



Försättsblad till skriftlig tentamen vid Linköpings Universitet

(fylls i av ansvarig)

Datum för tentamen	<i>2010-08-24</i>
Sal	<i>TER3</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>10</i>
Antal sidor på tentamen (inkl. försättsbladet)	<i>3</i>
Jour/Kursansvarig	<i>Lars Ahrenberg</i>
Telefon under skrivtid	<i>2422</i>
Besöker salen ca kl.	
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	
Antal exemplar i påsen	<i>5+2</i>

TENTAMEN

TDDD02 Språkteknologi för informationssökning
tisdag 24 augusti 2010

Inga hjälpmedel är tillåtna.

Varje fråga är värd 3 poäng. Maximal poäng är 30. 15 poäng ger säkert godkänt. Det går bra att besvara flera frågor på samma papper.

1. För att representera ett dokument görs ofta olika former av bearbetningar såsom stemming eller eliminering av stopword. Förklara vad dessa båda bearbetningar innebär och motivera varför de används.
2. Ett dokumentökningsystems prestanda mäts ofta med måtten precision och recall.
(a) Ge definitioner av dessa två mått, (b) Förklara varför båda måtten behövs.
3. Förklara (a) vad som menas med en term-dokumentmatris och (b) måttet $tf*idf$.
4. Ange ett reguljärt uttryck som känner igen prisangivelser upp till 10,000:-. Uttrycket ska kunna hantera öresangivelser med kolon som i 89:50 eller 25:- och tusental med eller utan kommatecken som i 5,500 eller 5500. Däremot kan du anta att det inte ingår något blanktecken i uttrycket.
5. Ange med ledning av frekvensangivelserna i tabellen nedan Maximum Likelihood-uppskattningar av (a) bigramsannolikheten $p(\text{inte}|\text{kan})$, (b) trigramsannolikheten $p(\text{komma}|\text{kan inte})$. c) Ange också en nackdel med Maximum Likelihood-uppskattningar.

inte	300
kan	250
komma	200
kan inte	40
inte kan	2
inte komma	5
kan inte komma	4
komma kan inte	1

6. Ange ett uttryck för sannolikheten av sekvensen ' $\langle s \rangle$ nu skiner solen igen' givet att vi känner till bigramsannolikheterna för orden i sekvensen. $\langle s \rangle$ anger meningsstart.

7. Förklara begreppen hyponymi och antonymi med exempel.
8. Låt w_1^N stå för en sekvens av N ord w_1, w_2, \dots, w_N och $t_1^N = t_1, t_2, \dots, t_N$ vara en motsvarande sekvens av ordklasstaggar. Med en sannolikhetsbaserad modell för ordklasstagging anges den bästa taggningen t^* med formeln

$$t^* = \underset{t_1^N}{\operatorname{argmax}} p(t_1^N | w_1^N)$$

- (a) Uttryck ovanstående formel i ord. (b) Skriv om den med hjälp av Bayes regel som produkten av två sannolikheter, (c) Ge något skäl varför produkten $\prod p(t_i|w_i)$ med $1 \leq i \leq N$ ger en dålig approximation av t^* .
9. Antag att vi ur texter vill extrahera så många korrekta par $\langle A, B \rangle$ som möjligt, där A står för ett företag och B för företagets vd. Exempel: Om en text innehåller meningen "Olle Olsson, vd på Informatia AB, kandiderar för en riksdagsplats." vill vi extrahera paret $\langle \text{Informatia AB}, \text{Olle Olsson} \rangle$. Beskriv hur man kan använda en s.k. bootstrapping-metod för att göra detta.
10. (a) Vad skiljer informationsutvinning (eng. information extraction) från informationssökning? (b) Ange minst fyra vanliga komponenter i ett informationsutvinningssystem.