



# Försättsblad till skriftlig tentamen vid Linköpings universitet

(fylls i av ansvarig)

<b>Datum för tentamen</b>	2014-03-17
<b>Sal</b>	G33
<b>Tid</b>	8-12
<b>Kurskod</b>	TDDD01
<b>Provkod</b>	TEN1
<b>Kursnamn/benämning</b>	Språkteknologi
<b>Institution</b>	<i>IDA</i>
<b>Antal uppgifter som ingår i tentamen</b>	9 + 2
<b>Antal sidor på tentamen (inkl. försättsbladet)</b>	4
<b>Jour/Kursansvarig</b>	Marco Kuhlmann
<b>Telefon under skrivtid</b>	4644
<b>Besöker salen ca kl.</b>	09:00
<b>Kursadministratör (namn + tfnr + mailadress)</b>	<i>Isa Kärman 013-285760, asa.karrman@liu.se</i>
<b>Tillåtna hjälpmedel</b>	inga
<b>Övrigt (exempel när resultat kan ses på webben, betygsgränser, visning, övriga salar tentan går i m.m.)</b>	
<b>Vilken typ av papper ska användas, rutigt eller linjerat</b>	valfritt
<b>Antal exemplar i påsen</b>	





# Tentamen

Marco Kuhlmann, Institutionen för datavetenskap, Linköpings universitet  
marco.kuhlmann@liu.se

17 mars 2014

Inga hjälpmedel är tillåtna. Maximal poäng finns angiven för varje fråga. Maximal poäng på hela tentamen är 24; 12 poäng ger säkert godkänt. Det går bra att besvara flera frågor på samma papper.

## Del A

Besvara alla frågor i denna del. Varje fråga ger 2 poäng.

1. Vid språkteknologiskt korpusarbete tokeniserar man helst texter så att skiljetecken som komma och punkt utgör egna token. (a) Varför gör man det? (b) Varför resonerar man annorlunda när det gäller förkortningspunkter?

2. I en svensk korpus finner vi följande frekvenser för några utvalda ord och ordsekvenser:

*med*: 32 900; *tanke*: 400; *på*: 24 500; *med tanke*: 260; *tanke på*: 270; *med tanke på*: 250

Vad är den Maximum Likelihood-uppskattade sannolikheten  $P(\text{på} \mid \text{med tanke})$  om vi använder (a) trigramsannolikheter, (b) en omskrivning till bigramsannolikheter?

3. (a) Vilka typer av sannolikheter ingår i en Hidden Markov-modell för ordklasstagning? (b) Vilka konkreta sannolikheter måste man ha skattat för att i en sådan modell kunna räkna ut den kombinerade sannolikheten för följande taggade mening?

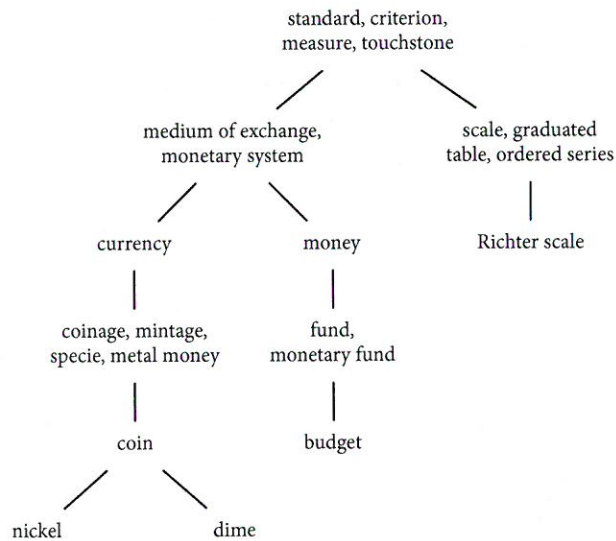
*Jag/PN äter/VB*

4. Nedanstående tabell visar reglerna i en probabilistisk kontextfri grammatik. Rita två olika frasstrukturträd (parseträd) enligt denna grammatik och ange deras sannolikheter.

$S \rightarrow NP VP$	1,00
$NP \rightarrow Lotta$	0,25
$NP \rightarrow cykeln$	0,75
$VP \rightarrow V NP$	1,00
$V \rightarrow lånar$	1,00



5. Nedanstående graf visar en liten del av WordNet. (a) Förklara vad noderna och bågarna representerar. (b) Bestäm avståndet (eng. *pathlength*) mellan "nickel" och "budget" och utifrån detta avstånd räkna ut den semantiska likheten mellan de två orden.



6. Ett automatiskt system för författaridentifikation ska hitta texter som är skrivna av författaren A. En utvärdering av systemet på en guldstandard ger följande resultat, där + betyder att texten är skriven av A och - betyder att texten inte är skriven av A. Räkna ut systemets precision och recall.

	system +	system -
guldstandard +	1	3
guldstandard -	0	19

7. En enkel metod för att tagga filmrecensioner med polariteter är att använda en Naive Bayes-klassificerare. Ange klassificerarens beslutsregel och förklara den.
8. (a) Vad innebär entitetsextraktion (eng. *named entity recognition*)? (b) Hur kan entitetsextraktion hanteras som ett taggningsproblem?
9. En central modul i ett frågebesvarande system är en analysator som bestämmer frågans svarstyp. Förklara vad som menas med en svarstyp och ge några exempel på möjliga svarstyper.



## Del B

Välj en fråga och besvara den utförligt. Varje fråga kan ge maximalt 6 poäng.

1. En probabilistisk parser ska räkna ut den mest sannolika syntaktiska analysen för en given mening. Förklara varför denna uppgift är beräkningsmässigt utmanande. Beskriv därefter två metoder för att bemöta denna utmaning.
2. Det finns för närvarande ett stort intresse inom både akademien och industrin i metoder för att analysera språk i sociala medier. Ange några skäl till detta. Diskutera därefter några av de utmaningar som språkteknologin ställs inför när den ska tillämpas på texter från Twitter och Facebook snarare än t.ex. tidningar och lexikon.

