

FORELÄSNING I SPRÅKTEKNOLOGI

STATISTISKA METODER 2

- MARKOVPROCESSE
- HMM - GÖMDA MARKOVMODELLER
- INFORMATIONSTEORI (ENTROPI)
- TILLÄMPNING: ORDPREDIKTION
 - UTVÄRDERING (PERPLEXITET)
 - MARKOVMODELL
 - ORDFREKVENSER, ORDKLASSSTATISTIK
 - HEURISTISKA FÖRBÄTTRINGAR
 - LAGRING AV DATA

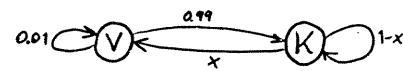
MARKOVPROCESS

EXEMPEL: STOKASTISK PROCESS X_1, X_2, \dots SOM LÄSER EN BOKSTAV I TAGET FRÅN ETT ORD OCH AVGÖR OM DEN ÄR EN VOKAL ELLER KONSONANT.

FREKVENSräKNING I KORPUS GER:

$$P[X_i = \text{vokal}] = 0.36 \text{ och } P[X_i = \text{vokal} | X_{i-1} = \text{vokal}] = 0.01$$

MODELLERA MED MARKOVKEDJA:



$$(P_{ij}) = \begin{pmatrix} 0.01 & 0.99 \\ x & 1-x \end{pmatrix}$$

$$P[X_i = V] = 0.36, P[X_i = K] = 0.64$$

$$(0.36 \ 0.64) \begin{pmatrix} 0.01 & 0.99 \\ x & 1-x \end{pmatrix} = (0.36 \ 0.64)$$

$$\Rightarrow 0.36 \cdot 0.01 + 0.64x = 0.36 \Rightarrow x = \frac{0.36(1-0.01)}{0.64} \approx 0.56$$

EN MARKOVPROCESS HAR INGET MINNE.

MODELLEN ÄR DÄRFÖR LÄMPLIG BARA OM SANNOLIGHETEN FÖR VOKAL INTET PÄVERKAS AV MER ÄN FÖREGÅENDE BOKSTAV, OCH DET GÖR DEN!

HMM - GÖMDA MARKOVMODELLER

GIVET: • TILLSTÅND q_1, \dots, q_n

- ÖVERGÅNGSSANNOLIKHETSMATRIS $P = (P_{ij})$
DÄR $P_{ij} = P[\text{ÖVERGÅNG FRÅN } q_i \text{ TILL } q_j]$

- STARTSANNOLIKHETER $v = (v_1, \dots, v_n)$

- SIGNALER $\sigma_1, \dots, \sigma_m$

- SIGNALMATRIS $A = (a_{ij})$
DÄR $a_{ij} = P[\sigma_j \text{ GENERERAS } | \text{TILLSTÅNDET ÄR } q_i]$

I EN HMM ÄR TILLSTÅNDEN OSYNLIGA OCH BARA SIGNALERNNA KAN OBSERVERAS.

EXEMPEL: TAGGNING AV EN MENING

TILLSTÅND \leftrightarrow TAGGOR, SIGNALER \leftrightarrow ORD

PROBLEM: HITTA TROLIGASTE PROCESSEN (TILLSTÅNDSFÖLJEN)
SOM GENERERAR MENINGEN

INFORMATIONSTEORI

INFORMATIONEN HOS EN SIGNAL σ SOM HAR SANNOLIKHETEN P ÄR $I_\sigma = \log_2 \frac{1}{P}$

EXEMPEL:

$$P[\sigma] = 0.25 \Rightarrow I_\sigma = \log_2 \frac{1}{0.25} = \log_2 4 = 2$$

$$P[\sigma] = 1 \Rightarrow I_\sigma = \log_2 \frac{1}{1} = \log_2 1 = 0$$

INFORMATIONS MÄNGDEN MÄTS I BITAR.

VANLIG SIGNAL \Rightarrow LITEN INFORMATIONS MÄNGD

ENTROPI

ENTROPIN ÄR VÄNTEVÄRDET FÖR INFORMATIONEN:

$$H[X] = E\left[\log_2 \frac{1}{P}\right] = \sum p_i \log_2 \frac{1}{p_i} = -\sum p_i \log p_i$$

DÄR X ÄR EN STOKASTISK VARIABEL, $P[X=\sigma_i] = p_i$

EXEMPEL:

LÄT $Y = \begin{cases} a & \text{MED SANNOLIKHET } 0.5 \\ b & \text{--- " --- } 0.25 \\ c & \text{--- " --- } 0.25 \end{cases}$

$$H[Y] = 0.5 \cdot \log_2 \frac{1}{0.5} + 0.25 \cdot \log_2 \frac{1}{0.25} + 0.25 \cdot \log_2 \frac{1}{0.25} = 0.5 + 0.25 \cdot 2 + 0.25 \cdot 2 = 1.5$$

ENTROPIN ÄR HÖGST DÅ ALLA n SIGNALER ÄR LIKA SANNOLIKA:

$$H_{\max} = \sum_{i=1}^n \frac{1}{n} \log_2 \frac{1}{1/n} = n \cdot \frac{1}{n} \log_2 n = \log_2 n$$

EXEMPEL:

$$n=3 \Rightarrow H_{\max} = \log_2 3 \approx 1.58$$

REDUNDANSEN ÄR DEN RELATIVA SKILLNADEN MELLAN H_{\max} OCH H : $R = \frac{H_{\max} - H}{H_{\max}} = 1 - \frac{H}{H_{\max}}$

$$\text{EXEMPEL: } R = 1 - \frac{1.5}{1.58} \approx 0.051 = 5.1\%$$

EXEMPEL: ENTROPIN FÖR SVENSKA

ENTROPIN FÖR SVENSKA BOKSTÄVER ÄR

$$H = \sum_{\alpha=A}^{\ddot{\alpha}} P[\alpha] \cdot \log_2 \frac{1}{P[\alpha]} = 4.34$$

Om alla bokstäver var lika vanliga vore entropin

$$H_{\max} = \log_2 29 \approx 4.86$$

$$\Rightarrow \text{REDUNDANSEN ÄR } 1 - \frac{4.34}{4.86} \approx 11\%$$

FÖRSÖK TA HÄNSYN TILL ATT BOKSTÄVERNA INTE FÖREKOMMER HERT SLUMPMÄSSIGT I SVENSKA ORD.

ENTROPIN FÖR SVENSKA ORD ÄR

$$H = \sum_{x \in \{\text{svenska ord}\}} P[x] \cdot \log_2 \frac{1}{P[x]} = 8.3$$

MEDELORDLÄNGDEN = 5.3 BOKSTÄVER

$$\text{ENTROPIN UTSLAGEN PÅ BOKSTÄVERNA ÄR } \frac{8.3}{5.3} \approx 1.6$$

$$\text{REDUNDANSEN BLIR DÄ } 1 - \frac{1.6}{4.86} \approx 67\%$$

ORDPREDIKTION

Exjobb av
JOHAN CARLBERGER
1997

PROBLEM: GIVET ORD w_1, w_2, \dots, w_{i-1}
GISSA ORDET w_i

IDÉ: Hitta X som maximerar $P[X=w_i | w_1, w_2, \dots, w_{i-1}]$

UTVÄRDERING AV EN ODPREDIKTIONSMETOD:

INFORMATIONSTEORETISK:

PERPLEXITETEN FÖR X ÄR

$$Perp[X] = 2^{-H[X]}$$

MOTSVARAR ANTALET ORD PREDIKTORN HAR ATT VÄLJA MELLAN.

OM ALLA N ORD ÄR LIKA SANNOLIKA BLIR

$$Perp[X] = 2^{\log_2 N} = N.$$

Mål: MINIMERA PERPLEXITETEN.

PRÄKTISK:

MÄT ANTALET TANGENTTRYCKNINGAR SOM BEHÖVS FÖR ATT GENERERA EN TESTTEXT MED ODPREDIKTORN.

Mål: MAXIMERA ANTALET SPARADE TANGENTTRYCKNINGAR

TESTA INTE ETT PROGRAM PÅ SAMMA TEXT SOM DET HAR TRÄNATS PÅ!

STATISTIKINSAMLING

SAMLA IN FÖLJANDE STATISTIK FRÅN EN STOR KORPUS:

ORDFREKVENSER $C_o(w)$

ORDBIGRAMSFREKVENSER $C_{oo}(w_1, w_2)$

TAGGFREKVENSER $C_t(t)$

TAGGBIGRAMSFREKVENSER $C_{t\bar{o}}(t_1, t_2)$

TAGGTREGRAMSFREKVENSER $C_{\bar{t}\bar{o}}(t_1, t_2, t_3)$

ORDTAGGPARSFREKVENSER $C_{ot}(w, t)$

ANTAL ORD I KORPUSSEN C

BERÄKNA SEDAN (OCH LÄGRA):

$$f_o(w) = \frac{C_o(w)}{C}, f_{oo}(w_1 | w_2) = \frac{C_{oo}(w_1, w_2)}{C_o(w_1)}$$

$$f_t(t) = \frac{C_t(t)}{C}, f_{\bar{o}o}(t_1 | t_2) = \frac{C_{\bar{o}o}(t_1, t_2)}{C_t(t_1)}, f_{\bar{t}\bar{o}}(t_1 | t_2, t_3) = \frac{C_{\bar{t}\bar{o}}(t_1, t_2, t_3)}{C_{t\bar{o}}(t_1, t_2)}$$

$$f_{ot}(t | w) = \frac{C_{ot}(w, t)}{C_o(w)}$$

TVÅ MARKOVMODELLER

TAGGMODELLEN

ANDRA ORDNINGENS MARKOVMODELL FÖR TAGGAR.

TILLSTÅNDEN BESKRIVS AV PAR AV TAGGAR (t_1, t_2) .

$$P_{(t_1, w_1), (t_2, w_2)} = P[t_3 | t_1, t_2] = q_0 f_{\bar{t}\bar{o}}(t_3 | t_1, t_2) + q_1 f_{\bar{t}\bar{o}}(t_3 | t_2) + q_2 f_{\bar{t}\bar{o}}(t_3)$$

DÄR q_0, q_1, q_2 ÄR LÄMPLIGT VALDA KONSTANTER

ORDMODELLEN

FÖRSTA ORDNINGENS MARKOVMODELL FÖR ORD.

TILLSTÅNDEN BESKRIVS AV ORDEN I KORPUSSEN.

$$\begin{aligned} P_{w_1, w_2} &= P[w_2 | w_1] = \\ &= r_0 f_{oo}(w_2 | w_1) + r_1 f_o(w_2) / (r_2 f_{oo}(w_2 | w_1) + r_3 f_o(w_2)) \cdot P_t(w_2) \end{aligned}$$

DÄR r_0, r_1, r_2, r_3 ÄR LÄMPLIGT VALDA KONSTANTER OCH

$$P_t(w_2) = f_{ot}(t_2 | w_2) \cdot P[t_2 | t_1, b_1]$$

LATMANNAHASHNING

HASHA BARA PÅ DOM TRE FÖRSTA BOKSTÄVERNA I SÖKNYCKELN. ANVÄND SEDAN BINÄRSÖKNING.

LÄMPLIGT FÖR SÖKNING MED FÅ DISKACESSER I STOR TEXT NÄR INDEXET INTE KAN LIGGA I PRIMÄRMINNET.

EXEMPEL: INDEXERA STORT SVENSK-ENGELSKT LEXIKON.

GIVET:

- STOR FIL L MED HEZA LEXIKONTEXTER.
- INDEX I DÄR VARJE RAD ÄR: SÖKORD POSITION I

SKAPA SÖKINDEX: SORTERA I MED SORT.

SKAPA EN INDEXARRAY $A[abc]$ SOM FÖR VARJE abc ANGER VAR I I FÖRSTA ORDET SOM BÖRJAR SÅ FINNS.

FÖRBEREDELSE INFÖR SÖKNINGAR:

LÄS IN A FRÅN FIL. ÖPPNA FILERNA I OCH L .

SÖKNINGSALGORITM (w) =

```

wprefix ← FÖRSTA 3 BOKST. I  $w$  | GÅ TILL POSITION  $i$  I  $w$ 
 $i \leftarrow A[wprefix]$  | WHILE TRUE DO
 $j \leftarrow A[wprefix+1]$  | LÄS  $w$  NÄSTA ORD FRÅN  $I$  IS
WHILE  $j-i > 1000$  DO | IF  $s=w$  THEN
     $m \leftarrow \lfloor \frac{i+j}{2} \rfloor$  | LÄS IN POSITIONEN  $i$  X
    GÅ TILL POSITION  $m$  I  $I$  | RETURN X
    LÄS IN NÄSTA ORD FRÅN  $I$  IS | IF  $s < w$  THEN
    IF  $s \leq w$  THEN  $i \leftarrow m$  | RETURN NOTFOUND
    ELSE  $j \leftarrow m$ 
  
```

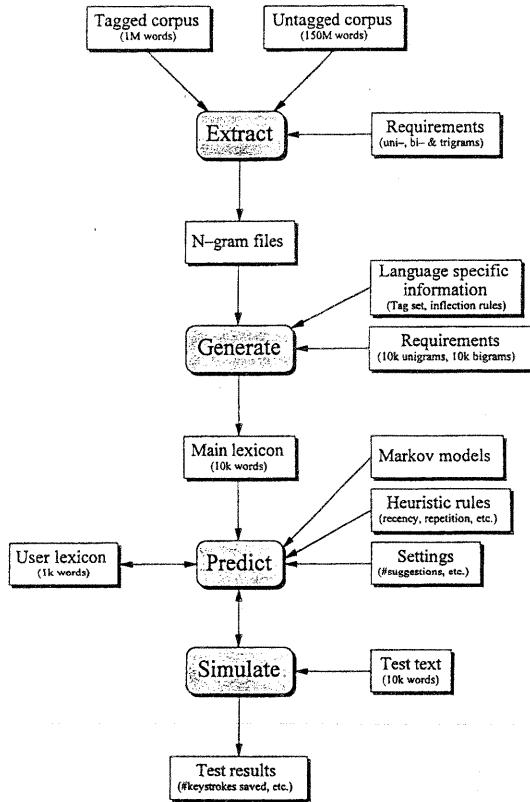


Figure 3. A system for word predictor construction.

Vitterbis algoritm

Anta att vi vet sannolikheten $P(x, y)$ för att ordklassen x ska följas av ordklassen y i en mening. Vi har också två speciella ordklasser $start$ för meningens början och $slut$ för meningens slut. $P(start, x)$ är då sannolikheten för att en mening inleds med ordklassen x och $P(y, slut)$ är sannolikheten för att en mening avslutas med ordklassen y .

Anta också att vi har ett lexikon som talar om vilka ordklasser varje ord kan ha. Exempelvis säger lexikonet att boken bara kan vara substantiv men att gula både kan vara adjektiv och substantiv.

Den trörligaste ordklassatagningen av en mening är den taggning som maximrar produkten av dom ingående sannolikheterna.

Indata består av en mening med n ord. Det finns m stycken ordklassataggar (inklusive $start$ och $slut$).

Vitterbis algoritmen löser detta problem i linjär tid i antalet ord n med hjälp av algoritmkonstruktionsmetoden dynamisk programmering.

Låt $M[i, j]$ vara lösningen (maximala produktssannolikheten) för delproblemet att tagga meningen från till ord i så att ord i får taggen j . $M[i, j]$ kan definieras rekursivt som

$$M[i, j] = \begin{cases} 1 & \text{om } i = 0 \text{ och } j = start, \\ 0 & \text{om } i = 0 \text{ och } j \neq start, \\ \max_{k \leq m} M[i-1, k] \cdot P[k, j] & \text{om } 1 \leq i \leq n \text{ och } j \notin \text{Lexikon}(w[i]), \\ \text{anmärks.} & \end{cases}$$

Vilket k -värdet som ger maximum i rekurrensen lagras i $father[i, j]$ så att algoritmen ska kunna följa den bästa taggningen tillbaka från $slut$ till $start$. Den bästa taggningen lagras sedan i $POS[1..n]$.

Indata: $P[1..m, 1..m]$, $w[1..n]$, Lexikon: ord → set of integer
for $j \leftarrow 1$ to m do $M[0, j] \leftarrow 0$
 $M[0, start] \leftarrow 1$

```

for  $i \leftarrow 1$  to  $n+1$  do
  if  $i \leq n$  then  $S \leftarrow \text{Lexikon}(w[i])$  else  $S \leftarrow \{slut\}$ 
  for  $j \leftarrow 1$  to  $m$  do
     $tmp \leftarrow M[i-1, k] \cdot P[k, j]$  else  $tmp = 0$ 
    foreach  $j \in S$  do
      if  $tmp > M[i, j]$  then  $M[i, j] \leftarrow tmp$ ;  $father[i, j] \leftarrow k$ 
    write  $w[i], POS[i]$ 
  
```

I värsta fallet kan varje ord ha m ordklasser. Dom tre nästlade forsslingorna går då $n+1, m$ respektive m varv, vilket ger tidskomplexiteten $O(nm^2)$.