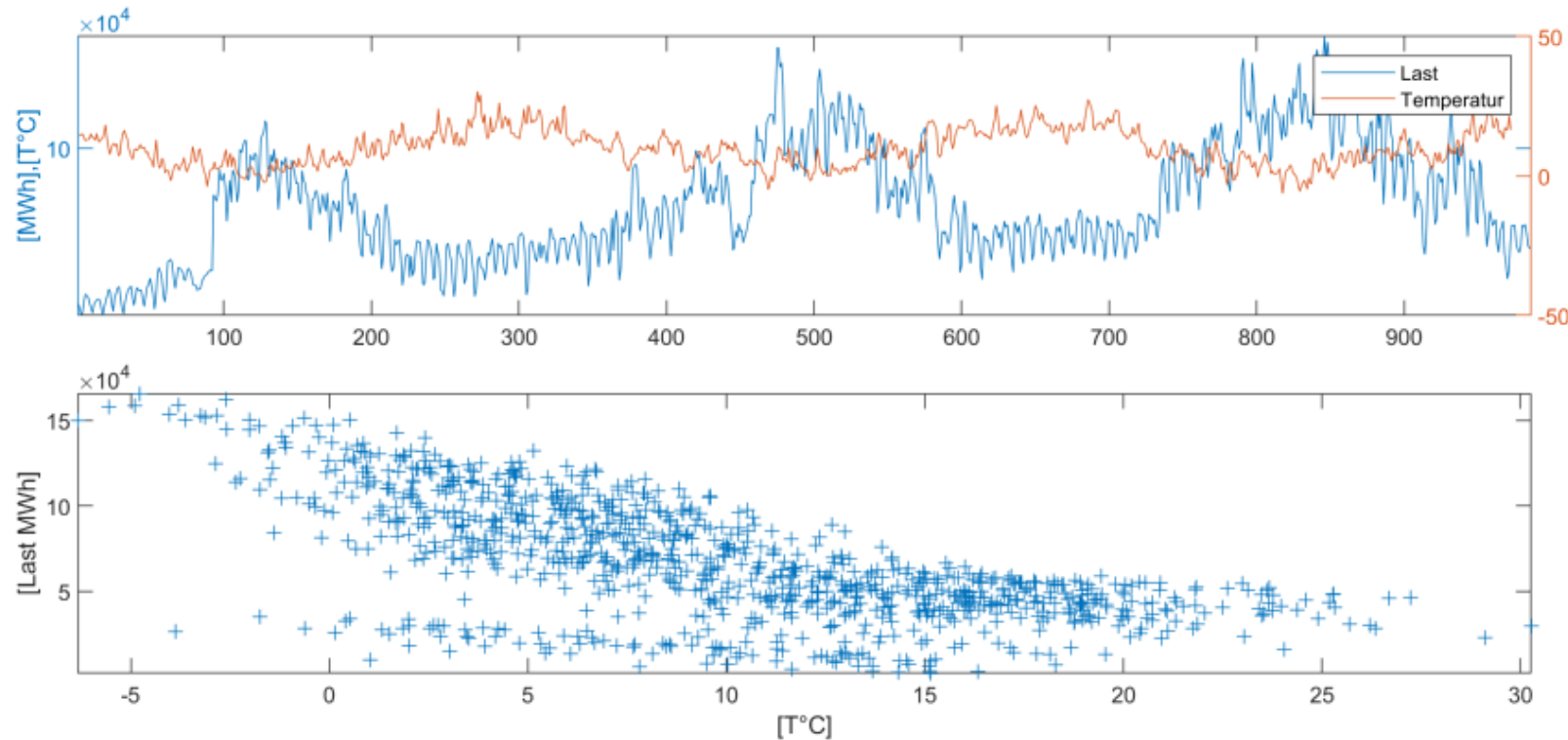
Lernziel:

1	Bedeutung des Bias-Variance Trade off	Notwendigkeit des Einsatzes von Validierungstechniken
2	Validation/ Cross-Validation-Techniken	Anwendung von Techniken zur Modellbeurteilung
3	Best Subset selection Stepwise selection	Prozess zur Modellauswahl

1	Bedeutung des Bias-Variance Trade off
2	Generierung von Schätzern für den Test-RMSE
3	Best Subset selection Stepwise selection

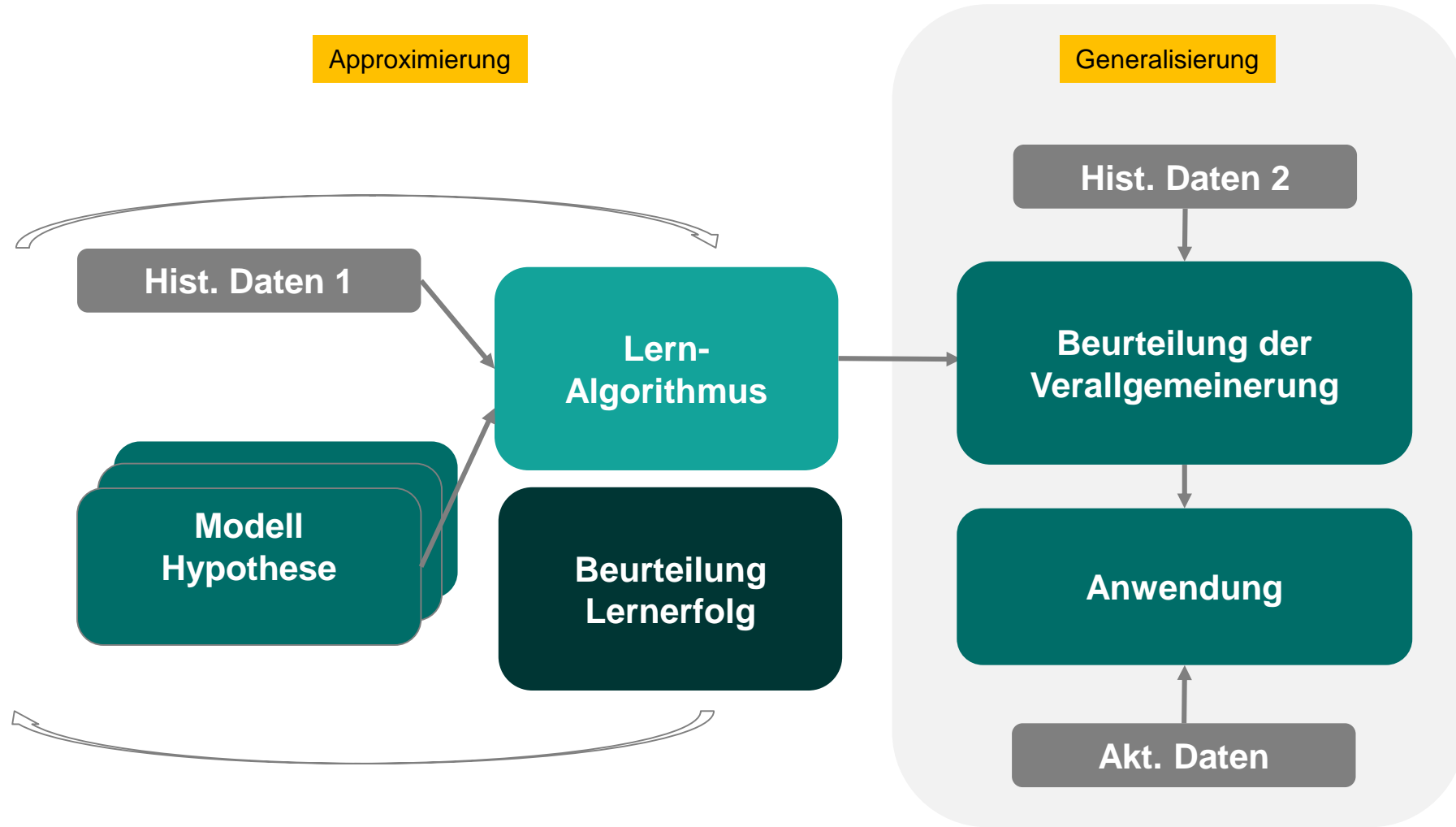
Literaturempfehlung und Quellennachweise der Abbildungen: James, Gareth; Witten, Daniela; Hastie, Trevor; Tibshirani, Robert. An Introduction to Statistical Learning: with Applications in R (Springer Texts in Statistics), Springer New York. Kindle-Version.

Essentielle Fragen im Rahmen der Modellkalibrierung

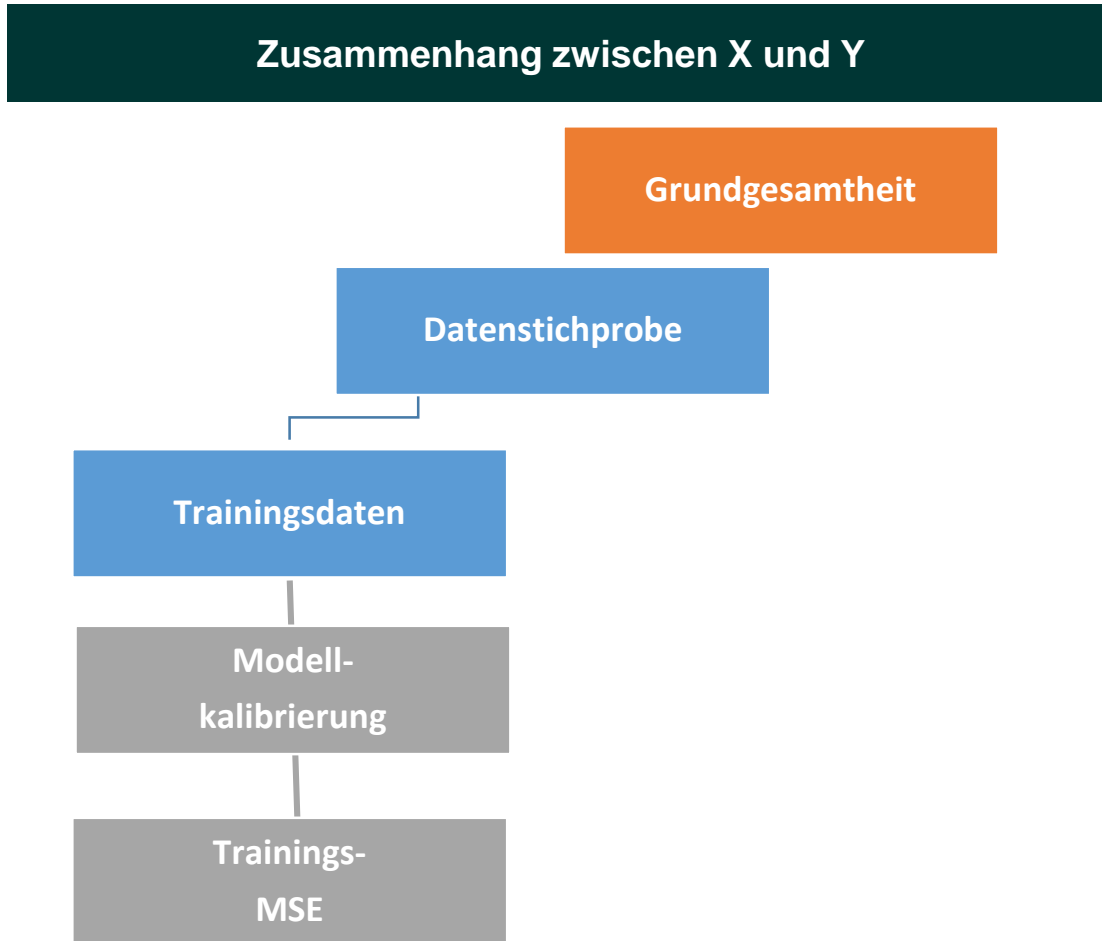


- Ziel 1: Auffinden eines Modells zur Approximation des Zusammenhangs der vorliegenden Daten → **Approximation**
- Ziel 2: Sicherstellung, dass das Modell den Zusammenhang bislang unbekannter Daten abbildet → **Generalisierung**

Eine wesentliche Aufgabe des Anwenders ist die richtige Modellauswahl



Die Approximation erfolgt durch Minimierung eines Zielkriteriums ausgewertet auf einer Datenstichprobe



Ausgangslage:

- der reale funktionale Zusammenhang zwischen Y und X ist unbekannt.

Vorgehensweise:

1. Bildung Modellhypothese,
2. Bildung Datenstichprobe,
3. Modellparametrierung.

- Zielkriterium zur Anpassung der Parameter:

■ Mean Squared Error MSE ($MSE = S_{\hat{e}\hat{e}}$)

$$MSE = \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{f}(x_i) \right)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

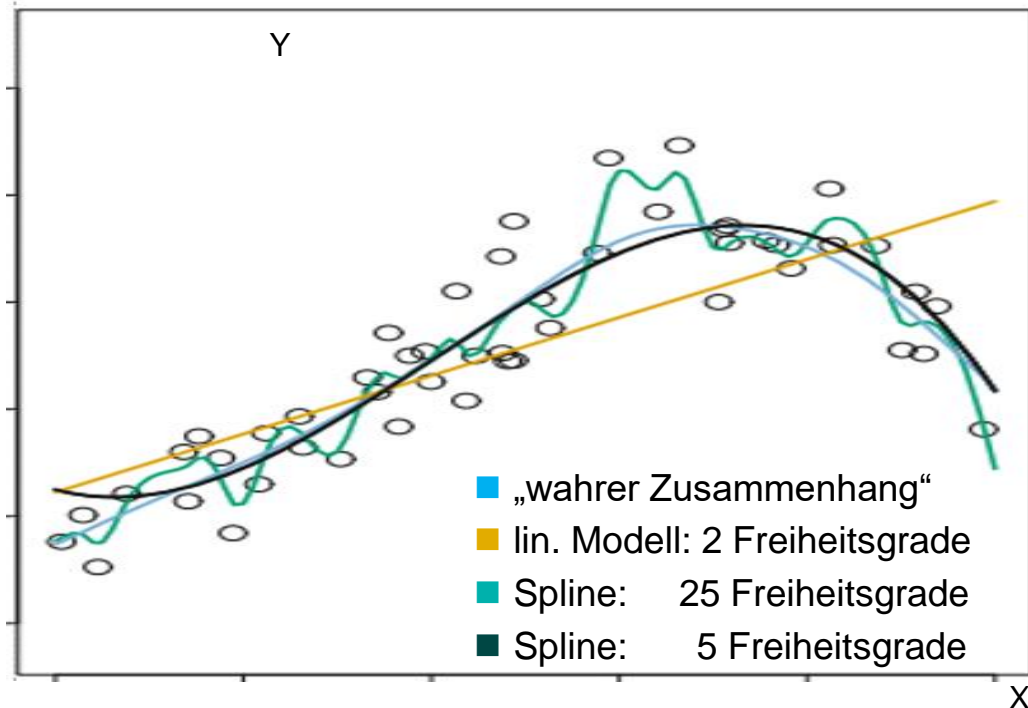
In Sample Fehler = Training MSE

\hat{y} ist der modellierte Wert bzw. Schätzer für y.

Der „In-Sample-Fehler“ misst die Abweichung des Modell auf den Daten anhand derer die Modellparameter optimal eingestellt werden.

Problem: durch hinzufügen weiterer Parameter wird der gesuchte Zusammenhang beliebig approximiert

Zusammenhang zwischen X und Y



- Modelle unterscheiden sich u.a. durch die Anzahl der verwendeten **Parameter**.
- Ein Modellparameter repräsentiert einen Freiheitsgrad.

- Generell gilt:

$$\text{Freiheitsgrade} \uparrow \rightarrow \text{MSE}(\text{Training}) \downarrow$$

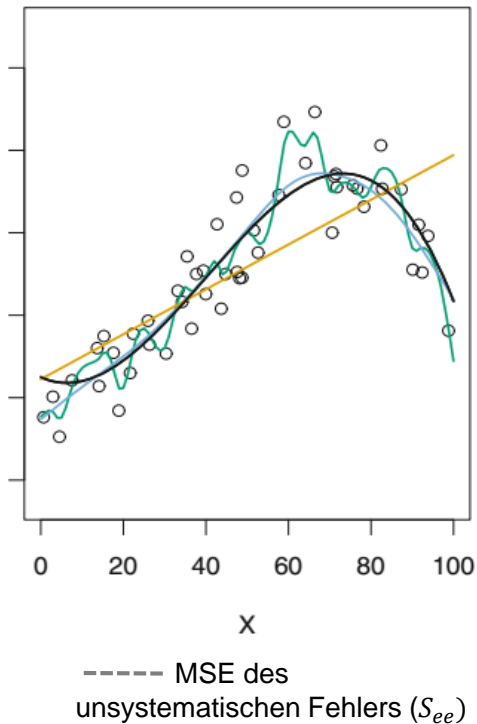
- Konsequenz für die Minimierung der Zielfunktion:
 - Verringerung Abweichung von \hat{Y} zu Y
 - Der MSE sinkt „zwangsläufig“ durch Vergrößerung der Regressorbasis p

Frage: wie schlägt sich so ein Modell bei „neuen“ bislang unbekannten Daten?

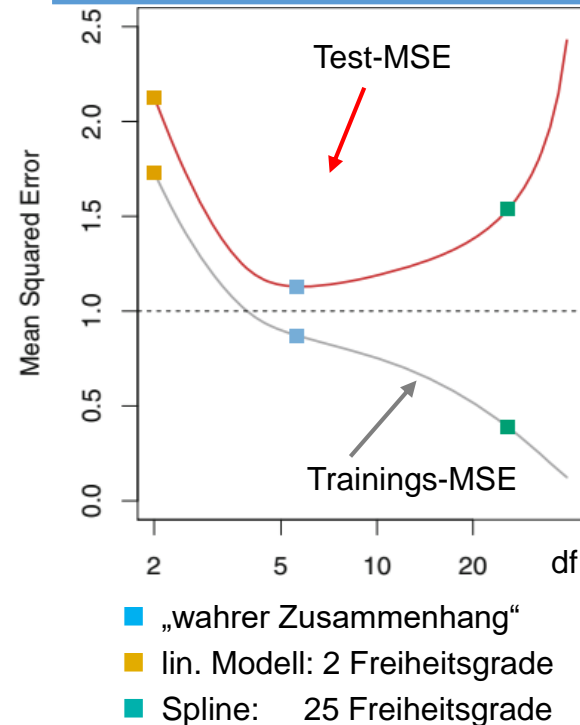
Beobachtung: MSE der Trainingsdaten und MSE der Testdaten verhalten ab einen gewissen Punkt gegensätzlich -> Zustand Overfitting

Verlauf des MSE bei Trainings- und Testdaten

MSE Trainingsdaten :=
Modellparametrierung



MSE Testdaten :=
Fit unbekannter Daten



Wirkungsweise Erhöhung der Freiheitsgrade:

- stetige Verringerung des Trainings-MSE
- Test-MSE verläuft i.d.R. u-förmig mit Minimum oberhalb des unsystematischen Fehlers (---)
- Test-MSE approximiert den „*Out of Sample Fehler*“:

$$\text{Out of Sample Fehler} = E \left[\left(\hat{f}(x) - f(x) \right)^2 \right]$$

Konsequenz:

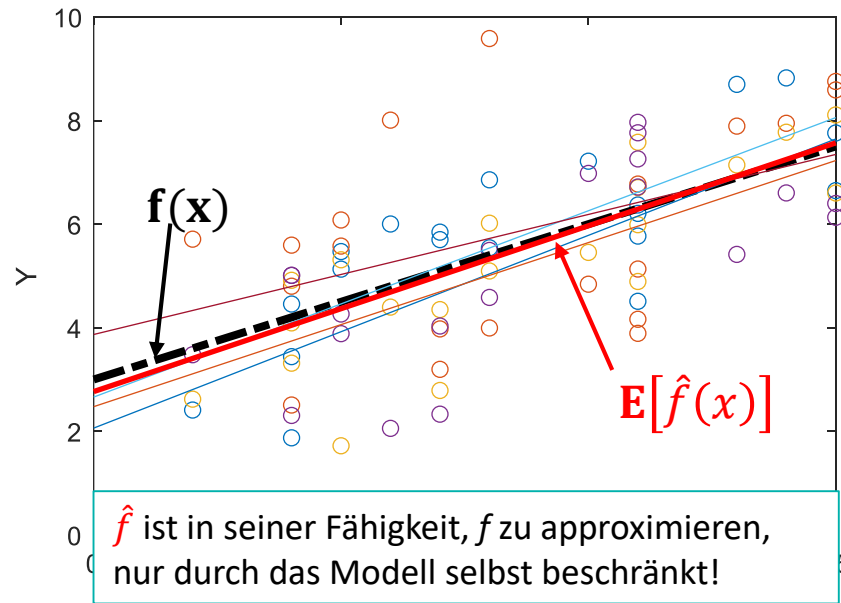
- Ab einer bestimmten Modellflexibilität setzt ein „Overfitting“ ein.
- Overfitting: Verschlechterung der Modellanpassung ggü. einfacheren Modellen.
- Begründung: es wird ein Teil des „unsystematischen“ Fehlers in dem Modell „erklärt“.

Bias und Varianz sind Gütekriterien für die Modellanpassung

BIAS (Verzerrung)

Der Bias misst die Abweichung zwischen **gemittelten** parametrisierten **Modell** $E[\hat{f}(x)]$ und dem **wahren Zusammenhang** $f(x)$

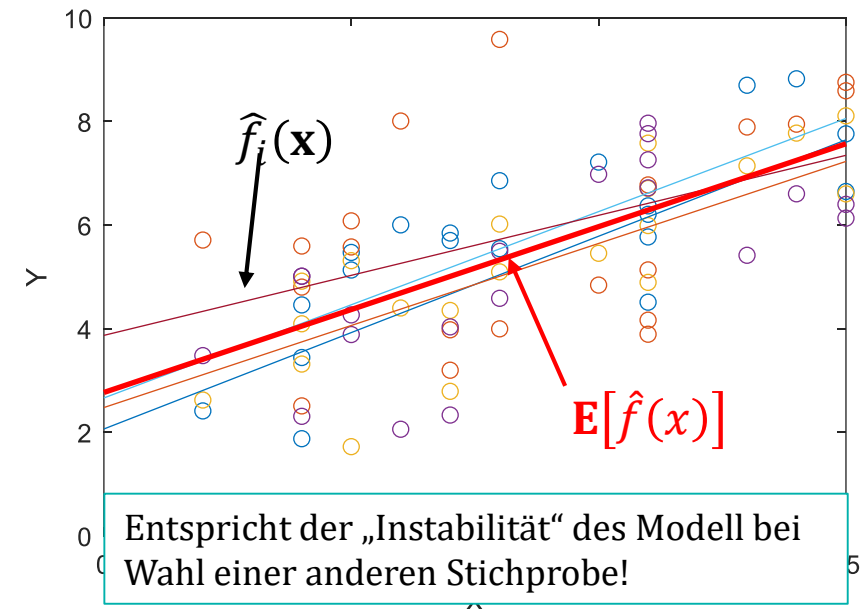
$$\text{Bias} = E[\hat{f}(x)] - f(x)$$



Varianz

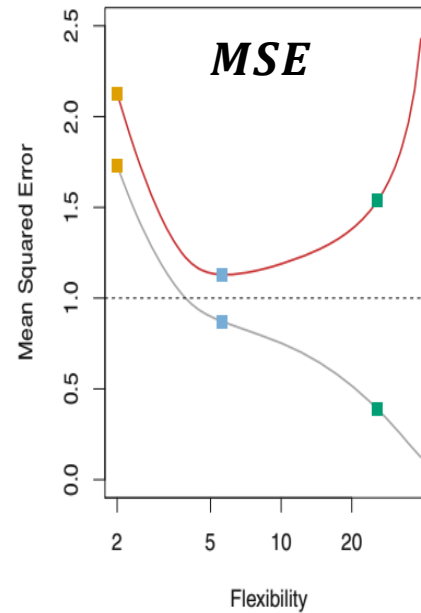
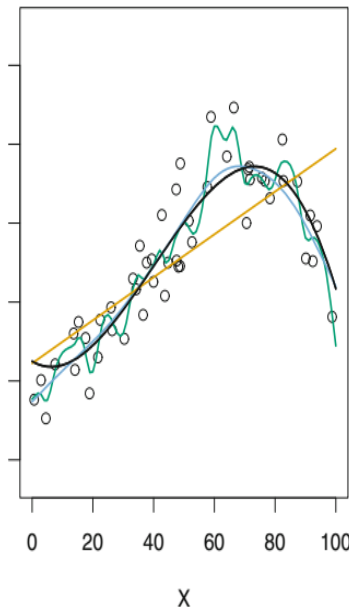
Die Varianz misst die **Variation** des parametrisierten Modells $\hat{f}(x)$ um das **gemittelte** Modell $E[\hat{f}(x)]$

$$\text{Varianz} = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$



Der Bias-Varianz Trade off hängt stark vom Problem und der gewählten Methodik ab

$$E[MSE] = \underbrace{E[\hat{f}(x)^2]}_{\text{Erwarteter Test-MSE}} - \underbrace{[E[\hat{f}(x)]]^2}_{\text{Varianz der Methodik}} + \underbrace{(E[\hat{f}(x)] - f(x))^2}_{BIAS^2} + \underbrace{Var(\epsilon)}_{\text{Unsystematischer Fehler}}$$



Bias-Varianz Trade off:

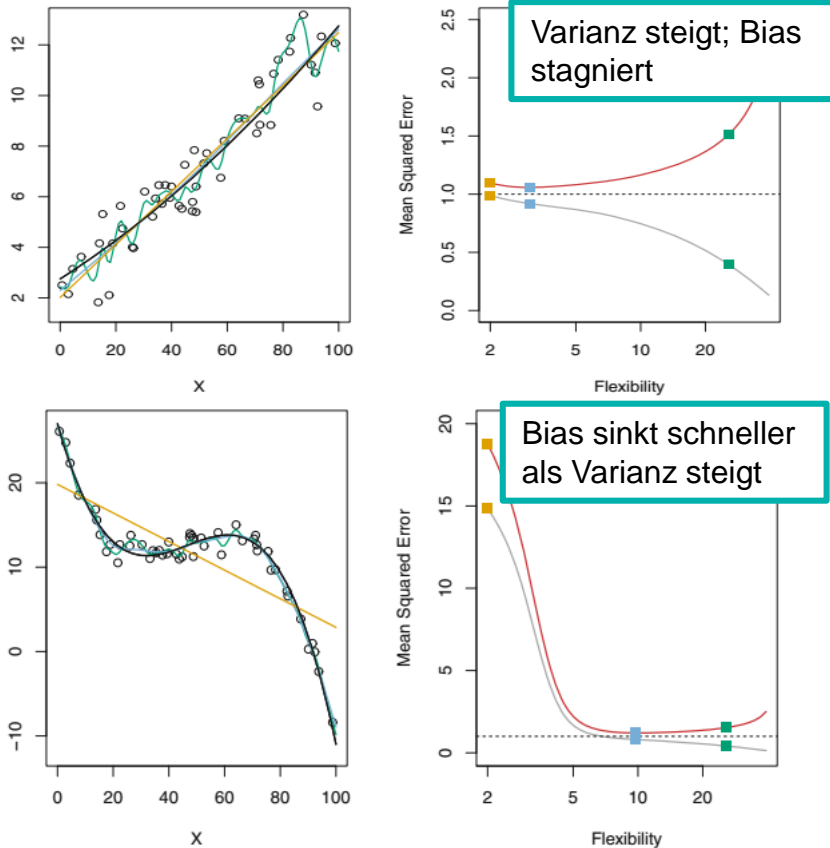
- Flexible Funktionen sind generell besser in der Lage den Bias zu verringern als unflexible Funktionen.
- Grundsätzlich haben flexiblere Methoden zur Abbildung der Funktion f eine höhere Varianz.

Nebenrechnung: $MSE = E[(\hat{f}(x) - f(x))^2] = E[\hat{f}(x)^2] + E[f(x)^2] - 2E[\hat{f}(x) \cdot f(x)] = E[\hat{f}(x)^2] - [E[\hat{f}(x)]]^2 + [E[\hat{f}(x)]]^2 + E[f(x)^2] - 2E[\hat{f}(x) \cdot f(x)]$

$E[(\hat{f}(x) - E[\hat{f}(x)])^2] = E[\hat{f}(x)^2] + [E[\hat{f}(x)]]^2 - 2E[\hat{f}(x) \cdot E[\hat{f}(x)]] = E[\hat{f}(x)^2] - [E[\hat{f}(x)]]^2$

Der Trainings-MSE und der Test-MSE verhalten ab einen gewissen Punkt gegensätzlich -> Zustand Overfitting II

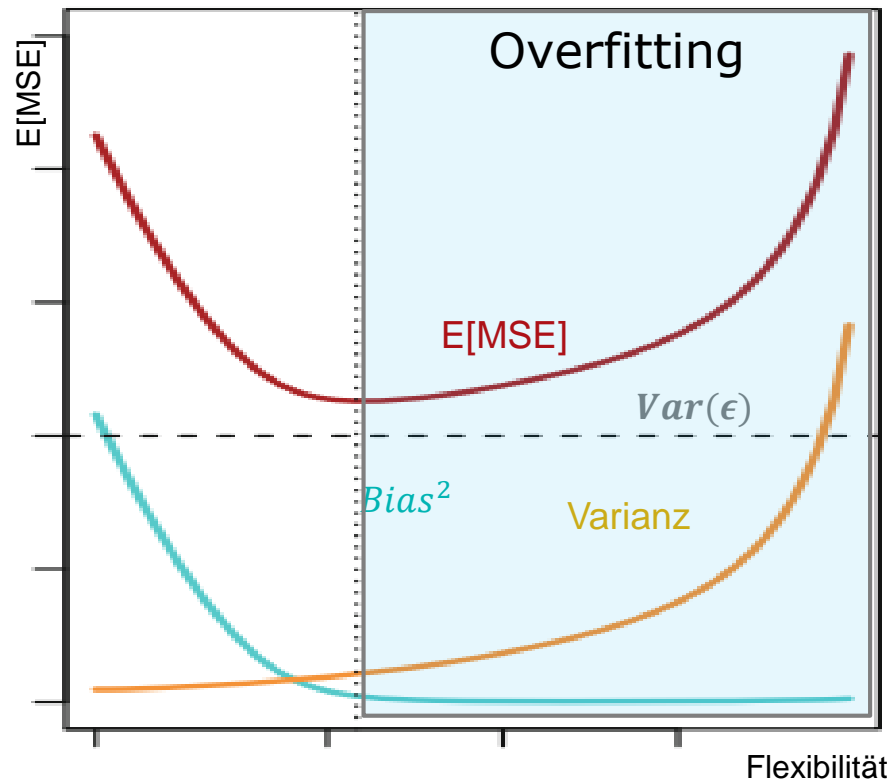
Grau MSE der Trainingsdaten
Rot MSE der Testdaten



- Die Änderungsrate von Bias und Varianz bestimmt die Veränderung des erwarteten MSE.
- Wenn wir die Flexibilität einer Klasse von Methoden erhöhen, neigt der Bias dazu, zunächst schneller zu sinken als die Varianz steigt.
- Infolgedessen sinkt der erwartete MSE (approximiert durch den Test-MSE).
- Allerdings hat die Erhöhung der Flexibilität ab einen bestimmten Punkt nur noch geringen Einfluss auf den Bias.
- Die Varianz erhöht sich ab diesem Punkt deutlich.

Gesucht wird das Minimum des Test MSE

Zusammenhang zwischen Bias, Varianz und $E[\text{MSE}]$



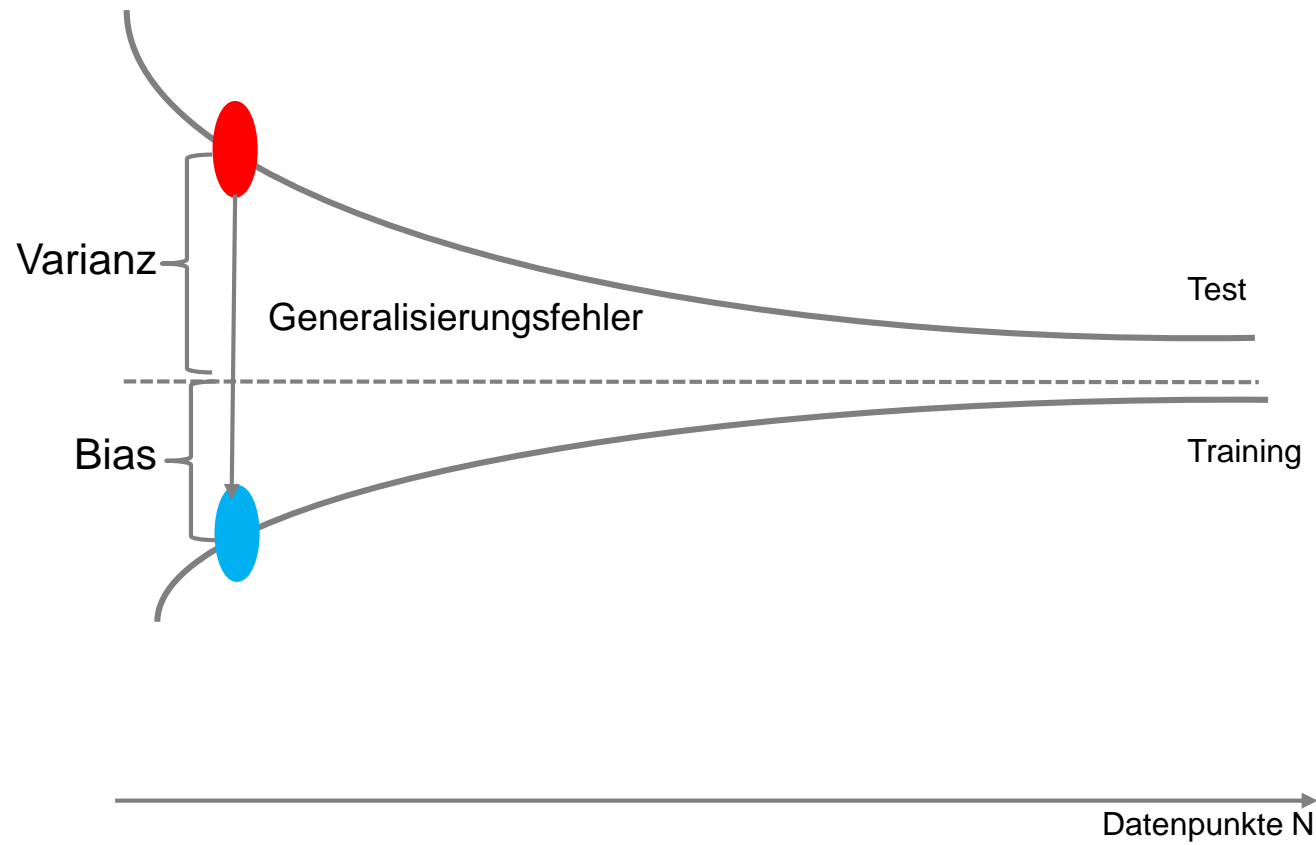
Erkenntnis:

- Der Erwartungswert des MSE kann nicht unterhalb des systematischen Fehlers liegen.
- Bei höherer Flexibilität wird dieser zunächst positiv und dann negativ (Overfitting) beeinflusst.

Zielsetzung:

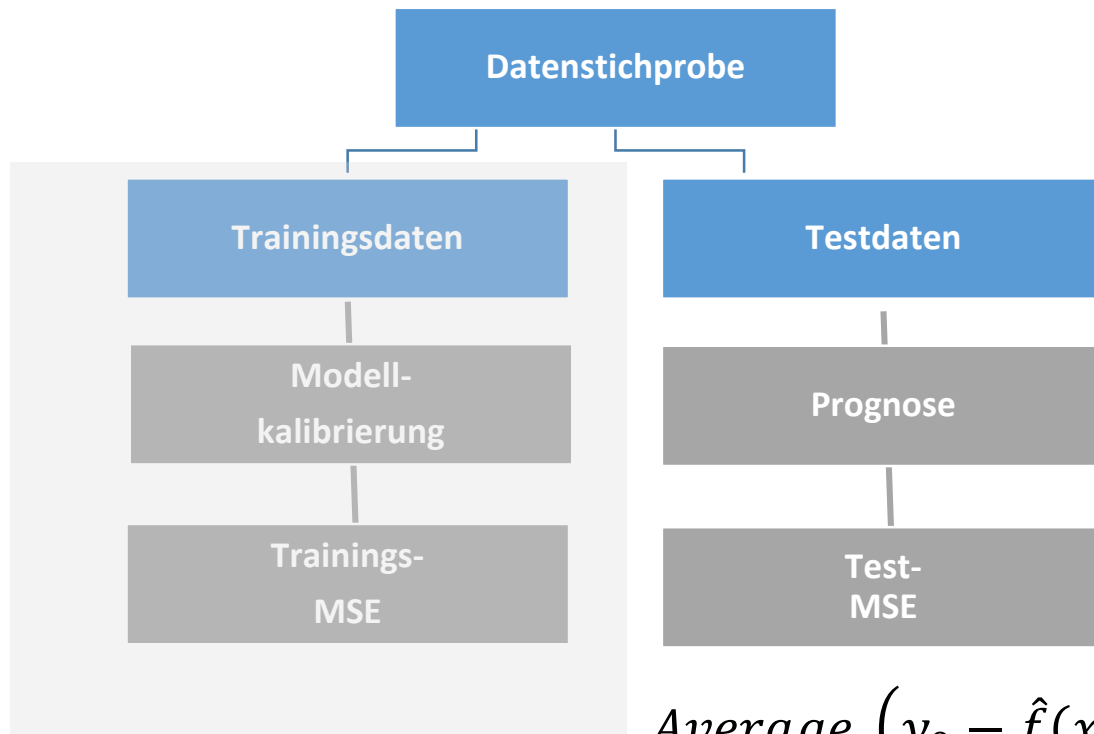
- Auffinden einer Funktion mit geringer Varianz und geringem Bias.
- In Realität ist der wahre funktionale Zusammenhang und damit auch die Varianz von f unbekannt, somit auch der Erwartungswert des Test-MSE.

Der Stichprobenumfang hat einen bedeutenden Einfluss auf die Varianz



Für den Einsatz des Modells ist nicht der Trainings-MSE, sondern die nachträgliche „Güte“ bei „unbekannten Daten“ (Testdaten) entscheidend

Zusammenhang zwischen X und Y



Erkenntnis:

- Für Modelleinsatz ist **nicht** entscheidend, wie gut das Modell die Daten wiedergibt, anhand dessen es kalibriert wurde! (Trainings-MSE)
- Entscheidend ist, wie gut das Modell mit „unbekannten“ Daten funktioniert (Test-MSE)
- Der Trainings-MSE ist kein guter Schätzer für den Test-MSE

Konsequenz:

- Der Trainings-MSE kann kein Kriterium für die Modellwahl darstellen.

Ziel bei Modellauswahl:

- **Modellwahl** anhand des geringsten **Test-MSE**

$$\text{Average } (y_0 - \hat{f}(x_0))^2$$

1	Bedeutung des Bias-Variance Trade off
2	Samplingtechniken
3	Best Subset selection Stepwise selection

Literaturempfehlung und Quellennachweise der Abbildungen: James, Gareth; Witten, Daniela; Hastie, Trevor; Tibshirani, Robert. An Introduction to Statistical Learning: with Applications in R (Springer Texts in Statistics), Springer New York. Kindle-Version.

Für die Schätzung des MSE der Testdaten gibt es zwei mögliche Vorgehensweise

Modellauswahl

Indirekte Schätzung des erwarteten MSE

Informationskriterien auf Basis Trainingsdaten

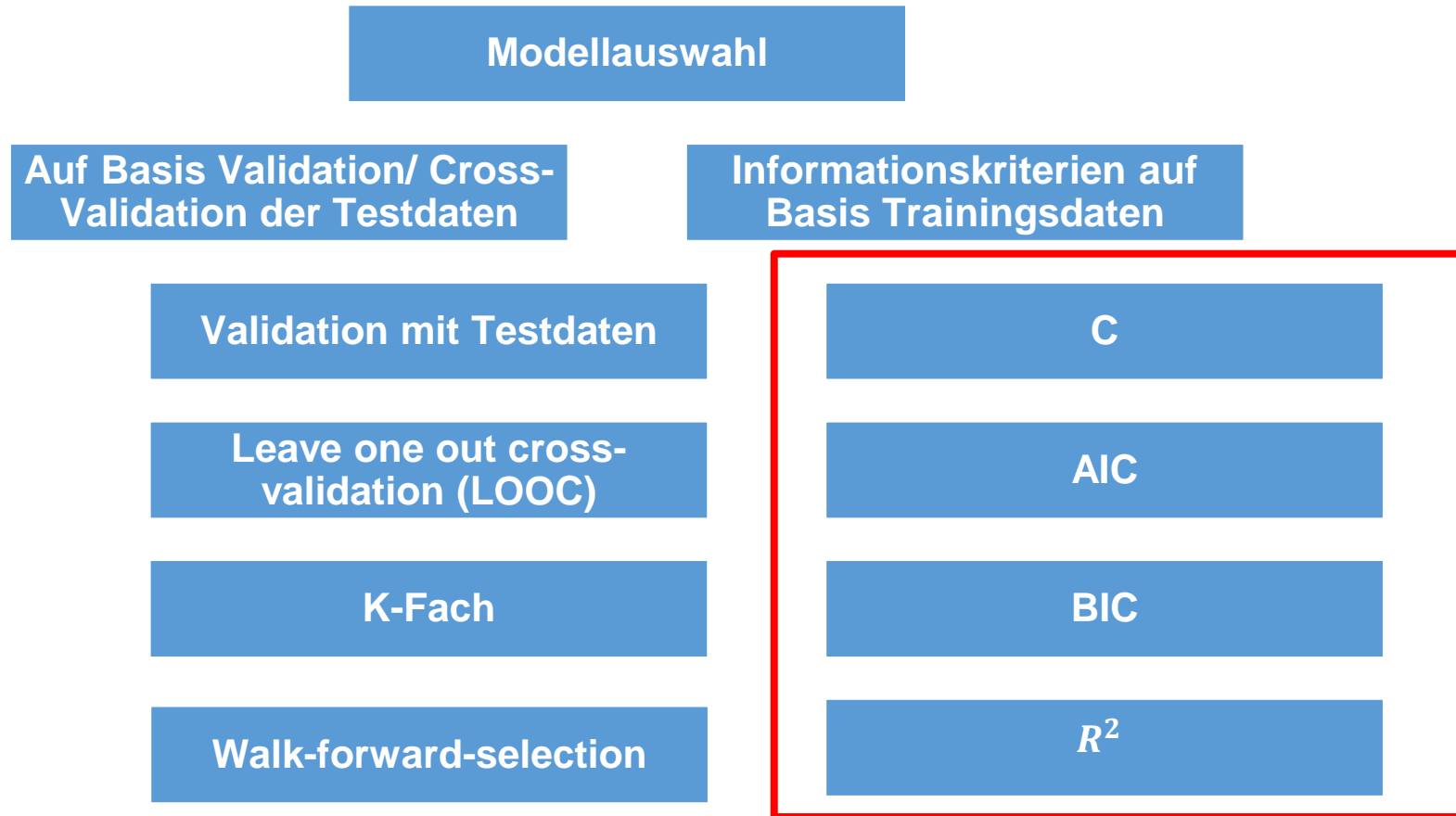
Der Trainings-MSE
wird mathematisch geeignet angepasst

Direkte Schätzung des erwarteten MSE

Auf Basis Validation/ Cross-Validation der Testdaten

Vom Trainingsdatensatz werden einzelne
Datenpunkte der Parametrierung vorenthalten
(Testdaten), um nach erfolgter Schätzung die
geschätzte Funktion auf diese anzuwenden.
Test-MSE

Die Modellauswahl kann durch direkter Betrachtung des Test-MSE oder Approximation des Test-MSE erfolgen



Der Einsatz von Informationskriterien kann bei bestimmter zugrunde liegender Verteilung hinreichende Rückschlüsse bieten

■ Zielsetzung von Informationskriterien

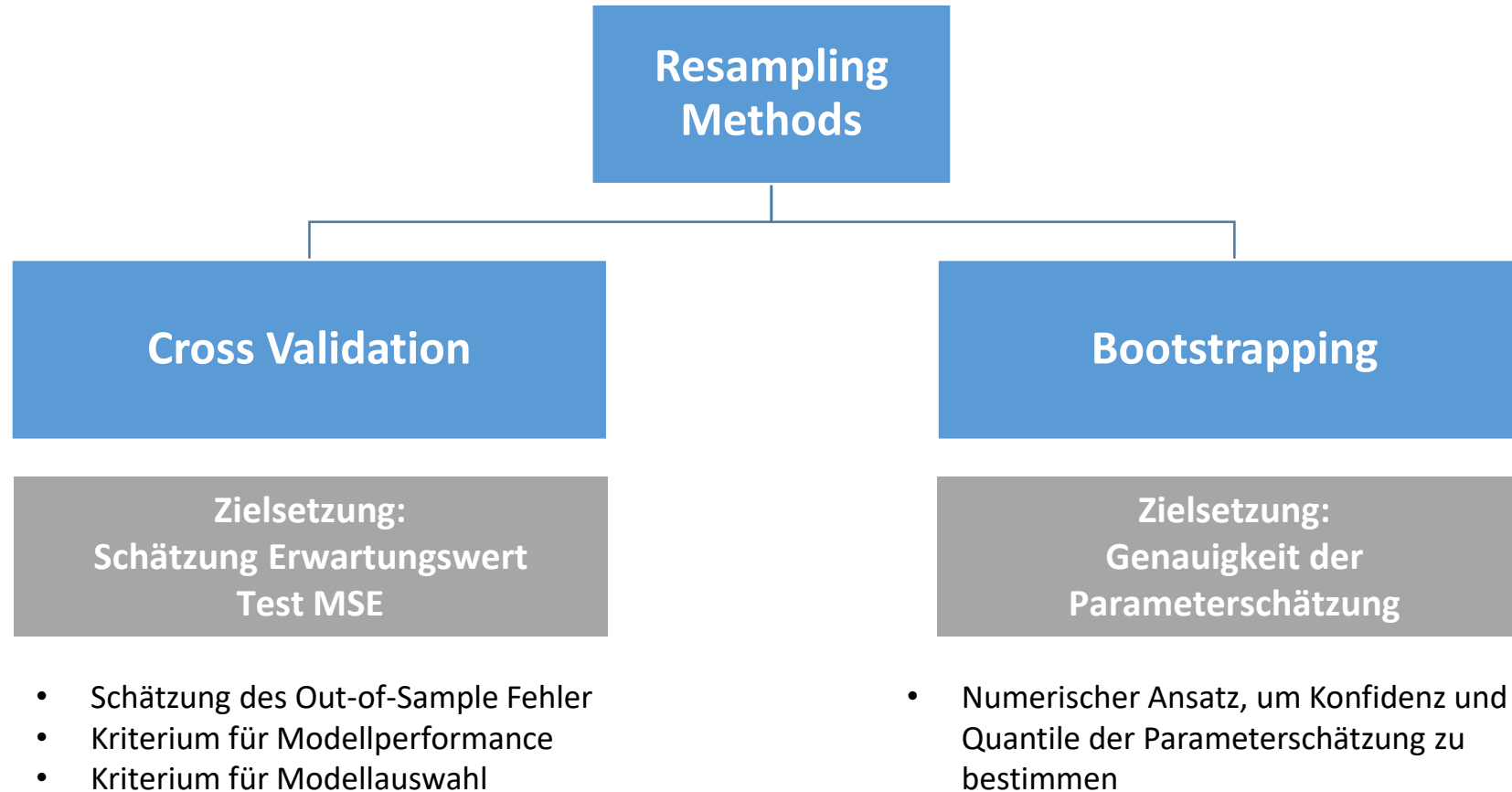
- Auffinden des optimalen Fits bzgl. Bias und Variance anhand eines Informationskriterium.
 - Als Anpassungsgüte wird i.d.R. die Trainings- MSE oder RSS verwendet ($RSS = MSE \cdot n$)
 - Ein Strafterm für die Anzahl der Parameter d berücksichtigt, um Overfitting zu verhindern.
- Es wird ein möglichst geringer Wert des Informationskriteriums angestrebt

Informationskriterien	Mathematische Formulierung
C_p	$C_p = MSE + 2d\hat{\sigma}^2/n$
Akaike Informationskriterium	$AIC = \frac{1}{\hat{\sigma}^2} (MSE + 2d\hat{\sigma}^2/n)$
Bayessches Informationskriterium	$BIC = MSE + \log(n)d\hat{\sigma}^2/n$

d = Anzahl Regressoren, $\hat{\sigma}^2$ Schätzer der Varianz des unsystematischen Fehlers ϵ , für $n > 7 \rightarrow \log(n) > 2$

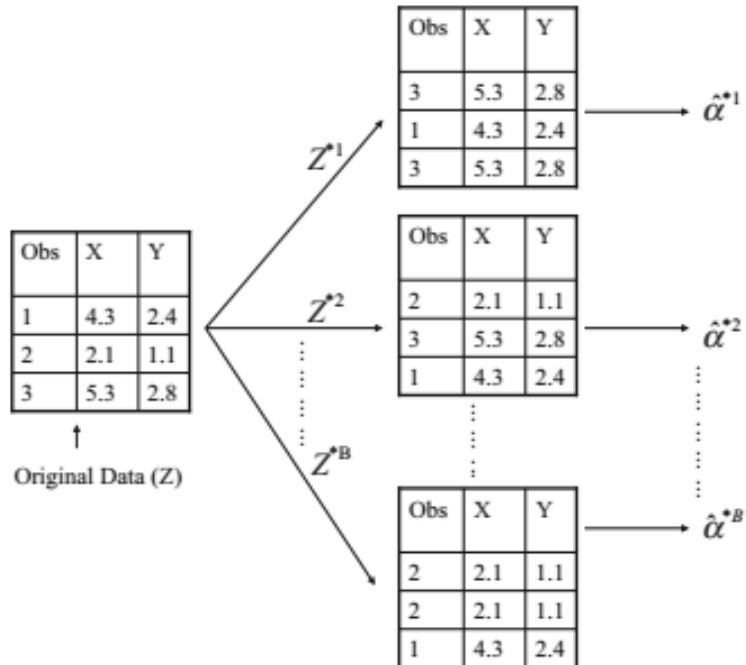
RSS = Residual squared sum; n = Anzahl Datenpunkte; d = Anzahl Parameter

Wie generieren wir Testdaten: Resampling Methoden stellen eine Möglichkeit dar



Bootstrapping als allgemeingültiger Ansatz zur Berechnung von Unsicherheit

Aufbau zufälliger Datensamples: Ziehen aus Stichprobe mit zurücklegen



Zielsetzung und Vorgehensweise

Zielsetzung:

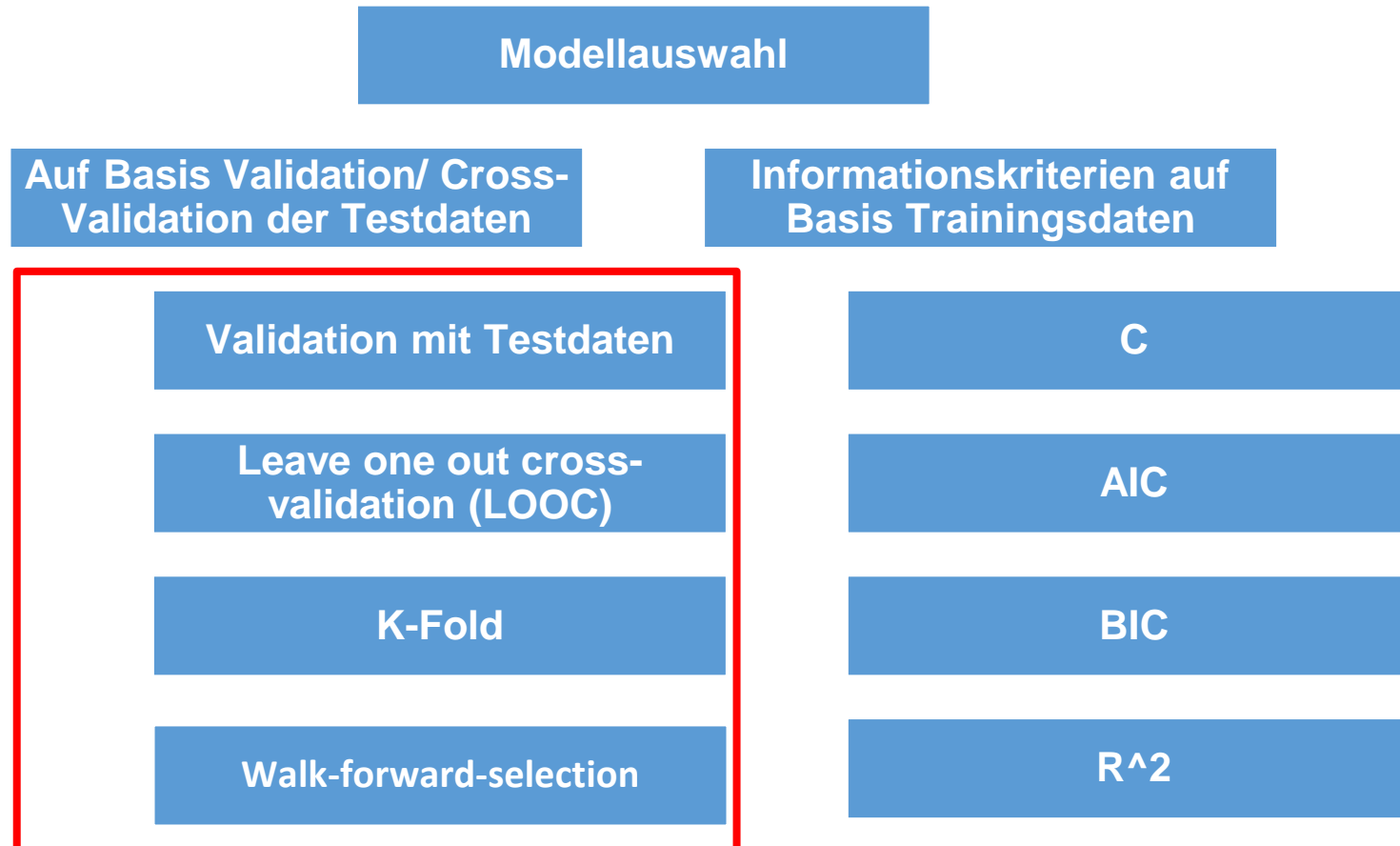
- Quantifizierung der Unsicherheit, die mit einer Parameterschätzung verbunden ist
 - Z.B. Standardabweichung der Parameterschätzung

Vorgehensweise

- Ziehen mit Zurücklegen aus der Datenstichprobe zur Generierung von n neuen Datensamples
- Parametrierung auf Basis eines jeden Datensamples

$$\text{std}(\hat{\boldsymbol{\varphi}}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^B (\hat{\boldsymbol{\varphi}}_r - \overline{\hat{\boldsymbol{\varphi}}})^2}$$

Die Modellauswahl kann durch direkte Betrachtung des Test-MSE oder Approximation des Test-MSE erfolgen



Durch die Aufteilung der Datenstichprobe in Trainings- und Testdaten sollen letztere als Schätzer für den Erwartungswert des MSE verwendet werden



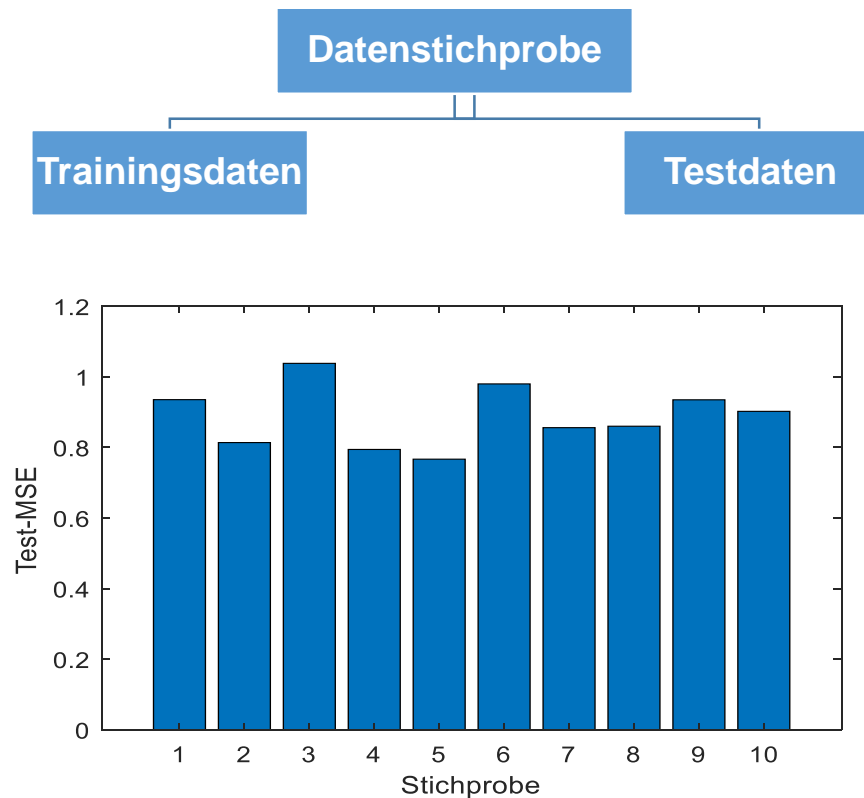
- Schätzung des Out-of-Sample Fehler

$E[\text{Test MSE}]$

$$\text{Average } \left(y_0 - \hat{f}(x_0) \right)^2$$

Der Validierungsset-Ansatz versucht anhand der Testdaten den Test-MSE zu schätzen; besitzt aber Nachteile insbesondere bei kleinen Stichproben

Ergebnis des Test-MSE nach 10-maliger Unterteilung von Trainings- und Testdaten



- Der Validierungsset-Ansatz unterteilt die Datenstichprobe in Trainings- und Testdaten

Funktion der Trainingsdaten:

- Schätzung funktionaler Zusammenhang

Funktion der Testdaten:

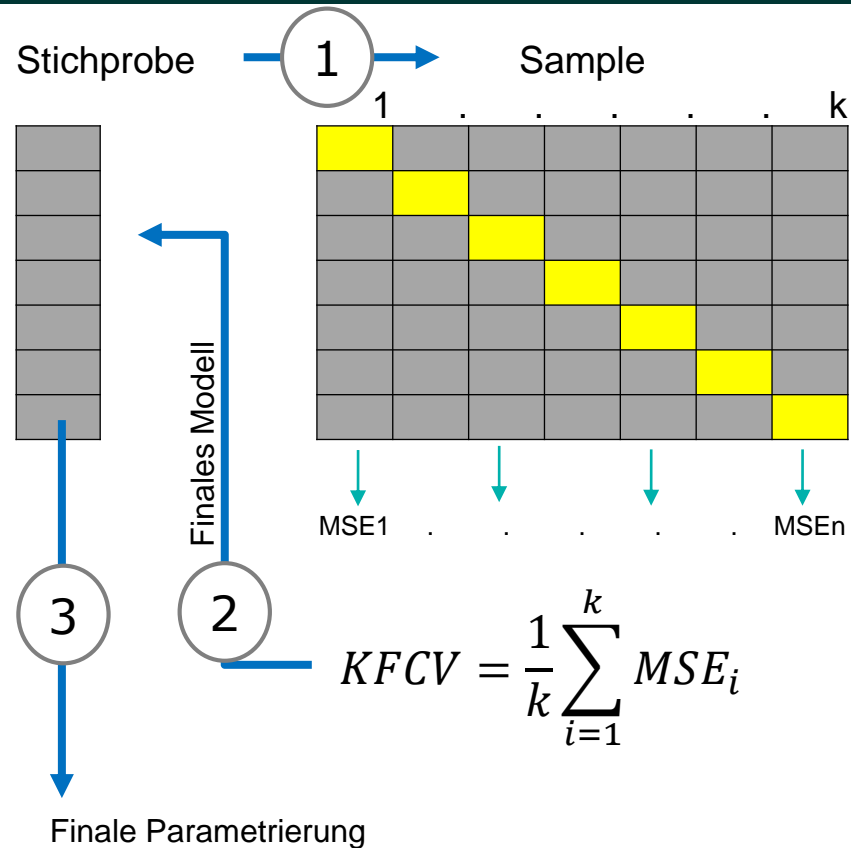
- Modellbewertung
- Schätzung des Test-MSE
- Modellauswahl

Nachteile:

- Schätzer für Test-MSE variiert in Abhängigkeit der zufällig gewählten Testdaten.
- Test-MSE wird bei kleiner Training-Stichprobe überschätzt.

Cross-Validierungs-Ansatz k-fold (KFCV) verringert den numerischen Aufwand ggü. LOOCV

Unterteilung der Datenstichprobe in k-Bereiche



Funktionsweise:

- Erstellung k disjunkter zufälliger Segmente der Datenstichprobe

1. Kalibrierung Modell i exkl. i-ten Segment
2. Prognose i-ten Segment mit Modell i

$$\hat{y}_i = f_i(x_i) \quad x_i \in \text{Segment}_i$$

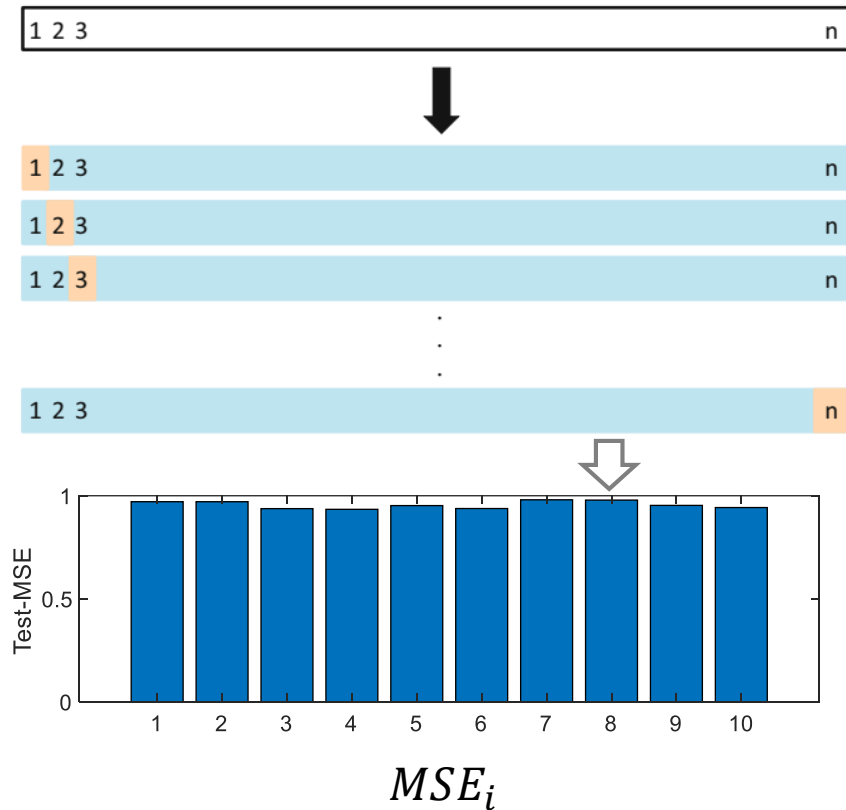
3. Berechnung des quadrierten Fehlers

$$MSE_i = \frac{1}{n_i} \sum_{i=1}^{n_i} (y_i - \hat{y}_i)^2$$

4. Berechnung des Cross Evaluation auf Basis aller durchlaufenden Testsets.
5. Wahl des Modells mit kleinsten KFCV

Leave one out Cross Validation LOOCV stellt der Parametrierung deutlich mehr Daten zur Verfügung

Unterteilung der Datenstichprobe



Funktionsweise:

1. Das Training Set =
Stichprobe exklusive einer Beobachtung (x_i, y_i)

2. Prognose des Testdatenpunktes x_i

$$\hat{y}_i = f_i(x_i)$$

3. Berechnung des quadrierten Fehlers

$$MSE_i = (y_i - \hat{y}_i)^2$$

■ Berechnung des Cross Evaluation auf Basis aller durchlaufenden Testsets.

$$LOOCV = \frac{1}{n} \sum_{i=1}^n MSE_i$$

Beurteilung:

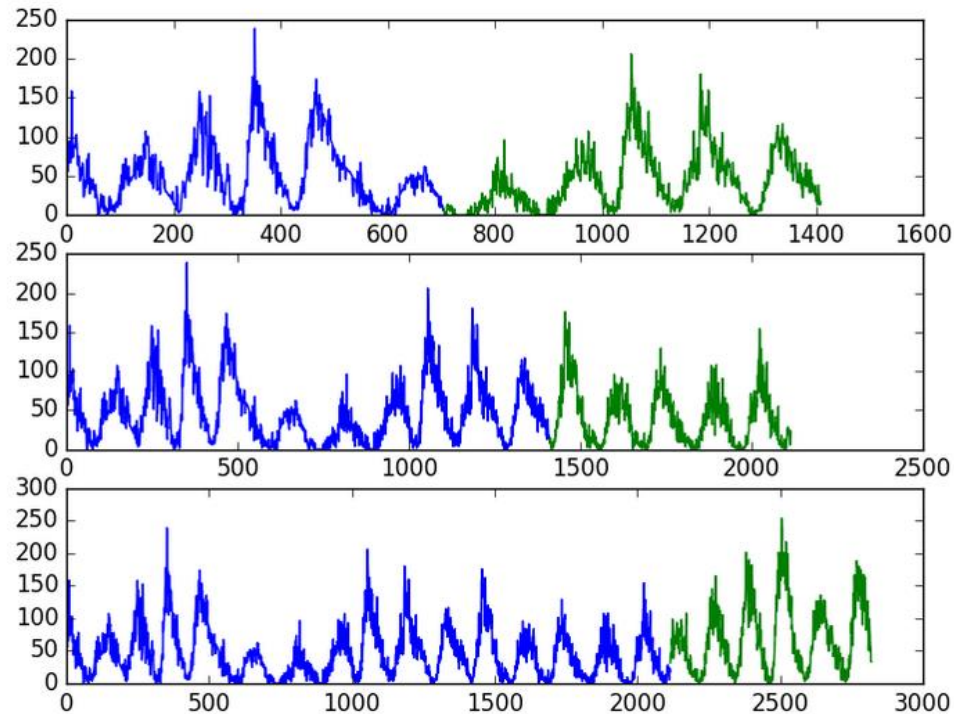
■ Vorteil: große Trainingsstichprobe, stabiler Schätzer (LOOCV) des Test-MSE

■ Nachteil: Performance

Für Zeitseriendaten sind die bisher besprochenen Unterteilungen ggf. nicht gut geeignet

Methoden für Zeitserien

TimeSeriesSplit: Verlängerung der Trainingsdaten bei konstant gehaltenem Testdatensatz



- Zeitseriendaten können zur Generierung von Training, Validierung und Testdaten nicht zufällig getrennt werden.
- Hierdurch würde die zeitliche Kopplung (autokorrelation) aufgehoben.
- Bisher betrachtete Lösungsmöglichkeit:
 - Teilgruppen müssen zeitliche Ordnung behalten (bei scikit learn mit dem Zusatz `shuffle='None'` erreichbar)
 - Sukzessive Verlängerung des Trainingsdatensatzes bei konstant gehaltenem Testdatensatz

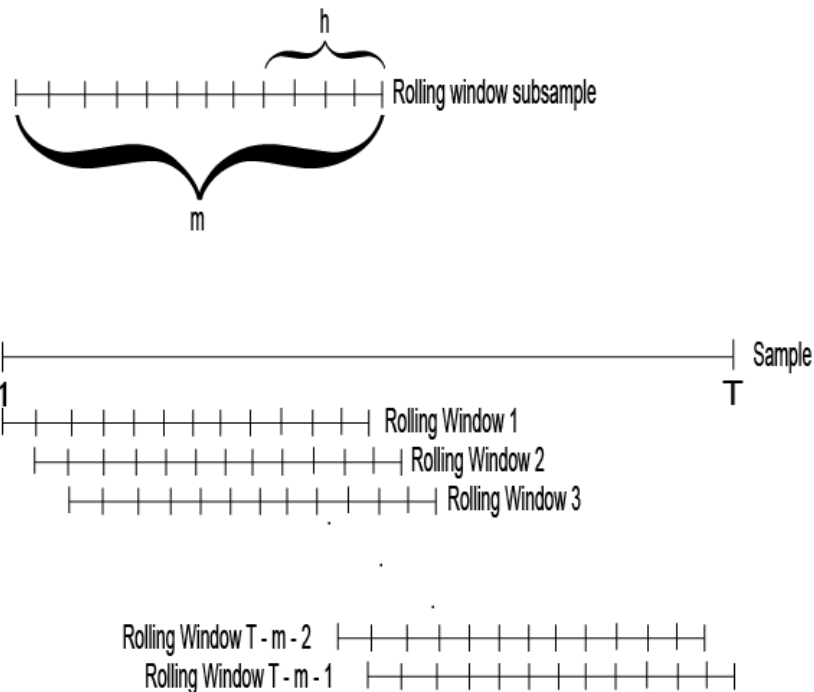
```
class sklearn.model_selection.TimeSeriesSplit(n_splits=5, *,
max_train_size=None, test_size=None, gap=0)
```

Provides train/test indices to split time series data samples that are observed at fixed time intervals, in train/test sets. In each split, test indices must be higher than before, and thus shuffling in cross validator is inappropriate.

- Alternative Walk-forward-Validation
 - Hierbei erfolgt eine kontinuierliches Update des Modells sobald neue daten verfügbar sind.

Die walk-forward-Evaluation ist ein iterativer Prozess und ein guter Prozess in der Praxis zur Aktualisierung eines Modells

walk-forward evaluation



rolling window size, m , forecast horizon, h .

1. Es wird ein Fenster (Samplesize) m definiert.
2. Die erste Modellparametrierung basiert auf: $[x_1, x_m]$
3. Erstellung einer Prognose \hat{y}_{m+h} für den Prognosehorizont h .
4. Berechnung des Fehlers: $e_1 = \hat{y}_{m+h} - y_{m+h}$.
5. Das Fenster wird um einen Zeitschritt verschoben und der Prozess wird wiederholt.
6. Berechnung des Test-RMSE aus den Einzelfehlern

$$RMSE = \sqrt{\sum_{i=1}^{T-m-h} e_i^2}$$

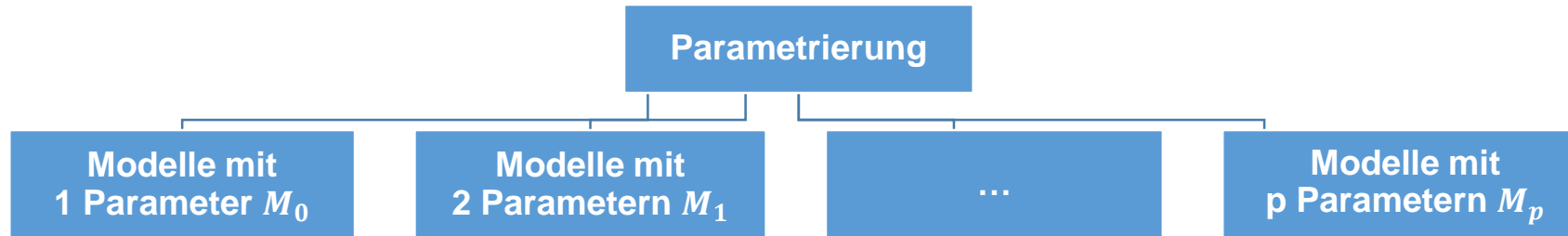
Notwendige Parameter:

- Minimale Anzahl an Beobachtungen zu Beginn m
- Prognose-Horizont h

1	Bedeutung des Bias-Variance Trade off
2	Generierung von Schätzern für den Test-RMSE
3	Best Subset selection Stepwise selection

Literaturempfehlung und Quellennachweise der Abbildungen: James, Gareth; Witten, Daniela; Hastie, Trevor; Tibshirani, Robert. An Introduction to Statistical Learning: with Applications in R (Springer Texts in Statistics), Springer New York. Kindle-Version.

Best Subset Selection ist ein standardisierter Suchalgorithmus zur Bestimmung des besten Modells im Fall einer multiplen Regression



1. Startpunkt: Modell M_0 enthält keinen Regressor \rightarrow *Prognose* = Mittelwert der Datenstichprobe
2. Für $k = 1, 2, \dots, p$
 - a) Parametrierung aller möglichen Modelle mit genau k Variablen (Regressoren)
 - b) Wahl des besten Modell M_k mit k Variablen entsprechend des kleinsten Trainings-MSE oder R^2
3. Wahl des besten Modells jeder Klasse aus M_0, \dots, M_p durch Verwendung Cross Validation.

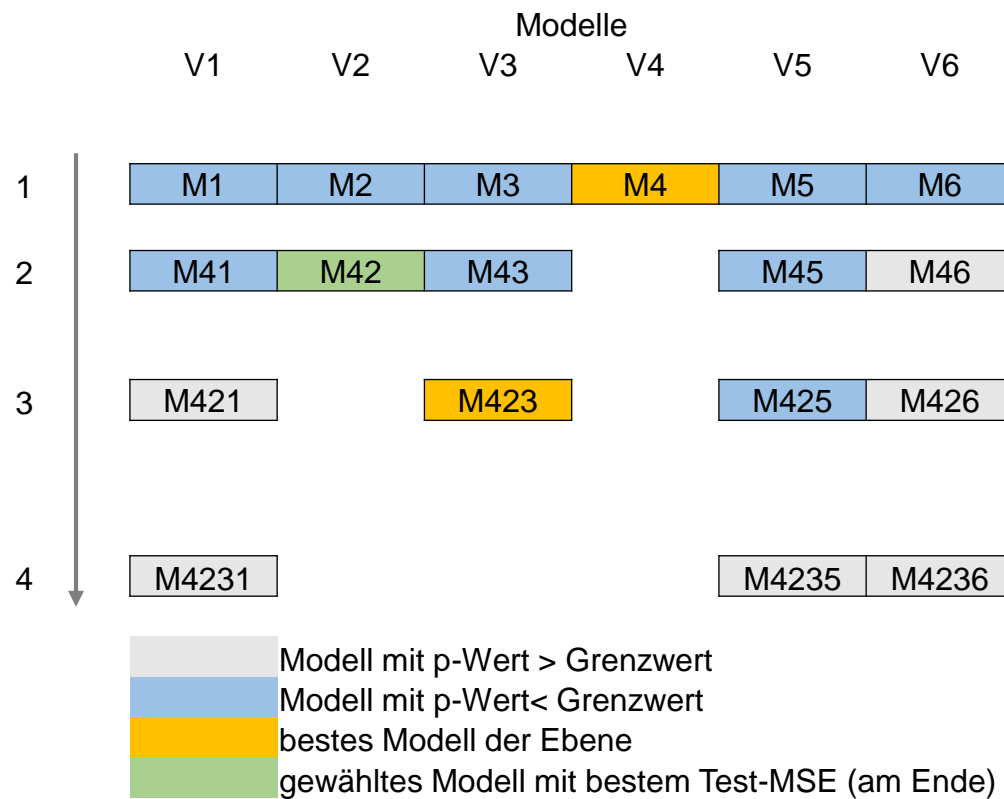
Erläuterung:

Schritt 2 identifiziert das beste Modell (auf Basis der Trainingsdaten) für jede Teilmengengröße, um das Problem von $\binom{p}{2}$ möglichen Modelle auf eines der $p + 1$ möglichen Modelle zu reduzieren.

In Schritt 3 wird Cross Validation, C_p , BIC oder R^2 verwendet, um zwischen M_0, M_1, \dots, M_p auszuwählen.

Forward Stepwise Selection (FSS) ist ein inkrementelles Suchverfahren des besten Modells mit p Variablen

Abbildung der Modellsuche



- **FSS** betrachtet schrittweise einen kleineren Satz von Modellen.
- Auf der 1.ten Ebene werden p -Modelle mit jeweils 1 Regressor miteinander verglichen
- Die Variable des besten Modells wird beibehalten.
- In jedem Schritt wird die Variable, die die größte zusätzliche Verbesserung der Passform bewirkt, dem Modell hinzugefügt.
- Zur Wahl stehen nur Modelle mit signifikanten Parametern
- Am Ende wird mit Hilfe des Test-MSE das beste Modell gewählt

Problem: Auftreten von lokalen Minima

Suchmethoden helfen dabei den Modellauswahlprozess in die Schätzung der Parameter zu integrieren

Modellhypothesen*

Modelle im CV

Klassische Modellauswahl

$$p \cdot \frac{p+1}{2}$$

Grundsätzliche Schwierigkeit
bei allen Ansätzen:

Es müssen alle in Frage
kommenden Modellhypothesen
separat geschätzt werden

Best subset Selection

$$\text{maximal} \\ p \cdot \frac{p+1}{2}$$

$$p+1$$

Bildung Modellhypothese
(**Vorauswahl** relevanter
Parameter)

Forward Stepwise Selection

$$\text{maximal} \\ p$$

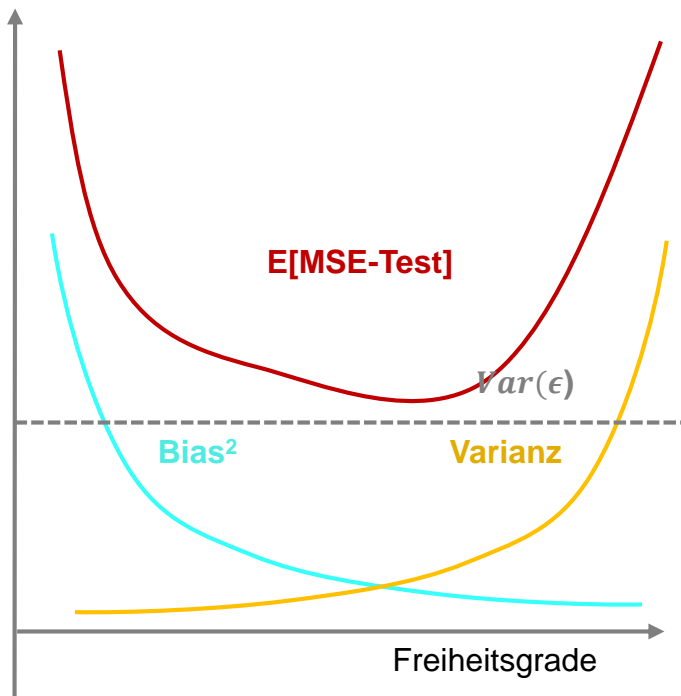
Risiko:

Modelle mit potentiell niedrigen
 $E[\text{Test-MSE}]$ bleiben
unberücksichtigt

*Annahme jeder Regressor wird nur maximal einmal in das Modell mitaufgenommen

Zusammenfassung

Im Rahmen der Modell-kalibrierung existiert ein trade-off zwischen Bias und Varianz



Es existieren direkte und indirekte Methoden zur Schätzung des erwarteten Test- MSE

Best Subject und Forward stepwise Selection verringern den Suchraum

