

# Teil IX

## Modellbasierte Optimierung

## 9 Zufallsfeld

Dieses Kapitel ist als Nachtrag zu Kapitel 3 gedacht.

- Wir werden modellbasierte Optimierung kennenlernen.
- Idee ist es, teure Auswertungen einer unbekannten Zielfunktion (black box) möglichst zu minimieren.
- Zunächst jedoch einige Grundlagen.

## 9.1 Zufallsfeld

- Ein *stochastischer Prozess* ist eine Familie von Zufallsvariablen  $\{Z_i \mid i \in I\}$  mit  $Z_i : (\Omega, \mathcal{A}, P) \mapsto (E, \mathcal{G})$ , wobei  $I \neq \emptyset$  eine Indexmenge,  $(\Omega, \mathcal{A}, P)$  ein Wahrscheinlichkeitsraum und  $(E, \mathcal{G})$  ein Zustandsraum ist [115].
- Die Indexmenge wird z.B. bei zeitlichen Prozessen als  $I \subseteq \mathbb{Z}$  oder  $I \subseteq \mathbb{N}$  im diskreten oder als  $I \subset \mathbb{R}$  im kontinuierlichen Fall festgelegt.
- Bei räumlichen Prozessen ist die Indexmenge häufig eine Teilmenge des  $\mathbb{R}^d$  ( $d \geq 1$ ) und der entsprechende Prozess wird auch als Zufallsfeld bezeichnet.
- Im Folgenden betrachten wir Zufallsfelder (ZF) der Form

$$\{Z(x) \mid x \in \mathbb{D} \subseteq \mathbb{R}^d, d \geq 1\}. \quad (1)$$

- Ein ZF ist ein *Gauß'sches ZF*, wenn für beliebige  $x_1, \dots, x_n \in \mathbb{D}$  der Vektor der Zufallsvariablen  $Z = (Z(x_1), \dots, Z(x_n))^T$  multivariat normalverteilt ist und sich damit vollständig durch  $E[Z(x)]$  und  $\text{Cov}[Z(x), Z(x')]$  charakterisieren lässt [109].

## 9.1 Zufallsfeld

- Ein ZF ist schwach stationär (mittelwert- und kovarianzstationär) wenn die beiden ersten Momente existieren und es gilt  $\forall x, x' \in \mathbb{D}$ :

$$E[Z(x)] = \mu \quad (2)$$

$$\begin{aligned} \kappa(x, x') &= \text{Cov}[Z(x), Z(x')] = \text{Var}[Z(x)] \text{Cor}[Z(x), Z(x')] \\ &= \sigma^2 \rho(x - x', \psi) \end{aligned} \quad (3)$$

Damit hängt die räumliche Variabilität des ZF nur von der Differenz  $x - x'$  und der Form der Korrelationsfunktion  $\rho$ , sowie deren Parametern  $\psi = (\psi_1, \dots, \psi_d)^T$  mit  $\psi_j > 0$  ( $j = 1, \dots, d$ ) ab [140].

- Annahme: Eigenschaften des Prozesses über alle  $x \in \mathbb{D}$  hinweg konstant.
- Bei Gaußprozess folgt aus schwacher auch die starke Stationarität, oder hier einfach Stationarität.
- Ist die Kovarianzfunktion

$$\kappa(\delta) = \kappa(x, x') \quad \forall x, x' \in \mathbb{D} : \|x - x'\| = \delta \quad (4)$$

nur vom Euklidischen Abstand  $\|\cdot\|$  zweier Punkte abhängig, so ist es ein isotropes ZF und damit rotationsinvariant.

- Damit sind stationäre isotrope ZFs bewegungsinvariant [173, 64].

## 9.1 Zufallsfeld: Kovarianz- und Korrelationsfunktion

Unter Berücksichtigung der schwachen Stationarität und Isotropie gilt für die Kovarianzfunktion  $\kappa$ :

$$(\kappa 1) \quad \kappa(x, x) = \sigma^2 \rho(x - x, \psi) = \sigma^2 = \text{Var}[Z(x)] \geq 0 \quad \forall x \in \mathbb{D}$$

$$(\kappa 2) \quad \kappa(x, x') = \kappa(x', x) \quad \forall x, x' \in \mathbb{D} \quad (\text{Symmetrie})$$

$$(\kappa 3) \quad |\kappa(x, x')| \leq \kappa(x, x) \quad \forall x, x' \in \mathbb{D} \quad (\text{Beschränktheit})$$

$$(\kappa 4) \quad \kappa \text{ ist positiv semidefinit.}$$

Da  $\rho(x - x', \psi) = \frac{\kappa(x, x')}{\kappa(x, x)}$  mit  $\rho(x - x, \psi) = \rho(0, \psi) = 1$ , sofern  $\kappa(x, x) > 0$ , gelten die Eigenschaften  $(\kappa 2)$ ,  $(\kappa 3)$  und  $(\kappa 4)$  entsprechend auch für die Korrelationsfunktion.

## 9.1 Zufallsfeld: Kovarianz- und Korrelationsfunktion

Mögliche Wahl der Korrelationsfunktion:

- Exponentielle Korrelationsfunktion:

$$\rho(x - x', \psi) = \exp(-\|x - x'\|_{\Delta}) \quad (5)$$

- Distanz zwischen  $x$  und  $x'$  wird über den Abstand

$\|x - x'\|_{\Delta} = \sqrt{(x - x')^T \Delta^{-1} (x - x')}$  mit Diagonalmatrix

$\Delta = \text{diag}(\psi_1^2, \dots, \psi_d^2)$  und  $\psi_1 = \dots = \psi_d$  gemessen.

- Da Parameter über alle Dimensionen gleich sind, hängt die Korrelationsfunktion und somit die Kovarianzfunktion nur vom Abstand der Punkte ab und ist damit isotrop.

- anisotrope Matérn 5/2 Korrelationsfunktion:

$$\rho(x - x', \psi) = \left[ 1 + \sqrt{5} \|x - x'\|_{\Delta} + \frac{5}{3} \|x - x'\|_{\Delta}^2 \right] \exp(-\sqrt{5} \|x - x'\|_{\Delta}) \quad (6)$$

mit  $\|x - x'\|_{\Delta} = \sqrt{(x - x')^T \Delta^{-1} (x - x')}$ ,  $\Delta = \text{diag}(\psi_1^2, \dots, \psi_d^2)$ .

- Der Skalierungsparameter  $\psi$  legt fest, wie schnell die Korrelation bzw. Kovarianz zwischen zwei Punkten mit steigender Distanz gegen Null läuft.
- Daher wird  $\psi$  oftmals als Reichweite interpretiert.

## 9.1 Zufallsfeld: Kovarianz- und Korrelationsfunktion

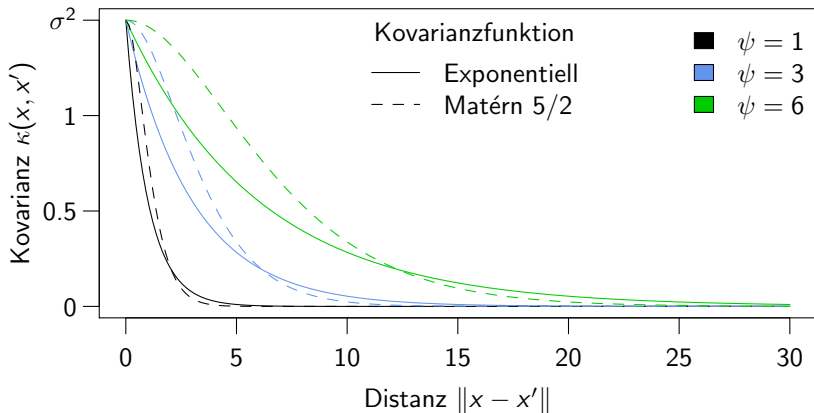


Abb. 9.1 : Matérn 5/2 und exponentielle Kovarianzfunktion mit  $\sigma^2 = 1.5$  und unterschiedlichen Reichweiten  $\psi$  für  $d = 1$ .

## 9.1 Zufallsfeld für verrauschte Beobachtungen

Zufallsfeld mit unkorreliertem Anteil:  $Z_\varepsilon(x) = Z(x) + \varepsilon(x)$

- $E[\varepsilon(x)] = 0$
- $\text{Cov}[\varepsilon(x), \varepsilon(x')] = \sigma_\varepsilon^2 \mathbb{1}_{x=x'}(x, x')$
- $\text{Cov}[Z(x), \varepsilon(x')] = 0$
- $\text{Cov}[Z_\varepsilon(x), Z_\varepsilon(x')] = \kappa_\varepsilon(x, x') = \sigma^2 \rho(x - x', \psi) + \sigma_\varepsilon^2 \mathbb{1}_{x=x'}(x, x')$

Die Varianz von  $Z_\varepsilon(x)$  ist somit um  $\sigma_\varepsilon^2$  größer als die von  $Z(x)$ .

Aus der nächsten Abbildung ist zu sehen, dass sich für  $Z(x)$  bzw.  $Z_\varepsilon(x)$  je nach Wahl der Parameter viele unterschiedliche Prozessverläufe ereignen können.



## 9.1 Zufallsfeld für verrauschte Beobachtungen

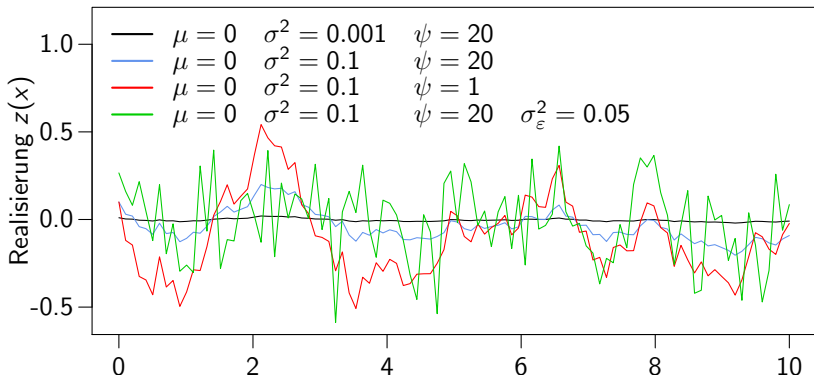


Abb. 9.2 : Realisierungen des Zufallsfeldes  $\{Z(x) \mid x \in [0, 10]\}$  für unterschiedliche Modellparameter mit exponentieller Kovarianzfunktion.

## 9.1 Zufallsfeld: Schätzung der Parameter

Um aus den Realisierungen  $z_\varepsilon(x_1), \dots, z_\varepsilon(x_n)$  die Parameter zu schätzen, wird

- das *Variogramm*

$$\begin{aligned}\nu(\delta) = \nu(x, x') &= \text{Var}[Z_\varepsilon(x) - Z_\varepsilon(x')] \\ &= \text{Var}[Z_\varepsilon(x)] + \text{Var}[Z_\varepsilon(x')] - 2\text{Cov}[Z_\varepsilon(x), Z_\varepsilon(x')] \\ &= 2\sigma^2[1 - \rho(x - x', \psi)] + 2\sigma_\varepsilon^2[1 - \mathbb{1}_{x=x'}(x, x')]\end{aligned}\tag{7}$$

- bzw. das *Semivariogramm*  $\nu_{\text{semi}}(\delta) = \frac{1}{2}\nu(\delta)$  herangezogen.
- Dabei entspricht  $\lim_{\delta \rightarrow \infty} \nu_{\text{semi}}(\delta)$  der Prozessvarianz  $\sigma^2 + \sigma_\varepsilon^2$  und  $\lim_{\delta \rightarrow 0} \nu_{\text{semi}}(\delta)$  der Fehlervarianz  $\sigma_\varepsilon^2$  (*Nuggetvarianz* in der Geostatistik).

## 9.1 Zufallsfeld: Schätzung der Parameter

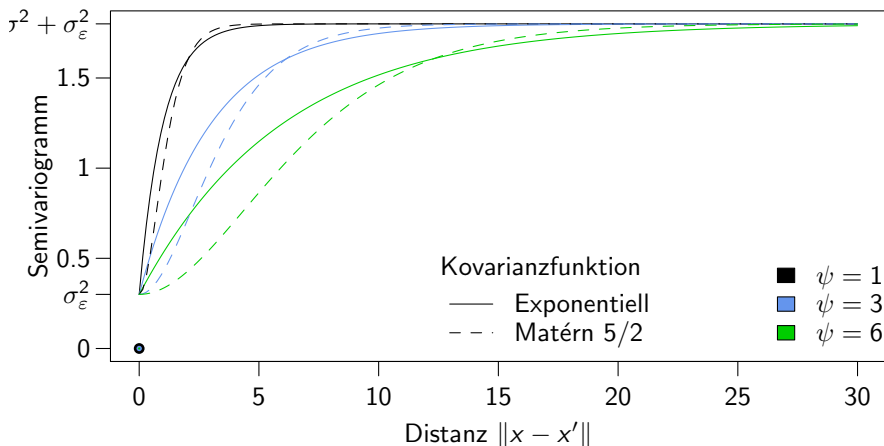


Abb. 9.3 : Theoretisches Semivariogramm  $\nu_{\text{semi}}(x, x')$  mit  $\sigma^2 = 1.5$  und  $\sigma_\varepsilon^2 = 0.3$ .

## 9.1 Zufallsfeld: Schätzung der Parameter

- Der Momentenschätzer für das Semivariogramm einer Punktdistanz  $\delta$  ist gegeben durch

$$\hat{\nu}_{\text{semi}}(\delta) = \frac{1}{2|N(\delta)|} \sum_{(x, x') \in N(\delta)} (z_{\varepsilon}(x) - z_{\varepsilon}(x'))^2, \quad (8)$$

wobei

$$N(\delta) = \{(x, x') \mid \|x - x'\| = \delta\} \quad (9)$$

die Menge der Punkte mit Abstand  $\delta$  ist und  $|N(\delta)|$  die Anzahl ihrer Elemente.

- Entsprechend ist  $\hat{\nu}(\delta) = 2\hat{\nu}_{\text{semi}}(\delta)$  der Momentenschätzer für das Variogramm.
- Damit die Anzahl der Elemente in  $N(\delta)$  nicht zu klein ausfällt, wird in der Praxis das Semivariogramm für disjunkte Distanzintervalle  $\delta = (\delta_l, \delta_u]$  betrachtet:

$$N(\delta) = N((\delta_l, \delta_u]) = \{(x, x') \mid \delta_l < \|x - x'\| \leq \delta_u\} \quad (10)$$

## 9.1 Zufallsfeld: Schätzung der Parameter

- Die gewichteten Kleinste-Quadrate-Schätzer für  $\sigma^2$ ,  $\sigma_\varepsilon^2$  und  $\psi$  werden durch das Minimieren der Verlustfunktion

$$L(\sigma^2, \sigma_\varepsilon^2, \psi) = \sum_{\delta} |N(\delta)| (\hat{\nu}(\delta) - \nu(\delta))^2 \quad (11)$$

ermittelt, wobei die Art der Kovarianzfunktion (z.B. Matérn 5/2) vorgegeben wird.

## 9.2 Kriging

- Das Kriging-Modell wurde von D. G. Krige, einem Bergbauingenieur aus Südafrika, 1951 vorgestellt und prägte die Forschung im Bereich der Geostatistik nachhaltig.
- Es liefert den besten linearen unverzerrten Schätzer für einen regionalisierten Gaußprozess gegeben eine Reihe von Beobachtungen.
- Sei  $Z(x)$  ein Gaußprozess mit  $E[Z(x)] = \mu$  und  $\text{Cov}[Z(x), Z(x')] = \kappa(x, x')$ .

Dann folgen die Zufallsvariablen  $Z = (Z(x_1) \cdots Z(x_n))^T$  zu den Realisierungen  $z = (z(x_1) \cdots z(x_n))^T$  an den Stellen  $x_1, \dots, x_n \in \mathbb{D} \subseteq \mathbb{R}^d$  der gemeinsamen Normalverteilung  $\mathcal{N}(\mu \mathbf{1}_n, K)$  mit Kovarianzmatrix

$$K = \begin{pmatrix} k(x_1)^T \\ \vdots \\ k(x_n)^T \end{pmatrix} \in \mathbb{R}^{n \times n} \quad (12)$$

bestehend aus den Kovarianzvektoren

$$k(x_i) = [\kappa(x_i, x_1) \cdots \kappa(x_i, x_n)]^T \in \mathbb{R}^n \quad (i = 1, \dots, n). \quad (13)$$

## 9.2 Kriging

- Da  $Z(x)$  und  $Z$  normalverteilt sind, folgt aus

$$\begin{pmatrix} Z(x) \\ Z \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mu \\ 1\mu \end{pmatrix}, \begin{pmatrix} \kappa(x, x) & k(x)^T \\ k(x) & K \end{pmatrix} \right) \quad (14)$$

die Verteilung des bedingten Prozesses

$$Z(x) | Z \sim \mathcal{N} \left( \mu + k(x)^T K^{-1} (Z - 1\mu), \kappa(x, x) - k(x)^T K^{-1} k(x) \right), \quad (15)$$

wobei wir davon ausgehen, dass  $\mu$  bekannt ist.

- Da dies in den meisten Fällen nicht zutrifft, muss  $\mu$  zunächst geschätzt werden, was typischerweise zu einer höheren Prognosevarianz führt.
- Es wird ab hier darauf verzichtet die Abhängigkeit der Kovarianzmatrix bzw. deren Vektoren von den Parametern  $\theta_{\kappa}$  hervorzuheben.

## 9.2 Kriging: Die Schätzung

- Die unbekannten Parameter der Verteilung  $\mu$  und  $\theta_\kappa$  werden mit der Maximum-Likelihood-Methode geschätzt, wobei die Likelihood

$$\begin{aligned} L(\mu, \theta_\kappa, Z) &= f(Z \mid \mu, \theta_\kappa) \\ &= (2\pi)^{-\frac{n}{2}} (\det K)^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (Z - \mathbf{1}_n \mu)^T K^{-1} (Z - \mathbf{1}_n \mu) \right) \end{aligned} \quad (16)$$

maximiert wird. Dabei lässt sich ein geschlossener Ausdruck

$$\hat{\mu} = \arg \max_{\mu \in \mathbb{R}} L(\mu, \theta_\kappa, Z) = \frac{\mathbf{1}_n^T K^{-1} Z}{\mathbf{1}_n^T K^{-1} \mathbf{1}_n} \quad (17)$$

für den Schätzer des Erwartungswerts  $\mu$  in Abhängigkeit von den übrigen Parametern herleiten.



## 9.2 Kriging: Die Schätzung

- Im nicht verrauschten Prozess mit  $K = \sigma^2 R$ , wobei  $R$  entsprechend zu  $K$  die Korrelationsmatrix von  $Z$  ist, lässt sich für die Prozessvarianz folgender Ausdruck ermitteln:

$$\hat{\sigma}^2 = \arg \max_{\sigma^2 \in \mathbb{R}_+} L(\hat{\mu}, \sigma^2, \psi, Z) = \frac{1}{n} (Z - 1_n \hat{\mu})^T R(\psi)^{-1} (Z - 1_n \hat{\mu}) \quad (18)$$

Damit die lineare Vorhersage  $\hat{\lambda}(x)^T Z$  für  $Z(x)$  an der Stelle  $x$  mit Gewichten  $\hat{\lambda}(x) \in \mathbb{R}^n$  erwartungstreu ist, muss

$$E[Z(x)] = \mu \stackrel{!}{=} \hat{\lambda}(x)^T 1_n \mu = E[\hat{\lambda}(x)^T Z] \quad (19)$$

gelten. Dies ist dann der Fall, wenn die Bedingung  $\hat{\lambda}(x)^T 1_n = 1$  erfüllt ist.

## 9.2 Kriging: Die Schätzung

- Um die Gewichte zu schätzen, wird der mittlere quadratische Fehler der Vorhersage

$$\begin{aligned} MSPE \left( \lambda(x)^T Z \right) &= E \left[ \left( Z(x) - \lambda(x)^T Z \right)^2 \right] \\ &= \kappa(x, x) - 2\lambda(x)^T k(x) + \lambda(x)^T K \lambda(x) \end{aligned} \quad (20)$$

unter der Nebenbedingung  $\lambda(x)^T 1_n = 1$  minimiert. Sei dabei  $c$  der Lagrange-Multiplikator und

$$L(\lambda(x), c) = \kappa(x, x) - 2\lambda(x)^T k(x) + \lambda(x)^T K \lambda(x) - c \left( \lambda(x)^T 1_n - 1 \right) \quad (21)$$

die zu minimierende Funktion mit den Ableitungen

$$\frac{\partial L(\lambda(x), c)}{\partial \lambda(x)} = -2k(x) + 2K\lambda(x) \quad (22)$$

und

$$\frac{\partial L(\lambda(x), c)}{\partial c} = -1_n^T \lambda(x) + 1. \quad (23)$$

## 9.2 Kriging: Die Schätzung

- Das Nullsetzen der Ableitungen führt zu

$$\hat{\lambda}(x) = K^{-1} \left( k(x) + \frac{c}{2} 1_n \right) \quad (24)$$

und

$$\hat{c} = 2 \frac{1 - 1_n^T K^{-1} k(x)}{1_n^T K^{-1} 1_n} \quad (25)$$

und das Einsetzen von (25) in (24) zur Lösung

$$\hat{\lambda}(x) = K^{-1} \left( k(x) + \frac{\hat{c}}{2} 1_n \right) = K^{-1} \left( k(x) + 1_n \frac{1 - 1_n^T K^{-1} k(x)}{1_n^T K^{-1} 1_n} \right). \quad (26)$$

## 9.2 Kriging: Die Schätzung

- Da die Hessematrix  $H(\hat{c}) = \frac{\partial^2 L(\lambda(x), c)}{\partial \lambda(x) \partial \lambda(x)^T}(\hat{c}) = 2K$  positiv definit ist, aufgrund der Eigenschaften der Kovarianzmatrix  $K$ , folgt, dass in

$$\hat{\lambda}(x) = \arg \min_{\{\lambda(x) \in \mathbb{R}^n \mid \lambda(x)^T \mathbf{1}_n = 1\}} MSPE \left( \lambda(x)^T Z \right) \quad (27)$$

tatsächlich das gesuchte Minimum ist.

- Daraus resultiert die Kriging-Vorhersage

$$\hat{Z}(x) = \hat{\lambda}(x)^T Z = \hat{\mu} + k(x)^T K^{-1} [Z - \mathbf{1}_n \hat{\mu}], \quad (28)$$

welche erwartungstreu ist, denn:

## 9.2 Kriging: Die Schätzung

•

$$\begin{aligned}
 E \left[ \widehat{Z}(x) \right] &= E \left( \widehat{\lambda}(x)^T Z \right) \\
 &= E(\widehat{\mu}) + k(x)^T K^{-1} [E(Z) - 1_n E(\widehat{\mu})] \\
 &= \mu + k(x)^T K^{-1} [1_n \mu - 1_n \mu] = \mu = E[Z(x)] \quad (29)
 \end{aligned}$$

aufgrund der Erwartungstreue des ML-Schätzers für  $\mu$

$$E(\widehat{\mu}) = E \left( \frac{1_n^T K^{-1} Z}{1_n^T K^{-1} 1_n} \right) = \frac{1_n^T K^{-1} E(Z)}{1_n^T K^{-1} 1_n} = \frac{1_n^T K^{-1} 1_n \mu}{1_n^T K^{-1} 1_n} = \mu \quad (30)$$

mit der Varianz

$$\text{Var}(\widehat{\mu}) = \frac{1}{1_n^T K^{-1} 1_n}. \quad (31)$$

## 9.2 Kriging: Die Schätzung

- Die Kriging-Varianz ist gegeben als

$$\begin{aligned}
 \hat{s}^2(x) &= \text{Var} \left( \hat{\lambda}(x)^T Z - Z(x) \right) \stackrel{(29)}{=} \text{MSPE} \left( \hat{\lambda}(x)^T Z \right) \\
 &\stackrel{(20)}{=} \kappa(x, x) - 2\hat{\lambda}(x)^T k(x) + \hat{\lambda}(x)^T K \hat{\lambda}(x) \\
 &= \kappa(x, x) - k(x)^T K^{-1} k(x) + \frac{(1 - 1_n^T K^{-1} k(x))^2}{1_n^T K^{-1} 1_n}, \quad (32)
 \end{aligned}$$

vgl. [57].

## 9.2 Kriging: Die Schätzung

Im Fall eines *unverrauschten* Prozesses  $Z(x)$  gegeben die *unverrauschten* Beobachtungen aus  $Z$  ist die Vorhersage gegeben als  $\hat{Z}(x_i) = z(x_i)$  mit  $\hat{s}^2(x_i) = 0$  für jedes bereits bekannte  $x_i \in \{x_1, \dots, x_n\}$ , da hier  $k(x_i)^T$  die  $i$ -te Zeile der Kovarianzmatrix  $K$  ist und damit  $k(x_i)^T K^{-1} = e_i^T$  mit  $e_i$   $i$ -ter Einheitsvektor.

- Dies gilt auch für den Fall, dass  $Z_\varepsilon(x) | Z_\varepsilon$  modelliert wird, sofern es nicht mehrere Realisationen  $z^{(1)}(x), \dots, z^{(l)}(x)$  für dasselbe  $x$  gibt.  
Dieses Modell wird auch als Kriging mit Nugget-Effekt bezeichnet.
- Im Zusammenhang mit der modellbasierten Optimierung ist die Vorhersage eines vielversprechenden Punktes anhand des Modells  $Z(x) | Z$  zum Standard bei Computerexperimenten geworden.  
Es wird angenommen, dass diese Computerexperimente deterministisch sind.

## 9.2 Kriging: Die Schätzung

Liegt dem Computerexperiment per Definition eine *stochastische* Komponente zugrunde, so erscheint es unplausibel ein Modell mit interpolierenden Eigenschaften zu verwenden.

- Falls die Beobachtungen Realisationen von  $Z_\varepsilon$  sind, wir jedoch an dem wahren Prozess  $Z(x)$  interessiert sind, verliert die zugehörige Krigingvorhersage ihre interpolierende Eigenschaft.
- In diesem Fall betrachten wir als Vorhersage den Erwartungswert von  $Z(x) | Z_\varepsilon$  (noisy Kriging-Modell).
- Hierbei enthält die Kovarianzmatrix  $K_\varepsilon$  als Einträge die Werte der Kovarianzfunktion  $\kappa_\varepsilon(x, x')$ , während  $k(x)$  immer noch der Kovarianzvektor des unverrauschten Prozesses aufgrund der Unkorreliertheit von  $Z(x)$  und  $\varepsilon(x')$  ist:

$$\begin{aligned}
 \text{Cov}[Z(x), Z_\varepsilon(x')] &= \text{Cov}[Z(x), Z(x') + \varepsilon(x')] \\
 &= \text{Cov}[Z(x), Z(x')] + \text{Cov}[Z(x), \varepsilon(x')] \\
 &= \text{Cov}[Z(x), Z(x')] = \kappa(x, x').
 \end{aligned} \tag{33}$$

Damit ist  $k(x_i)$  kein Vektor der Matrix  $K_\varepsilon = [\kappa_\varepsilon(x_i, x_j)]_{1 \leq i, j \leq n}$ , was direkt dazu führt, dass  $\hat{Z}(x_i) \neq z(x_i)$  und  $\hat{s}^2(x_i) > 0$  für  $i = 1, \dots, n$ .



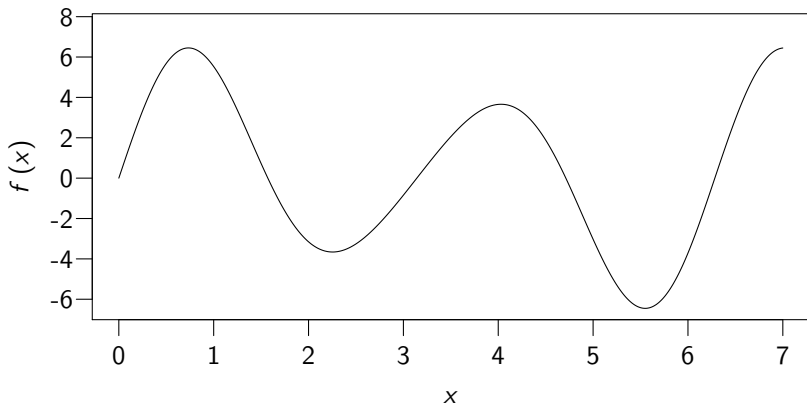
## 9.3 Modellbasierte Optimierung: Idee

Sei  $f : \mathbb{D} \mapsto \mathbb{R}$  eine Blackbox-Funktion, also eine Funktion

- über deren Eigenschaften nichts oder nur wenig bekannt ist und
- deren Auswertung teuer im Sinne von z.B. tatsächlich anfallenden Kosten und/oder zeitlichen Aufwand ist.

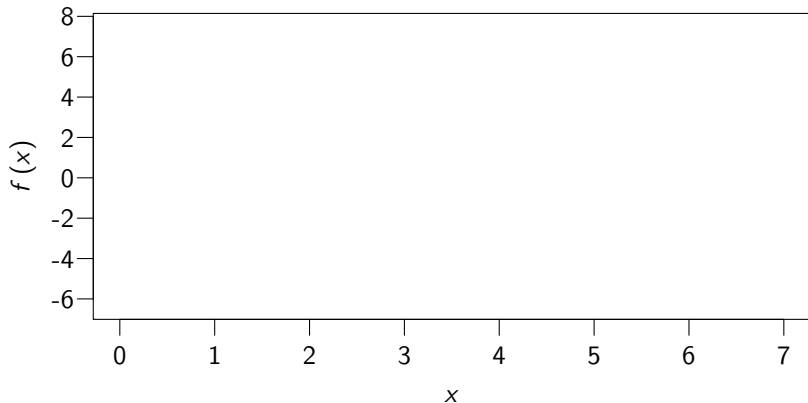
Um das Optimum bzw. die optimalen Parameter  $x \in \mathbb{D} \subseteq \mathbb{R}^d$  ( $d \in \mathbb{N}$ ) zu ermitteln, wird bei der modellbasierten Optimierung der unbekannte Funktionsverlauf von  $f$  modelliert und mit Hilfe des weit weniger teuren *Surrogatmodells*, z.B. Kriging, das Optimum gesucht.

## 9.3 Modellbasierte Optimierung: Idee



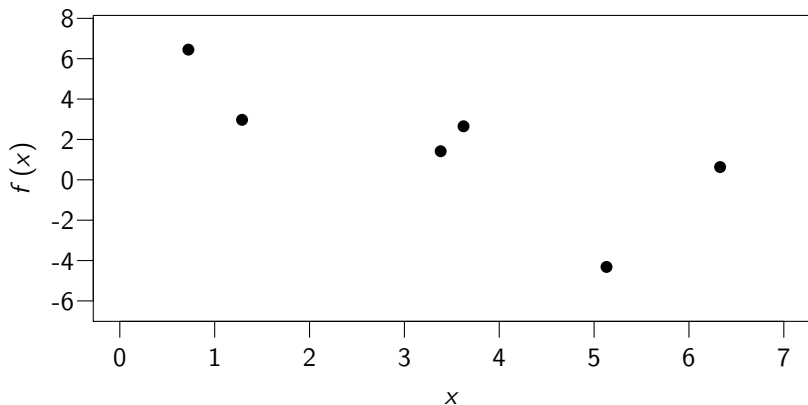
Sei  $f : \mathbb{D} \mapsto \mathbb{R}$  Funktion, deren Minimum wir finden wollen.

## 9.3 Modellbasierte Optimierung: Idee



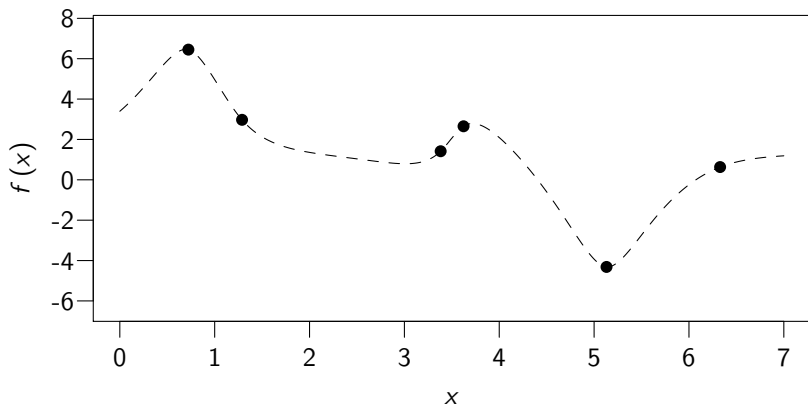
Sei  $f : \mathbb{D} \mapsto \mathbb{R}$  Funktion, die wir aber nicht kennen (Black-Box).

## 9.3 Modellbasierte Optimierung: Idee



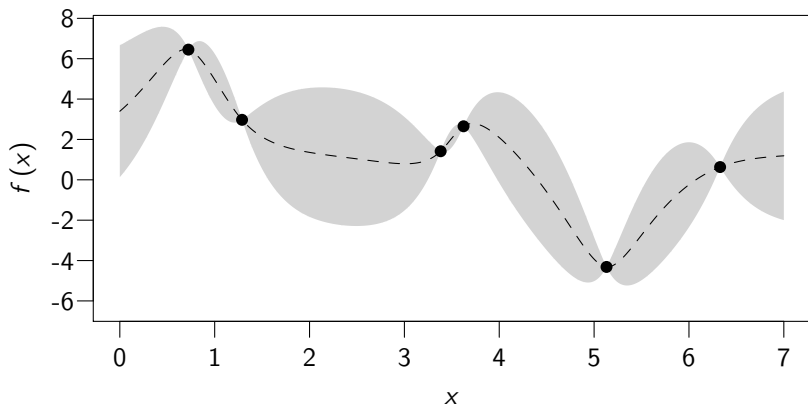
Werte  $f$  in  $x_1, \dots, x_n$  (hier  $n = 6$ ) eines raumfüllenden Designs aus:  $f(x_i) = z_i$

## 9.3 Modellbasierte Optimierung: Idee



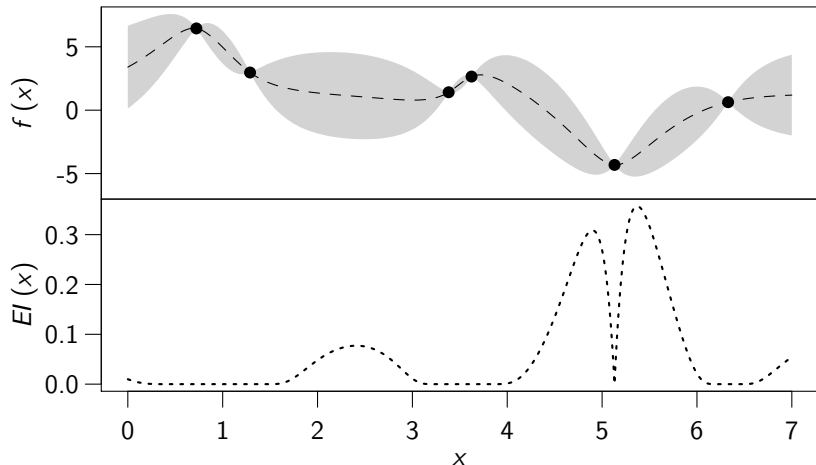
Bestimme den Kriging-Mittelwert gegeben  $z = (z_1, \dots, z_n)^T: \hat{Z}(x)$

## 9.3 Modellbasierte Optimierung: Idee



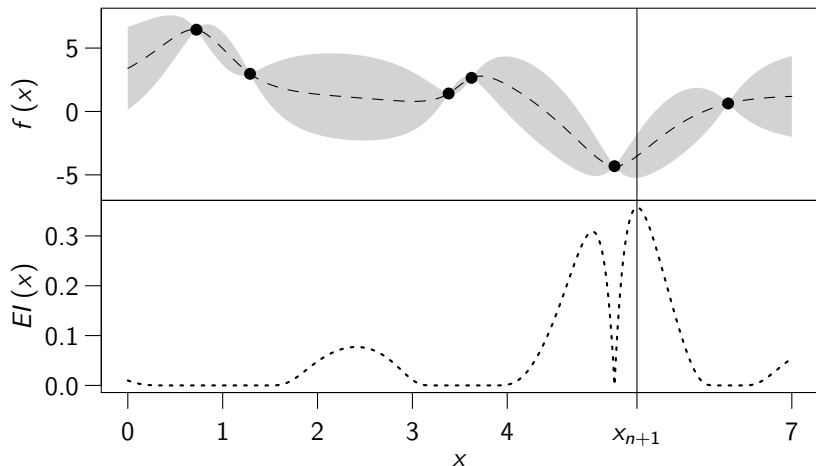
Bestimme den Kriging-Varianz gegeben  $z = (z_1, \dots, z_n)^T$ :  $\hat{s}^2(x)$

## 9.3 Modellbasierte Optimierung: Idee



Berechne Infill-Kriterium (hier: Expected Improvement), um den nächsten Auswertungspunkt zu bestimmen.

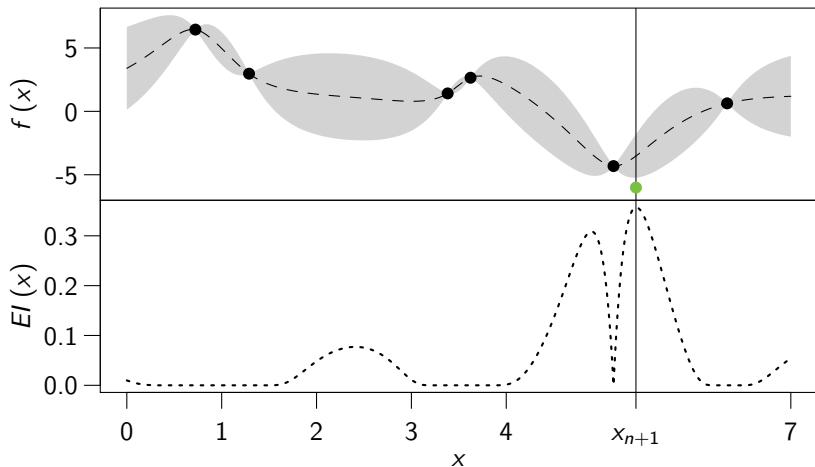
## 9.3 Modellbasierte Optimierung: Idee



Nächster Punkt für Auswertung:  $x_{n+1} = \arg \max_{x \in \mathbb{D}} EI(x)$

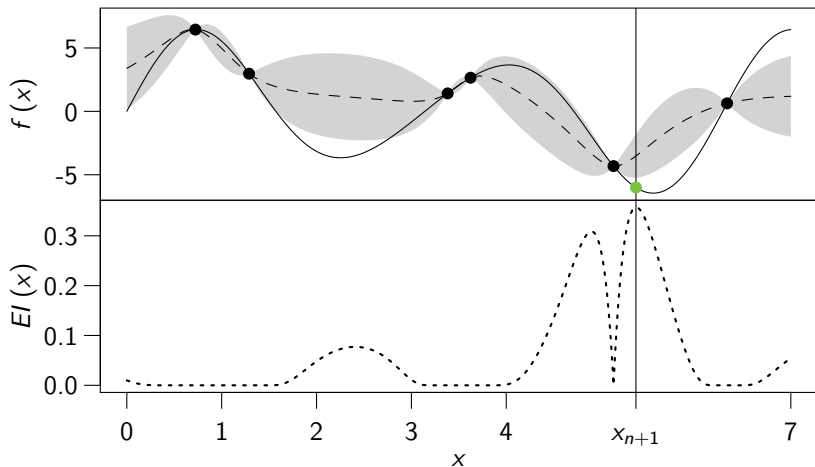


## 9.3 Modellbasierte Optimierung: Idee



Werte nächster Punkt  $x_{n+1}$  aus:  $z_{n+1} = f(x_{n+1})$

## 9.3 Modellbasierte Optimierung: Idee



Werte nächster Punkt  $x_{n+1}$  aus:  $z_{n+1} = f(x_{n+1})$

## 9.3 Modellbasierte Optimierung: Skizze

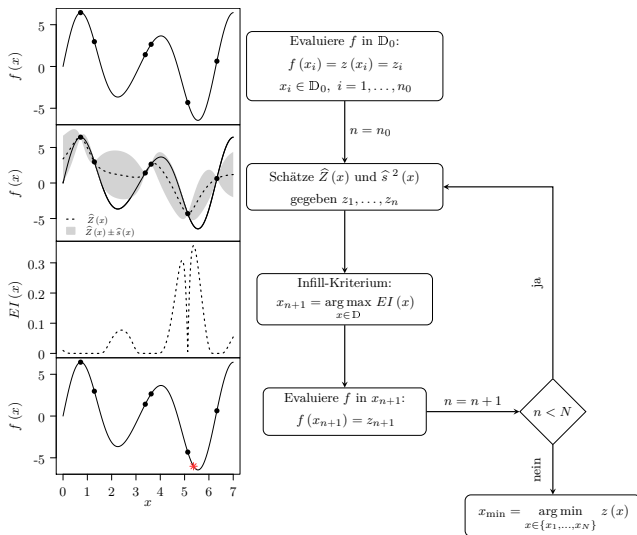


Abb. 9.4 : Beispiel für den Verlauf einer modellbasierten Optimierung (Modell: Kriging, Infill Kriterium: Expected Improvement (EI), Abbruchkriterium: Anzahl Iterationen).

## 9.3 Modellbasierte Optimierung: Schritte

Das Vorgehen gliedert sich dabei in fünf Teilschritte (vgl. Abbildung 4):

- 1 Die Funktion  $f$  wird in den Punkten eines üblicherweise raumfüllenden Versuchsplans  $\mathbb{D}_0 \subset \mathbb{D}$  ausgewertet.
- 2 Die Surrogatmodellparameter werden geschätzt.
- 3 Suche nach dem nächsten Auswertungspunkt.
- 4 Evaluiere die Zielfunktion.
- 5 Falls Abbruch, gib den bisher besten Auswertungspunkt zurück oder starte erneut bei 2.

## 9.3 Modellbasierte Optimierung: Schritte

Nun im Detail zu den Schritten 1–3:

- 1 Die Funktion  $f$  wird in den Punkten eines üblicherweise raumfüllenden Versuchsplans  $\mathbb{D}_0 \subset \mathbb{D}$  ausgewertet.
  - Damit wird eine Datengrundlage für die Schätzung der Surrogatmodellparameter in Schritt 2 geschaffen.
  - Mögliches raumfüllendes Design z.B. vollfaktorieller Plan oder Latin-Hypercube [106].
  - Die Anzahl  $n_0$  der Punkte im Initialdesign orientiert sich oft an der Dimension  $d$  des Parameterraumes. [81] empfehlen  $n_0 = 10d$  Initialauswertungen.

Nach Auswertung des Initialdesigns startet der Optimierungsprozess. Die Iteration wird beendet, wenn ein Abbruchkriterium (z.B. das Erreichen einer vorher festgelegten Anzahl  $N$  an Iterationsschritten) erfüllt ist.

## 9.3 Modellbasierte Optimierung: Schritte

- 2 Die Surrogatmodellparameter werden geschätzt.
- Wahl des Modells ist abhängig von den Eigenschaften der Funktion  $f$  und von der Beschaffenheit des Parameterraumes  $\mathbb{D}$ . Üblicherweise wird bei deterministischen Computerexperimenten mit reellwertigen Parametern  $x \in \mathbb{D} \subseteq \mathbb{R}^d$  das gewöhnliche Kriging-Modell  $Z(x) | Z$  verwendet. In seiner interpolierenden Art hat es bei deterministischen Anwendungen den Vorteil, dass es keine Unsicherheit bezüglich bereits ausgewerteter Punkte lässt.
  - Zahlreiche andere Möglichkeiten interpolierende Modelle zu konstruieren, z.B. Splines.
  - Meist wird hierbei eine Kombination aus der gewichteten Summe von Polynomen und der gewichteten Summe von bestimmten Basisfunktionen, welche um bereits beobachtete Punkte zentriert sind (z.B. dem Euklidischen Abstand), verwendet.
  - Eine ausführliche Diskussion unterschiedlicher Modelle ist bei [82] zu finden.

## 9.3 Modellbasierte Optimierung: verrauschte Beobachtungen

Ist die zu optimierende Zielfunktion  $f$  stochastisch, in dem Sinne, dass die Werte der Funktion verrauscht sind, so bietet sich das Kriging-Modell mit Nugget-Effekt  $Z_\varepsilon(x) | Z_\varepsilon$  oder das noisy Kriging-Modell  $Z(x) | Z_\varepsilon$  an.

- Die folgende Abbildung zeigt die Vorhersage  $\hat{Z}(x)$  und die Unsicherheit ausgedrückt als die ausgefüllte Fläche zwischen  $\hat{Z}(x) \pm \hat{s}(x)$ .
- Der Verlauf der Vorhersage mit  $Z_\varepsilon(x) | Z_\varepsilon$  und der Vorhersage mit  $Z(x) | Z_\varepsilon$  unterscheiden sich in der Hinsicht, dass die Vorhersage mit  $Z_\varepsilon(x) | Z_\varepsilon$  (blaue Kurve) Ausschläge zu den Beobachtungen aufweist und im Fall von nur einer Beobachtung in  $x$  diese auch interpoliert, was zu einem unstetigen Verlauf führt.
- Das R-Paket `m1rMBO` [9] bietet hier die Option `jitter`, die bei der Berechnung der Vorhersage eines bereits ausgewerteten Punktes  $x_i$  ( $i = 1, \dots, n$ ) statt  $\hat{Z}(x_i)$  einen benachbarten Wert  $\hat{Z}(x_i + c)$  zurückgibt, wobei  $c$  eine Größenordnung von  $10^{-12}$  hat. Das Ergebnis ist die gelb gestrichelte Kurve der Vorhersage.

## 9.3 Modellbasierte Optimierung: verrauschte Beobachtungen

- Ein weiterer Unterschied zwischen den beiden Modellen besteht bei der geschätzten Unsicherheit der Vorhersage. Im Vergleich zum Modell mit Nugget-Effekt ( $Z_\varepsilon(x) | Z_\varepsilon$ ) ist beim noisy Kriging-Modell ( $Z(x) | Z_\varepsilon$ ) die Unsicherheit etwas kleiner und fällt auch bei bereits beobachteten Punkten nicht auf Null (vgl. rote Fläche).
- Legende für folgende Grafik:
  - schwarze Linie: wahre Funktion;
  - schwarze Punkte: verrauschten Beobachtungen;
  - blaue Fläche:  $\hat{Z}(x) \pm \hat{s}(x)$  (Kriging mit Nugget-Effekt)  $Z_\varepsilon(x) | Z_\varepsilon$ ;
  - blau bzw. gelb (jitter): zugehörige Vorhersage;
  - rote Fläche:  $\hat{Z}(x) \pm \hat{s}(x)$  (noisy Kriging)  $Z(x) | Z_\varepsilon$ ;
  - rot: zugehörige Vorhersage.



## 9.3 Modellbasierte Optimierung: verrauschte Beobachtungen

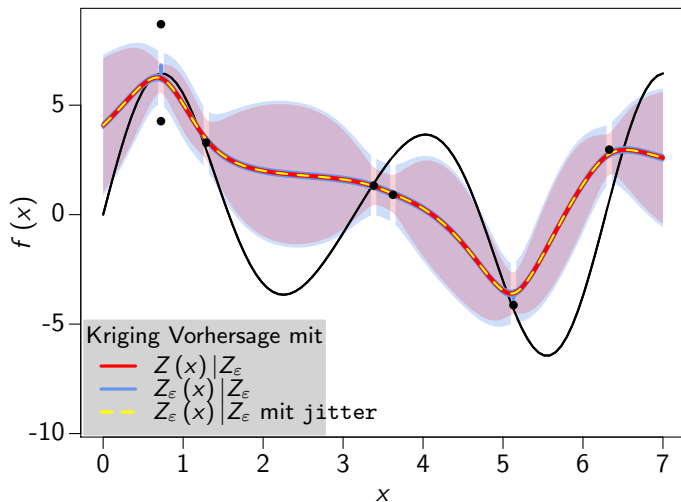


Abb. 9.5 : ausführliche Legende siehe vorige Folie.

## 9.3 Modellbasierte Optimierung: Schritte

### 🕒 Suche nach dem nächsten Auswertungspunkt.

- Da nun Verlauf der Modellfunktion bekannt ist, kann nächster Punkt mit „Bedacht“ gewählt werden:
- Suchstrategie als ausgewogene Kombination von Exploitation und Exploration,
- Optimum einerseits möglichst genau bestimmen und andererseits Chance erhöhen wirklich das globale und nicht nur ein lokales Optimum zu finden.
- Bei Exploitation kann sich die Punktwahl z.B. am Minimum der Modellvorhersage orientieren.

Stimmt das Minimum der Modellvorhersage mit einem bereits ausgewerteten Punkt überein, käme die Optimierung zum Stillstand. In [82] werden Alternativen diskutiert.

- Bei Exploration bietet das Kriging einen entscheidenden Vorteil, da hierbei neben der Vorhersage  $\hat{Z}(x)$  auch die Unsicherheit dieser Vorhersage  $\hat{s}^2(x)$  geschätzt wird. Das Abstandsmaß in der Kovarianzfunktion bewirkt dabei, dass  $\hat{s}^2(x)$  größere Werte annimmt, je weiter ein Punkt  $x$  von bereits beobachteten Punkten  $x_1, \dots, x_n$  entfernt ist. Auf diese Weise lassen sich die unbekannten Regionen im Parameterraum leicht identifizieren.

## 9.3 Modellbasierte Optimierung: Kriterien

Für die Kombination von Exploitation und Exploration wurden im Laufe der Zeit einige Kriterien (auch Infill-Kriterien genannt) vorgeschlagen;

- Das *Lower-Confidence-Bound*

$$LCB(x) = \hat{Z}(x) - c\hat{s}(x) \quad (34)$$

mit  $c > 0$  von [24] ist in seiner Form sehr intuitiv.

- In der zweiten Grafik von oben in Abbildung 4 markiert die untere Grenze der grauen Fläche gerade das Lower-Confidence-Bound mit  $c = 1$ .
- Als nächster Punkt wird vorgeschlagen:

$$x_{n+1} = \arg \min_{x \in \mathbb{D}} LCB(x). \quad (35)$$

- Die Wahl des Parameters  $c$  ist hierbei entscheidend. Betrachten wir z.B. Abbildung 4 mit  $c = 1$ , so wird  $x_{n+1} = 5.346$  gewählt. Mit  $c = 5$  erhalten wir  $x_{n+1} = 2.354$ , da in diesem Bereich die Unsicherheit besonders groß ist.

## 9.3 Modellbasierte Optimierung: Kriterien

- Das *Expected Improvement* hängt von keinen weiteren Parametern ab findet wohl am häufigsten Verwendung:

$$EI(x) = \left( z_{\min} - \hat{Z}(x) \right) \Phi \left( \frac{z_{\min} - \hat{Z}(x)}{\hat{s}(x)} \right) + \hat{s}(x) \varphi \left( \frac{z_{\min} - \hat{Z}(x)}{\hat{s}(x)} \right) \quad (36)$$

für  $\hat{s}^2(x) > 0$  und 0, sonst. Dabei ist  $z_{\min} = \min_{1 \leq i \leq n} z(x_i)$  das nach  $n$  ausgewerteten Punkten  $x_1, \dots, x_n$  realisierte Minimum der Funktion  $f$ .  $\Phi$  und  $\varphi$  bezeichnen die Verteilungsfunktion und die Dichte der Standardnormalverteilung.

- Das Konzept fand bereits 1978 Anwendung [113] und wurde später [81] im Zusammenhang mit effizienter globaler Optimierung (EGO) ein Standard.
- Durch Maximieren der Funktion erhalten wir als nächsten Punkt

$$x_{n+1} = \arg \max_{x \in \mathbb{D}} EI(x). \quad (37)$$

## 9.3 Modellbasierte Optimierung: Kriterien

- Im Fall einer stochastischen zu optimierenden Zielfunktion  $f$  ist das Minimum  $z_{\min}$  nicht exakt bekannt, da wir annehmen, dass alle Beobachtungen  $z(x_1) = f(x_1), \dots, z(x_n) = f(x_n)$  Realisationen des verrauschten Prozesses  $Z_\varepsilon(x)$  sind und daher mit einer Fehlervarianz  $\sigma_\varepsilon^2$  um den wahren Funktionswert schwanken.
- Vorgeschlagen wurde [75] statt  $z_{\min} = \min_{1 \leq i \leq n} z(x_i)$  die Vorhersage  $\hat{Z}(x^*)$  an der effektiv besten Stelle

$$x^* = \arg \min_{x \in \{x_1, \dots, x_n\}} \hat{Z}(x) + c\hat{s}(x) \quad (38)$$

zu wählen, wobei die Autoren  $c = 1$  empfehlen.

- Statt des ursprünglichen Expected Improvements  $El_{z_{\min}}(x)$  wird nun  $El_{\hat{Z}(x^*)}(x)$  verwendet.
- Zusätzlich wird die Fehlervarianz  $\sigma_\varepsilon^2$  durch einen Strafterm berücksichtigt, so dass das Augmented Expected Improvement für  $\hat{s}^2(x) > 0$  gegeben ist als

$$AEI(x) = El_{\hat{Z}(x^*)}(x) \left[ 1 - \sqrt{\frac{\sigma_\varepsilon^2}{\hat{s}^2(x) + \sigma_\varepsilon^2}} \right]. \quad (39)$$

## 9.3 Modellbasierte Optimierung: Augmented Expected Improvement

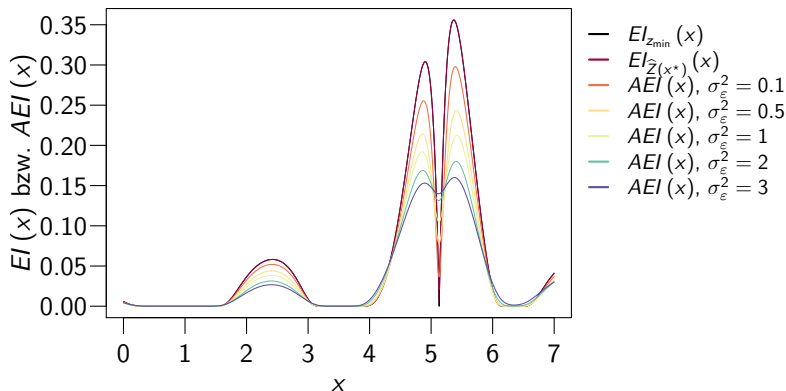


Abb. 9.6 : Beispiele für (Augmented) Expected Improvement.

## 9.3 Modellbasierte Optimierung: Software

Um in der Praxis die Stelle des Optimums  $x_{n+1}$  für ein gegebenes Infill-Kriterium zu finden, bietet das R-Paket `mlrMB0` [9] neben dem evolutionären Algorithmus u.a. auch den sogenannten Focus Search-Algorithmus [72] an.

Beim Focus Search wird der Parameterraum sukzessiv um das aktuelle Optimum verkleinert, wobei die Prozedur mehrmals wieder von vorne anfängt, was für eine ausgewogene Mischung aus Exploitation und Exploration sorgt. Der folgenden Pseudocode zeigt das genaue Vorgehen beim Minimieren eines Infill-Kriteriums  $c$ . Soll das Kriterium maximiert werden, wie dies beim Expected Improvement der Fall ist, so wird  $-c$  minimiert.

## 9.3 Modellbasierte Optimierung: Algorithmus

```

Input : infill criterion  $c: \mathbb{D} \rightarrow \mathbb{R}$ 
         control parameters  $n_{\text{restart}}, n_{\text{iters}}, n_{\text{points}}$ 
for  $u \in \{1, \dots, n_{\text{restart}}\}$  do
    Set  $\tilde{\mathbb{D}} = \mathbb{D}$ 
    for  $v \in \{1, \dots, n_{\text{iters}}\}$  do
        generate random design  $\mathcal{D} \subset \tilde{\mathbb{D}}$  of size  $n_{\text{points}}$ 
        compute  $x_{u,v}^* = (x_1^*, \dots, x_d^*)^T = \arg \min_{x \in \mathcal{D}} c(x)$ 
        # shrink  $\tilde{\mathbb{D}}$  by focusing on  $x_{u,v}^*$ :
        for each space dimension  $j \in \{1, \dots, d\}$  do
            if  $\tilde{\mathbb{D}}_j$  numeric with values  $[l_j, u_j]$  then
                 $l_j = \max\{l_j, x_j^* - \frac{1}{4}(u_j - l_j)\}$ 
                 $u_j = \min\{u_j, x_j^* + \frac{1}{4}(u_j - l_j)\}$ 
                 $\tilde{\mathbb{D}}_j = [l_j, u_j]$ 
            end
            if  $\tilde{\mathbb{D}}_j$  categorical with values  $\{x_{j1}, \dots, x_{js}\}$ ,  $s > 2$  then
                sample one category  $\bar{x}_j$  uniformly from  $\tilde{\mathbb{D}}_j \setminus \{x_j^*\}$ 
                 $\tilde{\mathbb{D}}_j = \tilde{\mathbb{D}}_j \setminus \{x_j^*\}$ 
            end
        end
    end
end

Result :  $x^* = \arg \min_{\substack{u \in \{1, \dots, n_{\text{restart}}\} \\ v \in \{1, \dots, n_{\text{iters}}\}}} c(x_{u,v}^*)$ 

```

Abb. 9.7 : Optimierung Infill-Kriterium: Focus Search [72].