



AGH

AKADEMIA GÓRNICZO-HUTNICZA IM. STANISŁAWA STASZICA W KRAKOWIE

Wydział Elektrotechniki, Automatyki, Informatyki i Inżynierii Biomedycznej

Praca dyplomowa

Gaussowskie aproksymacje filtrów optymalnych

English title

Autor:

Jakub Kłosiński

Kierunek studiów:

Automatyka i Robotyka

Opiekun pracy:

dr hab. inż. Piotr Bania

Kraków, 2022

Serdecznie dziękuję...

Streszczenie

W pracy zaprezentowano ...

Słowa kluczowe:

Abstract

The project presents ...

Keywords:

Spis treści

1. Algorytmy filtracji	9
1.1. Filtr Kalmana.....	11
1.2. Rozszerzony filtr Kalmana	13
1.3. Filtr Kalmana Gaussa-Hermite'a.....	15

1. Algorytmy filtracji

Termin *filtracja optymalna* odnosi się do zestawu metod, które mogą być używane do estymacji stanu systemu zmiennego w czasie. Stan systemu odnosi się do zbioru zmiennych, takich jak położenie, prędkość lub orientacja, które w pełni opisują badany system. Stan ten może być pośrednio obserwowany poprzez pomiary obciążone szumem, którego obecność oznacza, że obserwacje są niepewne; nawet w przypadku znajomości prawdziwego stanu systemu nie byłyby one jego deterministycznymi funkcjami, ale posiadałyby jedynie pewien rozkład możliwych wartości. Zmienność systemu w czasie jest modelowana jako system dynamiczny, który jest zaburzany pewnym szumem procesu. Szum ten jest używany do modelowania niepewności w dynamice systemu. W zdecydowanej większości przypadków zachowanie obiektu nie jest prawdziwie, wewnętrznie losowe, ale w celu przedstawienia niepewności modelu używany jest aparat matematyczny teorii prawdopodobieństwa.

Zadanie filtracji optymalnej można zaklasyfikować jako problem inwersji statystycznej, gdzie nieznaną wielkością jest wektorowy szereg czasowy $\{x_0, x_1, x_2, \dots\}$ obserwowany poprzez zbiór zaszumionych pomiarów $\{z_1, z_2, z_3, \dots\}$, przy obecności sygnałów sterujących $\{u_1, u_2, u_3, \dots\}$. Celem wspomnianej inwersji statystycznej jest oszacowanie ukrytych stanów $x_{0:T} = \{x_0, \dots, x_T\}$ na podstawie pomiarów $z_{1:T} = \{z_1, \dots, z_T\}$ i dostarczanych sygnałów sterujących $u_{1:T} = \{u_1, \dots, u_T\}$. W sensie statystyki bayesowskiej celem jest obliczenie rozkładu łącznego a posteriori wszystkich stanów przy znajomości wszystkich pomiarów i sygnałów sterujących. Zasadniczo jest to możliwe poprzez proste zastosowanie twierdzenia Bayesa (1.1).

$$p(x_{0:T}|z_{1:T}, u_{1:T}) = \frac{p(z_{1:T}|x_{0:T}, u_{1:T})p(x_{0:T}|u_{1:T})}{p(z_{1:T}|u_{1:T})} \quad (1.1)$$

Gdzie:

- $p(x_{0:T}|u_{1:T})$ to rozkład a priori zdefiniowany przez model dynamiczny,
- $p(z_{1:T}|x_{0:T}, u_{1:T})$ to wiarygodność (prawdopodobieństwo otrzymania danych wartości pomiarów pod warunkiem wartości stanu i sterowania)
- $p(z_{1:T}|u_{1:T})$ to stała normalizacyjna zdefiniowana jako $\int p(z_{1:T}|x_{0:T}, u_{1:T})p(x_{0:T}|u_{1:T}) dx_{0:T}$

Takie sformułowanie pełnego rozkładu a posteriori ma jednak poważną wadę w postaci konieczności ponownego obliczania całego rozkładu, kiedy tylko pojawi się nowy pomiar. Problem ten jest szczególnie widoczny przy dynamicznej estymacji stanu, kiedy pomiary są otrzymywane po kolei i celem jest

uzyskanie możliwie najlepszej estymaty po każdej takiej aktualizacji wartości mierzonej. Przy wzroście liczby kroków czasowych, wymiarowość pełnego rozkładu a posteriori również wzrasta, co z kolei pociąga za sobą wzrost potrzebnej mocy obliczeniowej.

Obliczenia stają się jednak znacznie prostsze, jeśli zamiast pełnego rozkładu a posteriori, obliczane są jedynie wybrane rozkłady brzegowe. Uproszczonym celem obliczeń może być zatem rozkład brzegowy stanu w kroku k przy założeniu znajomości historii pomiarów i wartości sterowania. Możliwe jest również zastosowanie twierdzenia Bayesa do wspomnianego rozkładu (1.2).

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}, \mathbf{u}_{1:t}) = \eta p(\mathbf{z}_t | \mathbf{x}_t, \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) \quad (1.2)$$

Gdzie η jest stałą normalizującą:

$$\eta = \frac{1}{p(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t})}$$

Możliwe jest przyjęcie szeregu założeń upraszczających:

- Żadne wartości pomiarów i sterowań przed krokiem t nie wpływają na przewidywanie pomiaru w kroku t przy założeniu znajomości stanu w kroku t (założenie Markowa):

$$p(\mathbf{z}_t | \mathbf{x}_t, \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) = p(\mathbf{z}_t | \mathbf{x}_t) \quad (1.3)$$

- Wprowadzenie zależności stanu w kroku t od stanu w kroku $t-1$ na podstawie twierdzenia o prawdopodobieństwie całkowitym:

$$p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) = \int p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) p(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) d\mathbf{x}_{t-1} \quad (1.4)$$

- Jedynie znajomość sterowania w kroku t może wpłynąć na przewidywanie stanu w kroku t przy założeniu znajomości stanu w kroku $t-1$. Żadne wartości pomiarów i sterowań przed krokiem t nie wpływają na to przewidywanie (założenie Markowa):

$$p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) = p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{u}_t) \quad (1.5)$$

- Znajomość wartości sterowania w kroku t nie wpływa na przewidywanie stanu w kroku $t-1$:

$$p(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) = p(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t-1}) \quad (1.6)$$

Wykorzystanie powyższych założeń pozwala na rekurencyjne obliczanie rozkładu z równania 1.2. W otrzymanym w ten sposób rekurencyjnym algorytmie filtru Bayesa można wyróżnić dwa zasadnicze kroki:

- Predykcję, polegającą na znajdowaniu przewidywanego rozkładu stanu systemu w kroku t na podstawie sterowania w kroku t i poprzedniego stanu (z kroku $t-1$). Rozkład szukany w kroku predykcji to $p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t})$,

- Korekcję, uwzględniającą pomiary z kroku t . Rozkład otrzymywany w tym kroku to $p(\mathbf{x}_t | \mathbf{z}_{1:t}, \mathbf{u}_{1:t})$.

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}, \mathbf{u}_{1:t}) = \underbrace{\eta p(\mathbf{z}_t | \mathbf{x}_t)}_{\text{Korekcja}} \underbrace{\int \underbrace{p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{u}_t) p(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t-1})}_{\text{Predykcja}} d\mathbf{x}_{t-1}}_{\text{Korekcja}} \quad (1.7)$$

gdzie η to stała normalizacyjna zdefiniowana jako $\int p(\mathbf{z}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) d\mathbf{x}_t$.

Równanie 1.7 jest fundamentem dla wielu algorytmów filtracji, wykorzystujących rekurencyjne wyznaczanie wspomnianych rozkładów stanu systemu w krokach predykcji i korekcji. Takie podejście wymaga zdefiniowania pierwotnego rozkładu obrazującego początkowe przekonanie o wartości stanu, a także dwóch modeli - jednego opisującego ewolucję systemu w czasie (model dynamiczny: $\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_{t-1})$) oraz drugiego pokazującego rozkład wartości pomiarów dla danego stanu systemu (model obserwacyjny: $\mathbf{z}_t \sim p(\mathbf{z}_t | \mathbf{x}_t)$).

1.1. Filtr Kalmana

Przy założeniu liniowości modeli dynamicznego oraz obserwacyjnego, a także addytywności szumów i normalnego charakteru ich rozkładów, można znaleźć rozwiązanie równania 1.7 w zwartej formie. Wspomniane modele wyglądają zatem następująco:

$$\begin{aligned} \mathbf{x}_t &= \mathbf{A}_{t-1} \mathbf{x}_{t-1} + \mathbf{B}_t \mathbf{u}_t + \mathbf{q}_{t-1} \\ \mathbf{z}_t &= \mathbf{H}_t \mathbf{x}_t + \mathbf{r}_t \end{aligned} \quad (1.8)$$

$\mathbf{q}_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{t-1})$ to szum procesu, natomiast $\mathbf{r}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t)$ jest szumem pomiaru. Macierz \mathbf{A}_{t-1} jest macierzą przejścia modelu dynamicznego, zaś przez \mathbf{H}_t została oznaczona macierz modelu obserwacji. Modele można również przedstawić w notacji probabilistycznej:

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{u}_t) &= \mathcal{N}(\mathbf{x}_t | \mathbf{A}_{t-1} \mathbf{x}_{t-1} + \mathbf{B}_t \mathbf{u}_t, \mathbf{Q}_{t-1}) \\ p(\mathbf{z}_t | \mathbf{x}_t) &= \mathcal{N}(\mathbf{z}_t | \mathbf{H}_t \mathbf{x}_t, \mathbf{R}_t) \end{aligned} \quad (1.9)$$

Działania wykonywane w krokach predykcji i korekcji nie powodują zmiany rodzaju rozkładu - wszystkie otrzymywane rozkłady są normalne:

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) &= \mathcal{N}(\mathbf{x}_t | \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) \\ p(\mathbf{x}_t | \mathbf{z}_{1:t}, \mathbf{u}_{1:t}) &= \mathcal{N}(\mathbf{x}_t | \mathbf{m}_t, \mathbf{P}_t) \\ p(\mathbf{z}_t | \mathbf{z}_{1:t-1}, \mathbf{u}_{1:t}) &= \mathcal{N}(\mathbf{z}_t | \mathbf{H}_t \bar{\mathbf{m}}_t, \mathbf{S}_t) \end{aligned} \quad (1.10)$$

Parametry powyższych rozkładów mogą zostać obliczone w krokach predykcji i korekcji filtru Kalmana:

- Krok predykcji:

$$\begin{aligned}\bar{\mathbf{m}}_t &= \mathbf{A}_{t-1} \mathbf{m}_{t-1} + \mathbf{B}_t \mathbf{u}_t \\ \bar{\mathbf{P}}_t &= \mathbf{A}_{t-1} \mathbf{P}_{t-1} \mathbf{A}_{t-1}^T + \mathbf{Q}_{t-1}\end{aligned}\quad (1.11)$$

- Krok korekcji:

$$\begin{aligned}\mathbf{v}_t &= \mathbf{z}_t - \mathbf{H}_t \bar{\mathbf{m}}_t \\ \mathbf{S}_t &= \mathbf{H}_t \bar{\mathbf{P}}_t \mathbf{H}_t^T + \mathbf{R}_t \\ \mathbf{K}_t &= \bar{\mathbf{P}}_t \mathbf{H}_t^T \mathbf{S}_t^{-1} \\ \mathbf{m}_t &= \bar{\mathbf{m}}_t + \mathbf{K}_t \mathbf{v}_t \\ \mathbf{P}_t &= (\mathbf{I} - \mathbf{K}_t \mathbf{H}_t) \bar{\mathbf{P}}_t\end{aligned}\quad (1.12)$$

Przewidywane parametry rozkładu $\bar{\mathbf{m}}_t$ i $\bar{\mathbf{P}}_t$ są obliczane przy użyciu modelu dynamicznego oraz sterowania dostarczanego do systemu. Sposób predykcji macierzy kowariancji $\bar{\mathbf{P}}_t$ bierze się z faktu, że zależność przyszłego stanu od stanu poprzedniego jest wyrażana poprzez macierz \mathbf{A}_{t-1} . W ten sposób przy obliczaniu niepewności uwzględniana jest również korelacja między zmiennymi stanu, wynikająca z modelu dynamicznego systemu. Do wyniku dodawana jest też macierz \mathbf{Q}_{t-1} , zatem po wykonaniu kroku predykcji wzrasta niepewność estymaty stanu.

Przewidywany stan jest korygowany poprzez uwzględnienie pomiarów w drugim etapie działania algorytmu. W zależności od podanej dokładności modelu dynamicznego oraz pomiarowego, algorytm podaje ostateczną estymatę bliższą przewidywaniom albo pomiarom. Macierz \mathbf{K}_t , nazywana macierzą wzmocnień Kalmana, precyzuje stopień zaufania do pomiarów i to na jej podstawie korygowane jest przewidywanie stanu. Uwzględnienie obserwacji jako kolejnego źródła informacji zmniejsza również niepewność oszacowania stanu.

Rozkład Gaussa jest w pełni określony przez wektor wartości średnich oraz macierz kowariancji, zatem obliczenia prowadzą do znalezienia tych dwóch charakterystyk rozkładu. Wektor wartości średnich zawiera optymalną estymatę stanu, natomiast diagonalne elementy macierzy kowariancji przedstawiają niepewność estymacji zmiennych stanu. Otrzymana estymata jest optymalna w każdym z najczęściej przyjmowanych sposobów, to znaczy *MAP* (*maximum a posteriori*), *MMSE* (*minimum mean squared error*) oraz przyjmując wartość bezwzględną błędu w funkcji kosztu (*Absolute error loss*). Wynika to z faktu, że moda, średnia arytmetyczna oraz mediana rozkładu normalnego pokrywają się.

Filtr Kalmana jest dość wydajnym obliczeniowo algorytmem. Dla najlepszych obecnie znanych algorytmów, złożoność obliczeniowa operacji odwracania macierzy jest w przybliżeniu rzędu $O(d^{2,8})$ dla macierzy rozmiaru $d \times d$. Każda iteracja algorytmu filtru Kalmana jest zatem ograniczona od dołu przez w przybliżeniu $O(k^{2,8})$, gdzie k jest rozmiarem wektora pomiarów \mathbf{z}_t . Wynika to z obserwacji, że każda iteracja algorytmu wiąże się z odwracaniem macierzy \mathbf{S}_t , rozmiaru $k \times k$. Kolejnym dolnym ograniczeniem złożoności filtru Kalmana jest $O(n^2)$, gdzie n to liczba zmiennych stanu, ze względu na mnożenie w ostatniej linii algorytmu. W wielu praktycznych zastosowaniach, wymiarowość przestrzeni pomiarów

jest znacznie mniejsza od przestrzeni stanu, i algorytm jest zdominowany przez operacje o złożoności $O(n^2)$.

1.2. Rozszerzony filtr Kalmana

Założenia o liniowych modelach dynamicznym oraz pomiarowym są rzadko spełnione w praktyce. Ten fakt, wraz z drugim założeniem o rozkładach jedynie normalnych, powoduje, że zwyczajny filtr Kalmana nadaje się tylko do najbardziej trywialnych rzeczywistych problemów. Istnieje kilka rozwiązań pozwalających na przewyższenie jednego z tych ograniczeń: założenia o liniowości. W tym wypadku zakładane jest jedynie, że wartością następnego stanu oraz pomiarami rządzą pewne (w ogólności nieliniowe) funkcje g_t i h_t :

$$\begin{aligned} \mathbf{x}_t &= \mathbf{g}_t(\mathbf{u}_t, \mathbf{x}_{t-1}) + \mathbf{q}_{t-1} \\ \mathbf{z}_t &= \mathbf{h}_t(\mathbf{x}_t) + \mathbf{r}_t \end{aligned} \quad (1.13)$$

Model przedstawiony w równaniu 1.13 uogólnia liniowy gaussowski model z równania 1.8, wykorzystywany w filtrze Kalmana. Funkcja g_t zastępuje macierze A_{t-1} oraz B_t , natomiast h_t występuje w miejsce macierzy H_t . W tym przypadku, przy dowolnych funkcjach g_t i h_t , otrzymywany rozkład nie jest już normalny. Wykonanie dokładnej aktualizacji estymaty stanu jest niemożliwe dla nieliniowych funkcji g_t i h_t , ponieważ algorytm filtru Bayesa z równania 1.7 nie posiada rozwiązania w zamkniętej formie.

Możliwe jest jednak poszukiwanie aproksymacji prawdziwego rozkładu stanu systemu. Jednym z pierwszych, podstawowych i najczęściej używanych rozwiązań jest rozszerzony filtr Kalmana (ang. *Extended Kalman Filter*, EKF). Algorytm ten również zakłada prostą reprezentację przekonania o stanie systemu za pomocą rozkładu normalnego, jednak w tym wypadku przekonanie to jest tylko przybliżeniem.

Główną ideą rozszerzonego filtru Kalmana jest linearyzacja, która przybliża g_t funkcją liniową styczną do g_t w miejscu średniej wartości rozkładu Gaussa. Poprzez projekcję rozkładu normalnego przez taką liniową aproksymację, wynikowy rozkład staje się normalny. W tym momencie cały mechanizm aktualizacji przekonania staje się taki sam jak w przypadku filtru Kalmana. Ten sam sposób może być zastosowany w przypadku funkcji h_t , zachowując w ten sposób gaussowską naturę rozkładu.

EKF wykorzystuje do linearyzacji metodę rozwinięcia Taylora pierwszego rzędu, która konstruuje liniowe przybliżenie funkcji poprzez jej wartość i pochodną cząstkową (równanie 1.14).

$$\mathbf{g}'_t(\mathbf{u}_t, \mathbf{x}_{t-1}) := \frac{\partial \mathbf{g}_t(\mathbf{u}_t, \mathbf{x}_{t-1})}{\partial \mathbf{x}_{t-1}} \quad (1.14)$$

Zarówno wartość funkcji g_t , jak i jej pochodna zależą od wartości argumentu funkcji. W rozszerzonym filtrze Kalmana jako argument wybiera się wartość stanu uznawaną za najbardziej prawdopodobną,

zatem funkcja g_t jest aproksymowana wokół \mathbf{m}_{t-1} (oraz \mathbf{u}_t):

$$\begin{aligned} g_t(\mathbf{u}_t, \mathbf{x}_{t-1}) &\approx g_t(\mathbf{u}_t, \mathbf{m}_{t-1}) + g'_t(\mathbf{u}_t, \mathbf{m}_{t-1})(\mathbf{x}_{t-1} - \mathbf{m}_{t-1}) \\ &= g_t(\mathbf{u}_t, \mathbf{m}_{t-1}) + \mathbf{G}_t(\mathbf{x}_{t-1} - \mathbf{m}_{t-1}) \end{aligned} \quad (1.15)$$

Macierz \mathbf{G}_t , często nazywana Jakobianem, jest macierzą rozmiaru $n \times n$, gdzie n to rozmiar wektora zmiennych stanu. Wartość Jakobianu zależy od \mathbf{u}_t oraz \mathbf{m}_{t-1} , zatem zmienia się dla różnych punktów w czasie.

EKF stosuje taką samą linearyzację dla funkcji h :

$$h'_t(\mathbf{x}_t) := \frac{\partial h_t(\mathbf{x}_t)}{\partial \mathbf{x}_t}$$

W tym przypadku rozwinięcie Taylora następuje w punkcie $\bar{\mathbf{m}}_t$, jako wartości najbardziej prawdopodobnej w momencie linearyzacji h :

$$\begin{aligned} h_t(\mathbf{x}_t) &\approx h_t(\bar{\mathbf{m}}_t) + h'_t(\bar{\mathbf{m}}_t)(\mathbf{x}_t - \bar{\mathbf{m}}_t) \\ &= h_t(\bar{\mathbf{m}}_t) + \mathbf{H}_t(\mathbf{x}_t - \bar{\mathbf{m}}_t) \end{aligned} \quad (1.16)$$

Modele przedstawione w notacji probabilistycznej wyglądają następująco:

$$\begin{aligned} p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{u}_t) &= \mathcal{N}(\mathbf{x}_t | g_t(\mathbf{u}_t, \mathbf{m}_{t-1}) + \mathbf{G}_t(\mathbf{x}_{t-1} - \mathbf{m}_{t-1}), \mathbf{Q}_{t-1}) \\ p(\mathbf{z}_t | \mathbf{x}_t) &= \mathcal{N}(\mathbf{z}_t | h_t(\bar{\mathbf{m}}_t) + \mathbf{H}_t(\mathbf{x}_t - \bar{\mathbf{m}}_t), \mathbf{R}_t) \end{aligned} \quad (1.17)$$

Podobnie jak w przypadku zwykłego filtra Kalmana, algorytm EKF, przedstawiony w równaniach 1.18 i 1.19, wyznacza potrzebne parametry w krokach predykcji oraz korekcji.

◦ Krok predykcji:

$$\begin{aligned} \bar{\mathbf{m}}_t &= g_t(\mathbf{u}_t, \mathbf{m}_{t-1}) \\ \bar{\mathbf{P}}_t &= \mathbf{G}_t \mathbf{P}_{t-1} \mathbf{G}_t^T + \mathbf{Q}_{t-1} \end{aligned} \quad (1.18)$$

◦ Krok korekcji:

$$\begin{aligned} \mathbf{v}_t &= \mathbf{z}_t - h_t(\bar{\mathbf{m}}_t) \\ \mathbf{S}_t &= \mathbf{H}_t \bar{\mathbf{P}}_t \mathbf{H}_t^T + \mathbf{R}_t \\ \mathbf{K}_t &= \bar{\mathbf{P}}_t \mathbf{H}_t^T \mathbf{S}_t^{-1} \\ \mathbf{m}_t &= \bar{\mathbf{m}}_t + \mathbf{K}_t \mathbf{v}_t \\ \mathbf{P}_t &= (\mathbf{I} - \mathbf{K}_t \mathbf{H}_t) \bar{\mathbf{P}}_t \end{aligned} \quad (1.19)$$

Bardziej ogólną sytuacją jest brak addytywności szumu w modelach dynamicznym i obserwacyjnym:

$$\begin{aligned} \mathbf{x}_t &= g_t(\mathbf{u}_t, \mathbf{x}_{t-1}, \mathbf{q}_{t-1}) \\ \mathbf{z}_t &= h_t(\mathbf{x}_t, \mathbf{r}_t) \end{aligned} \quad (1.20)$$

W takim przypadku możliwe jest obliczenie Jakobianów po zmiennych stanu oraz składowych szumu, oznaczonych odpowiednio G_{x_t} i G_{q_t} dla modelu dynamicznego, a także H_{x_t} i H_{r_t} dla obserwacji. Aproksymacja następuje, podobnie jak dla przypadku szumu addytywnego, wokół m_{t-1} i u_t (macierze G_{x_t} i G_{q_t}) oraz \bar{m}_t (dla macierzy H_{x_t} i H_{r_t}), a także wokół zerowych wartości składowych szumów. Algorytm rozszerzonego filtru Kalmana przyjmuje wówczas postać:

- Krok predykcji:

$$\begin{aligned}\bar{m}_t &= g_t(u_t, m_{t-1}, 0) \\ \bar{P}_t &= G_{x_t} P_{t-1} G_{x_t}^T + G_{q_t} Q_{t-1} G_{q_t}^T\end{aligned}\quad (1.21)$$

- Krok korekcji:

$$\begin{aligned}v_t &= z_t - h_t(\bar{m}_t, 0) \\ S_t &= H_{x_t} \bar{P}_t H_{x_t}^T + H_{r_t} \bar{R}_t H_{r_t}^T \\ K_t &= \bar{P}_t H_{x_t}^T S_t^{-1} \\ m_t &= \bar{m}_t + K_t v_t \\ P_t &= (I - K_t H_{x_t}) \bar{P}_t\end{aligned}\quad (1.22)$$

Algorytm rozszerzonego filtru Kalmana stał się najbardziej popularnym narzędziem wykorzystywanym do estymacji stanu systemów. Siła tego rozwiązania leży w jego prostocie oraz efektywności obliczeniowej. Tak jak w przypadku filtru Kalmana, każda iteracja potrzebuje czasu $O(k^{2,8} + n^2)$, gdzie k jest rozmiarem wektora pomiarów z_t , a n jest wymiarem przestrzeni stanów.

Ważnym ograniczeniem algorytmu EKF jest fakt, że korzysta on z linearyzacji ewolucji stanu oraz pomiarów przy pomocy metody Taylora rozwinięcia funkcji w szereg. Dokładność uzyskanej w ten sposób aproksymacji zależy od dwóch głównych czynników. Po pierwsze, jest to stopień nieliniowości funkcji, która jest linearyzowana. Jeśli funkcja ta jest w przybliżeniu liniowa, aproksymacja algorytmu będzie dobra, co przełoży się na odwzorowanie wynikowego rozkładu z wystarczającą dokładnością. Drugim czynnikiem wpływającym na skuteczność takiego sposobu linearyzacji jest stopień niepewności estymaty stanu. Jeśli niepewność jest duża, gęstość rozkładu jest mniej skupiona wokół średniej, a przez to bardziej wpływają na nią nieliniowości funkcji.

1.3. Filtr Kalmana Gaussa-Hermite'a

Innym sposobem otrzymania rozkładu normalnego jest metoda dopasowania rozkładów za pomocą momentów. Jeśli zmienna losowa $x \sim \mathcal{N}(m, P)$ jest transformowana w liniowy sposób na zmienną losową $y = g(x) + q$, $q \sim \mathcal{N}(0, Q)$, to gaussowska aproksymacja bazująca na momentach rozkładu łącznego x i y jest dana wzorem:

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} m \\ \mu_M \end{bmatrix}, \begin{bmatrix} P & C_M \\ C_M^T & S_M \end{bmatrix}\right)\quad (1.23)$$

Gdzie:

$$\begin{aligned}\mu_M &= \int g(x) \mathcal{N}(x|\mathbf{m}, \mathbf{P}) dx \\ S_M &= \int (g(x) - \mu_M)(g(x) - \mu_M)^T \mathcal{N}(x|\mathbf{m}, \mathbf{P}) dx + \mathbf{Q} \\ C_M &= \int (x - \mathbf{m})(g(x) - \mu_M)^T \mathcal{N}(x|\mathbf{m}, \mathbf{P}) dx\end{aligned}\quad (1.24)$$

Dopasowanie rozkładów za pomocą momentów jest także możliwe w przypadku szumu nieaddytywnego, czyli jeśli $\mathbf{y} = g(\mathbf{x}, \mathbf{q})$:

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m} \\ \mu_M \end{bmatrix}, \begin{bmatrix} \mathbf{P} & \mathbf{C}_M \\ \mathbf{C}_M^T & \mathbf{S}_M \end{bmatrix}\right)\quad (1.25)$$

Gdzie:

$$\begin{aligned}\mu_M &= \int g(\mathbf{x}, \mathbf{q}) \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{P}) \mathcal{N}(\mathbf{q}|\mathbf{0}, \mathbf{Q}) d\mathbf{x} \\ S_M &= \int (g(\mathbf{x}, \mathbf{q}) - \mu_M)(g(\mathbf{x}, \mathbf{q}) - \mu_M)^T \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{P}) \mathcal{N}(\mathbf{q}|\mathbf{0}, \mathbf{Q}) d\mathbf{x} d\mathbf{q} \\ C_M &= \int (\mathbf{x} - \mathbf{m})(g(\mathbf{x}, \mathbf{q}) - \mu_M)^T \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{P}) \mathcal{N}(\mathbf{q}|\mathbf{0}, \mathbf{Q}) d\mathbf{x} d\mathbf{q}\end{aligned}\quad (1.26)$$

W ten sposób możliwe jest aproksymowanie wynikowych rozkładów pojawiających się po nieliniowych transformacjach rozkładów normalnych poprzez rozkład Gaussa. Średnia \mathbf{m}_t oraz kowariancja \mathbf{P}_t rozkładu $p(\mathbf{x}_t|\mathbf{z}_{1:t}, \mathbf{u}_{1:t}) \simeq \mathcal{N}(\mathbf{x}|\mathbf{m}_t, \mathbf{P}_t)$ jest przybliżana przy użyciu metody dopasowania momentów. Dla przypadku szumu addytywnego uzyskany filtr Gaussa ma postać:

- Krok predykcji:

$$\begin{aligned}\bar{\mathbf{m}}_t &= \int \mathbf{f}(\mathbf{x}_{t-1}) \mathcal{N}(\mathbf{x}_{t-1}|\mathbf{m}_{t-1}, \mathbf{P}_{t-1}) d\mathbf{x}_{t-1} \\ \bar{\mathbf{P}}_t &= \int (\mathbf{f}(\mathbf{x}_{t-1}) - \bar{\mathbf{m}}_t)(\mathbf{f}(\mathbf{x}_{t-1}) - \bar{\mathbf{m}}_t)^T \mathcal{N}(\mathbf{x}_{t-1}|\mathbf{m}_{t-1}, \mathbf{P}_{t-1}) d\mathbf{x}_{t-1} + \mathbf{Q}_{t-1}\end{aligned}\quad (1.27)$$

- Krok korekcji:

$$\begin{aligned}\mu_t &= \int \mathbf{h}(\mathbf{x}_t) \mathcal{N}(\mathbf{x}_t, \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) d\mathbf{x}_t, \\ S_t &= \int (\mathbf{h}(\mathbf{x}_t) - \mu_t)(\mathbf{h}(\mathbf{x}_t) - \mu_t)^T \mathcal{N}(\mathbf{x}_t, \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) d\mathbf{x}_t + \mathbf{R}_t, \\ C_t &= \int (\mathbf{x}_t - \bar{\mathbf{m}}_t)((\mathbf{h}(\mathbf{x}_t) - \mu_t))^T \mathcal{N}(\mathbf{x}_t, \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) d\mathbf{x}_t\end{aligned}\quad (1.28)$$

$$\begin{aligned}\mathbf{K}_t &= \mathbf{C}_t \mathbf{S}_t^{-1} \\ \mathbf{P}_t &= \bar{\mathbf{P}}_t - \mathbf{K}_t \mathbf{S}_t \mathbf{K}_t^T \\ \mathbf{m}_t &= \bar{\mathbf{m}}_t + \mathbf{K}_t (\mathbf{z}_t - \mu_t)\end{aligned}\quad (1.29)$$

Możliwe jest rozszerzenie algorytmu na przypadek szumu nieaddytywnego. Równania filtru przyjmą wówczas postać:

○ Krok predykcji:

$$\begin{aligned}\bar{\mathbf{m}}_t &= \int \mathbf{f}(\mathbf{x}_{t-1}, \mathbf{q}_{t-1}) \mathcal{N}(\mathbf{x}_{t-1} | \mathbf{m}_{t-1}, \mathbf{P}_{t-1}) \mathcal{N}(\mathbf{q}_{t-1} | \mathbf{0}, \mathbf{Q}_{t-1}) d\mathbf{x}_{t-1} d\mathbf{q}_{t-1} \\ \bar{\mathbf{P}}_t &= \int (\mathbf{f}(\mathbf{x}_{t-1}, \mathbf{q}_{t-1}) - \bar{\mathbf{m}}_t)(\mathbf{f}(\mathbf{x}_{t-1}, \mathbf{q}_{t-1}) - \bar{\mathbf{m}}_t)^T \\ &\quad \times \mathcal{N}(\mathbf{x}_{t-1} | \mathbf{m}_{t-1}, \mathbf{P}_{t-1}) \mathcal{N}(\mathbf{q}_{t-1} | \mathbf{0}, \mathbf{Q}_{t-1}) d\mathbf{x}_{t-1} d\mathbf{q}_{t-1}\end{aligned}\quad (1.30)$$

○ Krok korekcji:

$$\begin{aligned}\boldsymbol{\mu}_t &= \int \mathbf{h}(\mathbf{x}_t, \mathbf{r}_t) \mathcal{N}(\mathbf{x}_t | \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) \mathcal{N}(\mathbf{r}_t | \mathbf{0}, \mathbf{R}_t) d\mathbf{x}_t d\mathbf{r}_t \\ \mathbf{S}_t &= \int (\mathbf{h}(\mathbf{x}_t, \mathbf{r}_t) - \boldsymbol{\mu}_t)(\mathbf{h}(\mathbf{x}_t, \mathbf{r}_t) - \boldsymbol{\mu}_t)^T \mathcal{N}(\mathbf{x}_t, \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) \mathcal{N}(\mathbf{r}_t | \mathbf{0}, \mathbf{R}_t) d\mathbf{x}_t d\mathbf{r}_t \\ \mathbf{C}_t &= \int (\mathbf{x}_t - \bar{\mathbf{m}}_t)(\mathbf{h}(\mathbf{x}_t, \mathbf{r}_t) - \boldsymbol{\mu}_t)^T \mathcal{N}(\mathbf{x}_t, \bar{\mathbf{m}}_t, \bar{\mathbf{P}}_t) \mathcal{N}(\mathbf{r}_t | \mathbf{0}, \mathbf{R}_t) d\mathbf{x}_t d\mathbf{r}_t \\ \mathbf{K}_t &= \mathbf{C}_t \mathbf{S}_t^{-1} \\ \mathbf{P}_t &= \bar{\mathbf{P}}_t - \mathbf{K}_t \mathbf{S}_t \mathbf{K}_t^T \\ \mathbf{m}_t &= \bar{\mathbf{m}}_t + \mathbf{K}_t (\mathbf{z}_t - \boldsymbol{\mu}_t)\end{aligned}\quad (1.31)$$

Powyższe ogólne równania filtru Gaussa są raczej teoretycznymi konstrukcjami, a nie praktycznymi algorytmami filtracji. Należy przyjąć pewną metodę rozwiązywania potrzebnych całek występujących w formie 1.32, aby uzyskać funkcjonalny algorytm.

$$\begin{aligned}\int \mathbf{g}(\mathbf{x}) \mathcal{N}(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x} \\ = \frac{1}{(2\pi)^{n/2} (\det \mathbf{P})^{1/2}} \int \mathbf{g}(\mathbf{x}) \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{P}^{-1}(\mathbf{x} - \mathbf{m})\right) d\mathbf{x}\end{aligned}\quad (1.32)$$

Jedną z takich numerycznych metod jest algorytm Gaussa-Hermite'a, który w swojej najprostszej formie odnosi się do przypadku jednowymiarowego ze standardową funkcją gęstości. Aproksymacja wygląda wówczas następująco:

$$\int g(x) \mathcal{N}(x | 0, 1) dx \approx \sum_{i=1}^p W_i g(\xi_i) \quad (1.33)$$

W_i to wagi, natomiast punkty ξ_i nazywane są węzłami lub punktami sigma. Jest nieskończenie wiele sposobów wyboru wag oraz węzłów. Przy rozwiązywaniu całek metodą Gaussa-Hermite'a, tak jak w przypadku innych kwadratur, są one wybierane w ten sposób, że dla funkcji podcałkowych będących wielomianami określonego stopnia, wynik jest dokładny. Okazuje się, że stopień ten jest maksymalizowany przy wyborze węzłów jako pierwiastków wielomianu Hermite'a. Dla wielomianu Hermite'a stopnia p , całkowanie jest dokładne dla wielomianów rzędu $2p - 1$ lub niższego.

Wielomian Hermite'a stopnia p jest definiowany następująco:

$$H_p(x) = (-1)^p e^{x^2/2} \frac{d^p}{dx^p} e^{-x^2/2} \quad (1.34)$$

Pierwsze wielomiany Hermite'a to:

$$\begin{aligned} H_0(x) &= 1, \\ H_1(x) &= x, \\ H_2(x) &= x^2 - 1, \\ H_3(x) &= x^3 - 3x, \\ H_4(x) &= x^4 - 6x^2 + 3, \end{aligned} \quad (1.35)$$

a kolejne mogą być obliczone z zależności rekurencyjnej:

$$H_{p+1}(x) = xH_p(x) - pH_{p-1}(x) \quad (1.36)$$

Dla każdego punktu sigma ξ_i można obliczyć odpowiadającą mu wagę W_i , korzystając z następującej zależności:

$$W_i = \frac{p!}{p^2 [H_{p-1}(\xi_i)]^2} \quad (1.37)$$

Całki z niestandardową funkcją gęstości $\mathcal{N}(x|m, P)$ mogą być obliczone poprzez zwykłą zmianę zmiennej:

$$\int g(x) \mathcal{N}(x|m, P) dx = \int g(P^{1/2}\xi + m) \mathcal{N}(\xi|0, 1) d\xi \quad (1.38)$$

W takim przypadku przybliżenie całki wygląda następująco:

$$\int g(x) \mathcal{N}(x|m, P) dx \approx \sum_{i=1}^p W_i g(P^{1/2}\xi_i + m) \quad (1.39)$$

Równanie 1.39 może być dalej uogólnione na przypadek wielowymiarowy, poprzez zdefiniowanie wektora nowych zmiennych ξ oraz zastosowanie rozkładu Choleskiego do macierzy kowariancji ($P = \sqrt{P}\sqrt{P}^T$):

$$x = m + \sqrt{P}\xi \quad (1.40)$$

Otrzymana w ten sposób całka nad wielowymiarową funkcją ze standardowym rozkładem normalnym jako funkcją wagową to:

$$\int g(x) \mathcal{N}(x|m, P) dx = \int g(m + \sqrt{P}\xi) \mathcal{N}(x|0, I) d\xi \quad (1.41)$$

Całka otrzymana w równaniu 1.41 może być przedstawiona jako całka iterowana i każdą z całek wchodzących w skład całki iterowanej można aproksymować z wykorzystaniem kwadratury Gaussa-Hermite'a:

$$\begin{aligned} & \int g(m + \sqrt{P}\xi) \mathcal{N}(x|0, I) d\xi \\ &= \int \cdots \int g(m + \sqrt{P}\xi) \mathcal{N}(\xi_1, 0, 1) \times \cdots \times \mathcal{N}(\xi_n, 0, 1) d\xi_1 \cdots d\xi_n \\ &\approx \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} g(m + \sqrt{P}\xi_{(i_1, \dots, i_n)}) \end{aligned} \quad (1.42)$$

Węzły powstają jako iloczyn kartezjański jednowymiarowych punktów sigma, $\xi_{(i_1, \dots, i_n)} = (\xi_1 \dots \xi_n)^T$, natomiast wielowymiarowe wagi są tworzone poprzez pomnożenie jednowymiarowych wag odpowiadających węzłom: $W_{(i_1, \dots, i_n)} = W_{i_1} \times \dots \times W_{i_n}$.

Całkowanie metodą Gaussa-Hermite-a jest dokładne dla jednomianów $x_1^{d_1} x_2^{d_2} \dots x_n^{d_n}$ i ich dowolnej kombinacji liniowej, gdzie każda potęga $d_i \leq 2p - 1$. Liczba węzłów (rozmiaru n) oraz wag potrzebnych do obliczenia całki n -wymiarowej przy zastosowaniu p węzłów jednowymiarowych jest równa p^n , zatem złożoność aproksymacji Gaussa-Hermite-a rośnie bardzo szybko wraz ze wzrostem wymiarowości i liczby p .

Zastosowanie metody Gaussa-Hermite-a do obliczenia całek z równań 1.27 i 1.28 daje w wyniku algorytm filtra Kalmana Gaussa-Hermite-a (ang. *Gauss-Hermite Kalman Filter*, GHKF) dla przypadku szumu addytywnego:

- o Krok predykcji:

$$\begin{aligned} \chi_{t-1}^{(i_1, \dots, i_n)} &= \mathbf{m}_{t-1} + \sqrt{\mathbf{P}_{t-1}} \boldsymbol{\xi}_{(i_1, \dots, i_n)} \\ \hat{\chi}_t^{(i_1, \dots, i_n)} &= \mathbf{f}(\chi_{t-1}^{(i_1, \dots, i_n)}) \\ \bar{\mathbf{m}}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} \hat{\chi}_t^{(i_1, \dots, i_n)} \\ \bar{\mathbf{P}}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} (\hat{\chi}_t^{(i_1, \dots, i_n)} - \bar{\mathbf{m}}_t)(\hat{\chi}_t^{(i_1, \dots, i_n)} - \bar{\mathbf{m}}_t)^T + \mathbf{Q}_{t-1} \end{aligned} \quad (1.43)$$

- o Krok korekcji:

$$\begin{aligned} \bar{\chi}_t^{(i_1, \dots, i_n)} &= \bar{\mathbf{m}}_t + \sqrt{\bar{\mathbf{P}}_t} \boldsymbol{\xi}_{(i_1, \dots, i_n)} \\ \hat{\Psi}_t^{(i_1, \dots, i_n)} &= \mathbf{h}(\bar{\chi}_t^{(i_1, \dots, i_n)}) \\ \boldsymbol{\mu}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} \hat{\Psi}_t^{(i_1, \dots, i_n)} \\ \mathbf{S}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} (\hat{\Psi}_t^{(i_1, \dots, i_n)} - \boldsymbol{\mu}_t)(\hat{\Psi}_t^{(i_1, \dots, i_n)} - \boldsymbol{\mu}_t)^T + \mathbf{R}_t \\ \mathbf{C}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} (\bar{\chi}_t^{(i_1, \dots, i_n)} - \bar{\mathbf{m}}_t)(\hat{\Psi}_t^{(i_1, \dots, i_n)} - \boldsymbol{\mu}_t)^T \\ \mathbf{K}_t &= \mathbf{C}_t \mathbf{S}_t^{-1} \\ \mathbf{P}_t &= \bar{\mathbf{P}}_t - \mathbf{K}_t \mathbf{S}_t \mathbf{K}_t^T \\ \mathbf{m}_t &= \bar{\mathbf{m}}_t + \mathbf{K}_t (\mathbf{z}_t - \boldsymbol{\mu}_t) \end{aligned} \quad (1.44)$$

Filtr dla szumu nieaddytywnego (do zweryfikowania):

- o Krok predykcji:

$\xi_{(i_1, \dots, i_n)}$ będzie tworzone na podstawie zwiększonej liczby wymiarów (suma wymiarów zmiennych stanu i wektora szumu \mathbf{q}). \mathbf{m}_{t-1} będzie miało także zwiększony wymiar, dla zmiennych

szumu będą tam zera. Macierz P_{t-1} będzie miała dołożoną na przekątnej macierz Q_{t-1} , a pozostałe miejsca macierzy (odpowiadające kowariancji zmiennych stanu i zmiennych szumu) będą uzupełnione zerami. Funkcja f będzie przyjmować współrzędne węzłów odpowiadające zmiennym stanu i szumom:

$$\begin{aligned}\chi_{t-1}^{(i_1, \dots, i_n)} &= \mathbf{m}_{t-1} + \sqrt{P_{t-1}} \boldsymbol{\xi}_{(i_1, \dots, i_n)} \\ \hat{\chi}_t^{(i_1, \dots, i_n)} &= f(\chi_{t-1}^{(i_1, \dots, i_n)}) \\ \bar{\mathbf{m}}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} \hat{\chi}_t^{(i_1, \dots, i_n)}\end{aligned}$$

Nie będzie już dodawana macierz Q :

$$\bar{P}_t = \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} (\hat{\chi}_t^{(i_1, \dots, i_n)} - \bar{\mathbf{m}}_t)(\hat{\chi}_t^{(i_1, \dots, i_n)} - \bar{\mathbf{m}}_t)^T \quad (1.45)$$

o Krok korekcji:

Tutaj podobnie jak w kroku predykcji, $\boldsymbol{\xi}_{(i_1, \dots, i_n)}$ tworzone będzie na podstawie zwiększonej liczby wymiarów (suma wymiarów zmiennych stanu i wektora szumu r). \bar{P}_t będzie miało dołożoną macierz R_t . $\bar{\mathbf{m}}_t$ będzie też uzupełnione zerami. Funkcja h przyjmie współrzędne węzłów odpowiadające i zmiennym stanu, i szumom:

$$\begin{aligned}\bar{\chi}_t^{(i_1, \dots, i_n)} &= \bar{\mathbf{m}}_t + \sqrt{\bar{P}_t} \boldsymbol{\xi}_{(i_1, \dots, i_n)} \\ \hat{\Psi}_t^{(i_1, \dots, i_n)} &= h(\bar{\chi}_t^{(i_1, \dots, i_n)}) \\ \boldsymbol{\mu}_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} \hat{\Psi}_t^{(i_1, \dots, i_n)}\end{aligned}$$

Bez dodawania macierzy R :

$$\begin{aligned}S_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} (\hat{\Psi}_t^{(i_1, \dots, i_n)} - \boldsymbol{\mu}_t)(\hat{\Psi}_t^{(i_1, \dots, i_n)} - \boldsymbol{\mu}_t)^T \\ C_t &= \sum_{i_1, \dots, i_n} W_{(i_1, \dots, i_n)} (\bar{\chi}_t^{(i_1, \dots, i_n)} - \bar{\mathbf{m}}_t)(\hat{\Psi}_t^{(i_1, \dots, i_n)} - \boldsymbol{\mu}_t)^T \\ K_t &= C_t S_t^{-1} \\ P_t &= \bar{P}_t - K_t S_t K_t^T \\ \mathbf{m}_t &= \bar{\mathbf{m}}_t + K_t (z_t - \boldsymbol{\mu}_t)\end{aligned} \quad (1.46)$$