

11-830 Computational Ethics for NLP

Intellectual Property, Copyright
And Plagiarism



Carnegie Mellon University

Language Technologies Institute

Intellectual Property

- Three areas of IP
 - Patents
 - Copyright
 - Trade Marks
- Plagiarism



Copyright

- Right to copy things
 - Given by the government in exchange for protection
 - But today is about not copying things
- When is a copy a copy?
 - Some copying is fine
 - Copying from your disk to play/display your bought copy
 - Copying small parts for review purposes
 - Parody
 - Fair Use
- Is Intent an issue?

Copyright and Intent

- I record a TV show from broadcast TV (legally)
- I upload my mkv file to dropbox
- Someone else does the same
- Dropbox notes the files are exactly the same
 - Deletes the other file and links to mine
- Other person distributes the file on a file sharing network
 - My file is being distributed illegal
 - Am I liable for copyright infringement

Copyright and Unknown Infringement

- “Men at Work” singing “Land Down Under”
- Contains a snippet of traditional Kookaburra song
 - Actually half of the song, and is still under copyright
- Only explicitly noticed when a quiz show asks a question
 - (20 years after the song was released)
- Found to be a copyright violation
 - We all were copyright pirates

Copyright and Maliciousness

- But people are violating copyright deliberately for profit
- Marvell vs CMU on noise reduction in hard disk access
 - Marvell asked about licensing algorithm (patent)
 - Thought it was too expensive
 - Implemented it anyway (naming it for the author)
 - Authors got \$750m
- Distributing music and video without permission for profit
- But it can also be hard to know if you are violating copyright
 - Don't use anything you find on the web

Detecting Copyright Violations

- Music and Video hash systems to find “same” files
 - Can it be robust to encoding technique and resolution
 - Music ID systems have been doing this for some time
 - Often works but can make errors too
 - Accidental (it sounds like it) or “Something” (the NASA case)
- What about text and NLP?



Finding Text Copyright Violations

- Specifically called Plagiarism
 - Student Homeworks
 - Scientific Papers
- Its a copy with the title/author changed
 - But you need a big database of papers (hence TurnItIn)
- But it sometimes is just a sentence, or paragraph, or diagram
 - But with proper attribution that's a **good** thing
 - Sometimes people generate the same text
 - Can you find the overlap in text with others
- N-gram overlap (sentence)
 - But need thresholds for amount of copying.
 - What if there are minor changes, synonyms, etc

EU Copyright Law “Article 13”

- Websites can no longer claim “just a carrier”
- Websites must delete identified copyright violations “quickly”
- Sounds good at first reading
- Consequences:
 - Small websites can afford to do this so only big websites can continue
 - Small websites cannot have user uploaded content (so no commenting ever)
 - Could there be centralized filtering
 - Could there be open source filtering
 - Could there be independent filterings
- Law passed (but as Article 17 not 13)

When is copying not copying

Harry Potter and the Philosopher's Stone.

Star Wars: A New Hope; synopsis

Harry Potter

~~Luke Skywalker~~ is an orphan living with his uncle and aunt on the remote wilderness of ~~Tatooine~~.

He is rescued from ~~aliens~~ by wise, bearded ~~Ben Kenobi~~, who turns out to be a ~~Jedi Knight~~.

~~Ben~~ reveals to ~~Luke~~ that ~~Luke's~~ father was also a ~~Jedi Knight~~, and was the best pilot he had ever seen.

~~Luke~~ is also instructed in how to use ~~the Jedi light sabre~~ as he too trains to become a ~~Jedi~~.

~~Luke~~ has many adventures in ~~the galaxy~~ and makes new friends such as ~~Han Solo~~ and ~~Princess Leia~~.

In the course of these adventures he distinguishes himself as a top ~~X-wing pilot~~ in the battle of ~~the Death Star~~, making the ~~direct hit~~ that secures ~~the Rebels~~ victory against the forces of evil, ~~Slytherin~~.

~~Luke~~ also sees off the threat of ~~Darth Vader~~, who we know murdered his ~~uncle and aunt~~.

In the finale, ~~Luke~~ and his new friends receive medals of valour.

All of this will be set to an orchestral score composed by John Williams.



Plagiarism Detection

- Plagiarism detection is a big field
 - Researchers, conferences, large software installations, businesses
 - CMU's Lightside text machine learning systems build for plagiarism detection
 - Bought by TurnItIn
- Plagiarism goes beyond the surface form
 - Its not just overlapping n-grams (or phrases/sentences)
 - Need to detect structure too
 - Were these two examples written independently
 - Even though they contain the same algorithm
- We now care about both surface form and deep form

Plagiarism Avoidance

- Automated Detection spawns Automated Avoidance
 - (some form of legal justification)
- Adversarial Classifiers
 - Find methods to avoid detection or be different enough
 - (Humans do this explicitly sometimes)
 - But perfect adversarial plagiarism avoidance systems are creative

Back to IP

- Copyright is about the surface form
 - The actual expression of the idea
 - West Side Story is not plagiarism of Romeo and Juliet
- Patents are about the deep form
 - The idea itself
 - (But as applied to a particular endeavor)

Detecting Patent Violations

- Given two patents are they about the same thing
 - You only have the surface text description of each
 - You need to derive the underlying idea
 - Compare the ideas (not the text)
- Patent search is a big field
 - The USPTO now uses Google technology to do this
 - Its hard due to the surface variation in text
 - It still uses surface level features, topics etc not “meaning”
- Patents though are still hard (algorithm equivalence)
 - Solving the Traveling Salesman in polynomial time vs
 - Solving 3-SAT in polynomial time vs
 - Solving the Traveling Mailman in polynomial time

Detection vs Creation

- Can machines create new works?
- If a language generation system generates a funny limerick
 - Who has the copyright?
- We know Monkey's can't own copyright
 - (But corporations can)
- We are building more sophisticated generative models
 - They are “creating” new works
 - You can be sure if these works are worth money there will be copyright



● From Wikimedia

Creation of Melodies

- Build generative model for songs
- From the major/minor scale
- Generate 4-4-4-4 bar structure
- Follow techniques for coherence in songs
- Generate them all
- There is a court case that rules against these as creations
- [Google JS Bach Counterpoint generator]
- If a tune plays and no one listens is it still copyright?
 - For Trade Mark laws there **must** be a trade
- But if this generates successful songs there will be copyright

IP and Plagiarism

- IP can be compared with NLP techniques
 - Copyright targets surface similarity
 - Patents target deep similarity
- Similarity is beyond similarity
 - Its a legal question not just a computational one
- But useful systems don't need to deal with all cases
 - But useful systems will have false positives
 - (and will be tuned for the funder of the system)