

IN1140: Introduksjon til språkteknologi

Forelesning #10

Samia Touileb

Universitetet i Oslo

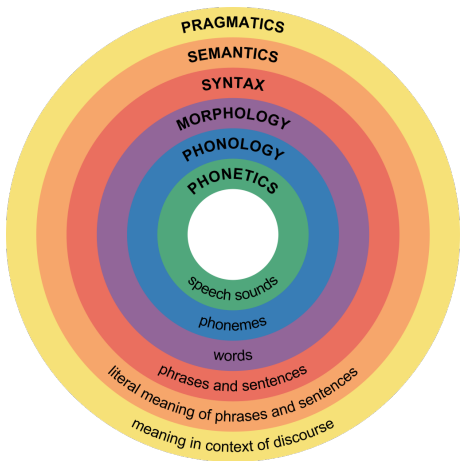
31. oktober 2019



Beregning av poeng og bestått/ikke-bestått for oblig'er

Vi kommer til å operere med en viss margin i beregning av poeng og bestått/ikke bestått per oblig. Som dere vet må dere akkumulere minst 100 poeng for oblig 1 (1a og 1b) for å bestå oppgaven og minst 50 poeng for hver av oblig 2 og oblig 3. Dersom noen mangler poeng på en oblig, vil dere ha mulighet til å hente inn disse på neste oblig. F.eks. dersom dere mangler 10 poeng på oblig 2 (har 40 poeng) må dere få 10 poeng ekstra (dvs minimum 60 poeng) på oblig 3.

November							
Week	Mon	Tue	Wed	Thu	Fri	Sat	Sun
44					1	2	3
45	4	5	6	7	<u>8</u>	9	10
46	11	12	13	14	15	16	17
47	18	19	20	21	22	23	24
48	25	26	27	28	29	30	





Forrige uke:

- ▶ Semantikk
- ▶ Studiet av betydning slik det uttrykkes gjennom språk
- ▶ Betydning til morfemer, ord, fraser og setninger
 - ▶ leksikal semantikk
 - ▶ setningssemantikk
 - ▶ (pragmatikk: hvordan konteksten påvirker betydning)



Tema for i dag:

- ▶ Hva slags oppgaver inngår i semantisk analyse?
- ▶ Med hva slags metoder kan disse oppgavene løses?



Flertydighet

- ▶ They rarely **serve** red meat, preferring to prepare seafood.
- ▶ He **served** as U.S. ambassador to Norway in 1976 and 1977.
- ▶ He might have **served** his time, come out and led an upstanding life.

Flertydighet

- ▶ Which of those flights **serve** breakfast?
- ▶ Does Air France **serve** Philadelphia?
- ▶ ?Does Air France **serve** breakfast and Philadelphia?

Ordbetydning ("word sense")

Flertydighet

- ▶ *The astronomer married the star* – **the star**
- ▶ *You are free to execute your laws, and your citizens, as you see fit* (Star Trek, Next Generation) – **execute**
- ▶ *Oh, flowers are common here, Miss Fairfax, as people are in London* (Oscar Wilde, The Importance of Being Earnest) – **common**



- ▶ *The dog bit the mailman* er ikke det samme som *The mailman bit the dog*
- ▶ Sammenligne med *The dog was bitten by the mailman*
- ▶ Semantiske roller beskriver "hvem som gjør hva mot hvem"
[The mailman]_{AGENT} bit [the dog]_{PATIENT}



- ▶ *Apple bought Cisco*
- ▶ *Apple acquired Cisco*
- ▶ *Cisco was taken over by Apple*



- ▶ Synonymer (samme betydning)
 - ▶ *ascend/rise*
 - ▶ *sweater/pullover*
- ▶ Antonymer (motsetninger)
 - ▶ *good/evil*
 - ▶ *ascend/descend*
- ▶ Hypernym/hyponym (mer generell, mer spesifikk)
 - ▶ *tree/birch*
 - ▶ *animate object / mammal / whale*
- ▶ Meronymi (del-helhet)
 - ▶ *knob/door*
 - ▶ *wheel/car*

- ▶ Lakoff & Johnson: "Metaphors We Live By". Noen av deres metaforer
 - ▶ ARGUMENT IS WAR
 - ▶ *He shot down all of my arguments*
- ▶ LOVE IS A JOURNEY
 - ▶ *They hit a bumpy stretch in their relationship*
- ▶ TIME IS MONEY
 - ▶ *You're wasting my time*
- ▶ hvorfor er metaforer problematiske for språkteknologi?



- ▶ *The chairman announced yesterday that they would have the problem solved within three days*
 - ▶ Når ble problemet løst dersom setningen forekom i en avis 10 mai, 2007?
- ▶ *Book me a flight on the 7:35 tomorrow*



- ▶ *Apple bought Cisco*
- ▶ *Apple acquired Cisco*
- ▶ *Cisco was taken over by Apple*

- ▶ Spørsmål: *Who bought Cisco?*
- ▶ Forventet svar: *Apple bought Cisco*
 - ▶ *Cisco's acquisition by Apple* → (entails) *Apple bought Cisco*



- ▶ *Apple acquired Cisco*
- ▶ *Apple did not acquire Cisco*
- ▶ *Apple failed to acquire Cisco*
- ▶ *Apple denied not acquiring Cisco*

- ▶ Automatisk negasjonsanalyse – internasjonal forskningskonkurranse
 - ▶ *SEM Shared Task on Negation Resolution
 - ▶ system som angir
 - ▶ negation **cue**
 - ▶ negation scope
 - ▶ negated *event*
 - ▶ There was **no** answer.

Metoder



- ▶ Klassifisering
 - ▶ Gitt et ord i en setning, og en liste av mulige betydninger, velg en betydning (den mest sannsynlige?).
 - ▶ Gitt et predikat i en setning, finn dets semantiske roller.
- ▶ Automatisk tilegnelse av semantisk informasjon fra rå tekst
 - ▶ Hvilke ord og fraser betyr det samme.
 - ▶ Distribusjonell semantikk:
<http://vectors.nlp1.eu/explore/embeddings/en/similar/>
<https://research.google.com/semantris/>

⇒ **Maskinlæring**

Semantiske ressurser



- ▶ Klassifisering forutsetter treningsdata.
- ▶ Leksikalske databaser (WordNet, FrameNet).
- ▶ Korpuser annotert med semantisk informasjon (PropBank).



- ▶ Manuelt konstruert database
- ▶ Betydningen til ord karakteriseres gjennom **relasjoner** til andre ord
- ▶ Semantiske konsepter karakteriseres gjennom relasjoner til andre konsepter
- ▶ Hva slags relasjoner kan det være snakk om?



- ▶ Mellom ord:
 - ▶ Synonymi (samme betydning).
 - ▶ Synonymi-relasjonen grupperer ord i synonymmengder, såkalte **synsets**.
- ▶ Mellom konsepter (=synsets)
 - ▶ Hypernyymi (mer generell, mer spesifikk).
 - ▶ Varierer noe, men antonymi og meronymi er også spesifisert for noen synsets.



- ▶ Elektronisk leksikon
 - ▶ Online grensesnitt
 - ▶ Lastes ned
 - ▶ Tilgjengelig på <http://wordnet.princeton.edu/>
 - ▶ Også tilgjengelig via NLTK

- ▶ Består av tre separate databaser:
 1. Substantiv (117798 lemmaer)
 2. Verb (11529 lemmaer)
 3. Adjektiv og adverb (22479 adjektiver, 4481 adverb)

- ▶ verbet *skim*: synonymer, definisjoner og eksempler

Verb

- **S:** (v) [plane](#), **skim** (travel on the surface of water)
- **S:** (v) [skim over](#), **skim** (move or pass swiftly and lightly over the surface of)
- **S:** (v) [scan](#), **skim**, [rake](#), [glance over](#), [run down](#) (examine hastily) *"She scanned the newspaper headlines while waiting for the taxi"*
- **S:** (v) **skim**, [skip](#), [skitter](#) (cause to skip over a surface) *"Skip a stone across the pond"*
- **S:** (v) **skim** (coat (a liquid) with a layer)
- **S:** (v) **skim**, [skim off](#), [cream off](#), [cream](#) (remove from the surface) *"skim cream from the surface of milk"*
- **S:** (v) **skim**, [skim over](#) (read superficially)



- ▶ Synsets er koblet sammen ved
 - ▶ Hyponym/hypernym relasjonen (hovedhierarkiet)
 - ▶ Meronymi: del-helhet relasjoner
 - ▶ Komponent/del (*leg – table, finger – hand*)
 - ▶ Medlem av en gruppe *tree – forest, student – class*
 - ▶ Materiale et objekt er laget av (*oxygen – water*)
- ▶ Ord er koblet sammen ved antonymi

- **S: (n) cat, true cat** (feline mammal usually having thick soft fur and no ability to roar: domestic cats; wildcats)
 - [direct hyponym](#) / [full hyponym](#)
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - **S: (n) feline, felid** (any of various lithe-bodied roundheaded fissiped mammals, many with retractile claws)
 - **S: (n) carnivore** (a terrestrial or aquatic flesh-eating mammal)
"terrestrial carnivores have four or five clawed digits on each limb"
 - **S: (n) placental, placental mammal, eutherian, eutherian mammal** (mammals having a placenta; all mammals except monotremes and marsupials)
 - **S: (n) mammal, mammalian** (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
 - **S: (n) vertebrate, craniate** (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
 - **S: (n) chordate** (any animal of the phylum Chordata having a notochord or spinal column)
 - **S: (n) animal, animate being, beast, brute, creature, fauna** (a living organism characterized by voluntary movement)

- **S: (n)** [organism](#), [being](#) (a living thing that has (or can develop) the ability to act or function independently)
 - **S: (n)** [living thing](#), [animate thing](#) (a living (or once living) entity)
 - **S: (n)** [whole](#), [unit](#) (an assemblage of parts that is regarded as a single entity) *"how big is that part compared to the whole?"; "the team is a unit"*
 - **S: (n)** [object](#), [physical object](#) (a tangible and visible entity; an entity that can cast a shadow) *"it was full of rackets, balls and other objects"*



- S: (n)
physical
entity (an entity that has physical existence)
 - S: (n)
entity (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

WordNet: substantiv eksempler



- **S, N1 cat, feline cat** (feline mammal usually having thick soft fur and no ability to roar; domestic cats, wildcats)
- **sheep** **sheep** / **fat** **sheep**
- **giant** **sheep** / **horned** **sheep** / **white** **sheep**
 - **S, N1 catfish** (a bony fish of various fish-like bodies roundheaded flattened moustached, many with electric eels)
 - **S, N1 cetacean** (a terrestrial or aquatic flesh-eating mammal)
 - **terrestrial cetaceans** (two four- or five-toed digits on each limb)
 - **S, N1 placental** **placental mammal** **subclass** **subclass** **mammal** (mammals having a placenta; all mammals except monotremes and marsupials)
 - **S, N1 cetacean** **monotremes** (any warm-blooded vertebrate having the skin nose or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
 - **S, N1 cetacean** **cetacean** (mammals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
 - **S, N1 chondate** (any animal of the phylum Chordata having a notochord or spinal column)
 - **S, N1 animal** **cetacean** **being** **being** **entity** **entity** **entity** (a being expansion characterized by resistatory movement)
 - **S, N1 organism** **being** (a being thing that has (or can develop) the ability to act or function independently)
 - **S, N1 being** **being** **entity** **entity** (a being or once being) **entity**)
 - **S, N1 atomic** **part** (an assemblage of parts that is regarded as a single entity) "The **big** is that part compared to the whole?" "The **beam** is a **unit**"
 - **S, N1 object** **physical object** (a tangible and mobile entity; an entity that can cast a shadow) "It was full of **spacules**, **clubs** and other objects"
 - **S, N1 physical entity** (an entity that has physical existence)
 - **S, N1 entity** (that which is perceived or known or inferred to have its own distinct existence (being or naming))



Noun

- **S: (n) mammal**, [mammalian](#) (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
 - [direct hyponym](#) / [full hyponym](#)
 - **S: (n) female mammal** (animals that nourish their young with milk)
 - **S: (n) tusker** (any mammal with prominent tusks (especially an elephant or wild boar))
 - **S: (n) prototherian** (primitive oviparous mammals found only in Australia and Tasmania and New Guinea)
 - **S: (n) metatherian** (primitive pouched mammals found mainly in Australia and the Americas)
 - **S: (n) placental**, [placental mammal](#), [eutherian](#), [eutherian mammal](#) (mammals having a placenta; all mammals except monotremes and marsupials)
 - **S: (n) fossorial mammal** (a burrowing mammal having limbs adapted for digging)
 - [part meronym](#)
 - **S: (n) coat**, [pelage](#) (growth of hair or wool or fur covering the body of an animal)
 - **S: (n) hair**, [pilus](#) (any of the cylindrical filaments characteristically growing from the epidermis of a mammal) "*there is a hair in my soup*"
 - [member holonym](#)
 - **S: (n) Mammalia**, [class Mammalia](#) (warm-blooded vertebrates characterized by mammary glands in the female)



- ▶ Semantikk for hendelser ('events', verb) er ganske forskjellig fra semantikk for entiteter (substantiver)
- ▶ Relasjoner
 - ▶ Troponymi:
 - ▶ *A verb expressing a specific manner of elaboration of another verb. X is a troponym of Y if to X is to Y in some manner*
 - ▶ *wade – walk*
 - ▶ Entailment
 - ▶ *A verb X entails Y if X cannot be done unless Y is, or has been, done*
 - ▶ *walking – stepping*

- **S:** (v) **jump**, **leap**, **bound**, **spring** (move forward by leaps and bounds) *"The horse bounded across the meadow"; "The child leapt across the puddle"; "Can you jump over the fence?"*
 - **direct troponym** / **full troponym**
 - **S:** (v) **pronk** (jump straight up) *"kangaroos pronk"*
 - **S:** (v) **bounce**, **resile**, **take a hop**, **spring**, **bound**, **rebound**, **recoil**, **reverberate**, **ricochet** (spring back; spring away from an impact) *"The rubber ball bounced"; "These particles do not resile but they unite after they collide"*
 - **S:** (v) **burst** (move suddenly, energetically, or violently) *"He burst out of the house into the cool night"*
 - **S:** (v) **bounce** (leap suddenly) *"He bounced to his feet"*
 - **S:** (v) **capriole** (perform a capriole, of horses in dressage)
 - **S:** (v) **galumph** (move around heavily and clumsily) *"the giant tortoises galumphed around in their pen"*
 - **S:** (v) **ski jump** (jump on skis)
 - **S:** (v) **saltate** (leap or skip, often in dancing) *"These fish swim with a saltating motion"*
 - **S:** (v) **vault** (bound vigorously)
 - **S:** (v) **leapfrog** (jump across) *"He leapfrogged his classmates"*
 - **S:** (v) **vault**, **overleap** (jump across or leap over (an obstacle))
 - **S:** (v) **curvet** (perform a leap where both hind legs come off the ground, of a horse)
 - **S:** (v) **hop**, **skip**, **hop-skip** (jump lightly)
 - **S:** (v) **caper** (jump about playfully)
 - **S:** (v) **hop** (make a jump forward or upward)
 - **direct hypernym** / **inherited hypernym** / **sister term**
 - **S:** (v) **move** (move so as to change position, perform a nontranslational motion) *"He moved his hand slightly to the right"*

Verb

- **S: (v) walk** (use one's feet to advance; advance by steps) "Walk, don't run!"; "We walked instead of driving"; "She walks with a slight limp"; "The patient cannot walk yet"; "Walk over to the cabinet"
 - [direct troponym](#) / [full troponym](#)
 - [verb group](#)
 - [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - **entailment**
 - **S: (v) step** (shift or move by taking a step) "step back"
 - [phrasal verb](#)
 - [antonym](#)
 - [derivationally related form](#)
 - [sentence frame](#)



- ▶ **Deskriptive** adjektiver
 - ▶ Organisert i grupperinger basert på likhet ("similarity"), relatert via antonymi
- ▶ **Relasjonelle** adjektiver, linket til sine substantiver
 - ▶ *nasal – nose*
- ▶ Liten gruppe adjektiver "reference-modifying"
 - ▶ *former, alleged*

Adjective

- **S: (adj) cheap**, **inexpensive** (relatively low in price or charging low prices) *"it would have been cheap at twice the price"; "inexpensive family restaurants"*
 - **similar to**
 - **S: (adj) bargain-priced**, **cut-rate**, **cut-price** (costing less than standard price) *"buying bargain-priced clothes for the children"; "cut-rate goods"*
 - **S: (adj) catchpenny** (designed to sell quickly without concern for quality) *"catchpenny ornaments"*
 - **S: (adj) dirt cheap** (very cheap) *"a dirt cheap property"*
 - **S: (adj) low-budget** (made on or suited to a limited budget) *"a low-budget movie"; "a low-budget menu"*
 - **S: (adj) low-cost**, **low-priced**, **affordable** (that you have the financial means for) *"low-cost housing"*
 - **S: (adj) nickel-and-dime** (low-paying) *"a nickel-and-dime job"*
 - **S: (adj) sixpenny**, **threepenny**, **twopenny**, **tuppenny**, **two-a-penny**, **twopenny-halfpenny** (of trifling worth)
 - **antonym**
 - **W: (adj) expensive** [Opposed to: **cheap**] (high in price or charging high prices) *"expensive clothes"; "an expensive shop"*



- ▶ Hvor mange betydninger har et ord?
 - ▶ Antall synsets ordet forekommer i
- ▶ Nærhet i betydning kan utledes fra nærhet i hierarkiet
 - ▶ Korteste stien via hyponym/hypernym-linkene mellom synsets



- ▶ Utgangspunkt for Word Sense Disambiguation
 - ▶ Merke forekomster av et ord med riktig betydning (=synset)
- ▶ Generaliserer over ord via hypernymi-relasjonen
 - ▶ Fra *cat* til *living being*
- ▶ Generalisere over synonymer
- ▶ ...

- ▶ Aspekt ved setningsbetydning: hvilke roller de forskjellige deltagerene inntar
 - ▶ *Nina* hevet *bilen* med *jekken*
 - ▶ *Nina* – deltageren som er ansvarlig for å utføre handlingen beskrevet av verbet
 - ▶ *bilen* – blir påvirket av handlingen
 - ▶ *jekken* – middelet som Gina bruker til å utføre handlingen
- ▶ Semantiske roller beskriver den semantiske relasjonen som argumenter har til handlingen beskrevet av verbet



- eksempelet: *Nina* *hevet* *bilen* *med* *jekken*
AGENT THEME INSTRUMENT



- ▶ Ikke full enighet rundt rolleinventaret
- ▶ Vanskelig å formulere formelle definisjoner av roller
- ▶ \Rightarrow generaliserte semantiske roller
 - ▶ PROTO-AGENT, PROTO-PATIENT
- ▶ Verbspesifikke roller
- ▶ Semantiske ressurser med informasjon om semantiske roller: **PropBank** og FrameNet

- ▶ Inneholder alle setningene i Penn Treebank
- ▶ Roller er (stort sett) verbspesifikke
 - ▶ Arg0, Arg1 = PROTO-AGENT, PROTO-PATIENT
 - ▶ Arg2 ... verbspesifikke

agree.01

Arg0 Agreeer

Arg1 Proposition

Arg2 Other entity agreeing

Ex1 [*Arg0 The group*] *agreed* [*Arg1 it wouldn't make an offer*]Ex2 [*ArgM-TMP Usually*] [*Arg0 John*] *agrees* [*Arg2 with Mary*] [*Arg1 on everything*]

- ▶ Applikasjon: Semantic Role Labeling
- ▶ Gitt et predikat i en setning, finn dets semantiske roller
- ▶ Gir oss en felles representasjon for:
 - ▶ [*Arg0*Big Fruit Co.] increased [*Arg1*the price of bananas]
 - ▶ [*Arg1*The price of bananas] was increased again by [*Arg0*Big Fruit Co.]
 - ▶ [*Arg1*The price of bananas] increased [*Arg2*5%]

“Big Fruit Co.” er alltid *AGENT* og “the price of bananas” er alltid *PATIENT*

- ▶ Roller er “frame-spesifikke”
- ▶ Frame air travel = *reservation, flight, travel, buy, price, cost, fare, rates, meal, plane*

Ex1 $[_{Arg1} \textit{The price of bananas}] \textit{increased} [_{Arg2} 5\%]$

Ex2 $[_{Arg1} \textit{The price of bananas}] \textit{rose} [_{Arg2} 5\%]$

Ex3 $\textit{There has been a} [_{Arg2} 5\%] \textit{rise} [_{Arg1} \textit{in the price of bananas}]$

Semantisk klassifisering (WSD)

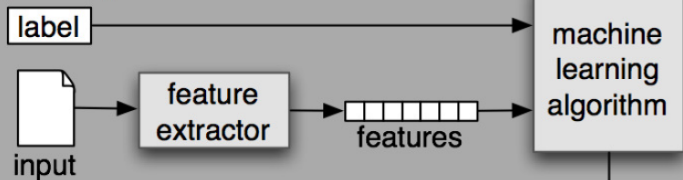


- ▶ Word Sense Disambiguation (WSD) – aktivt felt innenfor språkteknologi
 - ▶ gitt en setning med et spesifikt målord ("target word") og en liste med betydninger (f.eks. fra WordNet)
 - ▶ angi korrekt betydning for målordet i den setningen
- ▶ Klassifisering basert på et annotert datasett

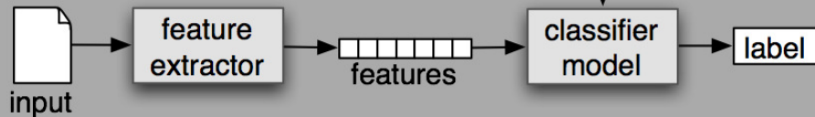


- ▶ Sentral metode innenfor maskinlæring
- ▶ Automatisk avgjøre hvilken kategori en observasjon tilhører
- ▶ Basert på **treningsdata**: observasjoner der kategorien er kjent
 - ▶ e-post \rightarrow {spam, ikke-spam}
 - ▶ pasient \rightarrow diagnose
- ▶ **Supervised** klassifisering: klassifisering som benytter treningsdata

(a) Training



(b) Prediction





- ▶ Første skritt består i å hente ut trekk (“features”) fra treningsdataene
- ▶ Eksempel: setninger merket med betydning
 - ▶ **SKIM** the pages for a clearer insight: **Reading**
 - ▶ She **SKIMS** through the novel which seems to fascinate them: **Reading**
 - ▶ Remove the vanilla pod, **SKIM** the jam, and let it cool: **Removing**
 - ▶ We **SKIMMED** across the surface of that sodding lake whilst all around us gathered the dark hosts of hell: **Self_Motion**
- ▶ Hvilke trekk (“features”) kan vi bruke for å skille mellom de forskjellige betydningene?



- ▶ **SKIM** the pages for a clearer insight: [Reading](#)
- ▶ She **SKIMS** through the novel which seems to fascinate them: [Reading](#)
- ▶ Remove the vanilla pod, **SKIM** the jam, and let it cool: [Removing](#)
- ▶ We **SKIMMED** across the surface of that sodding lake whilst all around us gathered the dark hosts of hell: [Self_Motion](#)

Henter ut alle ord (**ikke** ordnet):

- ▶ a, clearer, for, insight, pages, the: [Reading](#)
- ▶ fascinate, novel, seems, she, the, them, through, to, which: [Reading](#)
- ▶ and, cool, it, jam, let, pod, remove, the, the, vanilla: [Removing](#)
- ▶ across, all, around, dark, gathered, hell, hosts, lake, of, of, sodding, surface, that, the, the, us, we, whilst [Self_Motion](#)



- ▶ Konteksten til målordet kan representeres ved
 - ▶ ordformer
 - ▶ lemmaer
 - ▶ ordklassetagger
 - ▶ kombinasjon av disse
 - ▶ $[w_{i-2}, \text{POS}_{i-2}, w_{i-1}, \text{POS}_{i-1}, w_{i+1}, \text{POS}_{i+1}, w_{i+2}, \text{POS}_{i+2}]$
- ▶ Remove the vanilla pod, **SKIM** the jam, and let it cool: [Removing](#)
 - ▶ **trekkvektor**: [vanilla, JJ, pod, NN, the, DT, jam, NN]



- ▶ **SKIM** the pages for a clearer insight: [Reading](#)
- ▶ She **SKIMS** through the novel which seems to fascinate them: [Reading](#)
- ▶ Remove the vanilla pod, **SKIM** the jam, and let it cool: [Removing](#)
- ▶ We **SKIMMED** across the surface of that sodding lake whilst all around us gathered the dark hosts of hell: [Self_Motion](#)

- ▶ Verbets argumenter:
 - ▶ direkte_objekt: [Reading](#)
 - ▶ subjekt, pp_through: [Reading](#)
 - ▶ direkte_objekt: [Removing](#)
 - ▶ subjekt, pp_across: [Self_Motion](#)

- ▶ Kombinasjoner (argumenters hovedord)
 - ▶ direkte_objekt / pages: [Reading](#)
 - ▶ subjekt / she, pp_through / novel: [Reading](#)
 - ▶ direkte_objekt / jam: [Removing](#)
 - ▶ subjekt / we, pp_across / surface: [Self_Motion](#)



- ▶ N-gram av ord i nærheten av målordet
 - ▶ $n=1,2,3$
 - ▶ kan også bruke ordformer, lemmaer, ordklasser
- ▶ Eksempel:
 - ▶ **SKIM** the pages for a clearer insight: [Reading](#)
 - ▶ She **SKIMS** through the novel which seems to fascinate them: [Reading](#)
 - ▶ Remove the vanilla pod, **SKIM** the jam, and let it cool: [Removing](#)
 - ▶ We **SKIMMED** across the surface of that sodding lake whilst all around us gathered the dark hosts of hell: [Self_Motion](#)
- ▶ Trigram:
 - ▶ `_`, `_`, `_`, the, pages, for
 - ▶ `_`, `_`, She, through, the, novel
 - ▶ the, vanilla, pod, the, jam, and
 - ▶ `_`, `_`, We, across, the, surface



- ▶ Gitt treningsdataene og trekkvektorene, kan en rekke forskjellige maskinlæringsalgoritmer brukes til å trene en klassifiserer
- ▶ Her skal vi se på **Naive Bayes**-klassifisering
- ▶ Bruker informasjon om ord i konteksten for disambiguering av betydning
- ▶ Enkel metode, mye brukt i WSD



- ▶ Naive Bayes klassifiserer

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

- ▶ 2 sannsynligheter:

1. prior-sannsynligheten for betydningen $P(s)$

$$P(s_i) = \frac{\text{count}(s_i, w_j)}{\text{count}(w_j)}$$

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

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

Trekkvektorer (bag-of-words), der siste element angir betydningsklassen

- ▶ [a, clearer, for, insight, pages, the, Reading]
- ▶ [fascinate, novel, seems, she, the, them, through, to, which, Reading]
- ▶ [and, cool, it, jam, let, pod, remove, the, the, vanilla, Removing]
- ▶ [across, all, around, dark, gathered, hell, hosts, lake, of, of, sodding, surface, that, the, the, us, we, whilst, Self_Motion]

1. prior-sannsynligheten for betydningen $P(\textit{Reading})$

- ▶ $P(s_i) = \frac{\textit{count}(s_i, w_j)}{\textit{count}(w_j)}$
- ▶ $P(\textit{Reading}) = \frac{\textit{count}(\textit{Reading}, \textit{skim})}{\textit{count}(\textit{skim})}$
- ▶ $P(\textit{Reading}) = \frac{2}{4} = 0.5$

- ▶ [a, clearer, for, insight, **pages**, the, **Reading**]
- ▶ [fascinate, novel, seems, she, the, them, through, to, which, **Reading**]
- ▶ [and, cool, it, **jam**, let, pod, remove, the, the, vanilla, **Removing**]
- ▶ [across, all, around, dark, gathered, hell, hosts, lake, of, of, sodding, surface, that, the, the, us, we, whilst, **Self_Motion**]

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

- ▶ $P(f_j|s) = \frac{\text{count}(f_j,s)}{\text{count}(s)}$
- ▶ $P(\text{pages}|\text{Reading}) = \frac{1}{2} = 0.5$
- ▶ $P(\text{pages}|\text{Removing}) = \frac{0}{1} = 0$
- ▶ $P(\text{pages}|\text{Self_Motion}) = \frac{0}{1} = 0$
- ▶ ...
- ▶ $P(\text{jam}|\text{Reading}) = \frac{0}{2} = 0$
- ▶ $P(\text{jam}|\text{Removing}) = \frac{1}{1} = 1.0$
- ▶ $P(\text{jam}|\text{Self_Motion}) = \frac{0}{1} = 0$
- ▶ ...

- ▶ Vi har nå en NB-modell som vi kan bruke til å klassifisere en ny og usett setning:
 - ▶ I like to SKIM through the novel
 - ▶ [I, like, novel, the, through, to, ??]
- ▶ Vi bruker den velkjente formelen:

$$P(s) \prod_{j=1}^n P(f_j | s)$$

- ▶ Itererer gjennom alle betydningene og trekkene, ganger sammen
- ▶ Og velger den betydningen (s) som gir høyest sannsynlighet

$$P(\textit{Reading}) = \frac{2}{4}$$

$$P(I|\textit{Reading}) = \frac{0}{2}$$

$$P(\textit{like}|\textit{Reading}) = \frac{0}{2}$$

$$P(\textit{novel}|\textit{Reading}) = \frac{1}{2}$$

$$P(\textit{the}|\textit{Reading}) = \frac{2}{2}$$

$$P(\textit{through}|\textit{Reading}) = \frac{0}{2}$$

$$P(\textit{to}|\textit{Reading}) = \frac{1}{2}$$

$$P(\textit{Removing}) = \frac{1}{4}$$

$$P(I|\textit{Removing}) = \frac{0}{1}$$

$$P(\textit{like}|\textit{Removing}) = \frac{0}{1}$$

$$P(\textit{novel}|\textit{Removing}) = \frac{0}{1}$$

$$P(\textit{the}|\textit{Removing}) = \frac{2}{1}$$

$$P(\textit{through}|\textit{Removing}) = \frac{0}{1}$$

$$P(\textit{to}|\textit{Removing}) = \frac{0}{1}$$

$$P(\text{Self_motion}) = \frac{1}{4}$$

$$P(I|\text{Self_motion}) = \frac{0}{1}$$

$$P(\text{like}|\text{Self_motion}) = \frac{0}{1}$$

$$P(\text{novel}|\text{Self_motion}) = \frac{0}{1}$$

$$P(\text{the}|\text{Self_motion}) = \frac{2}{1}$$

$$P(\text{through}|\text{Self_motion}) = \frac{0}{1}$$

$$P(\text{to}|\text{Self_motion}) = \frac{0}{1}$$

- ▶ Som dere kan se både her og de to forrige sett av beregninger (forrige slide), vi kommer til å få 0 overalt pga sannsynlighetene som har verdien 0 (siden vi kommer til å multiplisere dem sammen ved å bruke Naive Bayes formelen).
- ▶ Vi må derfor bruke glatting (add-one smoothing – se neste slide).



$$P(\textit{Reading}) = \frac{2}{4}$$

$$P(\textit{I}|\textit{Reading}) = \frac{0+1}{2+3} = \frac{1}{5}$$

$$P(\textit{like}|\textit{Reading}) = \frac{0+1}{2+3} = \frac{1}{5}$$

$$P(\textit{novel}|\textit{Reading}) = \frac{1+1}{2+3} = \frac{2}{5}$$

$$P(\textit{the}|\textit{Reading}) = \frac{2+1}{2+3} = \frac{3}{5}$$

$$P(\textit{through}|\textit{Reading}) = \frac{0+1}{2+3} = \frac{1}{5}$$

$$P(\textit{to}|\textit{Reading}) = \frac{1+1}{2+3} = \frac{2}{5}$$



$$P(\textit{Reading}) \times \frac{\overbrace{1 \times 1 \times 2 \times 3 \times 1 \times 2}^{12}}{5^6} = 0,000384 = 3.84 \times 10^{-4}$$



$$P(\textit{Removing}) = \frac{1}{4}$$

$$P(\textit{I}|\textit{Removing}) = \frac{0 + 1}{1 + 3} = \frac{1}{4}$$

$$P(\textit{like}|\textit{Removing}) = \frac{0 + 1}{1 + 3} = \frac{1}{4}$$

$$P(\textit{novel}|\textit{Removing}) = \frac{0 + 1}{1 + 3} = \frac{1}{4}$$

$$P(\textit{the}|\textit{Removing}) = \frac{2 + 1}{1 + 3} = \frac{3}{4}$$

$$P(\textit{through}|\textit{Removing}) = \frac{0 + 1}{1 + 3} = \frac{1}{4}$$

$$P(\textit{to}|\textit{Removing}) = \frac{0 + 1}{1 + 3} = \frac{1}{4}$$



$$P(\textit{Removing}) \times \frac{\overbrace{1 \times 1 \times 1 \times 3 \times 1 \times 1}^3}{4^6} = 0,000183105 = 1.83 \times 10^{-4}$$



$$P(\textit{Self_motion}) = \frac{1}{4}$$

$$P(\textit{I}|\textit{Self_motion}) = \frac{0+1}{1+3} = \frac{1}{4}$$

$$P(\textit{like}|\textit{Self_motion}) = \frac{0+1}{1+3} = \frac{1}{4}$$

$$P(\textit{novel}|\textit{Self_motion}) = \frac{0+1}{1+3} = \frac{1}{4}$$

$$P(\textit{the}|\textit{Self_motion}) = \frac{2+1}{1+3} = \frac{3}{4}$$

$$P(\textit{through}|\textit{Self_motion}) = \frac{0+1}{1+3} = \frac{1}{4}$$

$$P(\textit{to}|\textit{Self_motion}) = \frac{0+1}{1+3} = \frac{1}{4}$$



$$P(\textit{Self_motion}) \times \frac{\overbrace{1 \times 1 \times 1 \times 3 \times 1 \times 1}^3}{4^6} = 0,000183105 = 1.83 \times 10^{-4}$$



$$P(\textit{Reading}) \times \frac{\overbrace{1 \times 1 \times 2 \times 3 \times 1 \times 2}^{12}}{5^6} = 0,000384 = 3.84 \times 10^{-4}$$

$$P(\textit{Removing}) \times \frac{\overbrace{1 \times 1 \times 1 \times 3 \times 1 \times 1}^3}{4^6} = 0,000183105 = 1.83 \times 10^{-4}$$

$$P(\textit{Self_motion}) \times \frac{\overbrace{1 \times 1 \times 1 \times 3 \times 1 \times 1}^3}{4^6} = 0,000183105 = 1.83 \times 10^{-4}$$

⇒ Den betydningen som gir høyest sannsynlighet for testsetningen vår [I, like, novel, the, through, to] er **Reading**



Merk at vi her har beregnet sannsynlighetene på en litt annen måte enn det vi gjorde under forelesning 8 om maskinlæring og klassifisering.

- ▶ Her har vi brukt:

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

Hvor vi teller antall ganger et ord f_j fra trekkvektoren forekommer med betydning s , delt på antall setninger klassifisert som betydning s .

- ▶ Når vi brukte Naive bayes for klassifisering brukte vi:

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

Hvor vi teller antall ganger et ord w_i forekommer i klassen c delt på antall ord typer som finnes i klassen c



Merk at det samme gjelder glatting!

- ▶ For Word Sense Disambiguation legger vi + 1 i telleren som vanlig. Men i nevneren legger vi til antall senses (betydninger) som finnes i treningsdataen vår. I eksemplet på slide 60, ble $P(I|Reading) = \frac{0+1}{2+3}$ fordi det finnes 3 betydninger i treningsdataen vår (altså *Reading*, *Removing*, og *Self_motion*).
- ▶ For klassifisering av klasser legger vi +1 i telleren, og + |V| i nevneren som representerer alle ord typer i hele treningskorpuset:

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



Treningsdata:

Reading	clearer insight pages
Reading	fascinate novel seems
Reading	like novel book pages
Reading	book action suspense
Reading	book pages thriller killer
Reading	book novel fascinate like
Removing	jam bread like breakfast
Removing	cool jam let pod remove vanilla
Removing	milk surface cream
Self_Motion	lake surface close novel surfboard
Self_Motion	dark gathered hell hosts lake sodding surface



Test data:

? like novel

Compute Naive Bayes
#####

$$P(\text{Reading})P(S|\text{Reading}) = \frac{6}{11} \times \frac{3 \times 4}{9^2} = 0,080 = 8.00 \times 10^{-2}$$

$$P(\text{Removing})P(S|\text{Removing}) = \frac{3}{11} \times \frac{2 \times 1}{6^2} = 0,0151 = 1.51 \times 10^{-2}$$

$$P(\text{Self}_m\text{otion})P(S|\text{Self}_m\text{otion}) = \frac{2}{11} \times \frac{1 \times 2}{5^2} = 0,0145 = 1.45 \times 10^{-2}$$



- ▶ Rekke oppgaver inngår i semantisk analyse
 - ▶ ordbetydningsdisambiguering (WSD)
 - ▶ semantiske roller
 - ▶ parafrasering
 - ▶ temporal analyse
 - ▶ entailment
 - ▶ negasjon
 - ▶ ...
- ▶ Sentral metode: [klassifisering](#)



- ▶ For klassifisering trenger vi treningsdata
- ▶ Semantiske ressurser
 - ▶ WordNet
 - ▶ leksikal database
 - ▶ innholdsord: substantiver, verb, adjektiver
 - ▶ bygget rundt leksikale relasjoner som synonymi, hyponymi, meronymi, etc.
 - ▶ PropBank/FrameNet
 - ▶ forskjellige ressurser for semantiske roller
 - ▶ korpus vs database
 - ▶ verbspesifikk vs ramme ("frame")
 - ▶ ...



- ▶ Nærmere kikk på betydningsdisambiguering
- ▶ Trekkrepresentasjon av treningsdata
 - ▶ ord
 - ▶ lemma
 - ▶ ordklasse
 - ▶ syntaktisk funksjon
 - ▶ etc.
- ▶ Naive Bayes-klassifisering
 - ▶ hvordan vi kan beregne den mest sannsynlige betydningen for et ord:
$$\hat{s} = \operatorname{argmax}_{s \in S} P(s | \vec{f})$$