

# Försättsblad till skriftlig tentamen vid Linköpings universitet



Datum för tentamen	2019-04-24
Sal (1)	U11(6)
Tid	14-18
Utb. kod	TDDE17
Modul	KTR1
Utb. kodnamn/benämning Modulnamn/benämning	Introduktion till språkteknologi Kontrollskrivning
Institution	IDA
Antal uppgifter som ingår i tentamen	5
Jour/Kursansvarig Ange vem som besöker salen	Marco Kuhlmann
Telefon under skrivtiden	013-284644
Besöker salen ca klockan	Besöker inte salen utan har endast jourtelefon. Se ovan.
Kursadministratör/kontaktperson (namn + tfnr + mailaddress)	Veronica Kindeland Gunnarsson 013-28 56 34 veronica.kindeland.gunnarsson@liu.se
Tillåtna hjälpmedel	Inga hjälpmedel tillåtna.
Övrigt	
Antal exemplar i påsen	

## Dugga 2019-04-24

Examinator: Marco Kuhlmann

Denna dugga består av 5 uppgifter som prövar din förståelse av de grundläggande begrepp och procedurer som behandlas på kursen. Dessa uppgifter kräver endast kompakta redogörelser, t.ex. en kort text, en uträkning eller ett diagram. Varje uppgift är värd 3 poäng. För att bli godkänd på duggan krävs minst 12 poäng totalt.

*Lycka till!*

## 01 Textklassificering

(3 poäng)

- a) Vid utvärderingen av en textklassificerare fick man ut nedanstående förväxlingsmatris. Den markerade cellen anger antalet gånger systemet klassade ett dokument som klass C medan det enligt guldstandarderna tillhörde klass A.

	A	B	C
A	58	6	1
B	5	11	2
C	0	7	43

Ange följande värden. Svara med bråk.

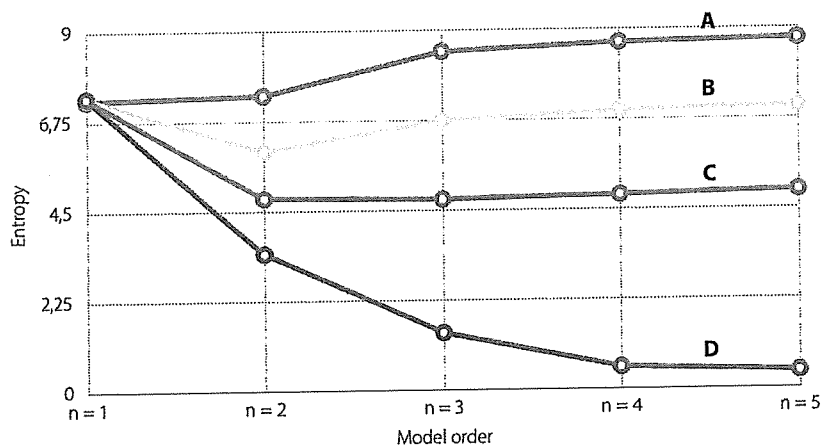
- i. precision med avseende på klass A    ii. recall med avseende på klass C
- b) Du ska bygga en textklassificerare som ska avgöra om ett inkommande dokument är en nyhet om Kina eller Japan. Träningsmängden innehåller 90% nyheter om Kina och 10% nyheter om Japan, och du räknar med samma fördelning även i framtiden. Förklara varför korrekthet (*accuracy*) är ett dåligt utvärderingsmått för din textklassificerare.
- c) Ange klassificeringsregeln för Naive Bayes-klassificeraren. Du kan använda antingen standardsannolikheter eller log-sannolikheter. Förklara din notation.

$$\hat{c} = \dots$$

Datamängden *Corpus of Contemporary American English* (COCA) består av ca. 560 miljoner token och innehåller 1 254 193 unika ord. Vi hittar följande frekvenser av unigram och bigram:

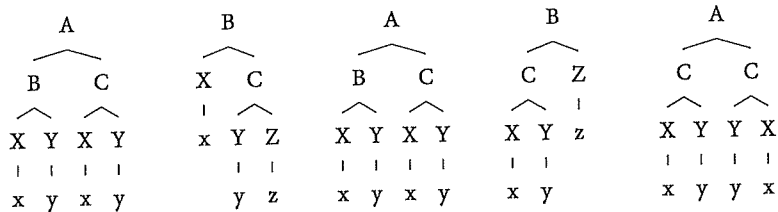
<i>snow</i>	<i>white</i>	<i>white snow</i>	<i>purple</i>	<i>purple snow</i>
38,186	256,091	122	11,218	0

- a) Skatta följande sannolikheter med hjälp av Maximum Likelihood-metoden (utan utjämning). Svara med bråk.
- $P(\textit{snow})$
  - $P(\textit{snow} \mid \textit{white})$
- b) Skatta bigramsannolikheten  $P(\textit{snow} \mid \textit{purple})$  med hjälp av Maximum Likelihood-metoden med addera- $k$  utjämning,  $k = 0.1$ . Svara med ett bråk.
- c) Vi använder Maximum Likelihood-metoden med addera- $k$  utjämning för att skatta  $n$ -gram-modeller på COCA, där  $n \in \{1, \dots, 5\}$  och  $k \in \{0, 0.01, 0.1, 1\}$ . Nedanstående graf visar entropin för varje modell när den utvärderas på samma data som den tränats på. Vilken kurva svarar mot vilket  $k$ -värde, och varför? Svara med en kort text.





- a) Du ska skatta en probabilistisk kontextfri grammatik utifrån den lilla trädbank som du ser här nedan. Skriv ner alla regler vars vänsterled är antingen B eller C och skatta deras sannolikheter med hjälp av maximum likelihood-metoden (utan smoothing). Svara med bråk.



- b) Du summerar alla regelsannolikheter i en viss probabilistisk kontextfri grammatik. Vilka (noll eller flera) av följande resultat kan du *inte* få, och varför? Svara med en kort text.

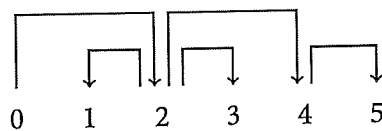
i. 0.42

ii. 1

iii. 4.2

iv. 42

- c) Ange två olika transitionssekvenser som gör att en dependensparser skapar nedanstående dependensträd:



- a) Välj en semantisk relation: synonym, antonym, hyponym, hyperonym?

blomma	är ... till	ros
människa	är ... till	barn
munter	är ... till	glad
ond	är ... till	god
maträtt	är ... till	ärtsoppa

- b) Rita ett partiellt WordNet-träd för följande synsets:

1 universitet	4 sjukhus
2 institution, inrättning	5 förskola
3 akademisk institution	6 medicinsk institution

- c) Vi läser av ordvektorer för ord på klingonska från nedanstående samförekomstmatis. (Målord svarar mot rader, kontextord svarar mot kolumner.)

	<i>HuSHa'</i>	<i>Ha'DIbaH</i>
<i>qa'vIn</i>	5	1
<i>qurgh</i>	5	5
<i>jonta'</i>	1	0
<i>Dargh</i>	1	4

Sortera de fyra målorden i avstigande grad av likhet till ordet *jonta'*. Antag att semantisk likhet kan mätas med hjälp av cosinusmättet.