

Luonnollisen kielen tilastollinen käsittely

T-61.281 (3 ov) L

Kevät 2004

Luennot:
ta Lagus

Timo Honkela ja Kris-

Laskuharjoitukset:

Vesa Siivola

1.	Markov-mallit	3
1.1	Näkyvät Markov-mallit	3
1.2	Piilo-Markov-Mallit	9
1.3	HMM-mallin käyttäminen	11
1.4	Parametrien estimointi	19
1.5	Soveltamistavoista	23

1. Markov-mallit

1.1 Näkyvät Markov-mallit

Olkoon satunnaismuuttuja X jolla ajanhetkillä $1 \dots t$ arvot $X = (X_1, \dots, X_t)$. Satunnaismuuttujan arvo jokin tiloista $S = s_1, \dots, s_N$.

Satunnaismuuttujan arvot peräkkäisinä ajanhetkinä ovat toisistaan riippuvaisia. Yleisessä tapauksessa (ei-Markov) tn:t riippuvat koko (äärettömästä) historiasta, eli siitä miten tämänhetkiseen tilaan tultiin.

Markov-ominaisuudet:

1. Äärellinen horisontti:

$$P(X_{t+1} = s_k | X_1, \dots, X_t) = P(X_{t+1} = s_k | X_t) \quad (1)$$

2. Aikariippumattomuus (stationaarisuus):

$$P(X_2 = s_k | X_1) = P(X_3 = s_k | X_2) = P(X_{t+1} = s_k | X_t) \quad (2)$$

Markov-ketjun esitys taulukkona

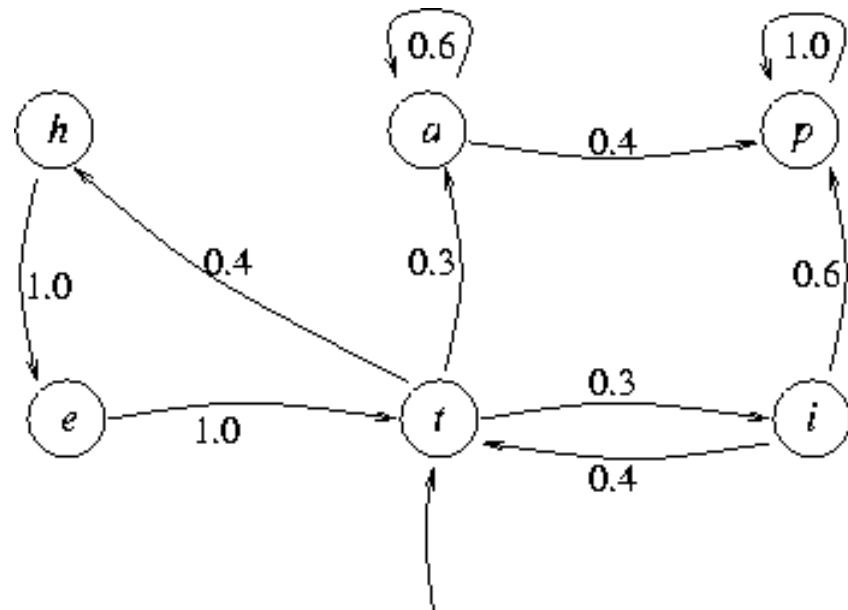
Tilojen joukko $S = \{a, e, i, h, t, p\}$

Tilasiirtymät: $a_{ij} = P(X_{t+1} = s_j | X_t = s_i)$ voidaan esittää 2-ulotteisena taulukkona:

$a_{i,j}$	a	e	i	h	t	p
a	0.6	0	0	0	0	0.4
e
i					0.6	0.4
h
t	0.3	0	0.3	0.4	0	0
p

Lisäksi määriteltävä tilojen aloitustodennäköisyydet: esim. ylimääräisen alkutilan s_0 avulla josta tilasiirtymät muihin tiloihin taulukossa.

Esittäminen probabilistisena äärellisenä tila-automaattina



- Äärellinen tila-automaatti eli Finite-State-Automaton, FSA: tilat piirretään ympyröinä, tilasiirtymät kaarina.
- Mikäli automaattia käytetään lukemaan jotain syötejonoa (ja esim. tarkastamaan onko se sallittu), tai tulostamaan kielessä sallittuja symbolijonoja, voidaan kaariin piirtää ko. tilasiirtymässä luettava/tulostettava symboli.
- epädeterministinen automaatti: yhdestä tilasta johtaa ulos monta kaarta (joissa sama symboli, ts. ei ko. tilassa tiedetä mihin tulisi mennä)
- Painotettu tai probabilistinen automaatti: jokaisella kaarella t_n , ja yksittäisestä tilasta ulosmenevien kaarien t_n :ien summa $=1$.
- Näkyvä Markov-ketju vastaa äärellistä probabilistista automaattia jossa tilajono on tunnettu. Tilajonoa voidaan siis pitää automaatin tulosteena.

Tilasekvenssin todennäköisyys:

$$P(X_1 \dots X_T) = P(X_1)P(X_2|X_1)P(X_3|X_1X_2) \dots P(X_T|X_1 \dots X_{T-1}) \quad (3)$$

$$= P(X_1)P(X_2|X_1)P(X_3|X_2) \dots P(X_T|X_{T-1}) \quad (4)$$

$$= \pi_{X_1} \prod_{t=1}^{T-1} a_{X_t X_{t+1}} \quad (5)$$

jossa $a_{X_t X_{t+1}}$ on siirtymätodennäköisyys tilasta hetkellä t tilaan hetkellä $t+1$ (katsotaan t:n:t taulukosta tai tilakoneen kaarista) ja π_{X_1} on tilan X_1 aloitustodennäköisyys.

N-grammimalli Markov-mallina:

- Esim. tri-grammimallissa $P(w_3|w_1w_2)$ merkitään tilojen joukko S on kaikkien 2-grammien joukko $S = \{w_1w_1, w_1w_2, \dots, w_1w_N, w_2w_1, \dots, w_2w_2, \dots, w_{N-1}w_{N-1}, w_{N-1}w_N\}$
- Rajallisen pituiset historiat voidaan aina esittää äärellisenä joukkona tiloja, jolloin 1. Markov-oletus täyttyy.

- Markov-mallit erotellaan historian pituuden m mukaan, ja puhutaan m :nnen asteen Markov-mallista.
- n grammimalli on siis $n - 1$:nnen asteen näkyvä Markov-malli.

1.2 Piilo-Markov-Mallit

Hidden Markov Models, HMMs, kätkeyty Markov-malli

Mallin rakennuspalat:

- Tilat: $S = \{s_1, \dots, s_N\}$
- Tilasiirtymätodennäköisyydet: $A = \{a_{ij}\}, 1 \leq i, j \leq N$
- Havainnot: $\{o_1 \dots o_t\}$
- Havaintotodennäköisyydet kussakin tilasiirtymässä: $B = \{b_{ijk}\}$

Ero näkyvään Markov-malliin: Vaikka havaintojono tunnetaan, tilasekvenssi ei yksikäsitteisesti tunnettu (kuitenkin tunnetaan tilasekvenssien tn-jakauma).

Esimerkki:

Tilat: Hyvä tuuli / Ärsyyntynyt

Havainnot: 'Upeaa!', 'Lähdetäänkö leffaan?', 'Miksei kukaan ole tyhjentänyt roskista?'

HMM:n variantteja

- mukana nollatransitioita: jotkin tilasiirtymät eivät emittoi mitään (merkitään esim. epsilonilla ϵ)
- arc-emission-HMM: havainnot emittoidaan kaarista: b_{ijk}
- state-emission-HMM: havainnot emittoidaan tiloista: b_{ik}

1.3 HMM-mallin käyttäminen

Output-sekvenssin tuottaminen

Jos HMM-malli on annettu, havaintojonojen tuottaminen siitä on triviaalia, seuraavasti:

t:=1

Arvotaan alkutila X_1 tilojen aloitustn:ien perusteella: $P(s_i) = \pi_i$

while (1)

..... Jos ollaan tilassa i , arvotaan tilasiirtymä $s_i \rightarrow s_j$ tn:llä a_{ij}

..... Arvotaan siirtymässä ij tulostettava symboli $o_t=k$ tn:llä b_{ijk}

end

Havaintojonon tn:n laskeminen

Jos on annettuna tietty havaintojono $O = o_1 \dots o_t$ ja malli $\mu = (A, B, \Pi)$, miten saadaan tehokkaasti laskettua havaintojonon O todennäköisyys mallissa, $P(O|\mu)$?

Suoraviivainen, tehoton ratkaisu: summataan havaintojonon tn kaikkien tilasekvenssien yli (edellyttää $(2T + 1)N^T$ kertolaskua, jossa T havaintojonon pituus)

$$P(O|X, \mu) = \prod_{t=1}^T P(o_t|X - t, X - t + 1, \mu) \quad (6)$$

$$P(O|\mu) = \sum_X P(O, X|\mu) = \sum_X P(O|X, \mu)P(X|\mu) \quad (7)$$

$$= \sum_{X_1 \dots X_{T+1}} \pi_{X_1} \prod_{t=1}^T a_{X_t X_{t+1}} \quad (8)$$

Huomattavasti tehokkaampi ratkaisu saadaan *dynaamisella ohjelmoinnilla* kun huomataan että sekvensseillä on yhteisiä alisekvenssejä, ja lasketaan t_n :t kunkin alisekvenssin osalta vain kerran, eli

Forward-algoritmi: (eteenpäinlaskenta)

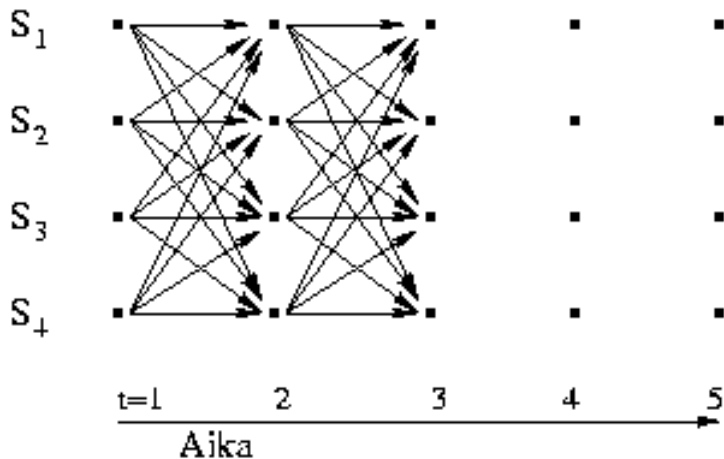
Kerätään tähänastisia t_n :iä tilakohtaisiin apuparametreihin α_i , kulkien hilaa ajassa eteenpäin:

$$\alpha_i(t) = P(o_1 o_2 \dots o_t - 1, X_t = i | \mu)$$

- Initialisoi tilojen alkutodennäköisyyksillä: $\alpha_i(1) = \pi_i$
- Induktioaskel: summaa i :n tulokaarista saapuva t_n -massa:
$$\alpha_i(t + 1) = \sum_{j=1}^N \alpha_j(t) a_{ij} b_{ij} o_t$$
- Lopuksi: $P(O | \mu) = \sum_{i=1}^N \alpha_i(T + 1)$

Forward-algoritmi

Tila



Apuparametri $\alpha_j(t)$ kussakin hilapisteessä:

Huom: **Backward-algoritmi** (taaksepäinlaskenta) samoin mutta ajassa päinvastaisena.

suuntaan.

Eteenpäin- ja taaksepäin-laskenta voidaan myös yhdistää missä kohdassa aikajanaa tahansa:

$$P(O|\mu) = \sum_{i=1}^N \alpha_i(t)\beta_i(t), \text{ jossa } 1 \leq t \leq T + 1 \quad (9)$$

Dekoodaus eli todennäköisimmän tilajonon etsiminen

Tehtävä: etsi tilajono joka parhaiten selittää havaintojonon O .

Eräs mahdollinen tulkinta: maksimoidaan oikeinarvattujen tilojen *lukumäärä*.
Kuitenkin tämä saattaa johtaa hyvin epätodennäköisiin tilasekvensseihin.

Siksi toinen tulkinta: etsitään *todennäköisin tilasekvenssi* joka on tuottanut havaintojonon.

Viterbi-algoritmi

(Tunnetaan myös nimillä DP alignment, Dynamic Time Warping, one-pass decoding)

- etsii havainnoille todennäköisimmän kokonaisen tilasekvenssin tehokkaasti
- Tallettaa jokaiseen hilapisteeseen siihenastisen todennäköisimmän po-

lun:

$$\delta_j(t) = \max_{X_1 \dots X_{t-1}} P(X_i \dots X_{t-1}, o_1 \dots o_{t-1}, X_t = j | \mu) \quad (10)$$

- Initialisoi tilojen alkutodennäköisyyksillä: $\delta_j(1) = \pi_j$
- Induktioaskel: valitse edellinen tila josta tulee suurin sekvenssin tn:
talleta tn: $\delta_j(t+1) = \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ij} o_t$
talleta pointteri ko. tilaan: $\phi_j(t+1) = \arg \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ij} o_t$
- Lopuksi: Lue hilaa lopusta taaksepäin seuraten ϕ_j :n talletamia kaaria ja kerää todennäköisin tilajono.
- Tasatilanteet voidaan ratkaista arpomalla.
- Keskeinen ero forward-algoritmiin: *td*dennäköisyyksien *summan* sijaan talletetaan *maksimi-tn*.
- Joskus yhden parhaan sekvenssin sijaan halutaan n parasta sekvenssiä (n -best-list). Tällöin talletetaan jokaisessa hilapisteessä $m < n$ parasta tn:ää ja tilaa.

- Viterbi tietyssä mielessä approksimaatio forward-algoritmille, joka ei laske *havainnon* tn :ää, ainoastaan selvittämään todennäköisimmän *tilasekvenssin* ja sen $tn:n$.

A*-dekodeausalgoritmi

Dekodeausalgoritmi joka soveltaa täydellistä forward-algoritmia (summaus), ja käy hilaa läpi pinon avulla prioriteettijärjestyksessä.

Kaksivaiheinen dekodeausalgoritmi (viterbi oli one-pass), jossa ensin nopeasti valitaan potentiaaliset seuraavat sanat, sitten lasketaan tarkemmin näiden $tn:t$ hilassa

Lisäksi on muita useampivaiheisia dekodeausalgoritmeja.

1.4 Parametrien estimointi

Jos annettuna tietty havaintojono O (opetusdata), etsitään HMM-mallin parametreille $\mu = A, B, \pi$ arvot jotka uskottavimmin selittävät havainnot. Sovelletaan MLE-estimointia:

$$\arg \max_{\mu} P(O_{\text{opetusdata}}|\mu) \quad (11)$$

Ei analyyttistä ratkaisua.

Baum-Welch eli forward-backward-algoritmi

- Iteratiivinen optimointimenetelmä, EM-algoritmin sovellus. Periaate: vuorotellaan seuraavia kahta askelta:
 1. Pidä mallin parametrit vakiona, laske $P(O|\mu)$ ja kerää tiedot siitä miten O :n tn-massa jakautui mallin hilaan (eli miten monta kertaa kuljettiin mitäkkin reittiä)
 2. Uudelleenestimoimoi mallin parametrit $\mu \rightarrow \hat{\mu}$ siten että O olisi mahdollisimman todennäköinen

- Taatusti joka kierroksella $P(O|\hat{\mu}) \geq P(O|\mu)$ (kuten EM yleensäkin)
- Malli voidaan initialisoida satunnaisesti tai jollain nopealla raa-alla menetelmällä
- Iteroidaan kunnes datan tn ei enää muutu merkittävästi.
- Löytää lokaalin maksimin, ei välttämättä globaalia (kustannusfunktio on suuren joukon parametrejä luultavasti epälineaarinen funktio).

Parametrien estimaatit sanallisesti:

$\hat{\pi}_i =$ odotettu frekvenssi tilassa i hetkellä $t = 1$

$\hat{a}_{ij} = \frac{\text{odotettu transitioiden lkm tilasta } i \text{ tilaan } j}{\text{odotettu transitioiden lkm tilasta } i}$

$\hat{b}_{ijk} = \frac{\text{odotettu transitioiden lkm tilasta } i \text{ tilaan } j \text{ kun havainto oli } k}{\text{odotettu transitioiden lkm tilasta } i \text{ tilaan } j}$

Käytännön implementointi sekä joitain ongelmia

Ongelma: Kerrotaan paljon hyvin pieniä tn:iä keskenään → ylivuoto-ongelmia.

Viterbissä vain kertolaskua ja max-operaatiota → voidaan toteuttaa kokonaan logaritmoituna, jolloin kertolaskut summia.

Baum-Welchissä: voidaan käyttää skaalauskerroimia joita kasvatetaan ajan t funktiona.

Ongelma: Parametrien suuri määrä

- tarvitaan paljon dataa TAI
- oletetaan että osa parametreista jaetaan verkon eri osien kesken ('parameter tying': sidottuja tiloja tai sidottuja tilasiirtymiä)
- oletetaan että verkossa on 'rakenteisia nollia', ts. rakenteellisesti mahdollomia tilasiirtymiä

Estimointialgoritmi löytää lokaalin maksimin, ei globaalia

- Parametrien fiksu initialisointi
- Erityisen tärkeää initialisoida havaintot:n:t $B = b_{ijk}$

1.5 Soveltamistavoista

Soveltaminen puheentunnistukseen

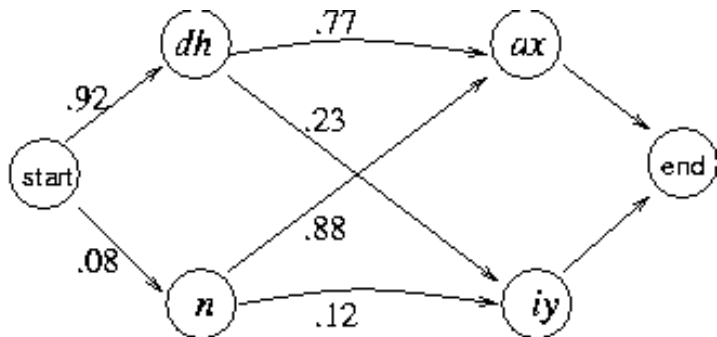
Yksinkertaistettu esimerkki: **englannin sanojen ääntämisen mallinnus** Jurafsky & Martin (2000) kuvat 7.5–7.8.

Havainnot fooneja. Sanaa vastaa joukko tiloja joita voidaan käyttää ko. sanan ääntämisessä.

Lasketaan kullekin sanalle oma HMM joka mallintaa sanan eri lausumistavat ja näiden tn:t.

Havaintojono: foonijono, esim. [dh iy]

Sana: 'the', jonka yksinkertaistettu tilamalli (kuva)



Hyvin yksinkertaistettu ääntämismalli sanalle 'the'

Jos tehtävänä on tunnistaa mikä yksittäinen sana lausuttiin, yhdistetään samamallit *rinnakkain* yhdeksi isoksi HMM:ksi

Lisätietoa: äärellisten automaattien yhdistäminen, ks. esim. Jurafsky&Martin.

Yhdistetty segmentointi + tunnistus

Viterbi-algoritmilla voidaan suorittaa samalla kertaa sekä segmentointi sanoiksi että sanojen tunnistus.

Havaintojono: foonijono, esim. [aa n iy dh ax]

Tuloste: sanajono 'I need the'

- Jokaiselle sanalle oma HMM
- Poimitaan sanaparien esiintymistn:t bi-grammimallista, ja käytetään näitä sanojenvälisiä tn:iä yhdistettäessä sanakohtaiset HMM:t *peräkkäin* (Jurafsky-Martin, kuva 7.8).
- Huom: näin tehdessä oletetaan että edellisen sanan ja seuraavan sanan lausuntatavat eivät riipu toisistaan. Oikeasti sanojen välillä riippuvuuksia lausumistavoissa, esim. savolainen viäntää, ja helsinkiläinen sihauttaa s-kirjainta sanasta riippumatta.
- Riippumattomuusoletusten tarkoitus helpottaa estimointia (vähemmän parametreja). Voidaan myös soveltaa pelkästään mallin initialisointiin.

Viterbin ongelma:

- Etsii todennäköisimmän tilajonon, mikä ei ole sama asia kuin todennäköisin sanajono.
- Kukin tilajono vastaa jotain sanojen lausuntavaihtoehtoa. Löytää siis todennäköisimmän *lausuntavaihtoehdon*.
- Kuitenkin sanajonon tn on summa sen sanajonon eri lausuntavaihtoehtojen tn:istä.
- Ei voida soveltaa kaikkien kielimallien kanssa (lähinnä bigrammin kanssa, jo trigrammien kanssa ongelmia)
- Eräs vaihtoehto: multi-pass-decoding: Viterbi nopeaan dekodaukseen joka tuottaa n-best-listan, sitten sovelletaan parempaa kielimallia näistä valitsemiseen.

Aidompi HMM:ien sovellus puheentunnistukseen:

- Havainnot puheen spektristä laskettuja piirrevektoreita.
- Tilat vastaavat fooneja
- Aika ei kulje fooni kerrallaan vaan tietynkestoisina segmentteinä.
- Erikestoiset foonit hoidetaan esim. siten että foonista on kaari myös siihen itseensä.

Muita HMM:ien sovelluksia

- sanaluokkien taggaus
- geenisekvenssien analyysi / taggaus