



Quality and accuracy

Introduction to GIS

Hans Skov-Petersen (hsp@ign.ku.dk)



**Guess what Chat GPT says
What are the major problems with present day GIS?**

Quality and accuracy of geodata (I)? Can you be a little more specific?

- **Accuracy** is the deviation between reality and a measurement
 - In principle, it **cannot be assessed!**
 - In practice, accuracy is assessed by **comparing with 'better' or repeated measurements** of 'ground truth'
- **Precision** describes the capabilities of the instruments and methods used.
 - It can be assessed as the **deviation between repeated measures** of the same object by the same method
 - With respect to the captured data (measurements) precision applies to the **number of digits** presented by a given (computer-) registration.
- **Completeness** is a measure of the extent or amount of a give phenomena that was captured by an automation process.
- **Resolution** is the granularity – the smallest objects or characteristics you can possibly reveal from a dataset.
- **Compatibility** describes the extent to which data from different sources can be used together due to diversity definitions and quality in
 - Space - e.g. different scales
 - Time
 - Semantics, characteristics (attributes)
- **Consistency** addresses potential differences in production method and quality – like compatibility - but across a single data set



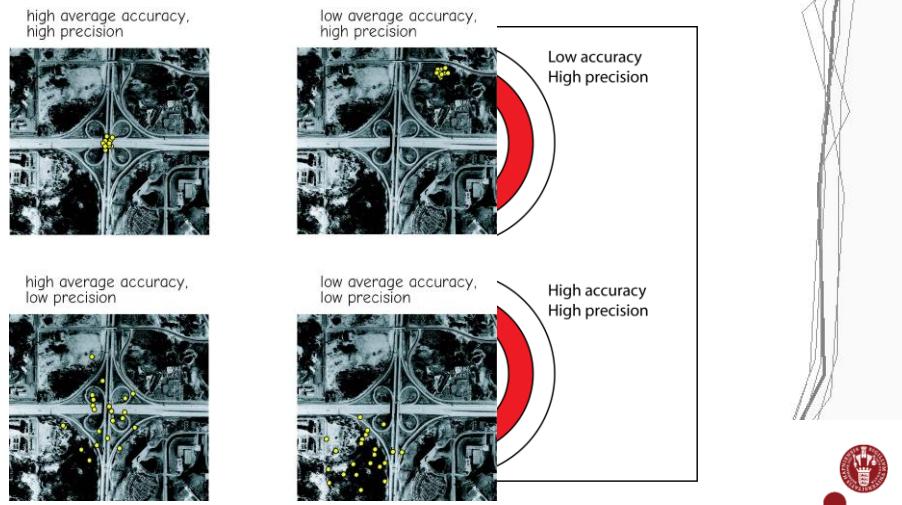
Quality and accuracy of geodata (II)?

- **Bias** is the systematic variation of data from reality
 - Data production method
 - Operator's preferences, education, knowledge
 - Instruments lacking calibration
 - Shifts in projection, datum etc.
- **Applicability/appropriateness/suitability** for specific analysis or cartographic presentations
- **Fitness for use.** 'Good quality' is not an objective norm.
 - The quality, accuracy, precision ect. of a data set can only be judged in relation to its use
 - ... and user
- **Trustworthiness.** The collective of all above: To which degree can you believe what you see – or can reveal from what you see - and make the right decision.

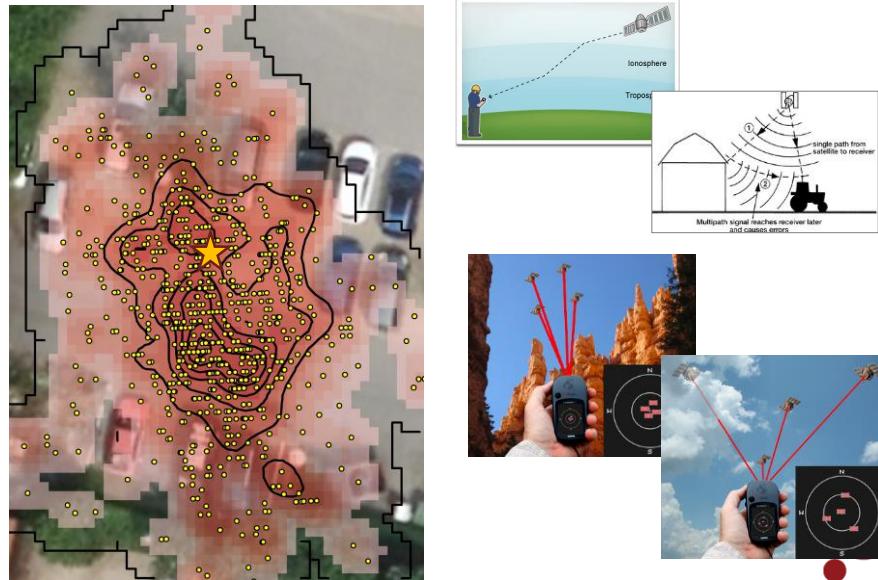


Accuracy and Precision

- Accuracy is how well data reflects reality
- Precision reflects the distribution of repeated measurements

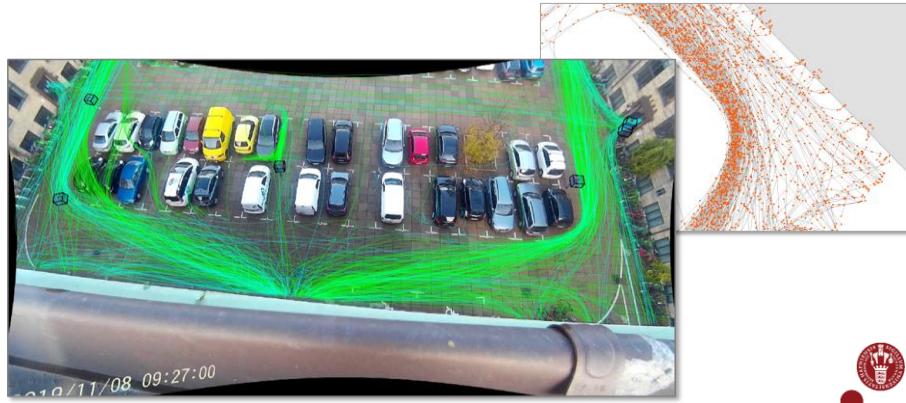


Assessment of precision ... of smart phone GPS location



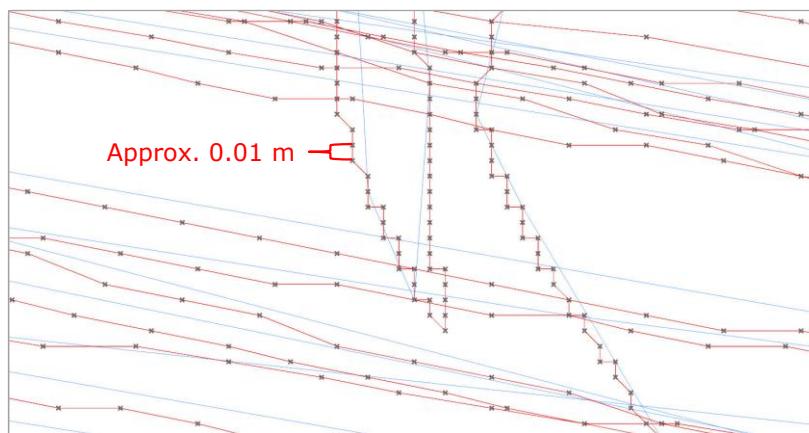
Assessment of precision A note on the number of decimals ... let's take a look at camera-base registration

- Very high spatial and temporal resolution and accuracy (compared to GPS)



Camera-base registration

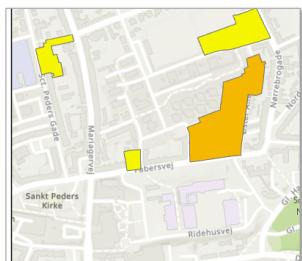
- Details at the edge of the reasonable...
- Why the edgy shapes?



Signalling precision

Another note on the number of decimals

... how do you label your legends?



Legend

- 96.656450 - 12278.849250
- 12278.849251 - 44619.656450
- 44619.656451 - 146254.479850
- 146254.479851 - 358114.871250
- 358114.871251 - 1075429.277250

Legend

- 96.6 - 12278.8
- 12278.8 - 44619.6
- 44619.6 - 146254.4
- 146254.4 - 358114.8
- 358114.8 - 1075429.2

Legend

- 0 - 10000
- 10 - 20000
- 20 - 30000
- 30 - 40000
- > 40000



Accuracy assessment (points)

According to the US 'National Standard for spatial Data Accuracy (NSSDA), Bolstad p 617:

- Identify/sample test points
- Identify their 'true' location (by higher order accuracy – measurements or existing, alternative data)
- Calculate mutual spatial deviations (distance)
- Calculate error statistics
 - Mean, STD and 95% confidence (assuming normal distribution)



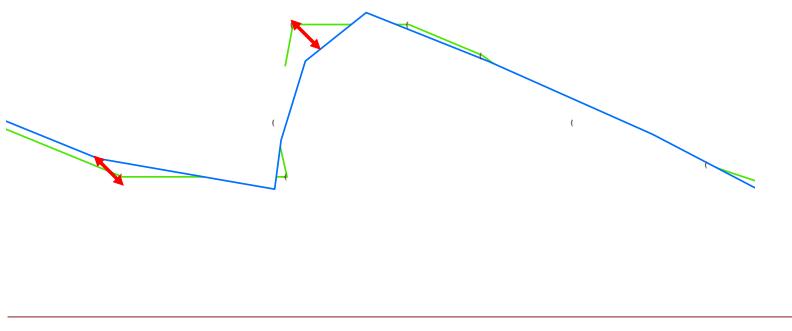
Standards (rules of thumbs) of accuracy

'National Standard for spatial Data Accuracy' (NSSDA), Bolstad p 617:

ID	x (true)	x (data)	x difference	(x difference) ²	y (true)	y (data)	y difference	(y difference) ²	sum x diff ² + y diff ²
1	12	10	2	4	288	292	-4	16	20
2	18	22	-4	16	234	226	6	36	52
3	7	12	-5	25	265	266	-1	1	26
4	34	34	0	0	243	240	3	9	9
5	15	19	-4	16	291	287	4	16	32
6	33	24	9	81	211	215	-4	16	97
7	28	29	-1	1	267	271	-4	16	17
8	7	12	-5	25	273	268	5	25	50
9	45	44	1	1	245	244	1	1	2
10	110	99	11	121	221	225	-4	16	137
11	54	65	-11	121	212	206	4	16	137
12	87	93	-6	36	284	278	6	36	72
13	23	22	1	1	261	259	2	4	5
14	19	24	-5	25	230	235	-5	25	50
15	76	80	-4	16	255	260	-5	25	41
16	97	108	-11	121	201	204	-3	9	130
17	38	43	-5	25	290	284	2	4	29
18	65	72	-7	49	277	282	-5	25	74
19	85	78	7	49	205	201	4	16	65
20	39	44	-5	25	282	278	4	16	41
21	94	90	4	16	246	251	-5	25	41
22	64	56	8	64	233	227	6	36	100

Assessment of accuracy (lines)

- In case of digitizing **linear data** - where vertexes cannot be mutually compared - the orthogonal distance from vertexes on one version of a line to another can be used



Assessment of accuracy (lines)

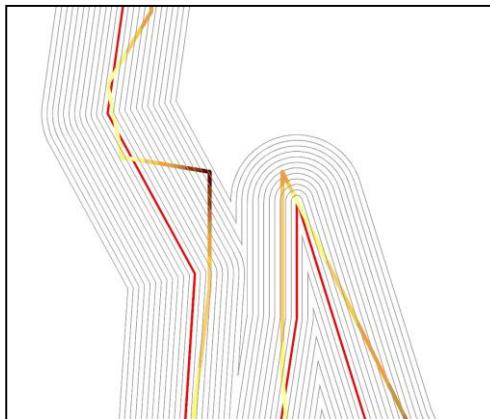
Or, you can assess the overlay of a buffer – of one line – vs another line ...



UNIVERSITY OF COPENHAGEN

Assessment of accuracy (lines)

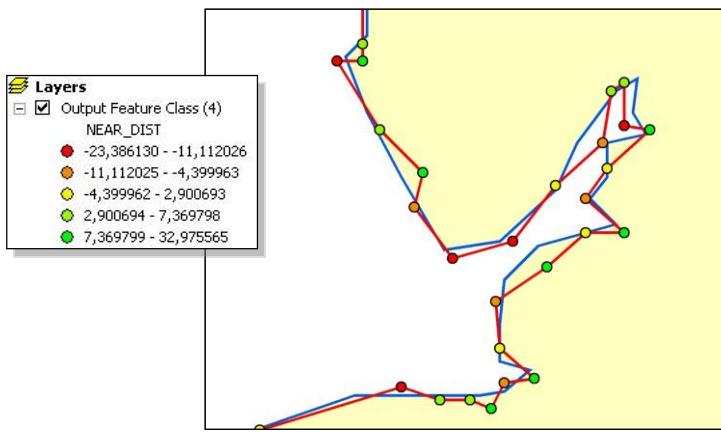
- Multi Ring Buffers (which will be the approach in the exercise)



UNIVERSITY OF COPENHAGEN

Assessment of accuracy (lines)

- A problem is that the deviation is two-sided – positive to one side, negative to the other
- If distances are assessed this way an assumption of normal distribution would be easier to justify



Accuracy assessment (polygons)

Let's say that you were to map the distribution of Rosa Rogusa
of a dune area...

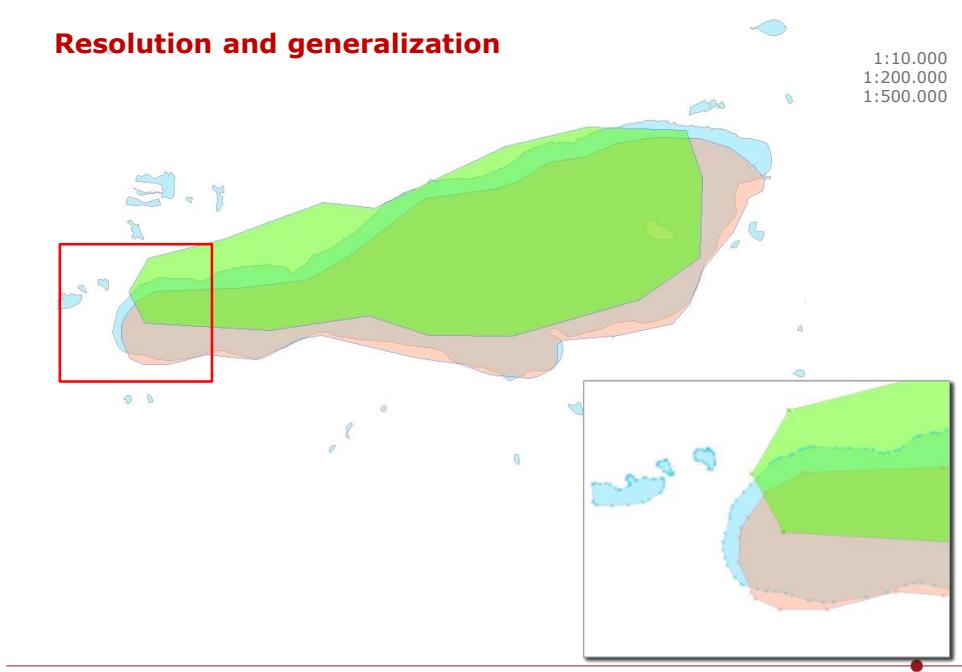
How would you approach and document the data you produce?



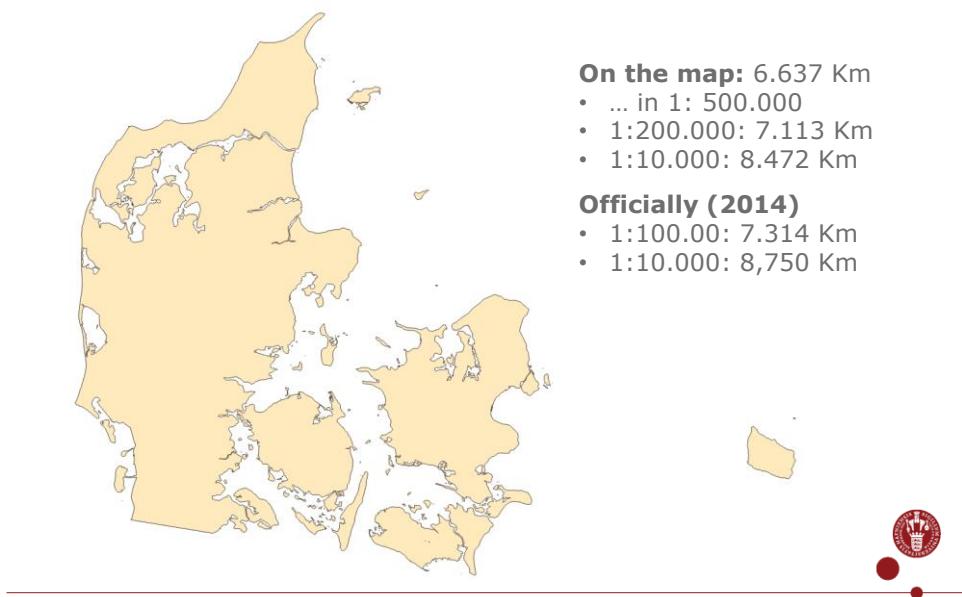
Pause



Resolution and generalization

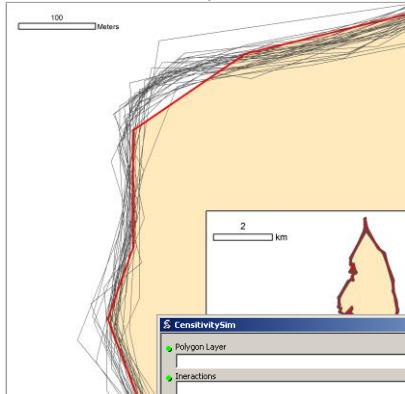


Resolution and the length of the Danish coastline

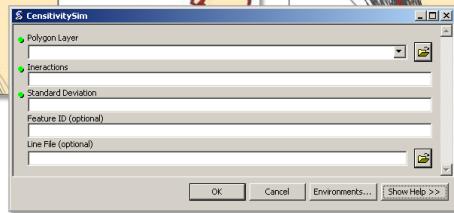


Assessment of accuracy of derived values (length and area)

34 students' attempts



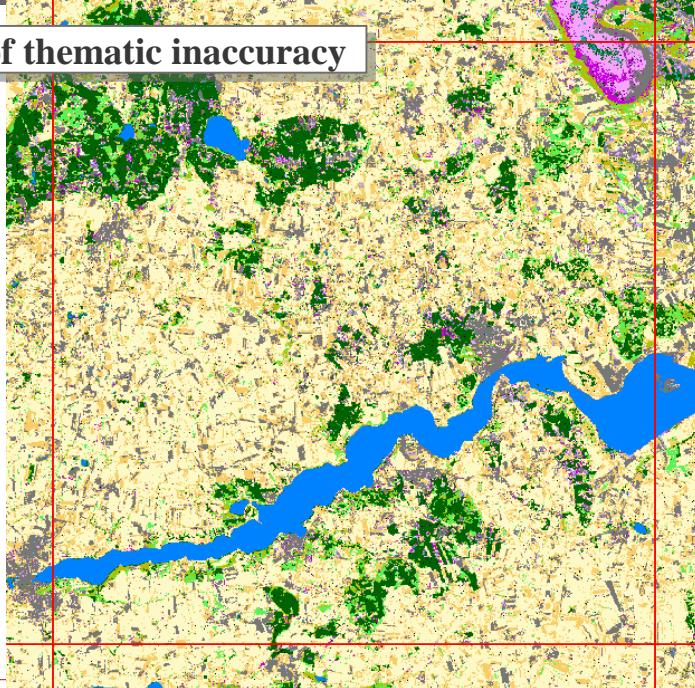
1000 automatically generated polygons



Assessment of thematic inaccuracy

The case of the
Area
Information
System (AIS)

Land Use
Land Cover



Thematic accuracy: Confusion matrix (1)

AIS Land use	AIS Land cover														Coniferous forest
	Grass land							Coniferous forest							
	X0	X1	X3	X5	X7	X8	X10	X11	X14	X15	X16	X18	Sum		
0	44													220	
1122	1		7	3	2	5	5	2	1	1	1	1	3	1.646	
1123	1		4	3	3			2	1	1	1	1	11	2.916	
1128													1	279	
1223			1	1	1	1	1	1	1				1	441	
1224			2	2	2	1	1	1				1	2	865	
1226													56		
Agriculture	2112	36	60	66	74	15	43	42	51	40	30	45	32.118		
	2300													100	
	2310													28	
	3110	3	1	1	4	4	2	4	9	19	14	7	3.931		
	3120	10	8	13	7	45	12	26	22	30	40	12	9.260		
	3130													22	
	3210		1	1	1	1	1		1	1	1	1	3	887	
	3220	1	1	5		1	2	5	2	2	3	2	1.012		
	3310													7	
	4110	1		1	2	2				1	1	3	3	219	
	4112			2	1			2					2	655	
	4120	2	96	2	2	2	9	30	12	5	4	4	4.685		
	4210	1				1	3	1						245	
	5120	3				1	9	1		4		1	1	622	
	5121													2	
	5126													21	
	5230			12			2							1.165	
	6000													116	
	Sum	100	100	100	100	100	100	0							

AIS codes:

- [Land use](#)
- [Land Cover](#)

Thematic inaccuracy: Confusion matrix (2)

AIS Land Use	AIS Land cover														Barren land	
	Barren land							Coniferous forest								
	X0	X1	X3	X5	X7	X8	X10	X11	X14	X15	X16	X18	Sum			
0	100													100		
1122		40	2	15	2	3	2	4	2	3	26	100				
1123		13	1	17			1	2	1	1	63	100				
1128		9	1	11			1	6	2	3	68	100				
1223	1	19	2	26	1	1	1	2	6	3	5	35	100			
1224		17	3	28	1	1	1	1	3	2	3	42	100			
1226		16	1	16	2	3	51	13	3	7	31	100				
Agriculture	2112	1	17	3	35		1	2	8	6	5	23	100			
	2300	2	1	40				3	4	2	4	49	100			
	2310	15	2	41				2	4	3	2	30	100			
	3110	2	14	1	1	1	13	22	18	28	100					
	3120	1	8	2	12	3	1	4	12	15	21	21	100			
	3130	1	6	6	2	2	3	25	23	21	11	100				
	3210	5	2	22	1	1	2	5	5	4	54	100				
	3220	6	6	7	1	2	8	12	7	14	37	100				
	3310	2	82	1	1	3	2				1	7	100			
	3330		12	1	3			3	4			76	100			
	4110	1	1	27	1			3	3	14	48	100				
	4112	2	4	34		1	5	3	3	1	46	100				
	4120	49	5	5	1	6	4	5	4	5	15	100				
	4210	5	9	2	33	7	3	2	5	1	5	28	100			
	5120	11	4	13	9	1	1	37	2	8	14	100				
	5121	20	3	22		2	2	3	9		40	100				
	5126	12	1	25	4	4	1	7	2	24	19	100				
	5230		95			1					1	1	100			
	6000		14	1	13		1	5	7	5	6	48	100			
	Sum	505	2.405	9.277	1.262	15.236	630	970	1.460	5.205	4.473	4.886	16.192	116		

AIS codes:

- [Land use](#)
- [Land Cover](#)

Completeness

Like in the case of accuracy/precision, completeness can only be assessed by either

- comparison with something regarded as being closer to the ‘truth’ or
- by repetition of *randomly selected portions* of data collection.

Completeness can be assessed in terms of e.g.

- *Counts* – the number of features within the test areas
- *Measures* – aggregated derived measures, e.g length and areas
- *Attributes* – amount of blank areas, confusion matrixes etc.



And so what?

We are at point in time where...

- Data production is hidden to its end-users
- Analytical techniques are well developed (and not a prime concern of most GIS-users)
- Storage and distribution of (massive amounts) of data is becoming increasingly important
- Dynamic – and often real-time – data are becoming more prominent
- Data and spatial analysis are used more and more in the (mass-) media
- ... and we don't know who we are talking to



How about all our non-professional users? We need to talk....

In the ol' days...

any GIS project would be constituted by 85% of data collection.

We don't need that anymore, which is so great.

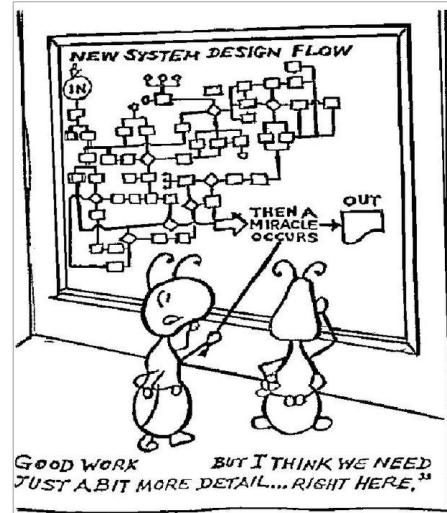
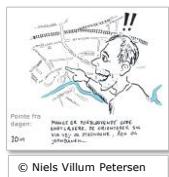
But... realizing that '**GIS is a mass-media**', things become much more challenging.

As professionals...

we have a fair idea of quality and what it means.
We know to which extent it is 'correct' and how we can have faith in it.

It cannot be expected for **the average, unknown target group** of mass-media communication.

We, as professionals, have an obligation to enhance communication of quality and trustworthiness of the geodata presented to the public



There are facts ... and alternative facts?

Talk to your neighbor.

How would you assess the **accuracy** and **precision** of measurements of the number of attendees at the two occasions?



How Kellyanne Conway ushered in the era of 'alternative facts'



Don't believe in any data you get...

They are nothing but representations
... of something
... for someone
... produced by someone
... for some purpose



**Thanks
That's it for today**



Hans Skov-Petersen
hsp@ign.ku.dk



AIS Land Use

Kode Beskrivelse			
1100 Befæstet overflade	1420 Sportsanlæg	4110 Eng	
1101 Bykerne	1421 Rekreativt område	4112 Vådområde	
1120 Lav bebyggelse	1422 Klippet græs	4120 Mose	
1121 Høj bebyggelse	2112 Landbrug	4130 Strandeng	
1122 Åben bebyggelse	2222 Gartneri	5120 So	
1210 Industri	2300 Græsarealer	5121 Vandlob > 8-12 m	
1221 Motorvej	2310 Græs i byområder	5123 Sø-rørskov	
1222 Motortrafikvej	2430 Blandet landbrug/natur	5126 Dambrug	
1223 Vej > 6 m	3100 Skov	5230 Hav	
1224 Vej 3-6 m	3110 Lovskov	6000 Uklassificeret	
1226 Jernbane	3120 Nåleskov		
1228 Bro	3130 Blandet skov		
1229 Dæmning	3210 Overdrev		
1240 Lufthavn	3220 Hede		
1242 Landingsbane	3250 Blandet natur		
1310 Råstofområde	3310 Sand/klit		
1340 Teknisk areal	3330 Anden overflade med ringe vegetation		
1341 Kirkegård			



AIS land Cover

Tabel 3. Oversigt over Land Cover Map og Land Cover Plus

Celle-verdi	Klasser	Kortkласse	Kort beskrivelse af klassen
		LCM LCM III LCP	
0	Ukendt klasse	•	Ikke-klassificerede billeddata
1	Åbent vand	•	Områder med permanent åbent vand
3	Ubevokset overflade	•	Inkluderer naturligt ubevokede overflader (f.eks. strande), ubevokset landstrugjord og betyggende områder
4	Vegetation påvirket af tidevand	1,3,7	Saltmarsk, alger og anden kystvegetation med tegn på oversvømmelse
5	Græsbevokset hedeområde	•	Tørre, semi-naturlige områder domineret af græsser, uden fjernelse af død; vegetativt materiale i større omfang
6	Permanent kort græs	7	Områder med permanent kort græs (g.a. hyppig græsning eller slæring)
7	Afgræsset eller stædt græs	•	Græsområder med lille akkumulering af død; vegetativt materiale – dog kan biomassen være enten høj eller lav
8	Engområde	•	Årligt våde, semi-naturlige områder domineret af græsser, med eller uden fjernelse af død; vegetativt materiale i større omfang
10	Busk- og græsbevokset hedeområde	•	Tørre, semi-naturlige områder domineret af en detaljeret mosaik af græsser og træagtige buske f.eks. lyng, revling, blæser, enebær
11	Buskbevokset hedeområde	•	Hovedsagelig tørre, semi-naturlige områder domineret af træagtige buske f.eks. lyng, revling, blæser, enebær
14	Busk- og skovområde	•	Terre eller våde områder domineret af buskbevoksning
15	Lavskov	•	Skove domineret af lavtræer
16	Nåleskov	•	Skove domineret af nåletræer
17	Urtebevokset moseområde	8	Sæsonbestemte våde, semi-naturlige områder med både græsser og andre urtebevoksninger
18	Sæsonbestemt arealdække	•	Områder, der skifter markant mellem et vegetationsdækket og ikke-vegetationsdækket stadium, inklusive dyrkede arealer



The you are here map

... of the Heisenberg institute

