Atmos. Meas. Tech. Discuss., 4, C904–C908, 2011 www.atmos-meas-tech-discuss.net/4/C904/2011/ © Author(s) 2011. This work is distributed under the Creative Commons Attribute 3.0 License.



**AMTD** 

4, C904–C908, 2011

Interactive Comment

# Interactive comment on "Near-surface profiles of aerosol number concentration and temperature over the Arctic Ocean" by A. Held et al.

### J. Piskozub (Referee)

piskozub@iopan.gda.pl

Received and published: 20 June 2011

The manuscript titled "Near-surface profiles of aerosol number concentration and temperature over the Arctic Ocean" by A. Held et al. is an interesting application of the gradient method of particle (and sensible heat) flux measurement. The paper presents new and important research results and is certainly worth publishing in Atmospheric Measurement Techniques.

In fact I must say I am surprised the authors have managed to measure Monin-Obukhov type of logarithmic profile on such a small altitude range. Surprised in the positive way. In the Petelski & Piskozub papers we used the ship mast of almost 20 m height and still we had difficulty convincing an anonymous reviewer that logarithmic gradients are possible to observe in nature. I believe the reason why this was possible with under





2 m height difference must have been the flat and smooth environment of sea-ice the authors were lucky enough to perform the measurements on.

The manuscript could be published almost as-is but I list some minor suggestions for the authors to consider hoping that may improve the paper.

1) I was surprised by the procedure of considering all height level permutations in order to determine the parameters of the logarithmic profile (first paragraph of Section 4.4). I do not see why it would be better than simply finding a best-fit for all the levels as we did. It seems to me we both use the same amount of information making the outcome equivalent. Am I wrong? I believe a comment on the reason of using the procedure would improve the paper allowing future users of the gradient method to do an educated choice between the variants.

2) The authors of the reviewed paper had the advantage of using eddy correlation at the same time as the gradient method. This allowed not only for the comparison of calculated fluxes but also made it possible to estimate independently friction velocity. At the time we made the measurements described in Petelski & Piskozub 2006, we did not have yet the possibility. Still we believed already then that simultaneous measurements with the gradient and eddy correlation methods could help establish whether the vor Karman constant is applicable also to particle flux (it's value was empirically established for heat fluxes and therefore its application for particle fluxes should be also checked experimentally). This was discussed in the Andreas comment to our paper and in our reply to it. We had seen some hints that the counterpart to van Karman constant for particle fluxes (let me call it Petelski constant) could be closer to 1.0 than 0.4. Would the authors care to comment whether their data can help constrain its value? I do not insist on including such a discussion in the manuscript (although I would not mind that). Commenting in the reply to this review would be enough if the authors do not feel their data could help constrain the Petelski constant in any meaningful way.

3) The third thing I would like to comment is using the statistical tests in the null hy-

# AMTD

4, C904–C908, 2011

Interactive Comment



Printer-friendly Version

Interactive Discussion



pothesis. First of all the phrase "the probabilities of acceptance of the null hypothesis" (line 10 in Section 4.1) is wrong. We never accept the null hypothesis. In fact we test how improbable it is to obtain our research hypothesis by accident, assuming the null hypothesis is true. If it's improbable enough (below the rather arbitrary threshold of 5% probability) we say we "rejected" the null hypothesis. However if the probability over the threshold we still do not accept the null hypothesis (as we never tested it in any way). We just say our research hypothesis "is not statistically significant".

However my comment goes further. I suggest not using the null hypothesis rejection and significance level analysis at all. The literature proposing this has long tradition. Cohen (1994) already said that after "4 decades of severe criticism, the ritual of null hypothesis significance testing - mechanical dichotomous decisions around a sacred .05 criterion - still persists." This methodology is criticized not only for the arbitrary threshold (the list of complains is too long to repeat here). Hunter (1997) in a paper which title itself tell it all argues that "The significance test as currently used is a disaster. Whereas most researchers falsely believe that the significance test has an error rate of 5%, empirical studies show the average error rate across psychology is 60% - 12 times higher than researchers think it to be". That is one of the reasons why Armstrong (2007) stated "I was unable to find empirical evidence to support the use of significance tests under any conditions" while Hubbard and Lindsay (2008) concluded "it is bad enough for researchers to misuse a measure that is useful: But it strains credulity to do so when that measure is seriously flawed in itself. And this paper has demonstrated - from a multitude of perspectives - that the p value is just that". Gigerenzer et al. (2004) actually compared using this methodology to rituals: "Elements of social rituals include (a) the repetition of the same action, (b) a focus on special numbers or colors, (c) fears about serious sanctions for rule violations, and (d) wishful thinking and delusions that virtually eliminate critical thinking [..]. The null ritual has each of these four characteristics: a repetitive sequence, a fixation on the 5% level, fear of sanctions by editors or advisers, and wishful thinking about the outcome (the p-value) combined with a lack of courage to ask guestions". To make it worse tests show that even 80% of

## AMTD

4, C904-C908, 2011

Interactive Comment



Printer-friendly Version

Interactive Discussion



scholars teaching statistics do not understand what significance testing actually means (Haller & Kraus 2002).

Most of the above examples of rejecting the "null ritual" come from social sciences and psychology. However at least two papers voiced the same concerns in the field of atmospheric (Nicholls 2000) and climate science (Ambaum 2010). However one may say: "OK, but what is the alternative?" There is more than one. Ambaum (2010) suggests Bayesian analysis which may be the future but the scientific world may not yet be ready for it (at least I'm not). The other proposition (one of the advices of Nicholls 2000) is using confidence intervals. This also is not a new proposal, Gardner and Altman proposed it in 1986 and later wrote a whole book promoting this approach (Altman et al. 2005).

In the case of the reviewed manuscript, the confidence interval approach would call for checking how many standard deviations ("sigmas") the values are from each other. If the distributions are normal, two sigmas correspondent to a 95% confidence interval, which actually implies what people expect from a 5% significance. I believe all the data presented in the paper would pass the 2 sigmas test. You may be surprised but I do not insist on implementing this suggestion. It is a matter of philosophy and I do not believe in coercion with respect to this matter, rather evangelizing (which I exactly what I did above).

Literature used in this section:

Cohen J (1994) The world is round (p<0.05). American Psychologist, 49, 997-1003

Hunter JE (1997) Needed: A ban on significance testing. Psychological Science, 8, 3-7

Amstrong JS (2007) Significance tests harm progress in forecasting. International Journal of Forecasting, 23, 321 - 327

Hubbard R, Lindsay RM (2008) Why P Values Are Not a Useful Measure of Evidence

4, C904–C908, 2011

Interactive Comment

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



in Statistical Significance Testing. Theory & Psychology, 18(1), 69-88

Gigerenzer G, Krauss S, Vitouch O (2004) The Null Ritual: What You Always Wanted to Know About SigniïňĄcance Testing but Were Afraid to Ask, Published in: D. Kaplan (Ed.). (2004). The Sage handbook of quantitative methodology for the social sciences (pp. 391–408). Thousand Oaks, CA: Sage.

Haller H, Krauss S (2002). Misinterpretations of signiïňĄcance: A problem students share with their teachers? Methods of Psychological ResearchâĂŤOnline [Online serial], 7 (1), 1–20

Nicholls N (2000) The Insignificance of Significance Testing., Bulletin of the American Meteorological Society, 81, 981-986

Ambaum MHP (2010) SigniïňĄcance Tests in Climate Science. Journal of Climate, 23, 5927-5932

Gardner MJ, Altman DG (1986) Confidence intervals rather than P values: estimation rather than hypothesis testing. British medical Journal, 292, 746-750.

Altman et al. (ed.) (2005) Statistics with confidence. Wiley-Blackwell; 2nd Edition, 240 pp.

There are some purely technical matters I would like also to mention:

data (abstract, line 5) is usually treated as plural of "datum" so I would prefer "were" to "was" height sensor "pointing normally" toward the ground (Section 2.1 line 34). I would prefer "pointing vertically"

Otherwise, the language of the manuscript is clear and easy to follow.

Regards,

Jacek Piskozub

# AMTD

4, C904–C908, 2011

Interactive Comment

Full Screen / Esc

Printer-friendly Version

Interactive Discussion



Interactive comment on Atmos. Meas. Tech. Discuss., 4, 3017, 2011.