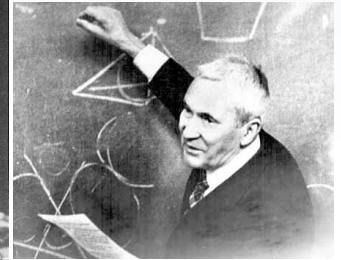
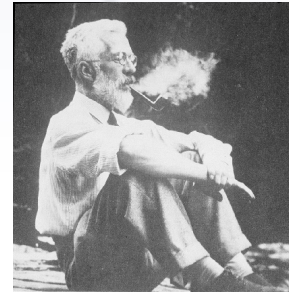
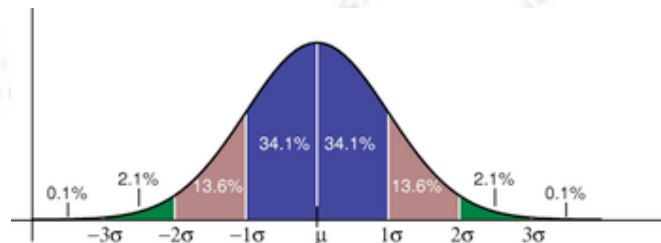


# Big Data Analysis

Data set: Housing Prices



Troels C. Petersen (NBI)



*“Statistics is merely a quantisation of common sense - Machine Learning is a sharpening of it!”*

# Data, goal, and misc.

## The data:

About **50.000 real estate sales**, including the final sales price along with several descriptive variables, many incomplete or missing.

## The goal:

**To determine the final sales price as accurately as possible.**

NOTE: “As accurately” is not a well determined measure, and we will discuss this.

## Miscellaneous:

While the dataset is on the border of “Big Data”, we have chosen it, as it fits all the ML methods well, and since its analysis can be **done in finite time**.

# Dataset variables - 90 in total

0 MI\_OBJ\_OIS\_PROPERTY\_ID  
1 MI\_OBJ\_OIS\_PROPERTY\_NUMBER  
2 MI\_OBJ\_OIS\_MOTHER\_ID  
3 MI\_OBJ\_OIS\_MUNICIPALITY\_NUMBER  
4 MI\_OBJ\_OIS\_POSTAL\_CODE  
5 MI\_OBJ\_OIS\_RENTED\_PLOT  
6 MI\_OBJ\_OIS\_OWNERSHIP\_CODE\_PROPERTY  
7 MI\_OBJ\_OIS\_OWNERSHIP\_CODE\_UNIT  
8 MI\_OBJ\_OIS\_PROPERTY\_APPLICATION\_CODE\_UNIT  
9 MI\_OBJ\_OIS\_PROPERTY\_APPLICATION\_CODE\_BUILDING  
10 MI\_OBJ\_OIS\_PROPERTY\_USE\_CODE  
11 MI\_OBJ\_OIS\_SALES\_PRICE  
12 MI\_OBJ\_OIS\_DATE\_OF\_SALES\_PRICE  
13 MI\_OBJ\_OIS\_PREVIOUS\_SALES\_PRICE\_FIRST  
14 MI\_OBJ\_OIS\_DATE\_OF\_PREVIOUS\_SALES\_PRICE\_FIRST  
15 MI\_OBJ\_OIS\_PREVIOUS\_SALES\_PRICE\_SECOND  
16 MI\_OBJ\_OIS\_DATE\_OF\_PREVIOUS\_SALES\_PRICE\_SECOND  
17 MI\_OBJ\_OIS\_PREVIOUS\_SALES\_PRICE\_THIRD  
18 MI\_OBJ\_OIS\_DATE\_OF\_PREVIOUS\_SALES\_PRICE\_THIRD  
19 MI\_OBJ\_OIS\_PREVIOUS\_SALES\_PRICE\_FOURTH  
20 MI\_OBJ\_OIS\_DATE\_OF\_PREVIOUS\_SALES\_PRICE\_FOURTH  
21 MI\_OBJ\_OIS\_TAXATION\_VALUE  
22 MI\_OBJ\_OIS\_TAXATION\_VALUE\_PLOT  
23 MI\_OBJ\_OIS\_TAXATION\_VALUE\_FARMHOUSE  
24 MI\_OBJ\_OIS\_DATE\_OF\_TAXATION\_VALUE  
25 MI\_OBJ\_OIS\_PROPERTY\_ADDRESS  
26 MI\_OBJ\_OIS\_HOUSE\_NUMBER  
27 MI\_OBJ\_OIS\_HOUSE\_LETTER  
28 MI\_OBJ\_OIS\_DOOR\_CODE  
29 MI\_OBJ\_OIS\_FLOOR\_NUMBER  
30 MI\_OBJ\_OIS\_MAX\_FLOOR\_NUMBER\_BUILDING  
31 MI\_OBJ\_OIS\_LAND\_ZONE  
32 MI\_OBJ\_OIS\_SIZE\_OF\_HOUSE  
33 MI\_OBJ\_OIS\_SIZE\_OF\_BUSINESS\_AREA  
34 MI\_OBJ\_OIS\_SIZE\_OF\_PLOT  
35 MI\_OBJ\_OIS\_SIZE\_OF\_INTEGRATED\_CARPORT  
36 MI\_OBJ\_OIS\_SIZE\_OF\_NOT\_INTEGRATED\_CARPORT  
37 MI\_OBJ\_OIS\_SIZE\_OF\_OUTDOOR\_LIVING\_ROOM  
38 MI\_OBJ\_OIS\_SIZE\_OF\_INTEGRATED\_OUTHOUSE  
39 MI\_OBJ\_OIS\_SIZE\_OF\_INTEGRATED\_GARAGE  
40 MI\_OBJ\_OIS\_SIZE\_OF\_LEGAL\_BASEMENT  
41 MI\_OBJ\_OIS\_SIZE\_OF\_BASEMENT  
42 MI\_OBJ\_OIS\_SIZE\_OF\_ATTIC  
43 MI\_OBJ\_OIS\_SIZE\_OF\_USED\_ATTIC  
44 MI\_OBJ\_OIS\_SIZE\_OF\_HOUSE\_EXCL\_UTILIZED\_ATTIC  
45 MI\_OBJ\_OIS\_SIZE\_OF\_BUSINESS\_AREA\_BUILDING  
46 MI\_OBJ\_OIS\_SIZE\_OF\_NOT\_INTEGRATED\_GARAGE  
47 MI\_OBJ\_OIS\_NUMBER\_OF\_FLOORS  
48 MI\_OBJ\_OIS\_CONSTRUCTION\_YEAR  
49 MI\_OBJ\_OIS\_CONSTRUCTION\_MATERIAL  
50 MI\_OBJ\_OIS\_REBUILD\_YEAR  
51 MI\_OBJ\_OIS\_ROOF\_MATERIAL  
52 MI\_KNN\_PROPERTY\_CONDITION  
53 MI\_KNN\_TOP\_FLOOR\_INDICATOR  
54 MI\_KNN\_GROUND\_FLOOR\_INDICATOR  
55 MI\_KNN\_GROUP\_VALID\_REGRESSION\_INPUT  
56 MI\_KNN\_GRP\_PERCENTILE\_MIN\_WEIGHTED\_SIZE\_OF\_HOUSE  
57 MI\_KNN\_GROUP\_PERCENTILE\_MIN\_SIZE\_OF\_PLOT  
58 MI\_KNN\_GROUP\_PERCENTILE\_MIN\_CONSTRUCTION\_YEAR  
59 MI\_KNN\_GROUP\_PERCENTILE\_MIN\_TAXATION\_VALUE  
60 MI\_KNN\_GROUP\_PERCENTILE\_MIN\_TAXATION\_VALUE\_PLOT  
61 MI\_KNN\_GRP\_PERCENTILE\_MAX\_WEIGHTED\_SIZE\_OF\_HOUSE  
62 MI\_KNN\_GROUP\_PERCENTILE\_MAX\_SIZE\_OF\_PLOT  
63 MI\_KNN\_GROUP\_PERCENTILE\_MAX\_TAXATION\_VALUE  
64 MI\_KNN\_GROUP\_PERCENTILE\_MAX\_TAXATION\_VALUE\_PLOT  
65 MI\_KNN\_M2\_P\_PREDIC  
66 MI\_KNN\_STD\_SALES\_PRICE\_NEIGHBORS  
67 MI\_KNN\_AVG\_GEO\_DISTANCE\_NEIGHBORS  
68 MI\_KNN\_AVG\_CONSTRUCTION\_YEAR\_NEIGHBORS  
69 MI\_KNN\_AVG\_WEIGHTED\_SIZE\_OF\_HOUSE\_NEIGHBORS  
70 MI\_KNN\_AVG\_SIZE\_OF\_PLOT\_NEIGHBORS  
71 MI\_KNN\_APARTMENTS\_NEIGHBORS\_INDICATOR  
72 MI\_KNN\_MATERIAL\_TYPE  
73 MI\_KNN\_APARTMENTS\_ACTUAL\_NUM\_OF\_NEIGHBORS  
74 MI\_KNN\_STATUS  
75 MI\_OBJ\_NUMBER\_OF\_EXTERNAL\_MATRS  
76 MI\_OBJ\_OIS\_SUM\_OF\_TAXATION\_VALUES  
77 MI\_OBJ\_OIS\_N\_COORDINATE  
78 MI\_OBJ\_OIS\_E\_COORDINATE  
79 C20\_1MONTH%  
80 C20\_3MONTH%  
81 C20\_6MONTH%  
82 C20\_12MONTH%  
83 SCHOOL\_DISTANCE\_1  
84 SCHOOL\_DISTANCE\_2  
85 SCHOOL\_DISTANCE\_3  
86 SUPERMARKET\_DISTANCE\_1  
87 SUPERMARKET\_DISTANCE\_2  
88 SUPERMARKET\_DISTANCE\_3  
89 KOEBESUM\_BELOEB

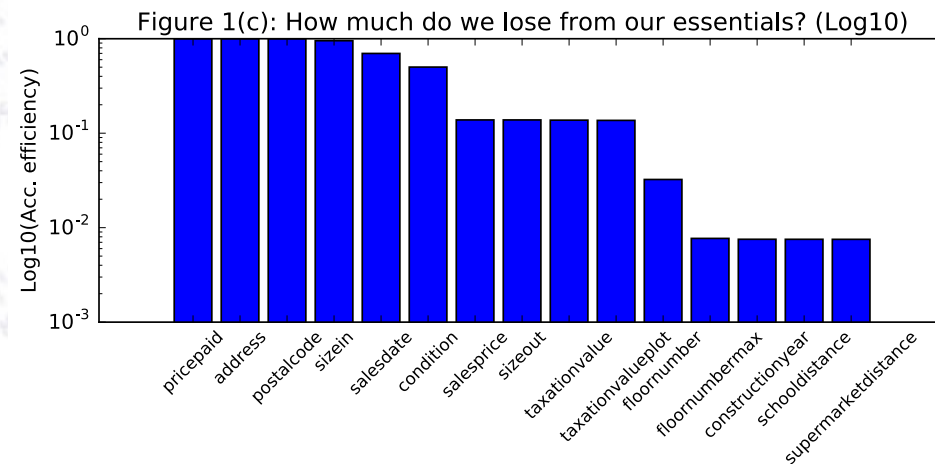
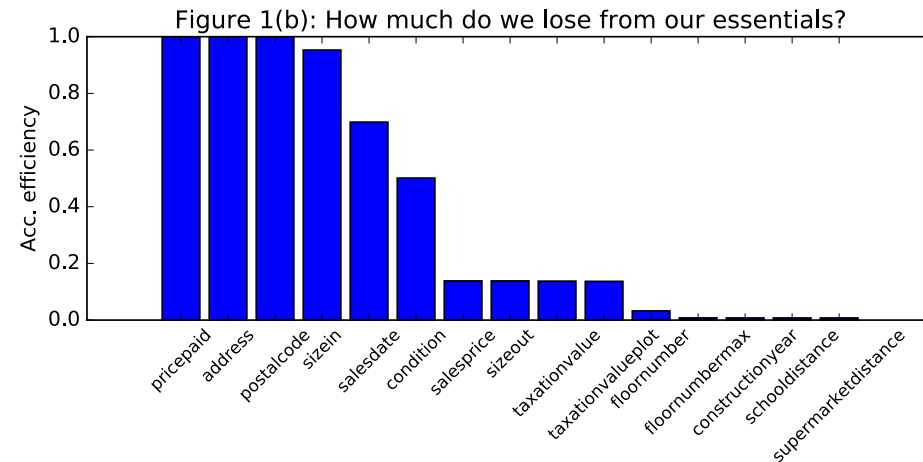
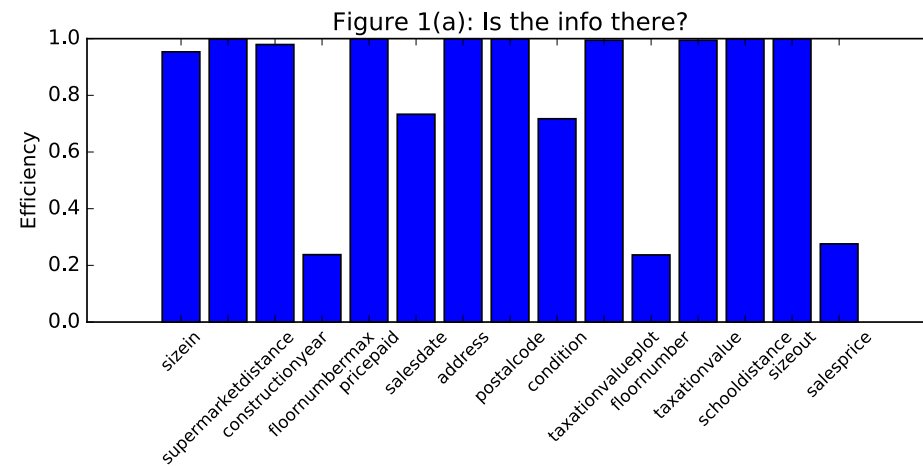


# Information available

While there are in principle 90 pieces of information on each property sale, it is in practice not the case! As it turns out, most entries are empty!!!

In the figure we consider the most crucial variables (see page before), and check what fraction of entries have information available here.

The conclusion is, that if we wanted all entries filled, we would only have < 1% of data remaining... not a great way forward!



# Information

## available

While there are in principle a lot of variables, not all of them are available. In practice, it is not the case that all information on each entry is available. In fact, it is in practice not the case that all information is available. In fact, out, most entries are empty.

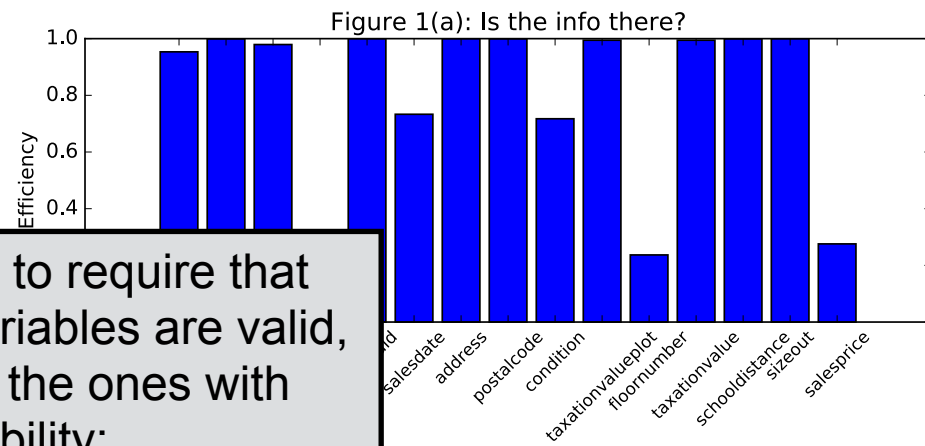
In the figure we consider the most crucial variables (see page 10). We check what fraction of the information available has been used.

The conclusion is, that if we require that all entries filled, we would have less than 1% of data remaining... not a great way forward!

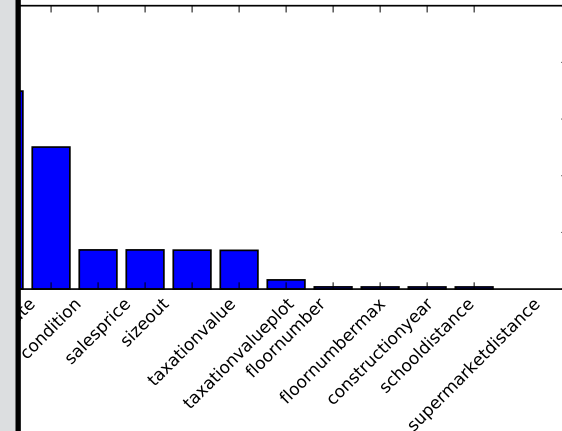
One could choose to require that e.g. the first six variables are valid, and then only add the ones with (almost) full availability:

- Price paid (of course)
- Address
- Postal Code
- Size inside
- Sales date
- Condition
- Size outside
- Taxation value
- Floor number
- School distance
- Supermarket distance

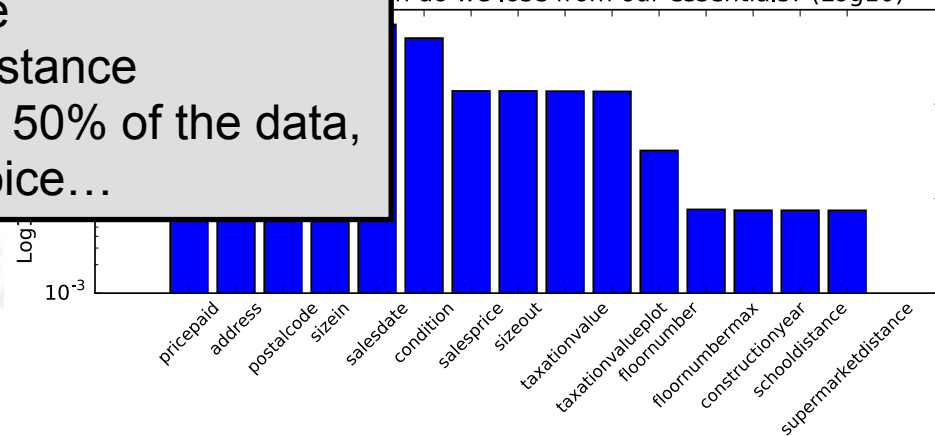
This leaves about 50% of the data, which is a fair choice...



How much do we lose from our essentials?

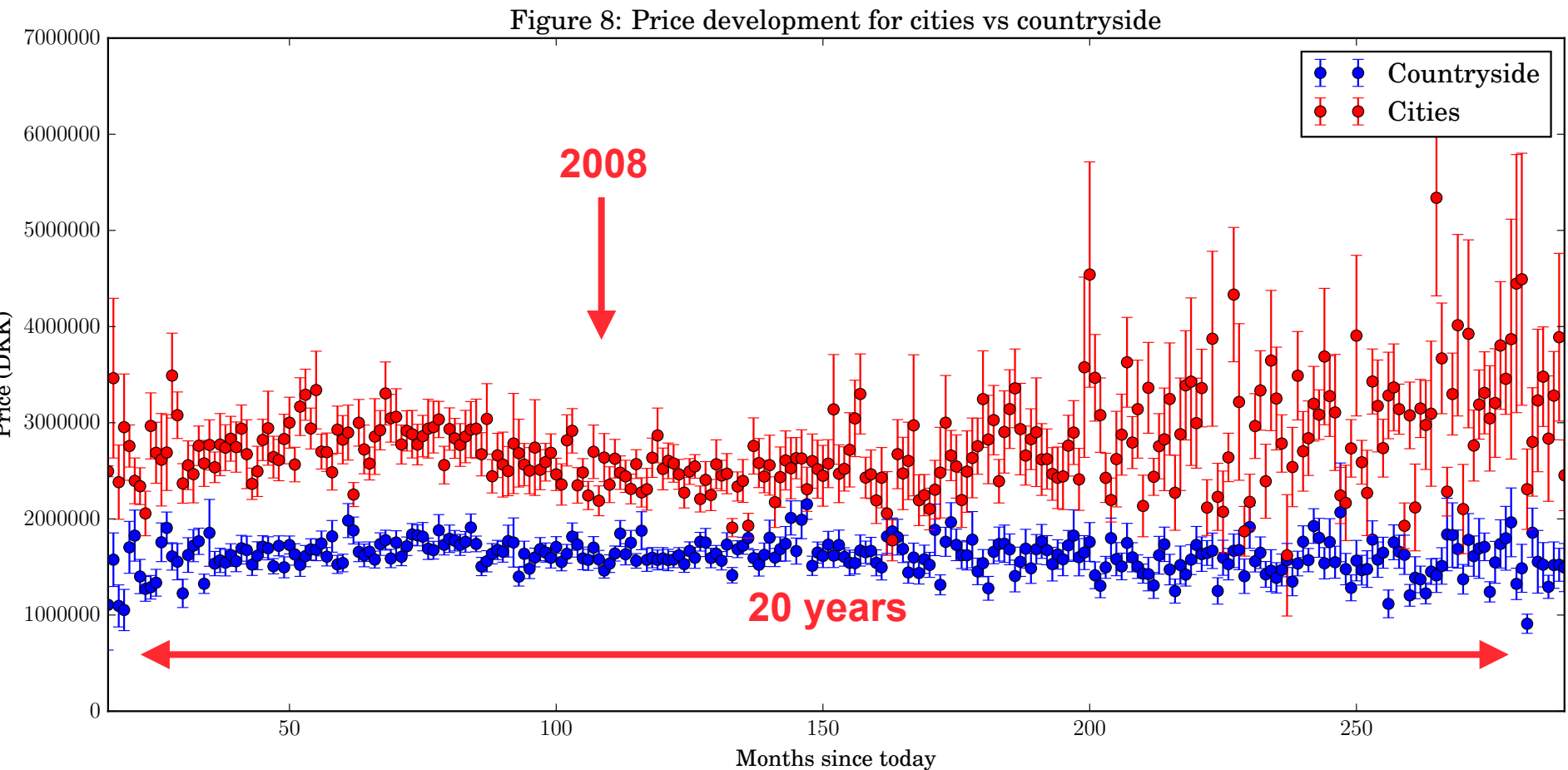


How much do we lose from our essentials? (Log10)



# Price vs. time

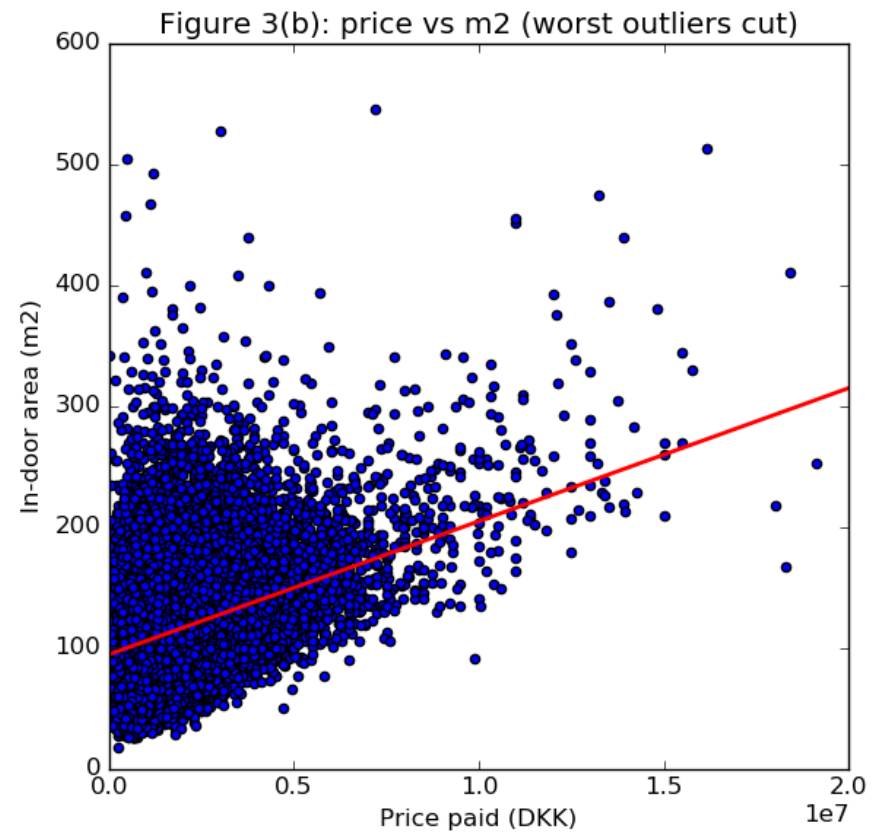
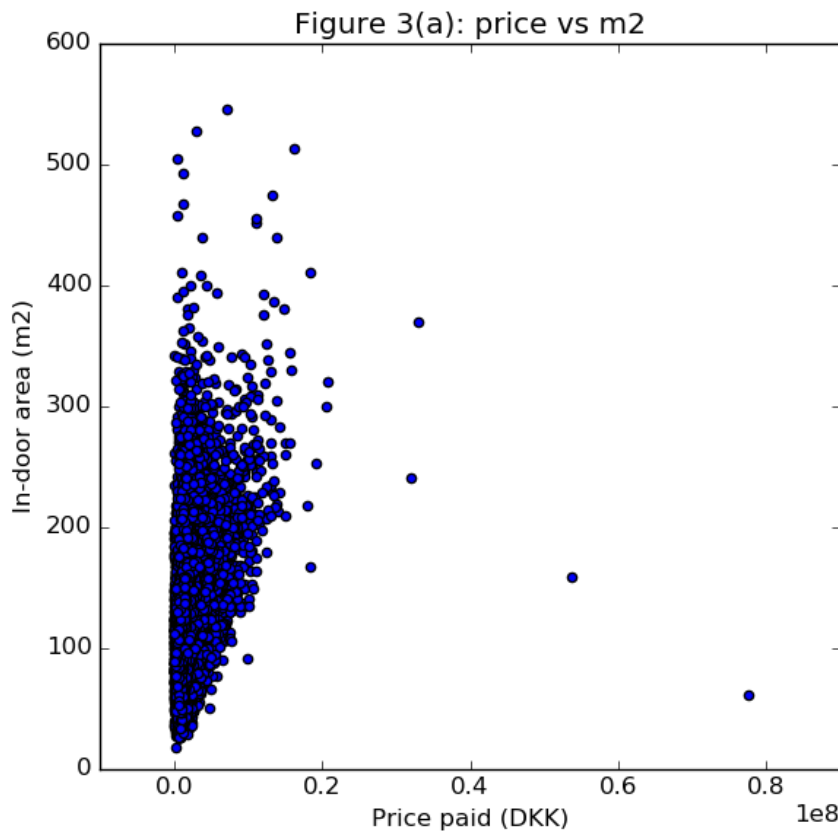
Just to gauge the data, we try to plot the average price over time:



Clearly, the data is corrected for inflation, but not much else, since 2008 doesn't clearly show up.

# Price per square meter

As a first step, one would estimate the price from the size, i.e. assume that the price per square meter was constant, and so we plot price vs. size:



As can be seen from the figure, this does not seem to be the case, and even after filtering away the worst outliers, we don't get any reasonable estimate!



# Price per square meter

Looking at the price/m<sup>2</sup>, most values are reasonable, but there are exceptions:

Figure 2(a): Is the info correct?

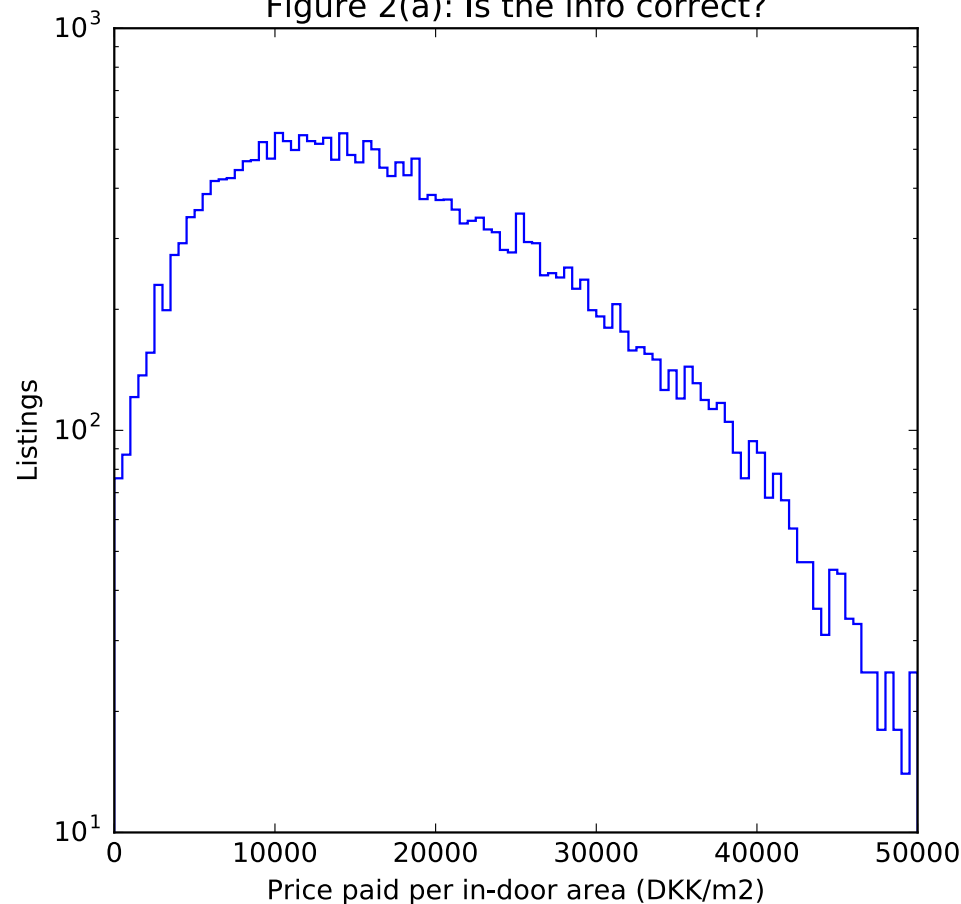
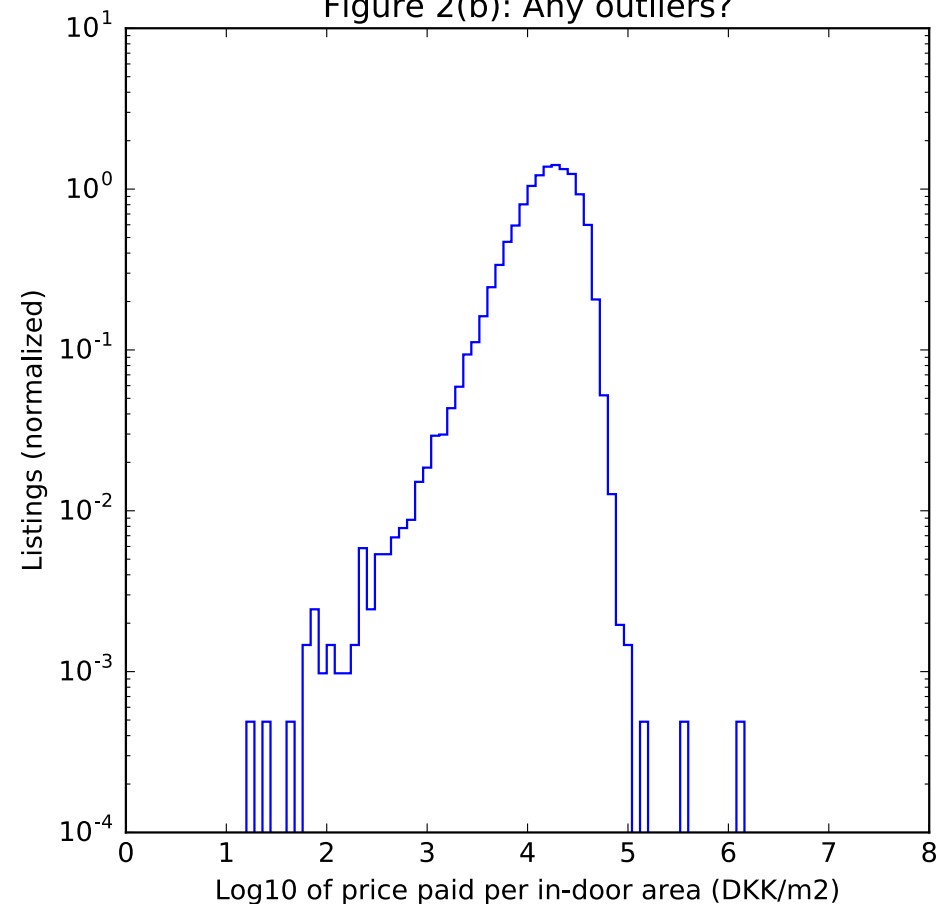


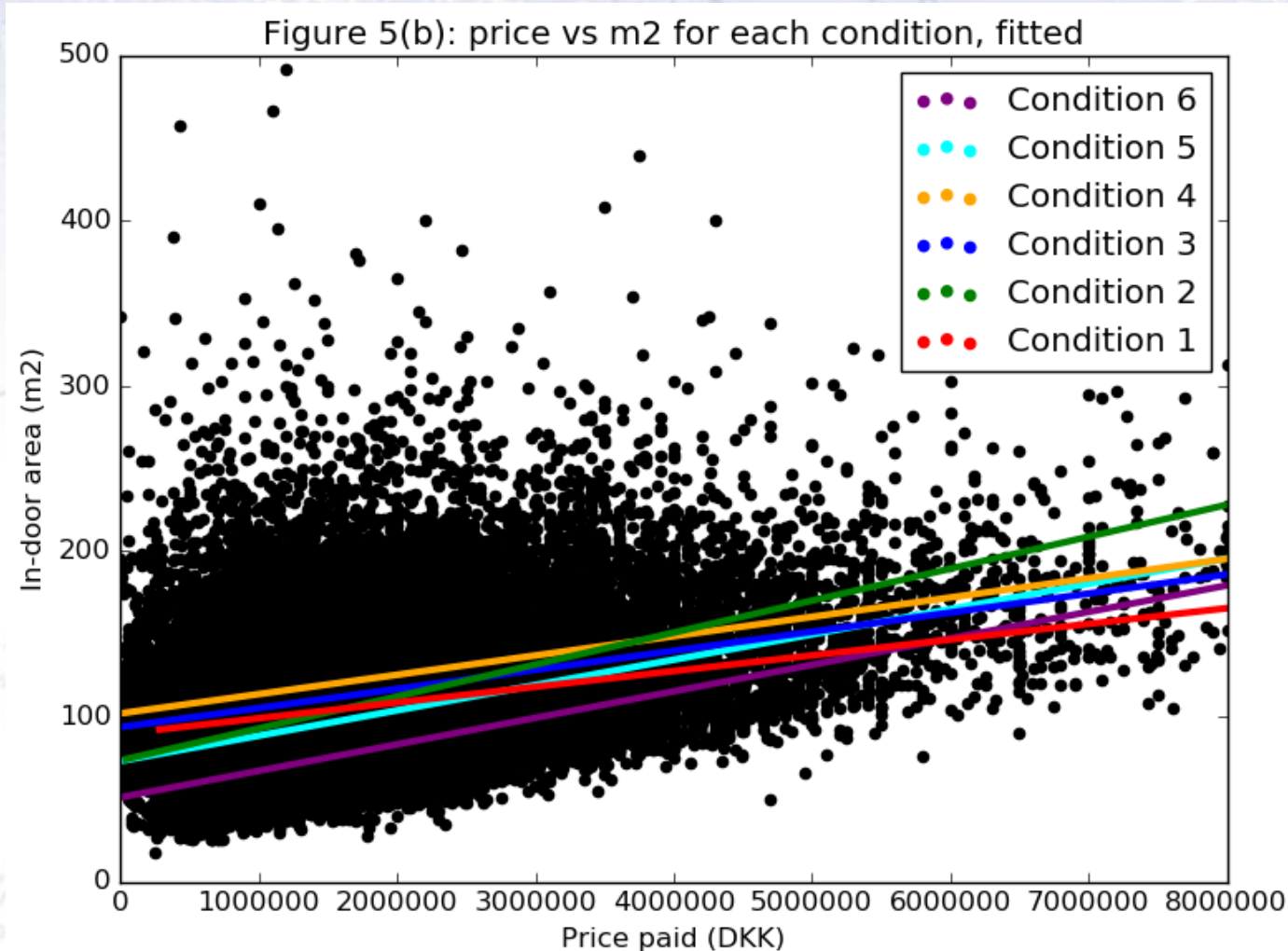
Figure 2(b): Any outliers?



I don't know who paid 1.000.000+ Kr./m<sup>2</sup>, but that is not a normal value!  
Similarly, < 100 Kr./m<sup>2</sup> seems odd, and also needs further investigation.

# Price per square meter

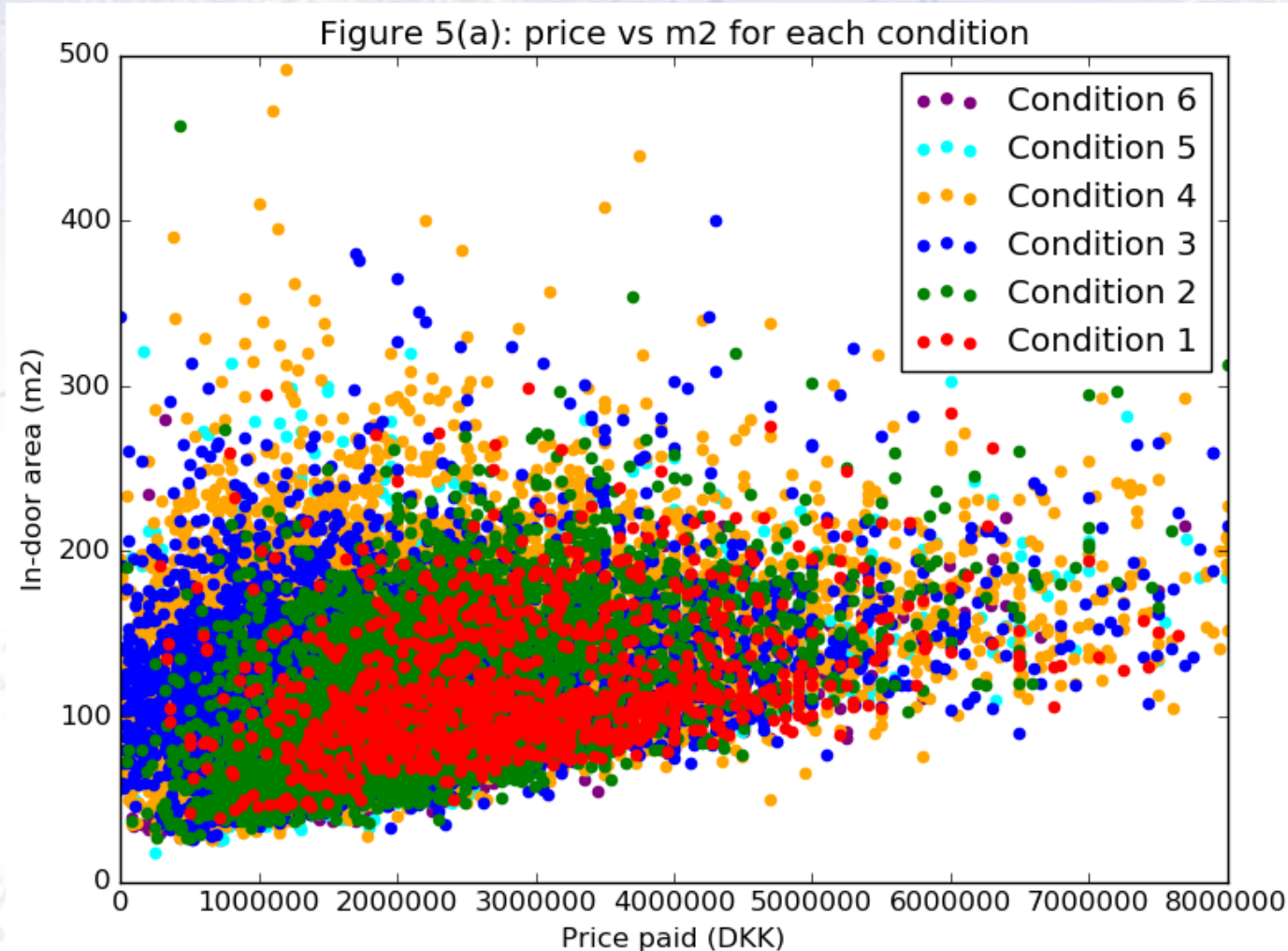
Dividing according to condition, one might expect a higher price/m<sup>2</sup>, but...



...the pattern is rather, that the basic price is higher!

# Price per square meter

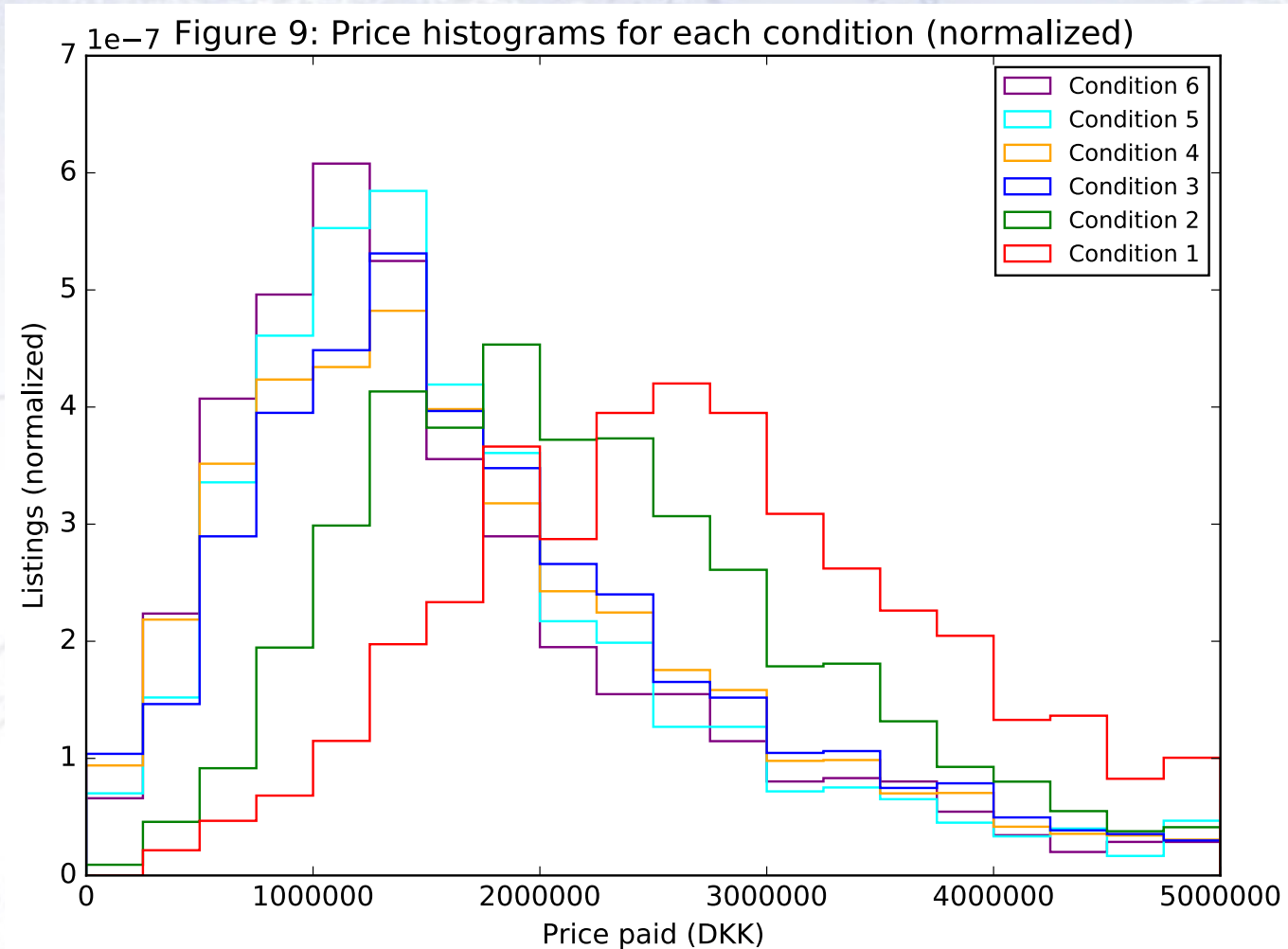
Dividing according to condition, one might expect a higher price/m<sup>2</sup>, but...



...the pattern is rather, that the basic price is higher! And condition 1 is best!!!

# Price per square meter

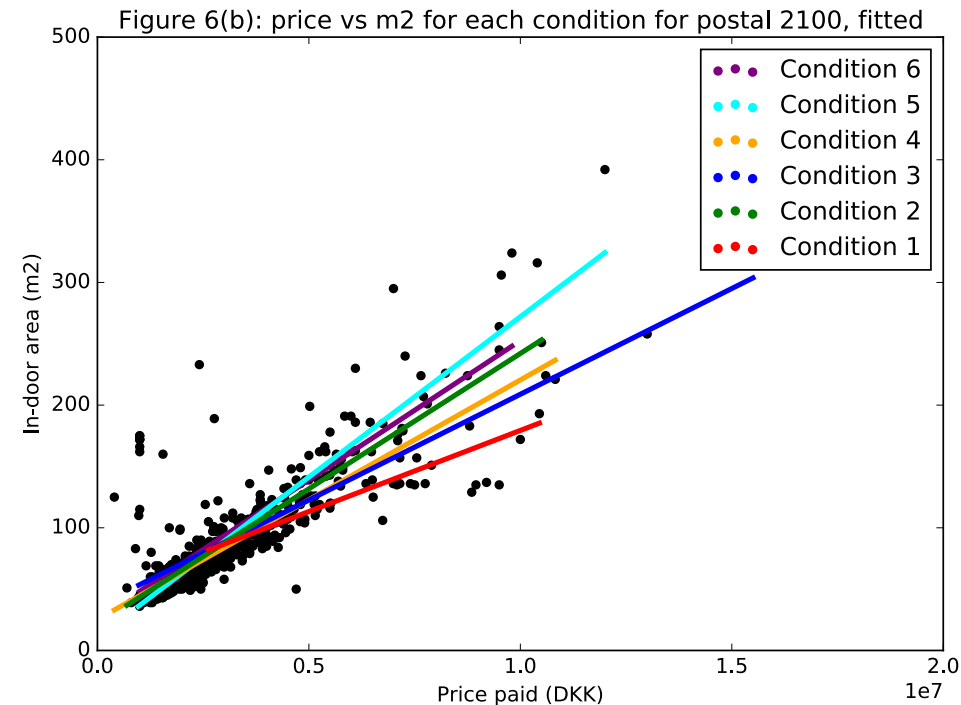
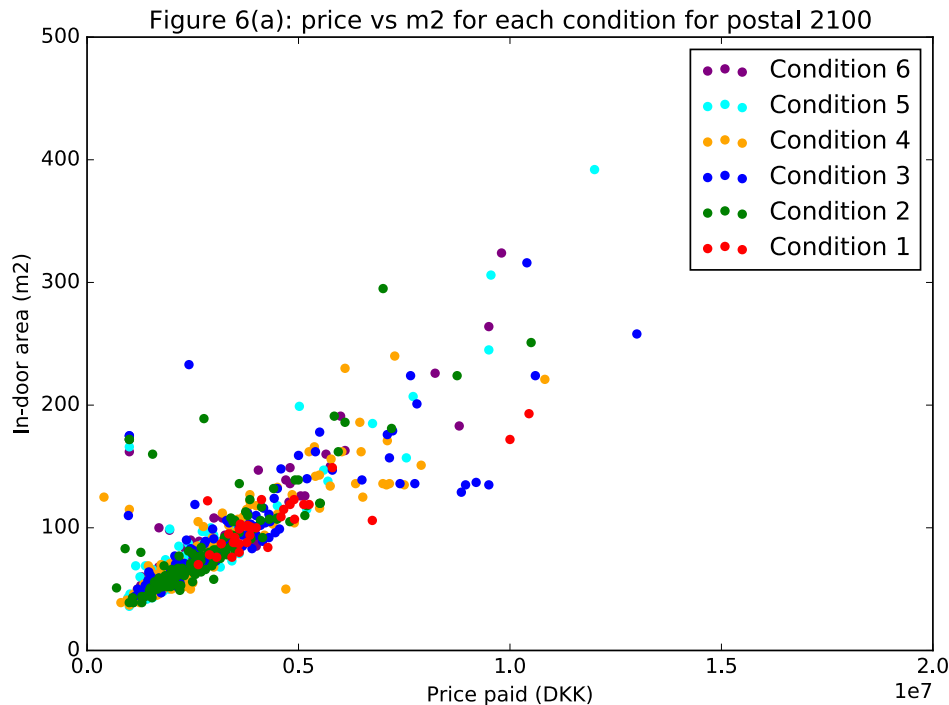
Dividing according to condition, one might expect a higher price/m<sup>2</sup>, but...



...the pattern is rather, that the basic price is higher! And condition 1 is best!!!

# Considering Østerbro only

If we restrict ourselves to Østerbro, the pattern suddenly becomes more clear:



The number of square meters suddenly become a much better indicator, and a condition suddenly also becomes a better variable.

So clearly, district/postal code is also a factor, as should be no surprise.

# Comparing districts

Now we consider the various postal codes (Østerbro, Nørrebro og Amager):

Figure 7(a): price vs m2 for some postal codes

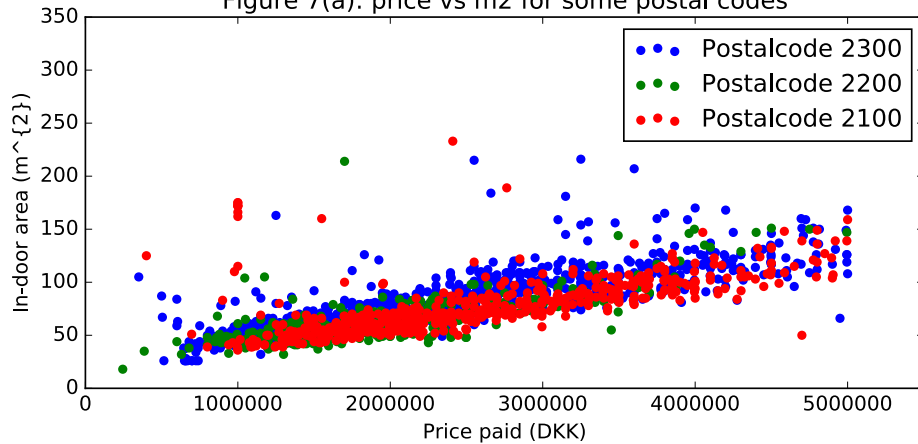


Figure 7(b): Histogram of prices for some postal codes

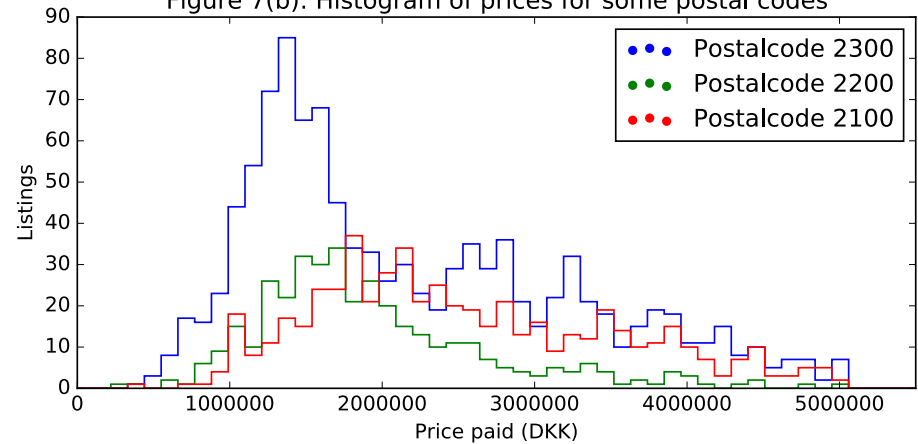


Figure 7(c): Histogram of sizes for some postal codes

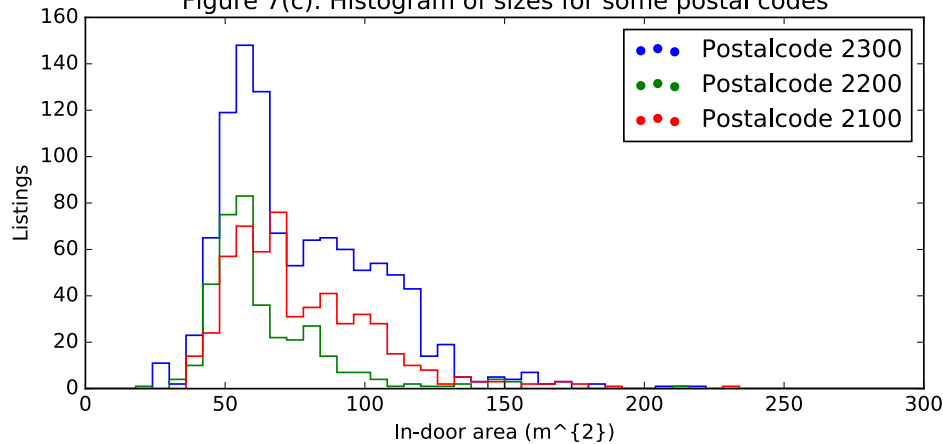
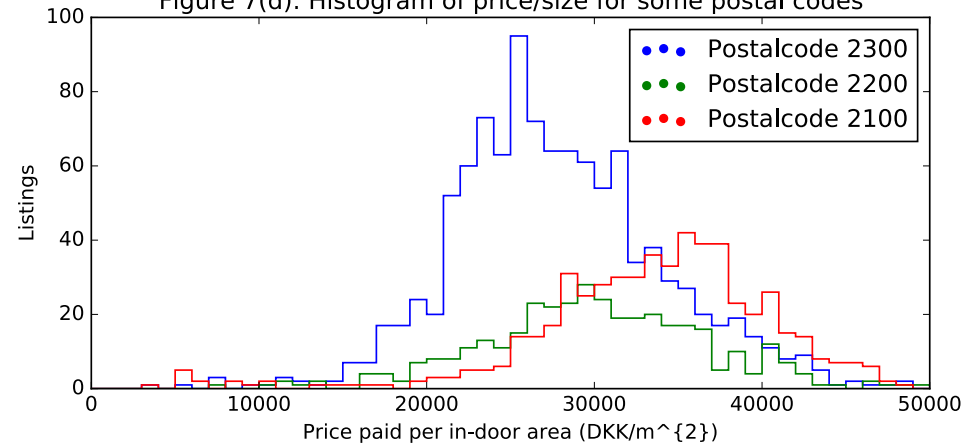


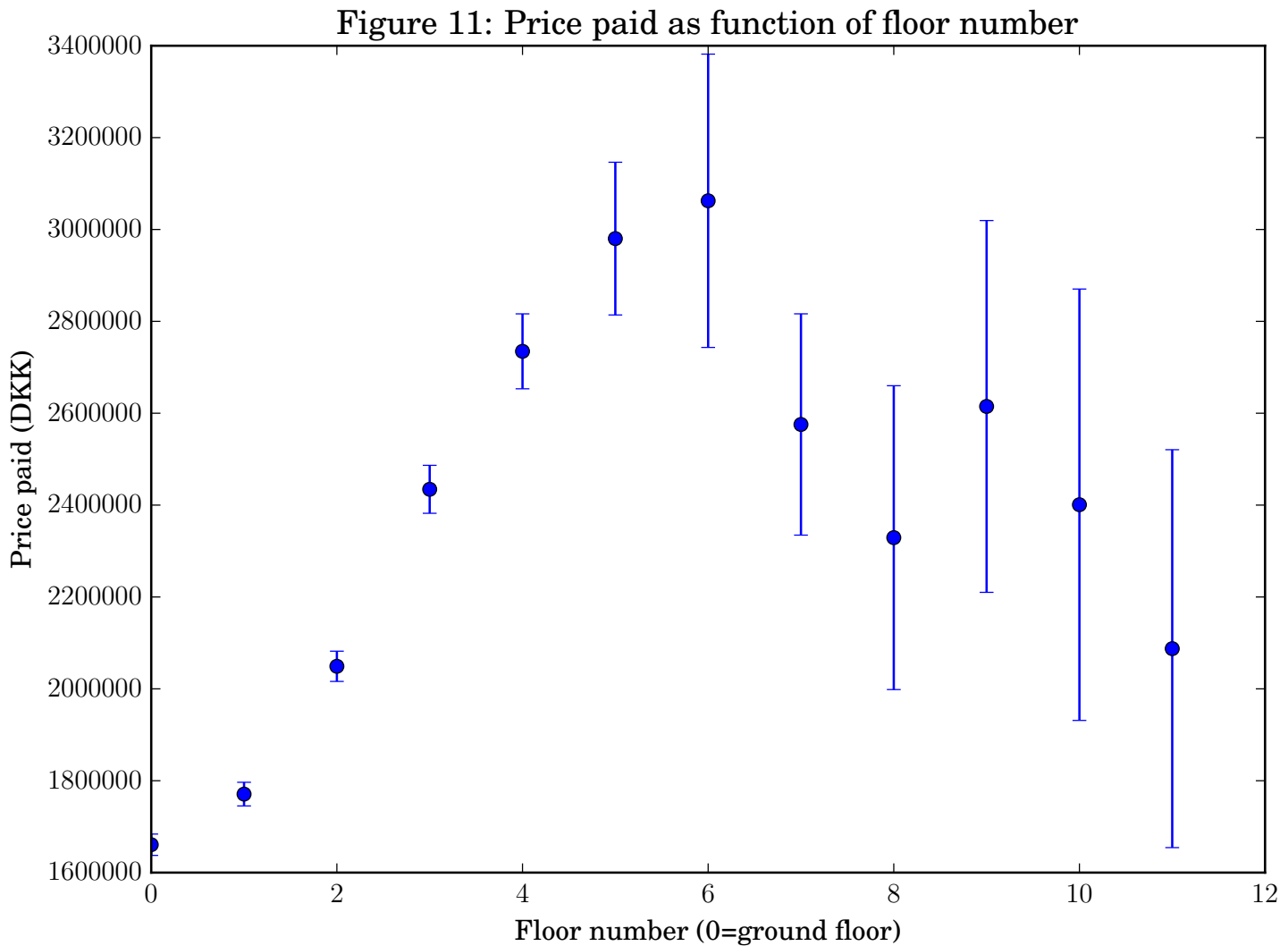
Figure 7(d): Histogram of price/size for some postal codes



Amager has small apartments and lower price/m<sup>2</sup>, and the linear model (price = price/m<sup>2</sup> \* size) holds OK for each district.

# Floor vs. price

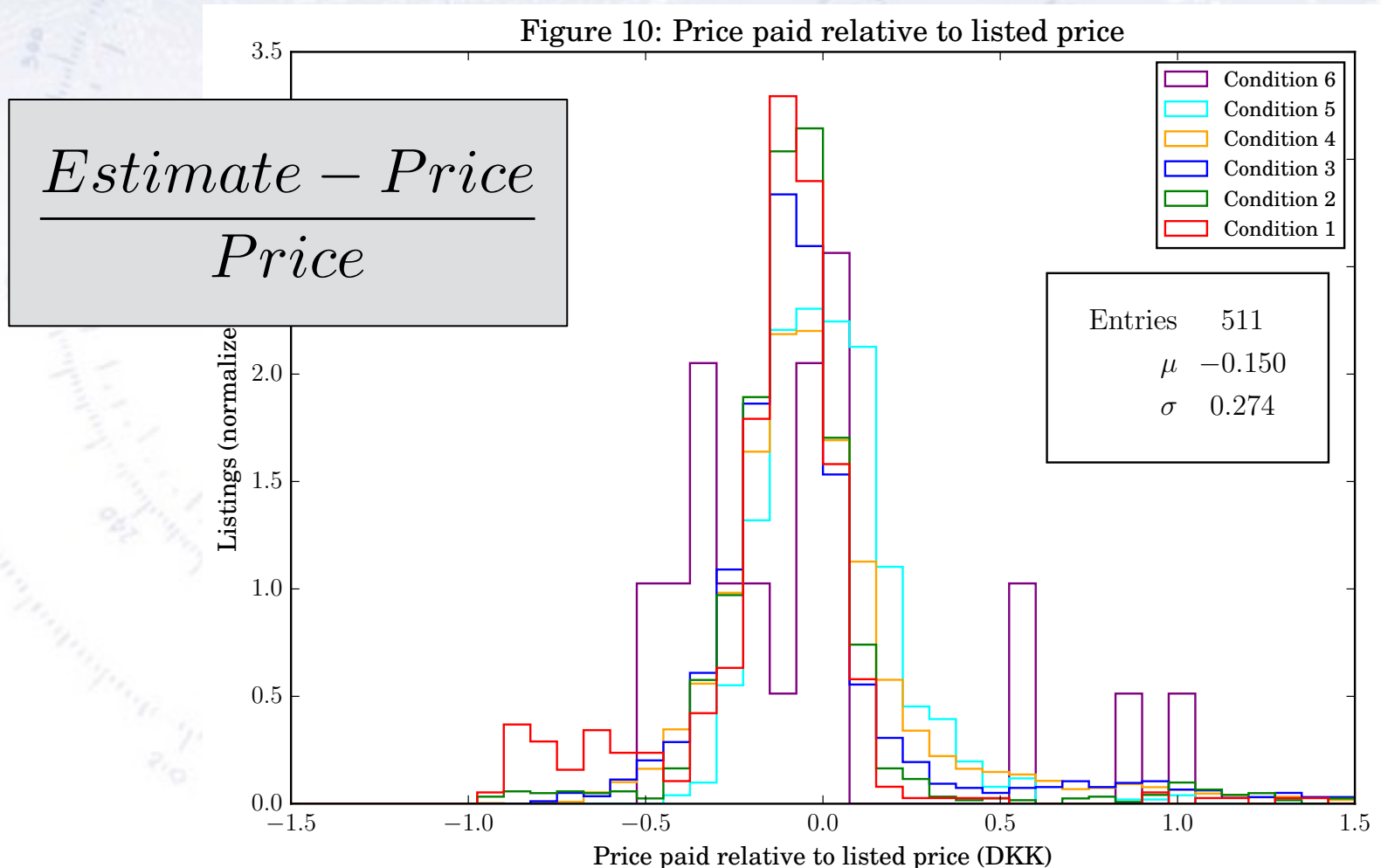
One can continue with all sorts of variables, such as e.g. floor:



# A “measure-of-goodness”

Q: How do we know, that we are improving our price estimates?

A: Well, consider how close the predictions are compared to actual price.

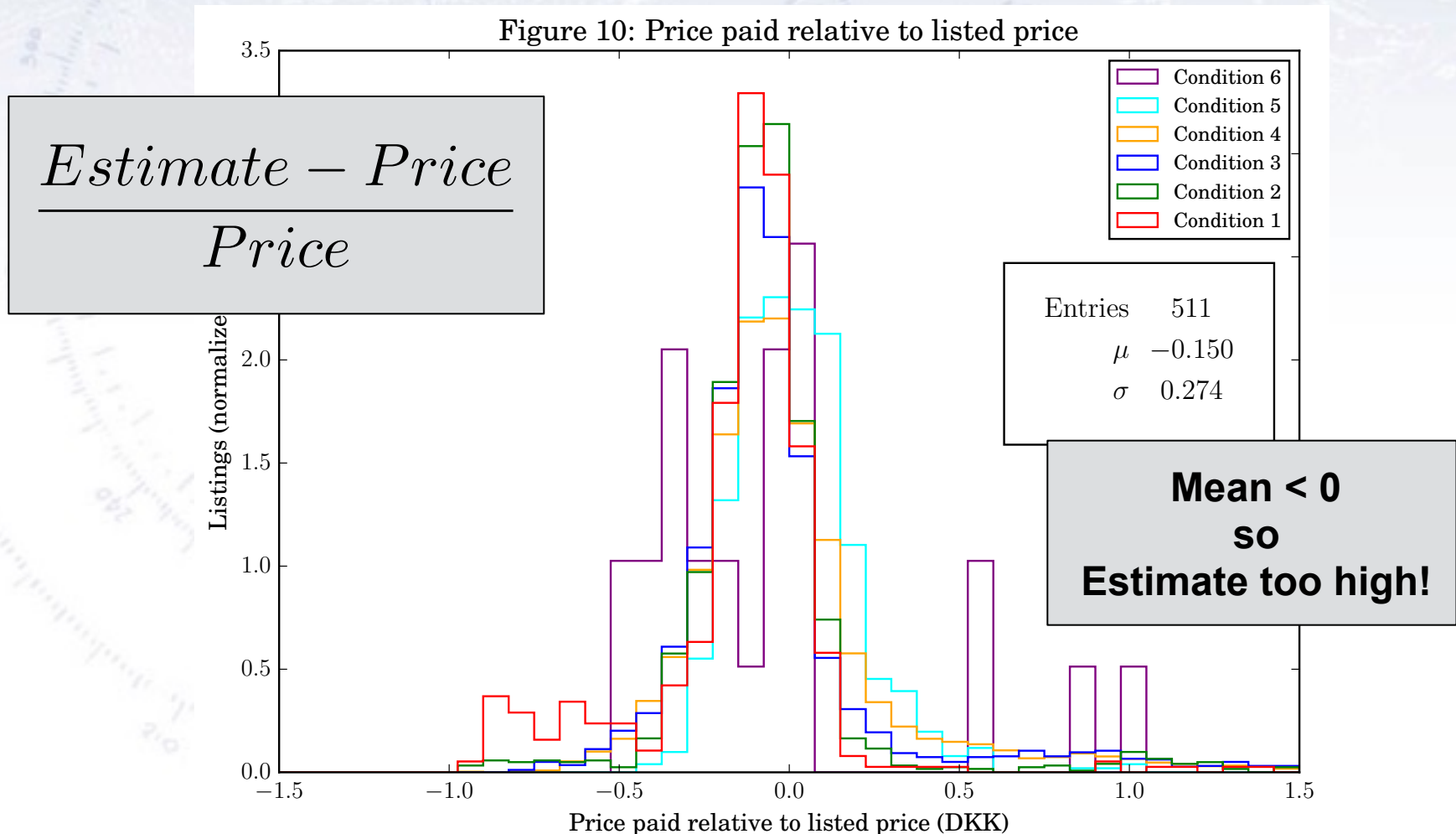




# A “measure-of-goodness”

Q: How do we know, that we are improving our price estimates?

A: Well, consider how close the predictions are compared to actual price.



# The path forward

Clearly, we could continue in this way, and produce a more and more refined model, which would give a rough estimate for most cases, but...

- The model gets more and more complicated to update or improve.
- There is no “system” by which the model can be improved.
- **The process is very manpower intensive.**

The solution is of course to use **MultiVariate Analysis (MVA)** on large datasets (which essentially is **Big Data analysis**), which in an automated and often very powerful way can combine many variables into one “optimal” prediction (or separation, if categorising).

