

Homework 1
Networked Life (NETS 112)
Fall 2018
Prof Michael Kearns

Posted September 26, 2018. Due in hard-copy format at the start of lecture on Thursday, October 11. Please don't forget to write your name and staple the pages together.

Collaboration of any kind is NOT permitted on the homework.

Your Name: ANSWER KEY

Grading:
Chris: 2 and 3
Adel: 1, 4 and 5
Karthik: 6 and 7

Problem 1. The following two websites:

<https://mathscinet.ams.org/mathscinet/freeTools.html?version=2>

<https://oracleofbacon.org/>

...each provide tools for computing shortest-path distances in collaboration networks (of mathematicians and actors, respectively). Use these tools to find the pair of mathematicians, and the pair of actors, with the **largest** shortest-path distance you can. The pairs you find must be in the same connected component --- i.e. there must be a finite distance between them.

You should provide screenshots documenting the longest distances you are able to find, and the paths found by the tools. You should provide a detailed description of the methodology and ideas you used to find your longest distances. Your methodology may include any material or research you like --- the sites themselves, information you find on the open Internet, systematic or random exploration, etc.

This problem will be graded on a combination of the actual distances you're able to find (larger is better), the creativity of your methodology, and your clear description of that methodology.

A prize of some kind (to be determined) will be awarded to the student(s) who find the largest distances.

For each website:

1/1 for having > 3 distance

2/4 for using only random methodology

OR

4/4 for doing anything more complex than random that makes sense (i.e. researching for factors that would separate people), even if best result was random.

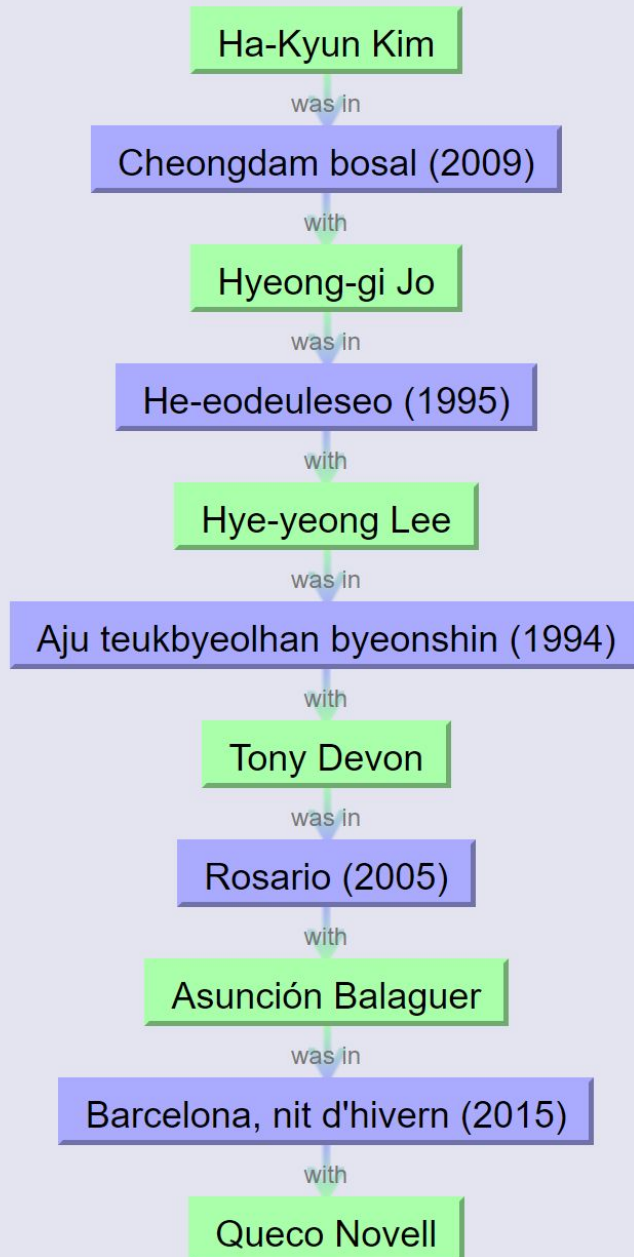
10 points total

Example Answer:

By looking at two actors who work in different countries and speak different languages. Google actors by searching by nationality, picking random ones that are not the most popular and then checking to see if there is a path.

Ha-Kyun Kim has a Queco Novell number of 5.

Find a different link



Queco Novell

to Ha-Kyun Kim

Find link

More options >>

Problem 2. Mark Zucchini, after taking NETS 112, is now interested in how content on Facebook is reshared throughout the network. Given the underlying graph of the social network and the timestamps at which the content is shared, Mark wants to see what the corresponding **cascade tree** looks like. For instance, here's the underlying social network of {A,B,C,D} and the timestamps at which some content was shared.

<i>Time</i> <i>t</i>	<i>people who shared</i> <i>at t</i>
1	A
2	B
3	C, D

Then the corresponding cascade tree will look as follows:

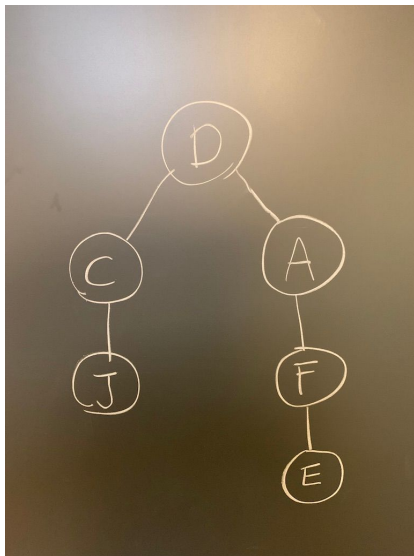
A first posts the content. After seeing the content shared by A, B reshares it. C could have seen the post from either A or B. However, we will **always break ties by attributing the cause of a reshare to the friend who shared most recently**. Hence, because by the time C shares at time $t=3$, B has shared the post more recently than A, we will attribute C's "infection" to B. Finally, we see that D reshares the content after seeing it from B; C and D share the post at the same time $t=3$, so the only friend that D could have had seen the post from is B. So the cascade tree, rooted at the person who first posts the content, illustrates who is responsible for other people resharing.

Also, Mark wants to calculate the **virality** of the cascade, where virality is defined as the **average distance across all pairs of nodes in the cascade tree**. In this case, the virality is $\frac{1}{n(n-1)} \sum_{i \neq j} d(i, j)$, where $d(i, j)$ is the distance between vertex i and j .

(a) Using the graph above and the following timestamps at which content was shared, draw the cascade tree, and calculate its virality.

<i>Time t</i>	<i>people who shared at</i> <i>t</i>
1	D
2	C
3	A
4	J
5	F
6	E

+2 for a correct cascade tree
 +1 for correct virality with work



Virality

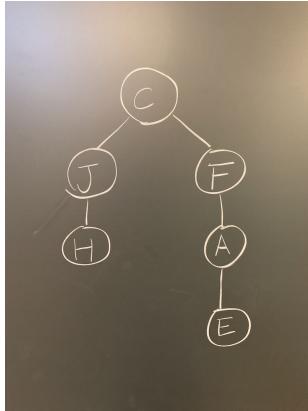
$$= (d_{jc} + d_{jd} + d_{ja} + d_{jf} + d_{je} + d_{cd} + d_{ca} + d_{cf} + d_{ce} + d_{da} + d_{df} + d_{de} + d_{af} + d_{ae} + d_{fe})/15$$

$$= 35/15$$

(b) Using the graph above and the following timestamps at which content was shared, draw the cascade tree, and calculate its virality.

Time t	people who shared at t
1	C
2	J
3	F
4	A
5	E
6	H

+2 for a correct cascade tree
 +1 for correct virality with work



Virality

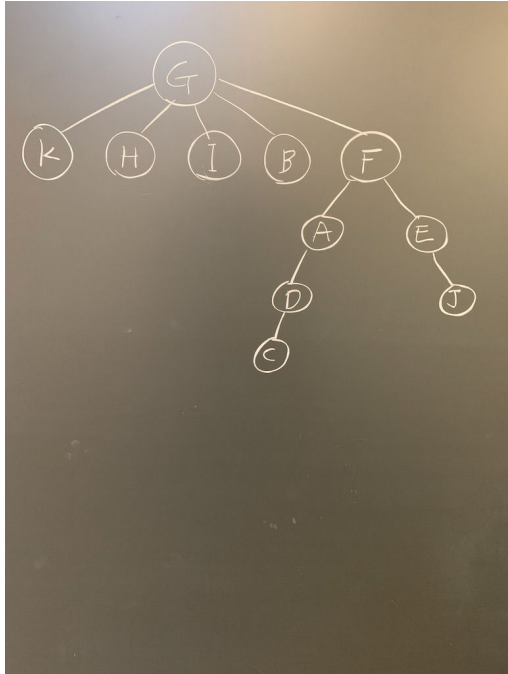
$$= (d_{hj} + d_{hc} + d_{hf} + d_{ha} + d_{he} + d_{jc} + d_{jf} + d_{ja} + d_{je} + d_{cf} + d_{ca} + d_{ce} + d_{fa} + d_{fe} + d_{ae})/15$$

$$= 35/15$$

(c) Using the graph above and the following timestamps at which content was shared, draw the cascade tree. You don't need to calculate the virality.

+4 for cascade tree

Time t	people who shared at t
1	G
2	K
3	H, I, B
4	F
5	A, E
6	J
7	D
8	C

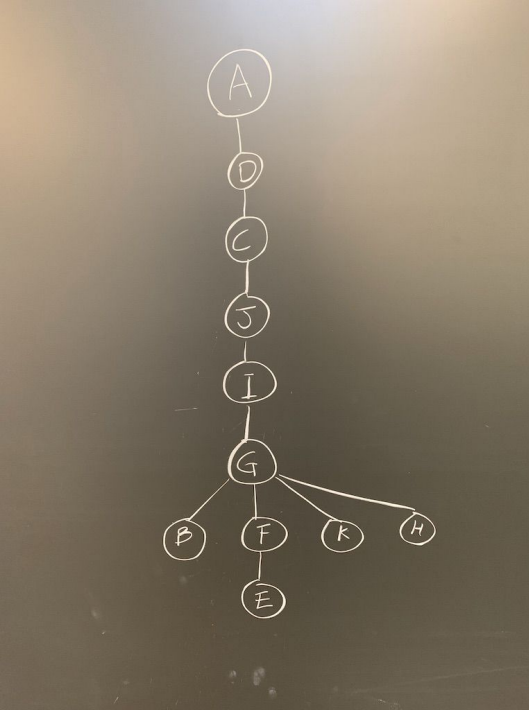


(d) Using the graph above and the following timestamps at which the content was shared, draw the cascade tree. You don't need to calculate the virality.

+4 for cascade tree

Time t people who shared at t

Time t	people who shared at t
1	A
2	D
3	C
4	J
5	I
6	G
7	K
8	F
9	B
10	H
11	E



(e) Based on the structure of the cascade trees, which of the cascades from (c) and (d) was most likely to be Taylor Swift posting a photo of her breakfast, and which was most likely to be a political petition? Justify your answers from the readings, ideally with specific citations.

4 points

+2 for correct tree identification: d political petition, c Taylor swift

+1 for correct explanation: e.g. skinny tree = grass-root petition / fat tree = celebrity posting

+1 for showing knowledge of concepts from readings

Problem 3. Suppose your first job out of college is as a “viral marketer”. You are given the task of designing a piece of content, and a plan for disseminating it on Facebook, with the goal that it achieve the **largest number of reshares possible**. You are able to control the following aspects of the campaign:

- The design of the content itself: It could be a photo (real or synthetic), a video an audio clip, a news article, etc. --- basically anything that it is possible for users to post on Facebook.
- The identities of the “seed” posters: Your firm has given you budget sufficient for you to convince/pay any 100 users of Facebook to post your content at any time of your choosing (not all necessarily simultaneously). So another component of your design is who these 100 users should be --- their identities, personalities, structural location in the Facebook network, etc. --- as well as the timing/schedule with which they will post your content.
- Other temporal features of the campaign --- e.g. particular dates or season of year in which you launch.
- Any other aspects of the campaign you want to design, except that in the end your only resource is the posting by the 100 seeds. In particular, you have **no ability to do any kind of external publicity or advertising outside of Facebook**. Your firm will judge your performance on the extent to which your content “goes viral”, purely on the basis of your 100 seeds and the cleverness of your campaign design.

Write an essay in which you describe your proposed campaign design as clearly as you can, along the dimensions suggested above or others you’d like to discuss. **You must justify as much of your design as possible by making reference to specific results in the contagion papers we have read.** Be as precise as you can in these references (e.g. paper title, quotations, page or figure numbers, etc.)

+4, +4, +4, +3 for each of the four components.

- Mention readings = must
- Mention only outside = partial credit
- Good readings to talk about: Structural Virality of Online Diffusion, Structural Diversity in Social Contagion, Can Cascades be Predictive
- Key points:
 - content is not that important.
 - Time between each reshare should be
 - Choose nodes with large degrees as seeds (e.g. celebrities) because most of adoptions are accounted for by the root nodes or by the immediate followers of the root nodes
 - Breadth vs depth (breadth helps initially)

Problem 4. Suppose there is a very large population of people for whom we can measure or obtain two numerical properties --- let's call them Property A and Property B --- and from these two properties, we'd like to predict some binary outcome y . For example, in class we discussed the problem of using high school GPAs (Property A) and SAT scores (Property B) to predict whether students admitted to Penn will graduate within four years or not (the outcome y). Importantly, we don't have Properties A and B and the associated y values for the entire population --- we only have a much smaller **training sample** from the overall population. The "machine learning" approach to this problem is to learn a "model" that does a good job of predicting y from A and B on the training sample --- and then hoping that this model will "generalize", in the sense that it will also do well at predicting y from A and B for the overall population.

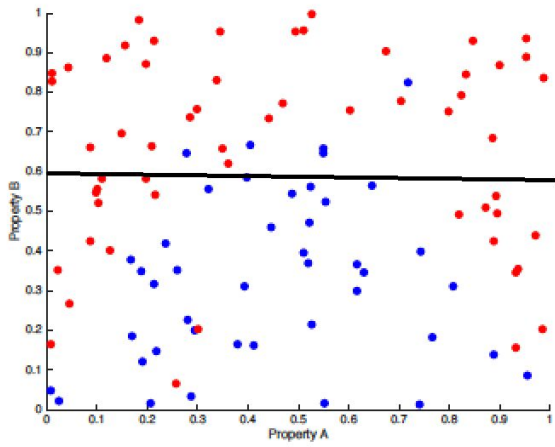
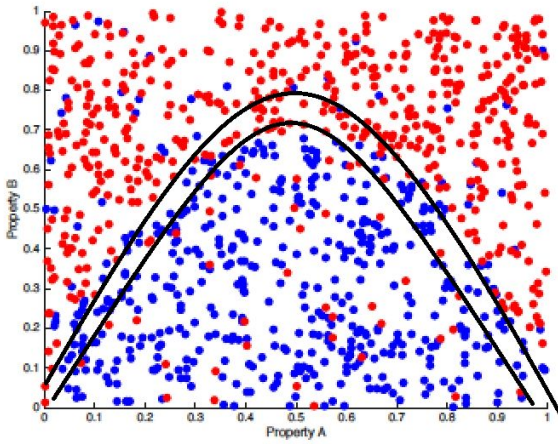
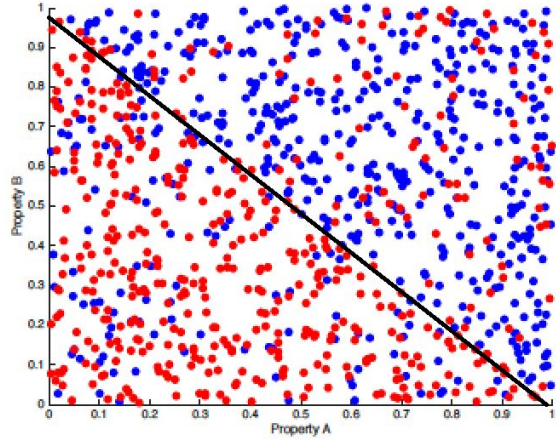
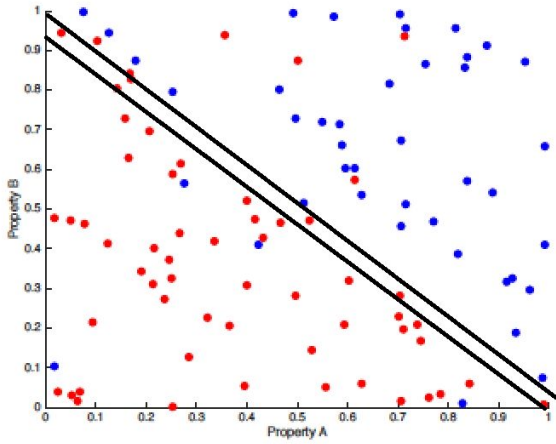
In each of the figures below, Property A is measured on the x axis, Property B on the y axis, and the dots show the (A,B) values of the training sample. There two possible outcome values are $y = \text{red}$ and $y = \text{blue}$, and the color of a dot indicates its outcome.

On each figure, you should draw a "model" that does a good/reasonable (though not necessarily perfect) job of separating red points from blue points, while remembering that the goal is not perfect fit to the training sample, but generalization to the unseen overall population. Your model should take the form of a curve in the plane whose shape/simplicity/complexity is up to you. In each case, you should write a brief justification of why you chose the curve you did, relying on our in-class discussion of overfitting. The following website might also be useful:
<https://en.wikipedia.org/wiki/Overfitting>.

+2 points for a shape that makes sense (i.e. isn't overfitted)

+2 for explanation (should mention sparsity for c vs d, general trend of blue vs red, not being too complex, overfitting w/ explanation, noise)

Examples below (two lines indicate an acceptable range of answers fit between)



On the fourth - if you draw a quadratic but justify it clearly and mention that there is not enough data to make a complex fit, full curve. If they just draw a quadratic and give a poor explanation then partial credit.

Problem 5. Consider the assigned paper “Navigation in a Small World” in the context of the network shown below.

(a) Add an edge to this network such that the shortest-path distance from A to B becomes 3, and the navigation algorithm from the paper also *could* take 3 hops from A to B.

+1 for shortest path, +2 for navigation

(b) Add an edge to this network such that the shortest-path distance from A to B becomes 2, but the navigation algorithm from the paper definitely takes 6 hops from A to B.

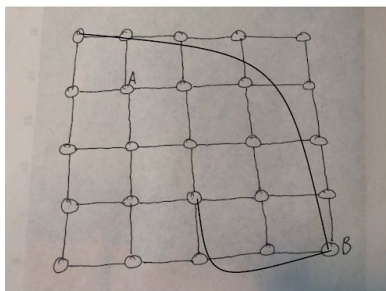
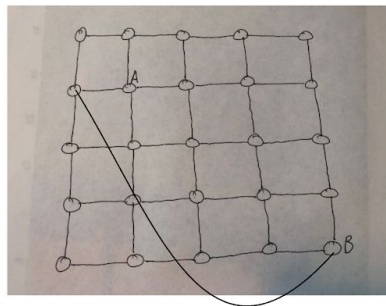
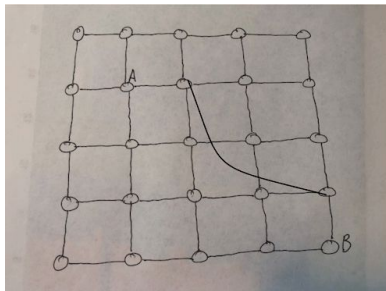
+1 for shortest path, +2 for navigation

(c) Add two edges to this network such that the shortest-path distance from A to B becomes 3, and the navigation algorithm from the paper *could* take 4 hops from A to B.

+1 for shortest path, +3 for navigation

Examples of acceptable edges below for A, B, and C, though there are more options.

Recall - the algorithm will always take steps towards the destination point so will not take edges that are “out of the way”.

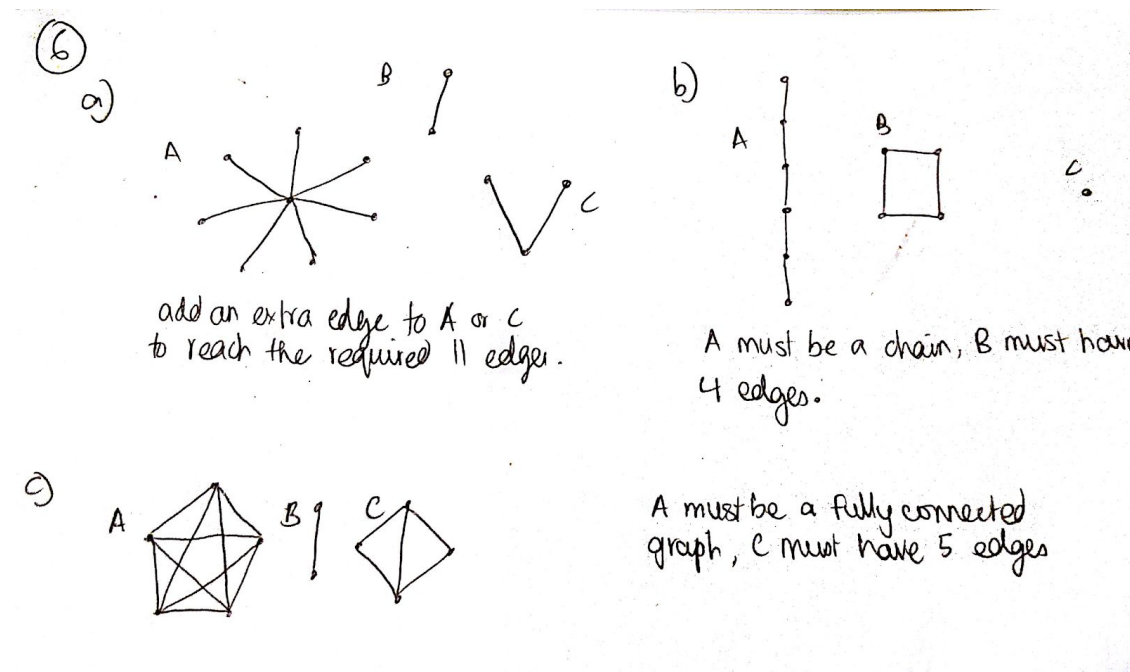


Problem 6. For each of the following parts, carefully draw an undirected network with the specified properties.

(a) Draw a network with three connected components of sizes 8, 3 and 2 vertices, in which the total number of edges is 11 and the worst-case diameter (longest shortest-path distance) of the largest component is 2.

(b) Draw a network with three connected components of sizes 6, 4 and 1 vertices, in which the total number of edges is 9 and the worst-case diameter of the largest component is 5.

(c) Draw a network with three connected components of sizes 5, 4 and 2 vertices, in which the total number of edges is 17 and the worst-case diameter of the largest component is 1.



EDIT: part c) contains an error. Component C should also be a fully connected graph with 4 vertices and 6 edges, as opposed to the 5 edges in the solution.

For each subpart:

+1 for the correct number of vertices in each component

+2 for the correct number of edges

+2 for the correct diameter

= 5 * 3 parts = 15 points

Problem 7. In this problem you are asked to compare and contrast the assigned readings “Navigation in a Small World” (Kleinberg) and “Identity and Search in Social Networks” (Watts, Dodd, Newman). For each of the following parts write a clear and concise response, supporting your answers with quotations or findings from the papers where possible.

(a) Both papers are trying to address a phenomenon that was implicit in the original Travers and Milgram paper, but was not explicitly examined there. What phenomenon is that?

Later papers call out the algorithmic components of navigation. Travers and Milgram focus on small diameter, which implicitly recognizes people’s ability to find short paths, as well as some of the criteria they used to find them. The two later papers formalize that short-path-finding as navigation algorithms and explore them.

(b) Discuss the similarities and differences between the approaches taken in the two papers. Your answer should address the models, methodologies and results of the papers.

Similarities should mention:

- Both are about searchability of networks and how people navigate with local information
- Both seek to explain small diameter using navigation

Differences should hit at the broad framework that Kleinberg has a simple clean model which is easy to prove properties for, but it’s unrealistic - it has a caricature of the world and only one dimension of relationships. In contrast, WDN pose a more representative model which factors in the complexity of the real world, but it’s almost intractable to analyze - they can’t prove any properties for it naturally and must rely on simulations to generate their results. They should also address both the modelling features and the results of the paper.

A sample answer that gets full credit might look like:

“The similarities are that both Kleinberg and WDN are seeking to identify navigation algorithms and network models that can allow people to search networks in a way that leads to short paths being found. Some of the differences between the papers are that 1) Kleinberg takes vertices as entities with no characteristics, whereas WDN give a vector of characteristics to each vertex. 2) Kleinberg’s model allows only for geographical basis for edges whereas WDN allow for more dimensions. 3) Kleinberg’s results are proven mathematically because his model is easy to analyze mathematically, whereas WDN must simulate networks to generate their results since their model is less tractable.”

(c) Which of the two papers do you find more compelling and convincing, and why?

+2 or +3 based on whether reasoning is provided.