



Försättsblad till skriftlig tentamen vid Linköpings Universitet

(fylls i av ansvarig)

Datum för tentamen	<i>2010-08-20</i>
Sal	<i>TER1</i>
Tid	<i>14-18</i>
Kurskod	<i>TDDD01</i>
Provkod	<i>TEN1</i>
Kursnamn/benämning	<i>Språkteknologi</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>Maria Holmqvist</i>
Telefon under skrivtid	<i>Ank. 1466 eller 0709-736834</i>
Besöker salen ca kl.	<i>15:00</i>
Kursadministratör (namn + telnr. + mailadress)	<i>Anna Grabska Eklund 013-282362, annek@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	

TENTAMEN

TDDD01 Språkteknologi
fredag 20 augusti 2010

Inga hjälpmedel är tillåtna.

Maximal poäng finns angiven för varje fråga. Maximal poäng på hela tentan är 26.
13 poäng ger säkert godkänt.

Det går bra att besvara flera frågor på samma papper.

1. Vad menas med att tokenisera en text? (2p)
2. Definiera måtten precision och recall, t.ex. för ett namnigenkänningsystem och förklara varför det inte räcker med bara ett av dessa mått för att beskriva systemets kvalitet. (3p)
3. (a) Vad menas med en trigrammodell för att predicera ord i kontext? (b) Vilket är det vanligaste sättet att uppskatta trigramsannolikheter? (2p)
4. Skriv ett reguljärt uttryck som kan känna igen prisangivelser upp till 10,000:-. Uttrycket ska klara av öresangivelser som i 99:90, 850:- och 75:00 men du kan anta att alla tecken är sammanskrivna, dvs att det inte finns något blanktecken inuti prisangivelsen. (3p)
5. Vad gör en Porter-stemmer och varför vill man använda en sådan? (2p)
6. I chartparsning används s.k. punkterade regler. Ett exempel är 'S → NP . VP'. Vad anger punkten? (2p)
7. Förklara principerna för en kaskadparser och ange minst två centrala komponenter i en sådan. (3p)
8. Språkteknologi använder sig ofta av rudimentära semantiska representationer (som många inte tycker förtjänar epitetet 'semantiska'), som t.ex. ordpåsar (eng. *bag-of-words*) eller mallar (eng. *templates*). Förklara vad dessa båda representationer står för och ange för var och en av dem någon typ av system som använder dem. (2p)
9. Ange två olika slags särdrag (eng. *features*) som är användbara i samband med ordbetydelsebestämning. (2p)

10. Förklara vad som menas med redigeringsavstånd och beskriv någon algoritm för att beräkna redigeringsavståndet mellan två ord. Tillämpa algoritmen på något exempel och ange slutligen vilka för- och nackdelar algoritmen har i samband med stavningskontroll. (3p)
11. En grundläggande modell för statistisk maskinöversättning är den brusiga kanalen (eng. The Noisy Channel). Beskriv denna modell och hur den är tänkt att ge oss en lösning för den bästa översättningen. (2p)