

Pattern Classification

All materials in these slides were taken from
Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000
with the permission of the authors and the publisher

Chapter 3:

Estymator największej wiarygodności i estymator Bayesa (part 1)

- Wstęp
- Estymator największej wiarygodności
 - Przykład dla przypadku szczególnego
 - Przypadek Gaussowski: nieznanne μ i σ
 - Obciążenie estymatora
- Dodatek: Sformułowanie zadania MW

- Wstęp

- Data availability in a Bayesian framework

- Możemy opracować optymalny klasyfikator Bayesa, znając:

- $P(\omega_i)$ (prawdopodobieństwa a'priori)
- $P(x | \omega_i)$ (warunkowe rozkłady w klasach)

Niestety, jedynie w wyjątkowych sytuacjach dysponujemy pełną informacją probabilistyczną

- Opracowywanie klasyfikatora w oparciu o ciąg uczący

- Łatwo oszacować prawdopodobieństwa a'priori
- Ciąg uczący jest zwykle zbyt mały do estymacji rozkładów warunkowych (duży wymiar przestrzeni cech)

- Informacja aprioryczna o problemie
- Normalny rozkład $P(x | \omega_i)$

$$P(x | \omega_i) \sim N(\mu_i, \Sigma_i)$$

- Rozkład charakteryzowany przez dwa parametry
- Techniki estymacji
 - Estymator Maksymalnej Wiarygodności (MW)
 - Estymator Bayesa
 - Wyniki końcowe są podobne, ale podejścia są różne

- Zestaw parametrów w zadaniu MW jest ustalony, ale wartości parametrów należy wyznaczyć
- Najlepsze oszacowania wartości parametrów uzyskujemy maksymalizując prawdopodobieństwo otrzymania danych, jakie zostały zaobserwowane
- W metodach Bayesowskich parametry rozpatrywane są jako zmienne losowe o znanych rozkładach
- W obu podejściach używamy $P(\omega_i | x)$ do konstrukcji reguł decyzyjnych klasyfikatora

- Estymator Maksymalnej Wiarygodności

- Ma dobre własności zbieżności dla rosnącej długości ciągu uczącego
- Łatwiejszy w wyznaczeniu od estymatorów innych metod

- Ogólna idea działania estymatora

- Załóżmy, że jest c klas obiektów

$$P(x | \omega_j) \sim N(\mu_j, \Sigma_j)$$

$$P(x | \omega_j) \equiv P(x | \omega_j, \theta_j) \text{ gdzie:}$$

$$\theta = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, \dots, \sigma_j^{11}, \sigma_j^{22}, \text{cov}(x_j^m, x_j^n) \dots)$$

- Wykorzystanie ciągu uczącego do estymacji $\theta = (\theta_1, \theta_2, \dots, \theta_c)$, gdzie θ_i ($i = 1, 2, \dots, c$) odpowiadają klasom
- Niech D zawiera n próbek, x_1, x_2, \dots, x_n

$$P(D | \theta) = \prod_{k=1}^n P(x_k | \theta) = F(\theta)$$

$P(D | \theta)$ – funkcja wiarygodności

- Estymator MW parametru θ maksymalizuje $P(D | \theta)$
“Jest taką wartością parametru θ , która najlepiej odpowiada zaobserwowanym pomiarom, zebranych w ciągu uczącym”

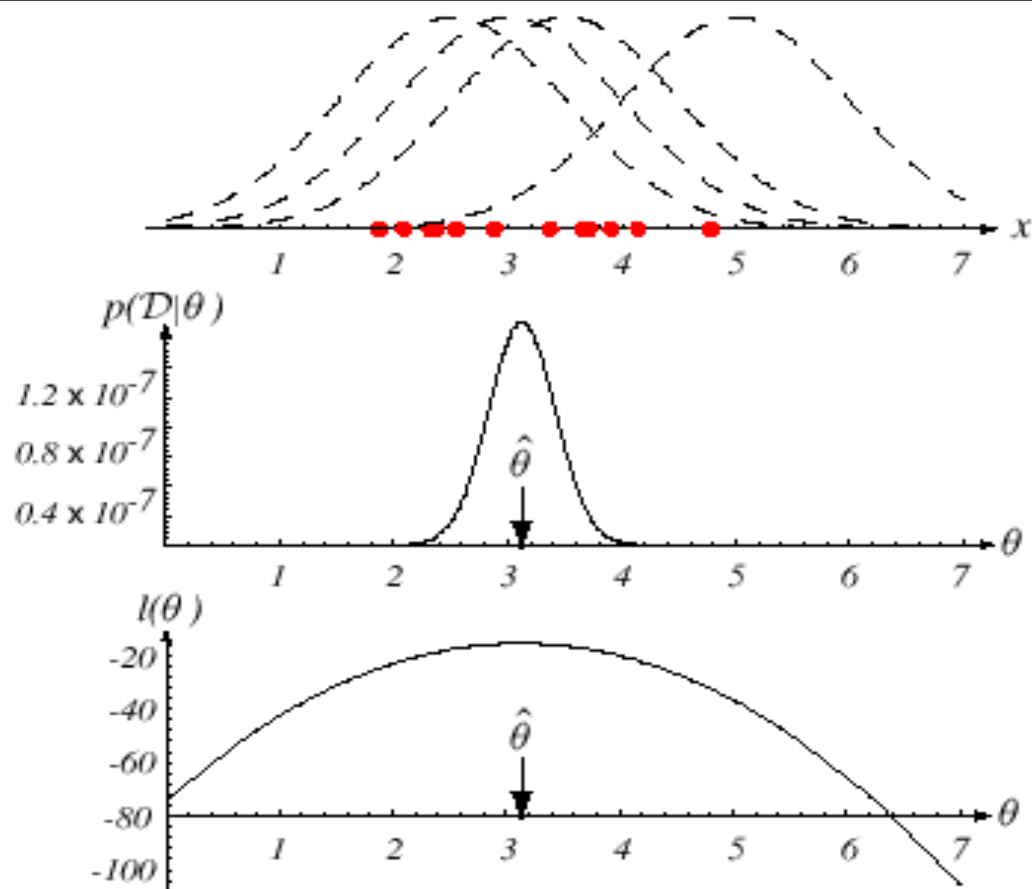


FIGURE 3.1. The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood $p(\mathcal{D}|\theta)$ as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked $\hat{\theta}$; it also maximizes the logarithm of the likelihood—that is, the log-likelihood $l(\theta)$, shown at the bottom. Note that even though they look similar, the likelihood $p(\mathcal{D}|\theta)$ is shown as a function of θ whereas the conditional density $p(x|\theta)$ is shown as a function of x . Furthermore, as a function of θ , the likelihood $p(\mathcal{D}|\theta)$ is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Estymator optymalny

- Niech $\theta = (\theta_1, \theta_2, \dots, \theta_p)^T$ oraz niech ∇_{θ} oznacza gradient

$$\nabla_{\theta} = \left[\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^T$$

- Oznaczmy $l(\theta)$ jako:

$$l(\theta) = \ln P(D | \theta)$$

- Sformułowanie problemu:

wyznaczyć wartość θ , która maksymalizuje:

$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$

Warunek konieczny optimum:

$$(\nabla_{\theta} l = \sum_{k=1}^n \nabla_{\theta} \ln P(x_k | \theta))$$

$$\nabla_{\theta} l = 0$$

- Szczególny przypadek: nieznanne μ
 - $P(x_i | \mu) \sim N(\mu, \Sigma)$
(Wartości cech obiektów są zmiennymi losowymi o wielowymiarowym rozkładzie normalnym)

$$\ln P(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^T \Sigma^{-1} (x_k - \mu)$$

$$\text{i } \nabla_{\theta=\mu} \ln P(x_k | \mu) = \Sigma^{-1} (x_k - \mu) = 0$$

$\theta = \mu$ a zatem:

- Estymator MW parametru μ musi spełniać warunek:

$$\sum_{k=1}^n \Sigma^{-1} (x_k - \hat{\mu}) = 0$$

- Po wymnożeniu przez Σ oraz po prostych przekształceniach otrzymujemy:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

Jest to średnia arytmetyczna próbek z ciągu uczącego

Wniosek:

Jeżeli $P(x_k | \omega_j)$ ($j = 1, 2, \dots, c$) jest rozkładem Gaussa w d -wymiarowej przestrzeni cech to można estymować

$\theta = (\theta_1, \theta_2, \dots, \theta_c)^T$ i dokonać optymalnej klasyfikacji (w sensie bayesowskim).

- Estymator MW:
 - Rozkład normalny: *nieznane μ oraz σ*
 $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$

$$l = \ln P(x_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln P(x_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln P(x_k | \theta)) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$

Sumowanie:

$$\left\{ \begin{array}{l} \sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \theta_1) = 0 \quad (1) \\ -\sum_{k=1}^n \frac{1}{\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 \quad (2) \end{array} \right.$$

Rozwiązanie układu daje wynik:

$$\hat{\theta}_1 = \mu = \sum_{k=1}^n \frac{x_k}{n} \quad ; \quad \hat{\theta}_2 = \sigma^2 = \frac{\sum_{k=1}^n (x_k - \mu)^2}{n}$$

- Obciążenie estymatora

- Estymator MW parametru σ^2 jest obciążony

$$E\left[\frac{1}{n}\sum(x_i - \bar{x})^2\right] = \frac{n-1}{n}\sigma^2 \neq \sigma^2$$

- Nieobciążony estymator macierzy Σ ma postać:

$$C = \frac{1}{n-1} \sum_{k=1}^{k=n} (x_k - \mu)(x_k - \hat{\mu})^T$$

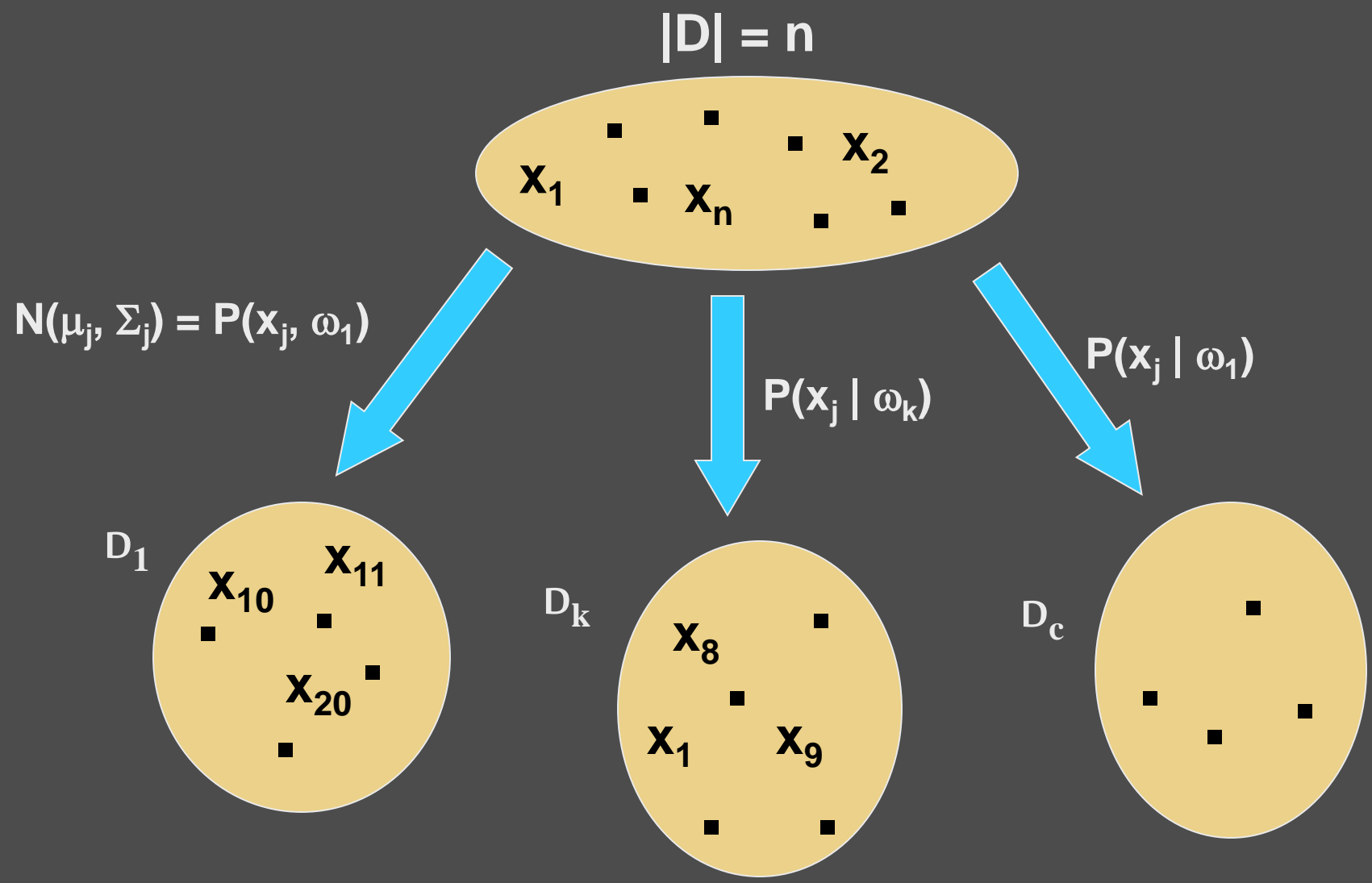
macierz kowariancji z próby

- Dodatek: Sformułowanie zadania estymacji MW

- Niech $D = \{x_1, x_2, \dots, x_n\}$

$$P(x_1, \dots, x_n | \theta) = \prod_{k=1}^n P(x_k | \theta); |D| = n$$

Należy wyznaczyć wartość $\hat{\theta}$, przy której ciąg uczący jest najbardziej reprezentatywny



$$\theta = (\theta_1, \theta_2, \dots, \theta_c)$$

Problem: wyznaczyć $\hat{\theta}$ takie, że:

$$\begin{aligned} \max_{\theta} P(D | \theta) &= \max P(x_1, \dots, x_n | \theta) \\ &= \max \prod_{k=1}^n P(x_k | \theta) \end{aligned}$$

Chapter 3:

Estymator największej wiarygodności i estymator Bayesa (part 1)

- Metoda Bayesa (MB)
 - Parametryczny estymator Bayesa: rozkład normalny
 - Parametryczny estymator Bayesa: przypadek ogólny
- Problemy z wielowymiarowością
- Złożoność obliczeniowa
- Analiza komponentów i analiza dyskryminacyjna
- Ukryte łańcuchy Markowa

- Estymator Bayesa (uczenie bayesowskie w rozpoznawaniu obrazów)
 - In MLE θ was supposed fix
 - In BE θ is a random variable
 - Wyznaczanie prawdopodobieństw a'posteriori $P(\omega_i | x)$
 - Wyznaczyć $P(\omega_i | x, D)$
Wzór Bayesa dla próbki D :

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D)P(\omega_i | D)}{\sum_{j=1}^c P(x | \omega_j, D)P(\omega_j | D)}$$

$$P(x, \mathbf{D} | \omega_i) = P(x | \mathbf{D}, \omega_i)P(\mathbf{D} | \omega_i)$$

$$P(x | \mathbf{D}) = \sum_j P(x, \omega_j | \mathbf{D})$$

$$P(\omega_i) = P(\omega_i | \mathbf{D})$$

Ostatecznie :

$$P(\omega_i | x, \mathbf{D}) = \frac{P(x | \omega_i, \mathbf{D}_i)P(\omega_i)}{\sum_{j=1}^c P(x | \omega_j, \mathbf{D})P(\omega_j)}$$

- Parametryczny estymator Bayesa: rozkład normalny

Cel: Estymacja θ z wykorzystaniem rozkładu a'posteriori $P(\theta \mid D)$

- Przypadek jednowymiarowy: $P(\mu \mid D)$
 μ jest jedynym nieznanym parametrem

$$P(\mathbf{x} \mid \mu) \sim \mathcal{N}(\mu, \sigma^2)$$

$$P(\mu) \sim \mathcal{N}(\mu_0, \sigma_0^2)$$

$$P(\mu | \mathbf{D}) = \frac{P(\mathbf{D} | \mu)P(\mu)}{\int P(\mathbf{D} | \mu)P(\mu)d\mu} \quad (1)$$

$$= \alpha \prod_{k=1}^k P(x_k | \mu)P(\mu)$$

- Estymacja rozkładu

$$P(\mu | \mathbf{D}) \sim N(\mu_n, \sigma_n^2) \quad (2)$$

Porównanie (1) i (2):

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

$$\text{and } \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

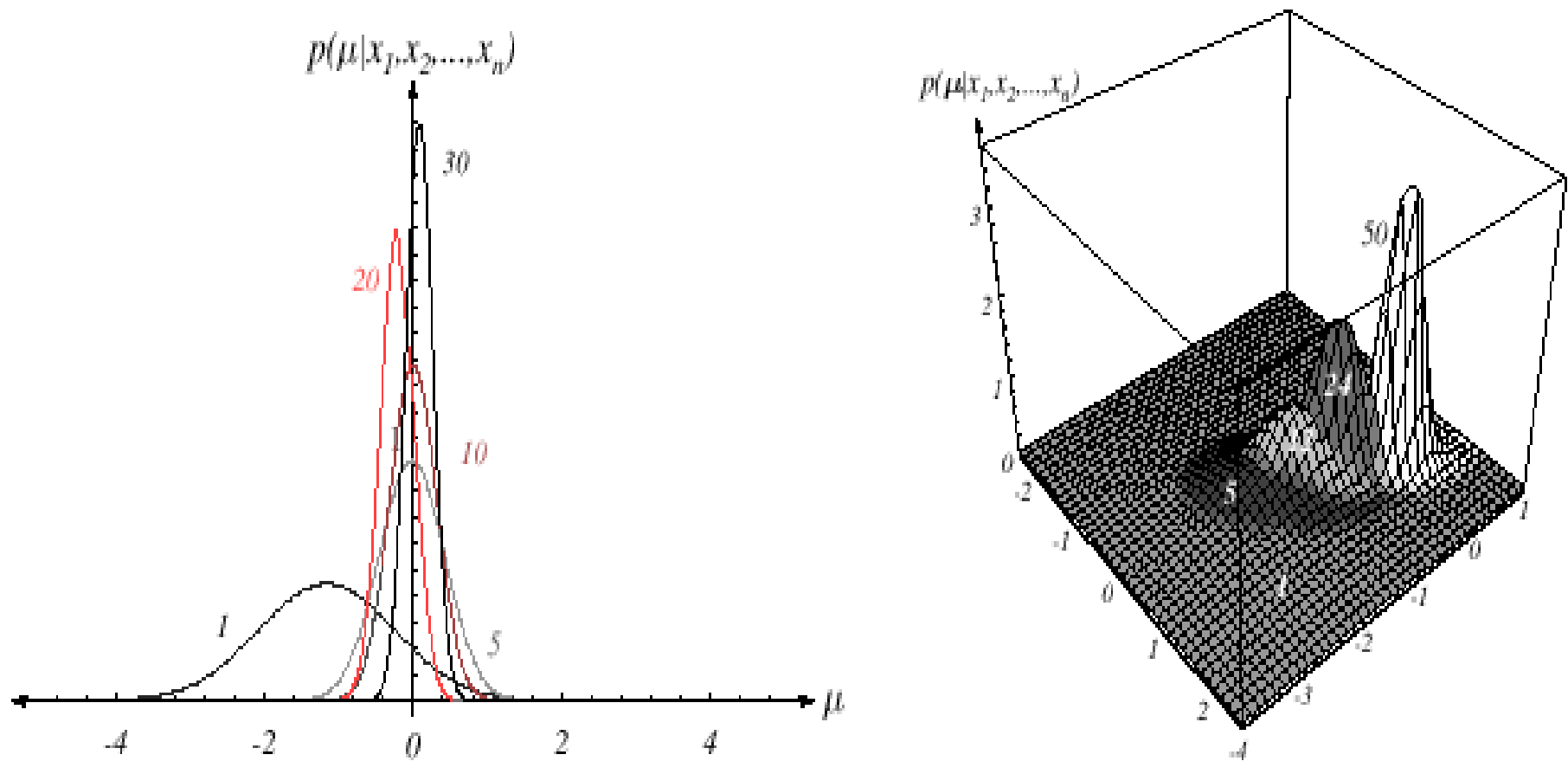


FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Przypadek jednowymiarowy $P(x | D)$
 - $P(\mu | D)$ jest już wyznaczone
 - $P(x | D)$ pozostało do wyznaczenia!

$P(x | D) = \int P(x | \mu)P(\mu | D)d\mu$ jest rozkładem normalnym

A więc: $P(x | D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$

(szukany rozkład w klasach $P(x | D_j, \omega_j)$)

Mając $P(x | D_j, \omega_j)$ wraz z $P(\omega_j)$ i wzorem Bayesa, otrzymujemy regułę decyzyjną dla klasyfikatora Bayesa:

$$\max_{\omega_j} [P(\omega_j | x, D)] \equiv \max_{\omega_j} [P(x | \omega_j, D_j)P(\omega_j)]$$

- Parametryczny estymator Bayesa: przypadek ogólny
 - Wyznaczenie $P(x | D)$ może być przeprowadzone w każdej sytuacji, w której nieznaną rozkład daje się parametryzować. Szczegółowe założenia są następujące:
 - Postać $P(x | \theta)$ jest znana z dokładnością do wartości θ , którą należy wyznaczyć
 - Wiedza o θ jest zawarta w znanym rozkładzie a’piori $P(\theta)$
 - Pozostała część wiedzy o θ jest zawarta w zbiorze D n zmiennych losowych x_1, x_2, \dots, x_n

Podstawowe zadanie:

“Wyznacz rozkład a’posteriori $P(\theta | \mathbf{D})$ ” a następnie “wyznacz $P(x | \mathbf{D})$ ”

Ze wzoru Bayesa mamy:

$$P(\theta | \mathbf{D}) = \frac{P(\mathbf{D} | \theta)P(\theta)}{\int P(\mathbf{D} | \theta)P(\theta)d\theta},$$

Z warunku niezależności otrzymujemy:

$$P(\mathbf{D} | \theta) = \prod_{k=1}^{k=n} P(x_k | \theta)$$

- Problemy z wielowymiarowością
- Zadania z 50 lub 100 binarnymi cechami
 - Dokładność klasyfikacji zależy od wymiaru i ilości danych uczących
 - Przypadek wielowymiarowy z rozkładem normalnym i dwiema klasami o tej samej kowariancji

$$P(\text{blad}) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{-\frac{u^2}{2}} du$$

$$\text{gdzie: } r^2 = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)$$

$$\lim_{r \rightarrow \infty} P(\text{blad}) = 0$$

- Jeżeli cechy są niezależne, to:

$$\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2)$$

$$r^2 = \sum_{i=1}^{i=d} \left(\frac{\mu_{i1} - \mu_{i2}}{\sigma_i} \right)^2$$

- Najbardziej użyteczne cechy to takie, których różnica między średnią jest duża w porównaniu do odchylenia standardowego
- Często obserwuje się w praktyce, że – począwszy od pewnego momentu – uwzględnianie kolejnych cech prowadzi do pogorszenia jakości klasyfikacji!.

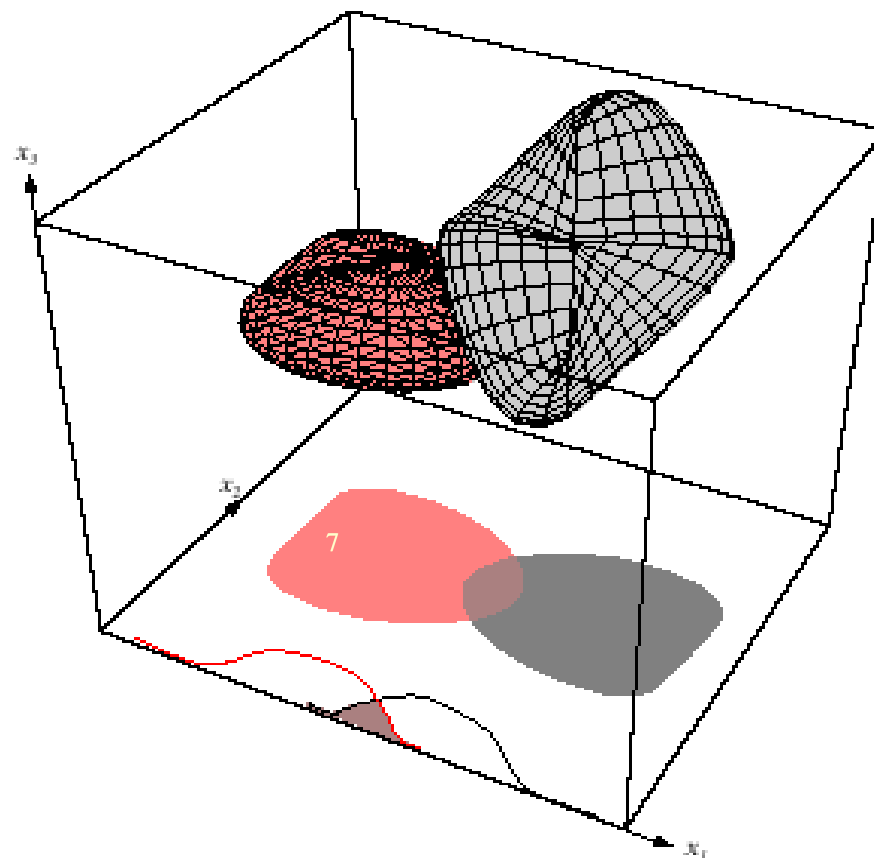


FIGURE 3.3. Two three-dimensional distributions have nonoverlapping densities, and thus in three dimensions the Bayes error vanishes. When projected to a subspace—here, the two-dimensional $x_1 - x_2$ subspace or a one-dimensional x_1 subspace—there can be greater overlap of the projected distributions, and hence greater Bayes error. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- Złożoność obliczeniowa

- Notacja “wielkie o”

$$f(x) = O(h(x))$$

Jeżeli:

$$\exists(c_0, x_0) \in \mathbb{R}^2; |f(x)| \leq c_0 |h(x)|$$

(górną granicą $f(x)$ rośnie nie szybciej niż $h(x)$
wystarczająco x o wystarczająco dużym
wymiarze!)

$$f(x) = 2 + 3x + 4x^2$$

$$g(x) = x^2$$

$$f(x) = O(x^2)$$

$$f(x) = O(x^2); f(x) = O(x^3); f(x) = O(x^4)$$

- Notacja “wielkie theta”

$$f(x) = \theta(h(x))$$

Jeżeli:

$$\begin{aligned} &\exists(x_0, c_1, c_2) \in \mathfrak{R}^3; \forall x > x_0 \\ &0 \leq c_1 g(x) \leq f(x) \leq c_2 g(x) \end{aligned}$$

$$f(x) = \theta(x^2) \text{ ale } f(x) \neq \theta(x^3)$$

- Złożoność estymatora MW
 - Rozkłady normalne a priori, d-wymiarowy klasyfikator, na każdą z c klas przypada n próbek ciągu uczącego
 - Dla każdej klasy wyznaczamy funkcję dyskryminującą:

$$g(x) = -\frac{1}{2} (x - \hat{\mu})^T \underbrace{\Sigma^{-1}}_{O(n \cdot d^2)} (x - \hat{\mu}) - \frac{\overbrace{d}^{O(1)}}{2} \ln 2\pi - \underbrace{\frac{1}{2} \ln |\hat{\Sigma}|}_{O(d^2 \cdot n)} + \underbrace{\ln P(\omega)}_{O(n)}$$

Złożoność ogółem = $O(d^2 \cdot n)$

Złożoność przy c klasach = $O(cd^2 \cdot n) \cong O(d^2 \cdot n)$

- Koszt obliczeniowy jest znaczący kiedy d i n są duże

- Analiza komponentów i analiza dyskryminacyjna
 - Łączenie cech w celu redukcji wymiaru przestrzeni cech
 - Najprościej: kombinacje liniowe cech
 - Projekcja danych z przestrzeni wysokowymiarowych do przestrzeni o mniejszej liczbie wymiarów
 - Dwa klasyczne podejścia do znajdowania „optymalnej” transformacji liniowej
 - PCA (Principal Component Analysis) “Projekcja dająca najlepszą reprezentację danych w sensie średniokwadratowym”
 - MDA (Multiple Discriminant Analysis) “Projekcja najlepiej separująca dane w sensie średniokwadratowym”

- Ukryte łańcuchy Markowa:
 - Łańcuchy Markowa
 - Cel: wyznaczyć sekwencję decyzji
 - Procesy dynamiczne, na stany w czasie t wpływają stany w czasie $t-1$
 - Zastosowania: rozpoznawanie i tagowanie mowy, rozpoznawanie gestów, sekwencjonowanie DNA,
 - Proces bez pamięci:

$$\omega^T = \{\omega(1), \omega(2), \omega(3), \dots, \omega(T)\}$$
 – sekwencja stanów, np.

$$\omega^6 = \{\omega_1, \omega_4, \omega_2, \omega_2, \omega_1, \omega_4\}$$
 - Proces może „odwiedzić” ten sam stan w różnych krokach, przy czym pewne stany w ogóle nie muszą być odwiedzane

- Model Markowa pierwszego rzędu
 - Prawdopodobieństwa przejść do kolejnych stanów:

$$P(\omega_j(t + 1) \mid \omega_i(t)) = a_{ij}$$

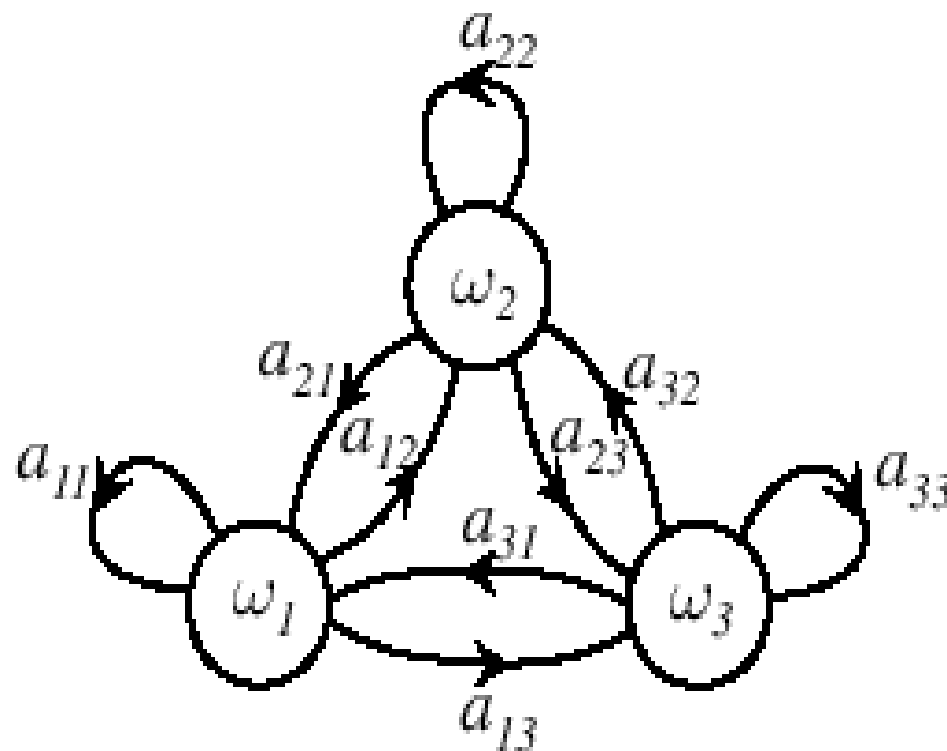


FIGURE 3.8. The discrete states, ω_j , in a basic Markov model are represented by nodes, and the transition probabilities, a_{ij} , are represented by links. In a first-order discrete-time Markov model, at any step t the full system is in a particular state $\omega(t)$. The state at step $t + 1$ is a random function that depends solely on the state at step t and the transition probabilities. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

$$\theta = (\mathbf{a}_{ij}, \omega^T)$$

$$P(\omega^T | \theta) = a_{14} \cdot a_{42} \cdot a_{22} \cdot a_{21} \cdot a_{14} \cdot$$

$$P(\omega(1) = \omega_i)$$

Przykład: rozpoznawanie mowy

“wypowiadanie kolejnych słów”

“obraz” słowa to reprezentacja przy użyciu fonemów: /p/ /a/ /tt/ /er/ /n/ // (// = stan ciszy)

Przejścia: /p/ do /a/, /a/ do /tt/, /tt/ do /er/, /er/ do /n/ i /n/ do stanu ciszy

Chapter 3:

Estymator największej wiarygodności i estymator Bayesa (part 3)

- Ukryte modele Markowa

- Ukryty model Markowa
 - Związek stanów procesu ze stanami ukrytymi
 $\sum b_{jk} = 1$ dla wszystkich j , dla których $b_{jk} = P(V_k(t) | \omega_j(t))$.
 - 3 zadania związane z tym modelem
 - The evaluation problem
 - Zadanie dekodowania
 - Zadanie uczenia

- The evaluation problem

Prawdopodobieństwo, że w trakcie działania procesu zaobserwowana zostanie sekwencja stanów V^T :

$$P(V^T) = \sum_{r=1}^{r_{\max}} P(V^T | \omega_r^T) P(\omega_r^T)$$

gdzie r jest numerem sekwencji T stanów ukrytych

$$\omega_r^T = \{\omega(1), \omega(2), \dots, \omega(T)\}$$

$$(1) \quad P(V^T | \omega_r^T) = \prod_{t=1}^{t=T} P(v(t) | \omega(t))$$

$$(2) \quad P(\omega_r^T) = \prod_{t=1}^{t=T} P(\omega(t) | \omega(t-1))$$

Z zależności (1) i (2) wynika:

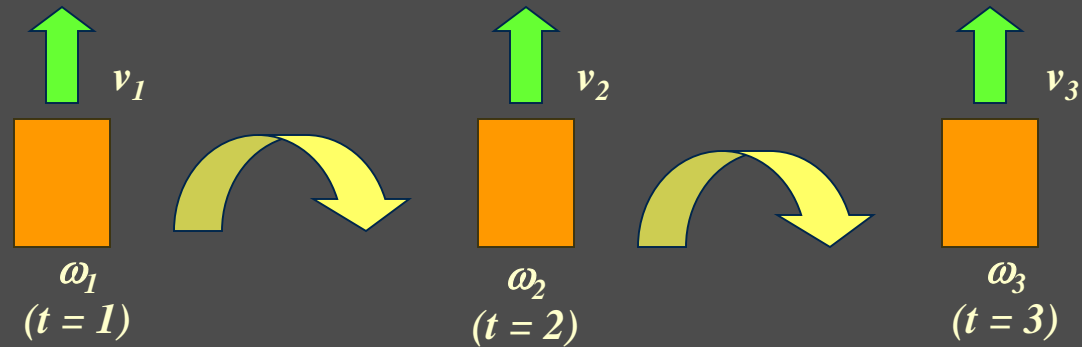
$$P(V^T) = \sum_{r=1}^{r_{\max}} \prod_{t=1}^{t=T} P(v(t) | \omega(t)) P(\omega(t) | \omega(t-1))$$

Interpretacja: Prawdopodobieństwo zaobserwowania konkretnej sekwencji T stanów procesu V^T jest sumą po wszystkich r_{\max} możliwych sekwencji stanów ukrytych, iloczynów prawdopodobieństw warunkowych przejścia do poszczególnych stanów i prawdopodobieństw zaobserwowania mierzalnych stanów w sekwencji.

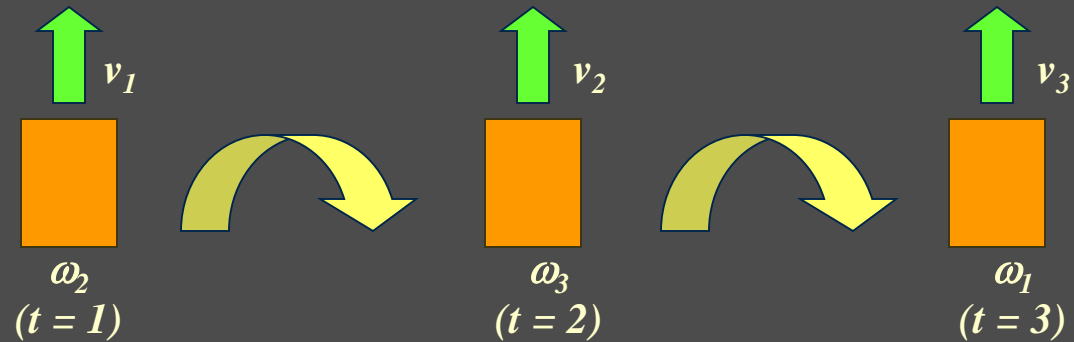
Przykład: Niech $\omega_1, \omega_2, \omega_3$ oznaczają ukryte stany; v_1, v_2, v_3 stany mierzalne a $V^3 = \{v_1, v_2, v_3\}$ jest sekwencją stanów widocznych.

$$P(\{v_1, v_2, v_3\}) = P(\omega_1) P(v_1 | \omega_1) P(\omega_2 | \omega_1) P(v_2 | \omega_2) P(\omega_3 | \omega_2) P(v_3 | \omega_3) + \dots + \text{(liczba wszystkich składników sumy to } (3^3=27)\text{)!}$$

Pierwszy przypadek:



Drugi przypadek:



$$P(\{v_1, v_2, v_3\}) = P(\omega_2) \cdot P(v_1 | \omega_2) \cdot P(\omega_3 | \omega_2) \cdot P(v_2 | \omega_3) \cdot P(\omega_1 | \omega_3) \cdot P(v_3 | \omega_1) + \dots +$$

Zatem:

$$P(\{v_1, v_2, v_3\}) = \sum_{\text{możliwe sekwencje stanów ukrytych}} \prod_{t=1}^{t=3} P(v(t) | \omega(t)) P(\omega(t) | \omega(t-1))$$

- **Zadanie dekodowania (optymalna sekwencja stanów)**

Dla sekwencji stanów mierzalnych V^T , the decoding problem polega na znalezieniu najbardziej prawdopodobnej sekwencji stanów ukrytych.

$\hat{\omega}(1), \hat{\omega}(2), \dots, \hat{\omega}(T)$, która spełnia

$$\hat{\omega}(1), \hat{\omega}(2), \dots, \hat{\omega}(T) = \arg \max_{\omega(1), \omega(2), \dots, \omega(T)} P[\omega(1), \omega(2), \dots, \omega(T), v(1), v(2), \dots, V(T) | \lambda]$$

gdzie: $\lambda = [\pi, A, B]$

$\pi = P(\omega(1) = \omega)$ (prawdopodobieństwo stanu początkowego)

$A = a_{ij} = P(\omega(t+1) = j | \omega(t) = i)$

$B = b_{jk} = P(v(t) = k | \omega(t) = j)$

Zauważmy, że brak jest tu sumowania, ponieważ poszukiwana jest jedna, najlepsza sekwencja

W poprzednim przykładzie, te obliczenia odpowiadają wyborowi najlepszej ścieżki spośród:

$$\begin{aligned} &\{\omega_1(t=1), \omega_2(t=2), \omega_3(t=3)\}, \{\omega_2(t=1), \omega_3(t=2), \omega_1(t=3)\} \\ &\{\omega_3(t=1), \omega_1(t=2), \omega_2(t=3)\}, \{\omega_3(t=1), \omega_2(t=2), \omega_1(t=3)\} \\ &\quad \{\omega_2(t=1), \omega_1(t=2), \omega_3(t=3)\} \end{aligned}$$

- Zadanie uczenia (estymacja parametrów)

Zadanie polega na zaproponowaniu metody estymacji parametrów $\lambda = [\pi, A, B]$ optymalizującej pewne kryterium. Należy znaleźć najlepszy model

$$\hat{\lambda} = [\hat{\pi}, \hat{A}, \hat{B}]$$

maksymalizujący prawdopodobieństwo zaobserwowania otrzymanej sekwencji :

$$\max_{\lambda} P(V^T | \lambda)$$

Do znalezienia lokalnego optimum można użyć metod iteracyjnych, np. metody Bauma-Welcha lub metody gradientowej