

IN1140: Introduksjon til språkteknologi

Forelesning #8

Samia Touileb

Universitetet i Oslo

05. oktober 2020



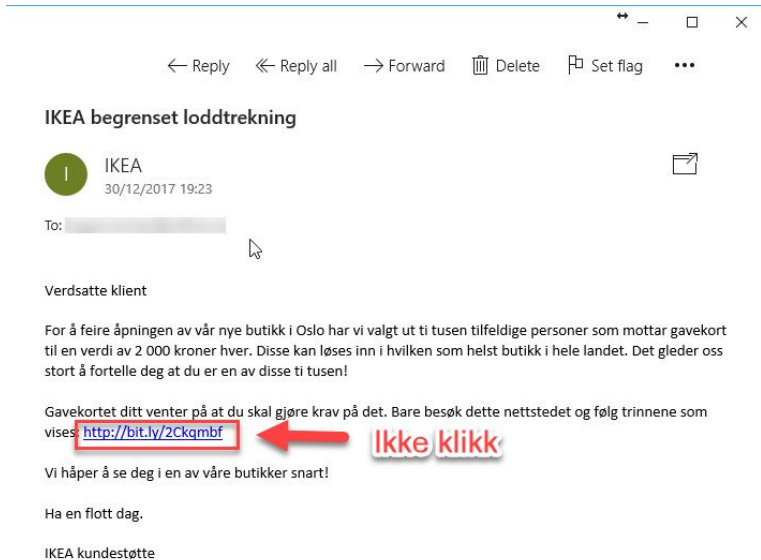
I dag skal vi se på:

- ▶ Maskinlæring og Klassifisering
 - ▶ Naive Bayes klassifisering

- ▶ Obligatorisk påmelding til fysiske gruppetimer hver uke.
- ▶ Petter har ekstratime om Python på tirsdag.

Maskinlæring og Klassifisering

Spam? Ikke spam?



Positive eller Negative?

- ▶ En varm og velopplagt bok. +
- ▶ Forfatteren burde ha fått bedre manushjelp av forlaget sitt. -
- ▶ Oppskrytt dameroman. -
- ▶ Denne boken er et hån mot underholdningssjangeren. -
- ▶ Dette er kanskje årets beste barnebok. +

Hva handler artikkelen om?

BBC: Sjansene for en brexitavtale mindre etter nattens forhandlinger

Det pågikk intense brexitforhandlinger mellom Storbritannia og EU i natt. Men det ble ikke noe gjennombrudd. Sjansene for en avtale er blitt mindre, sier britiske regjeringsskilder til BBC.



Veggmaleri av Boris Johnson på en bygning i London
FOTO: HANNAH MCKAY



Øystein Heggen
Journalist

Publisert i dag kl. 10:30
Oppdatert for én time siden

Kategori?

- ▶ Politikk
- ▶ Sport
- ▶ Underholdning
- ▶ Kultur
- ▶ ...

https://www.nrk.no/urix/forhandlinger-om-brexit-i-natt-_uten-losning-1.14744008

Automatisk tildeling av kategorier, tema, sjangere ...

- ▶ Spam-gjenkjenning.
- ▶ Sentiment analyse.
- ▶ Identifikasjon av forfattere/forfatterskap.
- ▶ Alder- og kjønnsidentifikasjon.
- ▶ Språkgjenkjenning.
- ▶ ...

- ▶ Input:
 - ▶ et dokument d
 - ▶ et forhåndsdefinert sett med klasser $C = \{c_1, c_2, \dots, c_n\}$
- ▶ Output:
 - ▶ tildele en klasse c

Manuell klassifisering?

- ▶ Regelbasert klassifisering basert f.eks. på ord:
 - ▶ Spam: inneholder “dårlig språk” AND “nevner ukjent person” AND “Klikk her!”
- ▶ Regler definert av eksperter → gode resultater.
- ▶ **Men** vanskelig å lage og vedlikeholde slike regler.

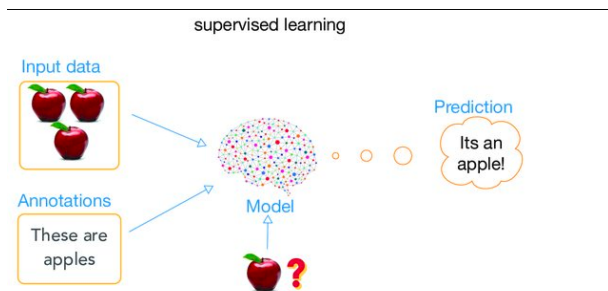
- ▶ Sentral metode innenfor maskinlæring.
- ▶ Automatisk avgjøre hvilken kategori en observasjon tilhører.
- ▶ Basert på annotert **treningsdata**: observasjoner der kategorien er kjent.
- ▶ **Supervised** klassifisering: klassifisering som benytter annotert treningsdata.

► Input:

- et dokument d
- et forhåndsdefinert sett med klasser $C = \{c_1, c_2, \dots, c_n\}$
- et treningsett m av manuelt annoterte dokumenter $(d_1, c_1), \dots, (d_m, c_m)$

► Output:

- en trent klassifiserer $y: d \rightarrow c$

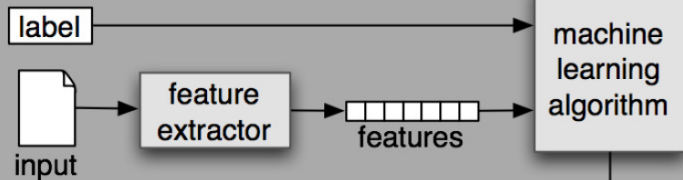


https://www.researchgate.net/publication/329533120/figure/fig1/AS:702267594399761015444445050584/Supervised-learning-and-unsupervised-learning-Supervised-learning-uses-annotation_W640.jpg

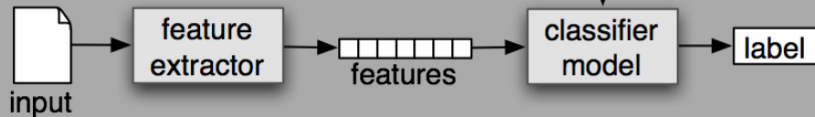
3 splits:

- ▶ **Train**: data for å trene modellen.
- ▶ **Dev**: data for å evaluere modellen underveis.
- ▶ **Test**: data for å evaluere den beste modellen etter trening og evaluering på dev.

(a) Training



(b) Prediction

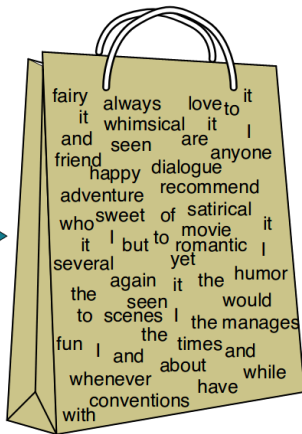


- ▶ Første skritt består i å finne og å hente ut **trekk** (“features”) fra treningsdataene.
- ▶ Eksempel:
 - ▶ En varm og velopplagt bok. **POS**
 - ▶ Forfatteren burde ha fått bedre manushjelp av forlaget sitt. **NEG**
 - ▶ Oppskrytt dameroman. **NEG**
 - ▶ Denne boken er et hån mot underholdningssjangeren. **NEG**
 - ▶ Dette er kanskje årets beste barnebok. **POS**
- ▶ Hvilke **trekk** (“features”) kan vi bruke for å skille mellom de positive og negative anmeldelser?

Trekk: Bag-of-words



I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

<https://web.stanford.edu/~jurafsky/slp3/4.pdf>

- ▶ Ordklasser.
- ▶ Bi-grams, tri-grams, n-grams av ord eller karakterer.
- ▶ Top N mest frekvente ord.
- ▶ Syntaktiske funksjoner.
- ▶ ...

- ▶ Gitt treningsdataene og trekkvektorene, kan en rekke forskjellige maskinlæringsalgoritmer brukes til å trene en klassifiserer.
- ▶ Her skal vi se på **Naive Bayes**-klassifisering:
 - ▶ Statistisk klassifiserer.
 - ▶ Bruker informasjon om ord i konteksten.
 - ▶ Enkel og naiv metode basert på Bayes regelen.

- ▶ **Hovedantagelse:**

for å finne en klasse \hat{c} (hentet fra alle mulige klasser C) for en trekkvektor \vec{f} må vi beregne den mest sannsynlige klassen, gitt vektoren

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c | \vec{f})$$

- ▶ Men det er problematisk å trene direkte: "sparse data"-problemet (alltid finnes både ord og sekvenser vi ikke har sett.)
- ▶ Kan bruke **Bayes teorem**

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- Betinget sannsynlighet:

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

- Produktsetningen:

$$P(A, B) = P(A|B)P(B) = P(B|A)P(A)$$

- Bayes regel:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- Omformulering ved Bayes teorem:

$$\hat{c} = \operatorname{argmax}_{c \in C} \frac{P(\vec{f}|c)P(c)}{P(\vec{f})}$$

- Men fremdeles ikke nok data for denne beregningen...
- Bryte opp trekkvektoren og se på individuelle trekk i kombinasjon med klasser.
- En **uavhengighetsantagelse**: trekk er **uavhengige** av andre trekk

$$P(\vec{f}|c) \approx \prod_{j=1}^n P(f_j|c)$$

- ▶ Naive Bayes klassifiserer

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c) \prod_{j=1}^n P(f_j|c)$$

- ▶ vi **trener** klassifisereren ved å beregne sannsynligheter fra et korpus (MLE – Maximum Likelihood Estimation)

► To sannsynligheter:

1. **prior**-sannsynligheten for klassen $P(c)$

$$P(c) = \frac{N_c}{N_{doc}}$$

N_c = antall dokumenter i treningsdataen som er i klassen c

N_{doc} = total antall dokumenter

2. sannsynligheten for **individuelle trekk** $P(f_j|c)$

$$P(f_j|c) = \frac{count(f_j, c)}{count(c)}$$

- For å beregne $P(f_j|c)$ kan vi anta at et trekk er ett ord som finnes i dokumentets bag of words, og kan derfor heller beregne følgende:

$$P(w_i|c) = \frac{\text{count}(w_i, c)}{\sum_{w \in V} \text{count}(w, c)}$$

- Med andre ord:

Antall ganger ordet w_i forekommer i alle dokumenter av klasse c , delt på antall ord som finnes i alle dokumenter av klasse c .

V her representerer hele vokabularet, altså mengden av **alle** ord (**types**) fra **alle** klasser.

- ▶ La oss anta at vi prøver å finne sannsynligheten for ordet “fantastisk” gitt klassen “positiv”.
- ▶ Ordet “fantastisk” blir ikke brukt i noen av dokumentene som er klassifisert som positive.

$$P(\textit{fantastisk}|\textit{positiv}) = \frac{\textit{count}(\textit{fantastisk}, \textit{positiv})}{\sum_{w \in V} \textit{count}(w, \textit{positive})} = 0$$

- ▶ Dette er problematisk! Siden **alle** sannsynligheter blir multiplisert med hverandre, å få 0 her vil lede til at sannsynligheten for **hele** klassen vil bli lik 0.

- ▶ Enkleste måte å løse dette på: bruke smoothing/glatting.
- ▶ Bruke **add-one** (Laplace) smoothing/glatting.

$$P(w_i|c) = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)} = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V} (\text{count}(w, c)) + |V|)}$$

- ▶ Hva om det finnes ord i testdataen som ikke finnes i vårt vokabular V ?
 - ▶ Ignorer dem, og slett dem fra testdataen, og ikke beregn sannsynligheter for dem.

	Kat	Dokument
Train	+	en varm og velopplagt bok
	+	dette er kanskje årets beste barnebok
	–	forfatteren burde ha fått bedre manushjelp av forlaget sitt
	–	oppskrytt dameroman
	–	denne boken er et hån mot underholdningssjangeren
Test	?	boken og forfatteren er oppskrytt

- ▶ Begynne med å regne ut **prior**-sannsynligheten for klassene $P(-)$ og $P(+)$.
- ▶ Husk at vi beregner $P(c)$ slik:

$$P(c) = \frac{N_c}{N_{doc}}$$

N_c = antall dokumenter i treningsdataen som er i klassen c

N_{doc} = total antall dokumenter

- ▶ $P(-)$ og $P(+)$ er da:

$$P(-) = \frac{3}{5}$$

$$P(+) = \frac{2}{5}$$

- ▶ Vi beregner sannsynlighetene for hvert ord i testsetningen, altså: “boken”, “og”, “forfatteren”, “er”, “oppskrytt”. Vi bruker smoothing/glatting.
- ▶ Vi må da beregne sannsynligheten for **individuelle trekk** $P(f_j|c)$ altså :

$$P(w_i|c) = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V} (\text{count}(w, c)) + |V|)}$$

$$P(\text{boken}|-) = \frac{1+1}{18+28}$$

$$P(\text{og}|-) = \frac{0+1}{18+28}$$

$$P(\text{forfatteren}|-) = \frac{1+1}{18+28}$$

$$P(\text{er}|-) = \frac{1+1}{18+28}$$

$$P(\text{oppskrytt}|-) = \frac{1+1}{18+28}$$

$$P(\text{boken}+) = \frac{0+1}{11+28}$$

$$P(\text{og}+) = \frac{1+1}{11+28}$$

$$P(\text{forfatteren}+) = \frac{0+1}{11+28}$$

$$P(\text{er}+) = \frac{1+1}{11+28}$$

$$P(\text{oppskrytt}+) = \frac{0+1}{11+28}$$

Eksempel – NB for sentiment klassifisering forts.

$$P(-)P(S|-) = \frac{3}{5} \times \frac{\overbrace{2 \times 1 \times 2 \times 2 \times 2}^{16}}{46^5} = 0,000000047 = 4.7 \times 10^{-8}$$

$$P(+)P(S|+) = \frac{2}{5} \times \frac{\overbrace{1 \times 2 \times 1 \times 2 \times 1}^4}{39^5} = 0,000000018 = 1.8 \times 10^{-8}$$

$$P(-)P(S|-) > P(+)P(S|+)$$

\implies Den blir derfor klassifisert som *negativ*.

- ▶ Forekomsten av ord er viktigere enn frekvensen: telle hvert ord i hvert dokument kun én gang.
- ▶ Håndtere negasjon:
 - ▶ I really like this movie. +
 - ▶ I really did't like this movie. -
 - ▶ Don't dismiss this film, doesn't really get us bored. +
 - ▶ Veldig enkel og naiv måte å løse dette på er å legge til prefikset NOT_ før hvert ord som blir negert:
 - ▶ didn't like this movie , but I liked the actors .
 - ▶ didn't NOT_like NOT_this NOT_movie , but I liked the actors.
- ▶ Ikke nok treningsdata?
 - ▶ Bruk et sentimentleksikon.
 - ▶ Legge til et trekk til Naive Bayes. Telle hver gang et ord er fra det positive eller negative leksikonet.

- ▶ Hvordan kan vi vite om vår Naive Bayes klassifiserer riktig?
- ▶ **Husk!** Vi har annotert treningsdata + annotert testsett.
- ▶ Testsettet brukes med Naive Bayes uten klassene.
- ▶ Klassene fra testsettet brukes som **gullstandard** (gold standard) for å sjekke om Naive Bayes har angitt riktig klasse.
- ▶ Men hvordan kan vi evaluere?

- ▶ Lage **contingency** tabell.
- ▶ Regne ut **accuracy**, **precision**, **recall**, og F_1 .

		<i>gold standard labels</i>		
		gold positive	gold negative	
<i>system output labels</i>	system positive	true positive	false positive	$\text{precision} = \frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		$\text{recall} = \frac{tp}{tp+fn}$		$\text{accuracy} = \frac{tp+tn}{tp+fp+tn+fn}$

$$F_1 = \frac{2PR}{P + R}$$

tp = true positive, fp = false positive, tn = true negative, fn = false negative,
 P = precision, R = Recall

Se filene:

- ▶ `doc_classification.py`
- ▶ `doc_classification_BOW.py`