



Försättsblad till skriftlig tentamen vid Linköpings Universitet

(fylls i av ansvarig)

Datum för tentamen	<i>2010-04-08</i>
Sal	<i>TER2</i>
Tid	<i>14-18</i>
Kurskod	<i>TDDD02</i>
Provkod	<i>TEN1</i>
Kursnamn/benämning	<i>Språkteknologi för informationssökning</i>
Institution	<i>IDA</i>
Antal uppgifter som ingår i tentamen	<i>11</i>
Antal sidor på tentamen (inkl. försättsbladet)	<i>3</i>
Jour/Kursansvarig	<i>Lars Ahrenberg</i>
Telefon under skrivtid	<i>0703182422</i>
Besöker salen ca kl.	<i>-</i>
Kursadministratör (namn + tfnr + mailadress)	<i>Elisabeth Qvarnström 013-281706, eliqv@ida.liu.se</i>
Tillåtna hjälpmedel	<i>inga</i>
Övrigt (exempel när resultat kan ses på webben, betygsgränser, visning, övriga salar tentan går i m.m.)	
Vilken typ av papper ska användas, rutigt eller linjerat	<i>valfritt</i>
Antal exemplar i påsen	<i>12+2</i>

TENTAMEN

TDDD02 Språkteknologi för informationssökning
torsdag 8 april 2010 kl. 14-18

Inga hjälpmedel är tillåtna.

Varje fråga är värd 3 poäng. Maximal poäng är 33. 16,5 poäng ger säkert godkänt. Det går bra att besvara flera frågor på samma papper.

1. I dokument sökningssystem används ofta någon form av textnormalisering, t.ex. via *stemming* eller *lemmatisering*. Förklara vad dessa bearbetningar går ut på och ge exempel som visar skillnaden mellan dem.
2. Förklara vad som menas med termviktning enligt tf-idf-modellen, dvs förklara vad tf står för, vad idf står för och hur de kombineras vid termviktning.
3. Ange en reguljär substitution eller en ändlig transduktor (eng. *Finite-state transducer*) som kan utföra mappningar som de följande. Din lösning ska generalisera till andra substantiv som slutar på -el, -er eller -en.

sedel → sedel N SG

sedlar → sedel N PL

utter → utter N SG

uttrar → utter N PL

öken → öken N SG

öknar → öken N PL

4. System för namnigenkänning utvärderas vanligen med både precision och recall. (a) Förklara varför det inte är tillräckligt att bara ange hur många namn systemet hittat, (b) Ange ett mått, med definition, som väger ihop precision och recall.
5. I en text på sammanlagt 2000 ord finns 150 som börjar med stor bokstav. 80 av dessa inleder en mening och 90 är egennamn. Det finns inga meningar och inga egennamn som inte börjar med stor bokstav. Det finns 30 ordbigram, d.v.s. båda de ingående orden är egennamn. Använd Maximum Likelihood-uppskattning för att uppskatta följande sannolikheter: (a) sannolikheten för att ett ord som börjar med stor bokstav är ett egennamn, (b) sannolikheten för att ett egennamn följs av ett egennamn, (c) sannolikheten för att ett ord som inte står först börjar med stor bokstav.

6. Ange (a) ett uttryck för sannolikheten av ordsekvensen *men han kom inte* givet att vi känner till bigramsannolikheterna för orden i meningen, (b) samma fråga men under förutsättning att vi känner till de ingående ordens trigramsannolikheter.

7. System för ordklasstaggning utnyttjar ords lokala kontext för att välja rätt tagg. Ange ett vanligt regelformat för sådan disambiguering och formulera någon eller några explicita regler som kan appliceras på ordet *visa* i dessa två meningar:

Eva ville visa oss sina leksaker.

Jag ville sjunga en visa om kärleken.

8. För att hitta ord med specifika semantiska relationer till varandra kan man använda mönster kodade som reguljära uttryck. Tre exempel på sådana mönster ges nedan. Ange för vart och ett av dem vilken eller vilka semantiska relationer man är ute efter. Svaret ska ange den relation som ordet X har till ordet Y.

(a) ... andra X så som Y, Z och W ...

(b) ... är X eller Y spelar ingen roll ...

(c) ... en/ett X är en/ett Y som ...

9. I modellen den brusiga kanalen (eng. *Noisy Channel*) används beslutsregeln

$$S^* = \underset{S}{\operatorname{argmax}} p(S|O)$$

(a) Förklara denna regel i ord. (b) Beräkningen av $p(S|O)$ baseras vanligtvis på en omskrivning som leder till två olika sannolikhetsfördelningar. Visa hur omskrivningen ser ut och ange vad de resulterande sannolikheterna kallas.

10. Ange fyra vanliga komponenter i ett informationsextraktionssystem och förklara vad de gör. Du kan förutsätta att indata till den första komponenten är en tokeniserad ordklasstaggad text som är relevant för systemets domän.

11. I extraktionsmetoden för textsammanfattning används olika kriterier för att bestämma om en mening ska inkluderas i sammanfattningen eller inte. Nämn tre sådana kriterier.