

Analysis of soil at India

Facundo Muñoz

October 9, 2015

Contents

1 Data description	1
2 Cost-based distances	6
3 Analysis of Calcium	7
4 Analysis of Copper	15
5 Analysis of Ferrum	23
6 Analysis of Potassium	31
7 Analysis of Magnesium	39
8 Analysis of Zinc	47
9 Conclusions	55

geoRcb package v.1.7.6

Here we compare the outcome of a classical kriging against a cost-based kriging which takes into account the presence of a barrier.

1 Data description

```
'data.frame': 70 obs. of 11 variables:  
$ sample: Factor w/ 70 levels "JIN10","JIN100",...: 28 45 66 1 16 22 25 26 27 29 ...  
$ Area  : Factor w/ 2 levels "Inside","Veranda": 1 1 1 1 1 1 1 1 1 1 ...  
$ Side   : Factor w/ 2 levels "Fireplace","Storage": 2 2 2 2 1 2 1 1 2 2 ...  
$ x      : num  13 12.5 13.5 13.5 14 13.5 14 14 14 ...  
$ y      : num  -11.5 -12 -12 -13 -11 -11.5 -10 -10.5 -13.5 -12.5 ...  
$ Ca     : num  2.95 3.4 4.3 5.7 3.97 4.5 3.3 3.2 2.81 5.07 ...  
$ Cu     : int  15 13 15 13 14 17 20 18 14 14 ...  
$ Fe     : num  1.21 1.18 1.3 1.36 1.04 1.27 1.23 1.43 1.41 1.27 ...  
$ K      : num  0.28 0.32 0.31 0.21 0.27 0.27 0.18 0.33 0.21 0.27 ...  
$ Mg     : num  0.54 0.51 0.57 0.55 0.5 0.53 0.49 0.56 0.45 0.6 ...  
$ Zn     : int  46 31 36 28 32 33 36 41 36 35 ...
```

sample	Area	Side	x	y
JIN10 : 1	Inside :47	Fireplace:36	Min. : 7.50	Min. :-13.50
JIN100 : 1	Veranda:23	Storage :34	1st Qu.: 9.50	1st Qu.:-12.00
JIN101 : 1			Median :11.50	Median :-11.00
JIN102 : 1			Mean :11.53	Mean :-10.99
JIN103 : 1			3rd Qu.:13.50	3rd Qu.:-10.00
JIN106 : 1			Max. :16.00	Max. :-8.00
(Other):64				
Ca	Cu	Fe	K	
Min. :0.660	Min. : 7.00	Min. :0.810	Min. :0.1400	
1st Qu.:2.490	1st Qu.:11.00	1st Qu.:1.050	1st Qu.:0.2300	
Median :2.945	Median :12.00	Median :1.190	Median :0.2600	
Mean :3.106	Mean :12.39	Mean :1.175	Mean :0.2679	
3rd Qu.:3.697	3rd Qu.:14.00	3rd Qu.:1.300	3rd Qu.:0.2900	
Max. :5.700	Max. :20.00	Max. :1.500	Max. :0.6300	
Mg	Zn			
Min. :0.220	Min. :17.00			
1st Qu.:0.390	1st Qu.:26.00			
Median :0.460	Median :32.00			
Mean :0.453	Mean :31.77			
3rd Qu.:0.520	3rd Qu.:35.75			
Max. :0.660	Max. :62.00			

Figures 1 and 2 display the raw data, and an exploratory smoothed surface.

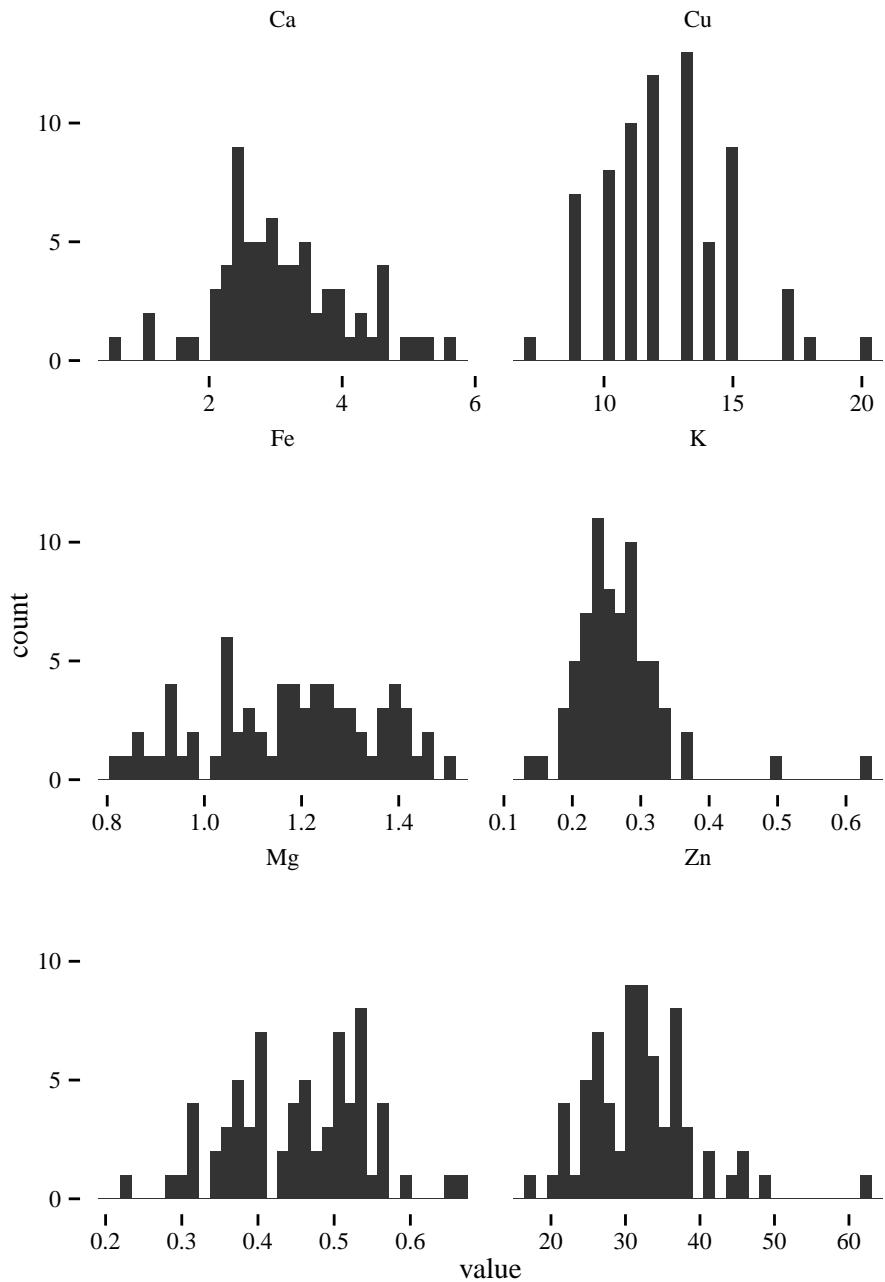


Figure 1: Histograms of measured variables.

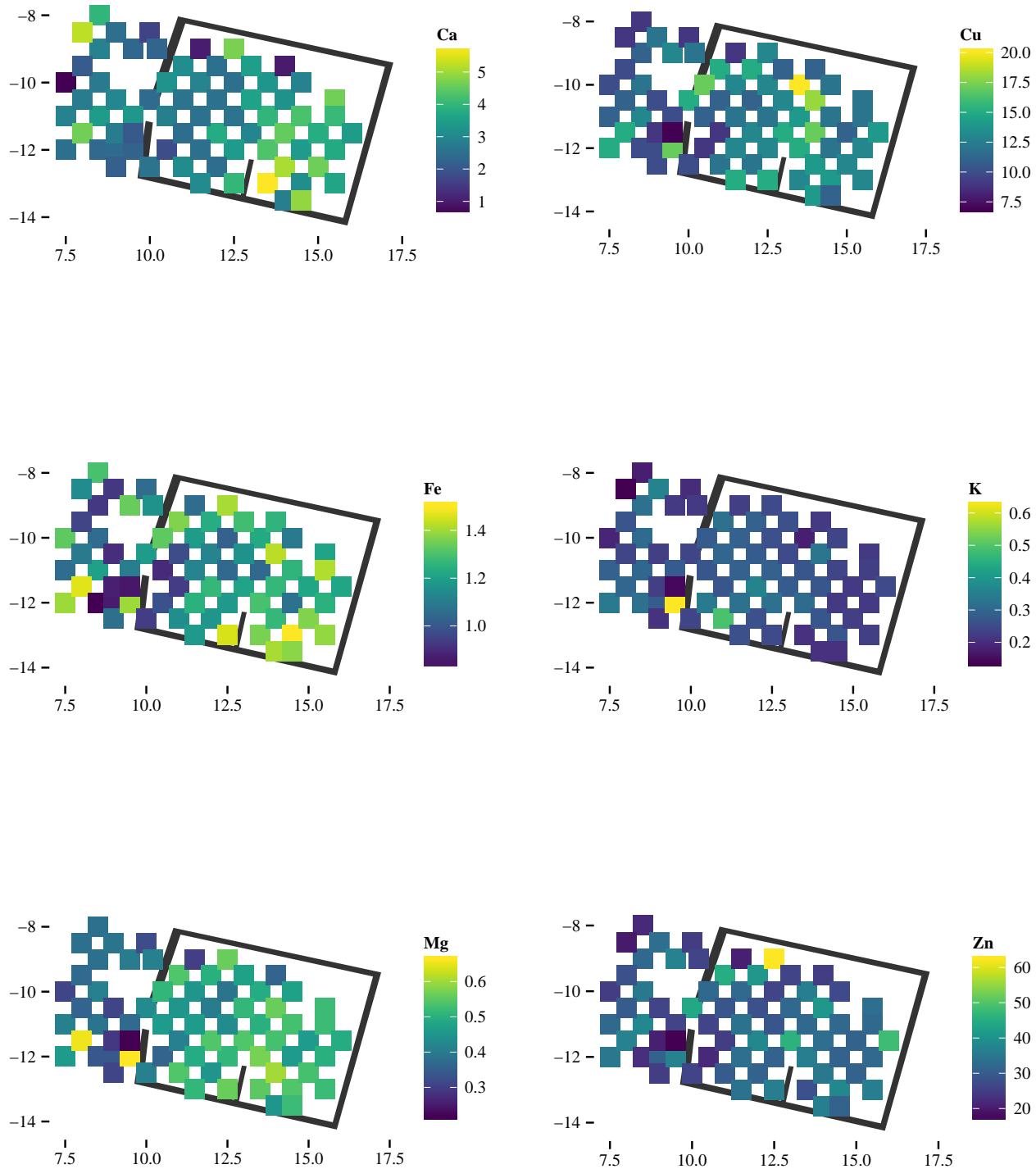


Figure 2: Measurement locations and observed values

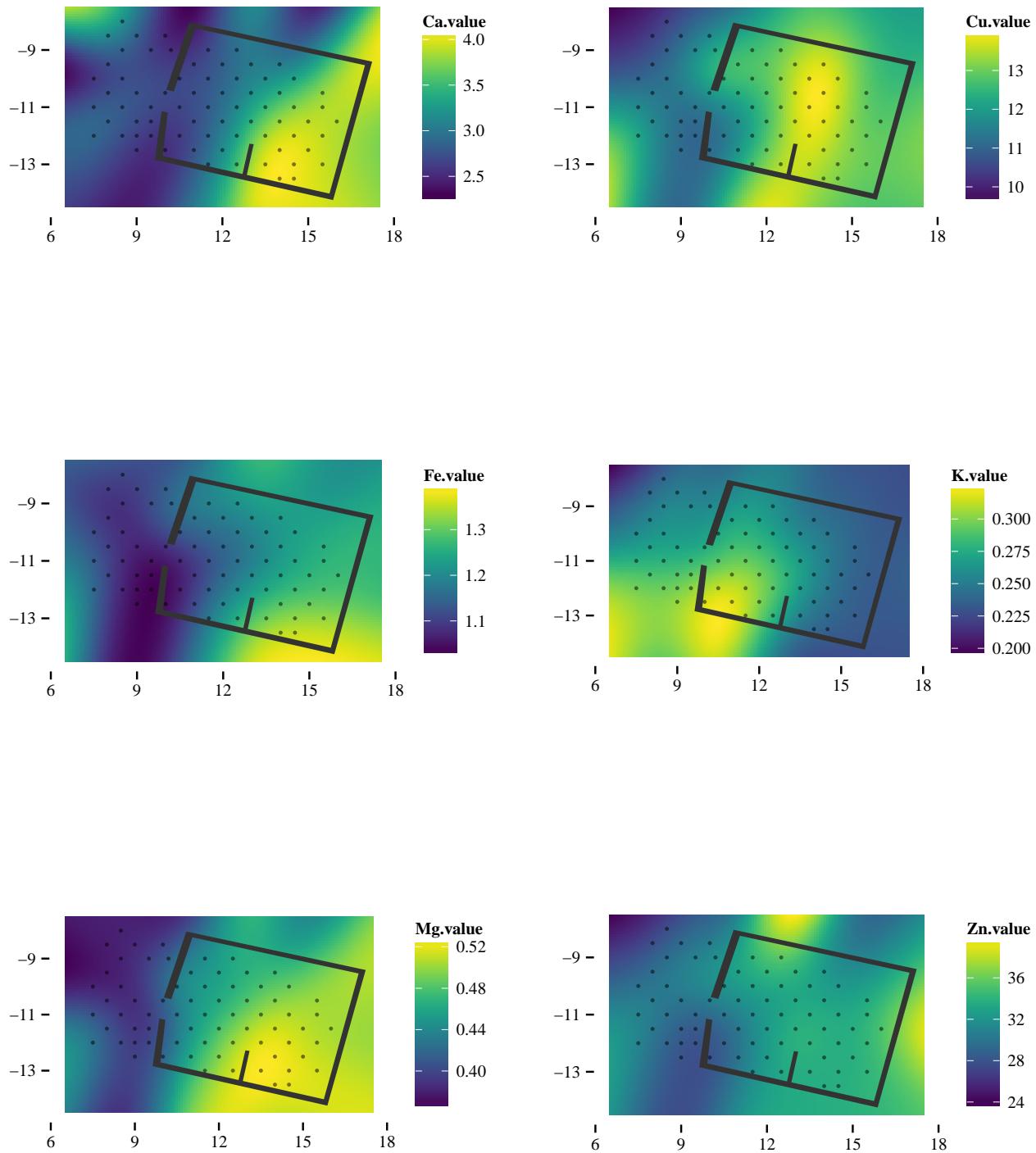


Figure 3: Exploratory kernel smoothing of the measurements

2 Cost-based distances

Here we set up the cost-based surface, and compute some cost-based maps, for verifications purposes.

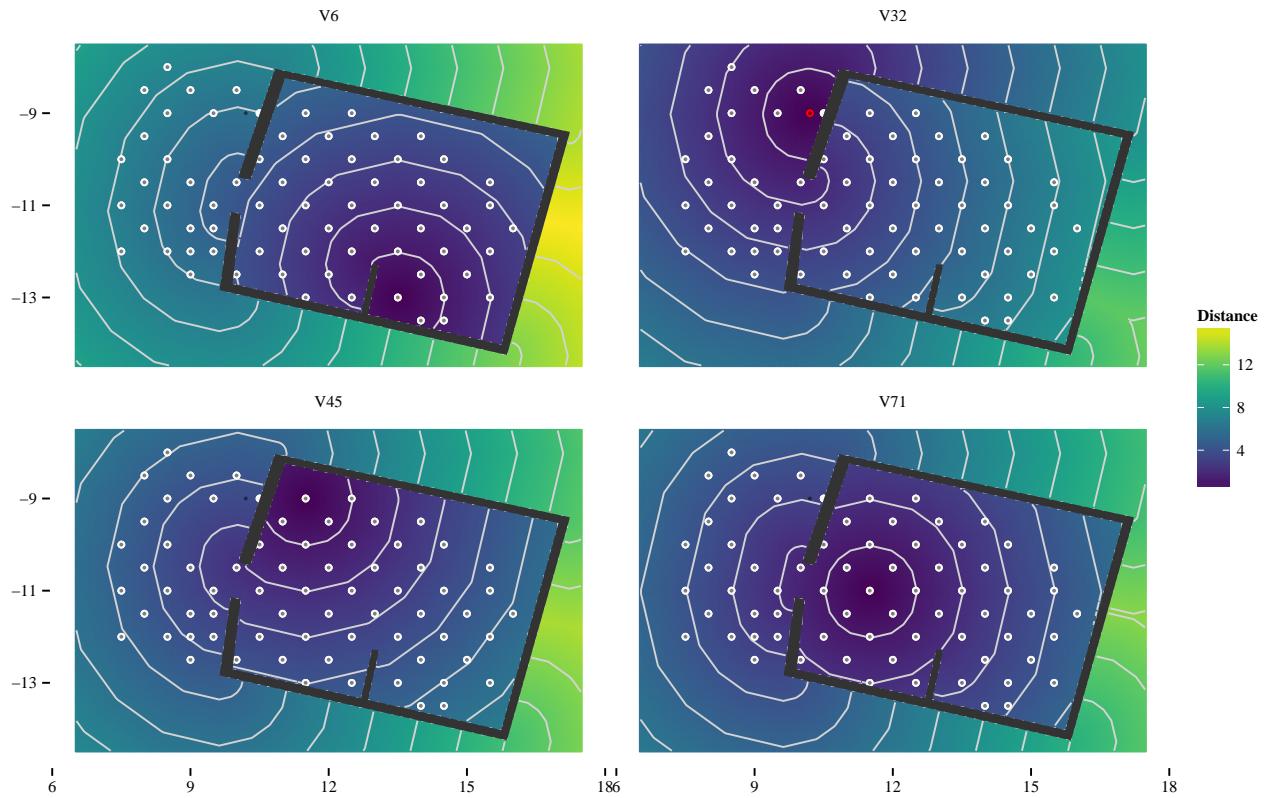


Figure 4: Some cost-based maps to selected observations.

3 Analysis of Calcium

3.1 Euclidean kriging

The variogram model is Exponential. We choose to estimate the nugget effect, which may account for measurement error, for example.

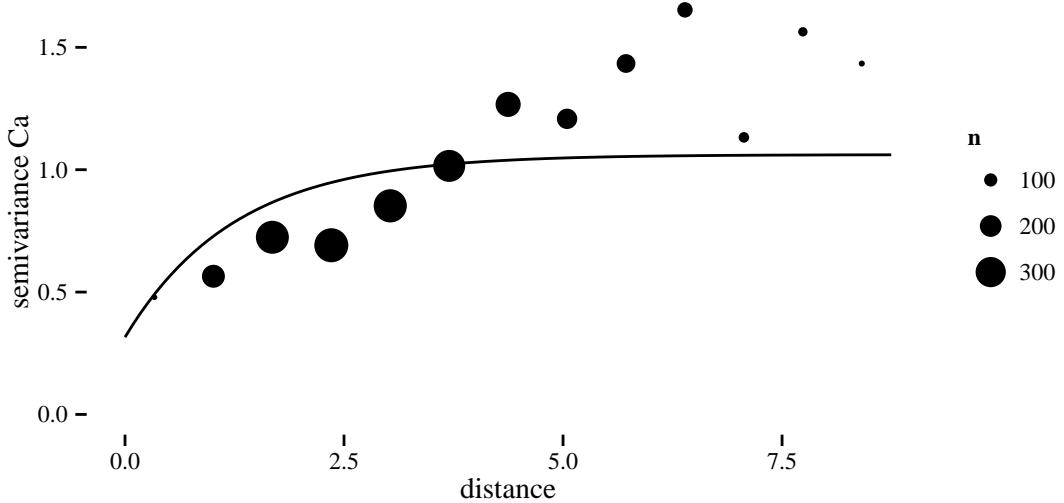


Figure 5: Empirical variogram and fitted model.

3.2 Cost-based kriging

3.3 Comparison of method outcomes

	Euclidean	Cost_based
Intercept	3.12	3.17
Nugget	0.32	0.60
Partial sill	0.75	0.85
phi	1.25	6.53
Pract. range	3.75	19.56
Log-likelihood	-89.25	-89.82

In the scatter plot, the horizontal patterns correspond to predictions on observed values. Otherwise, the differences are negligible.

Near the observations, the cost-based approach has a larger prediction error due to its increased estimation of the nugget (i.e. short-range variance). In the main area, the prediction errors are practically the same with both approaches. Behind the walls, the Euclidean prediction error is unrealistically low.

3.4 Leave-one-out Cross Validation (LOOCV)

method	rmse(error)
std	0.93

method	rmse(error)
cst	0.96

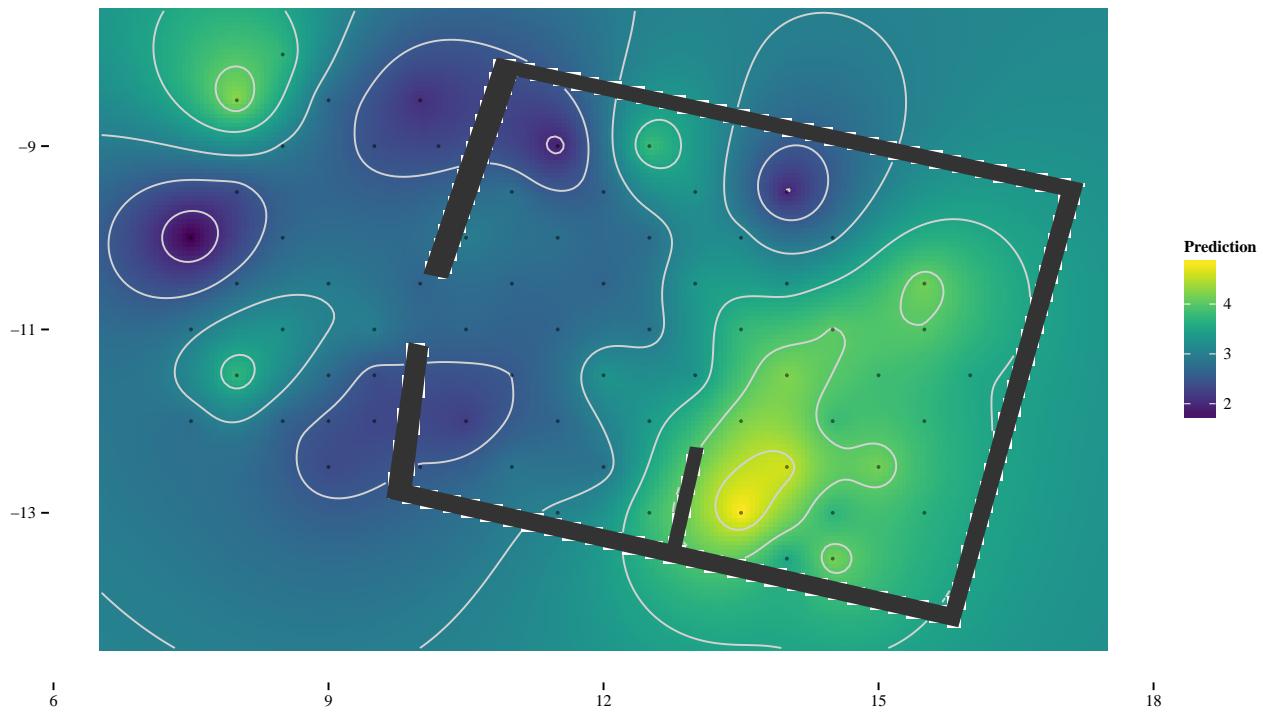


Figure 6: Euclidean kriging prediction

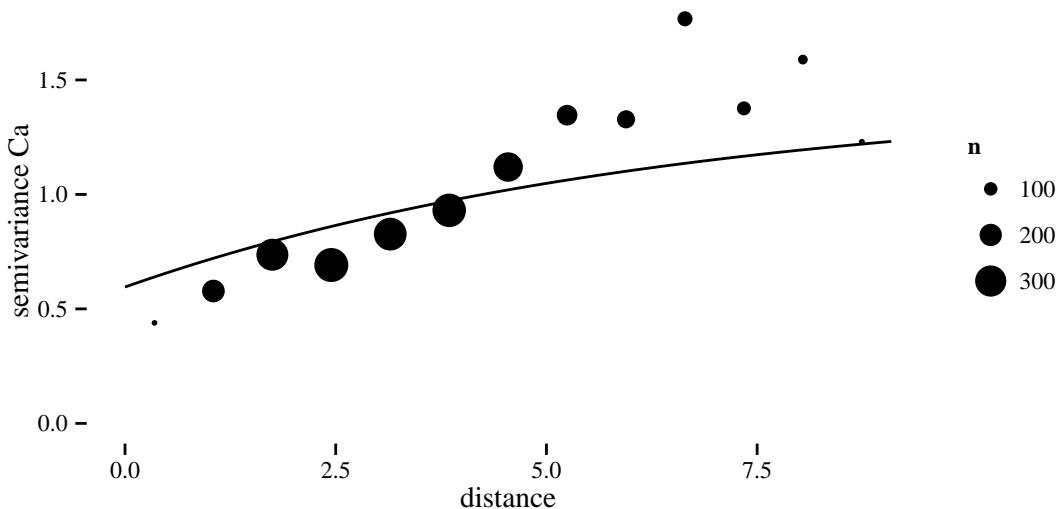


Figure 7: Empirical cost-based variogram and fitted model.

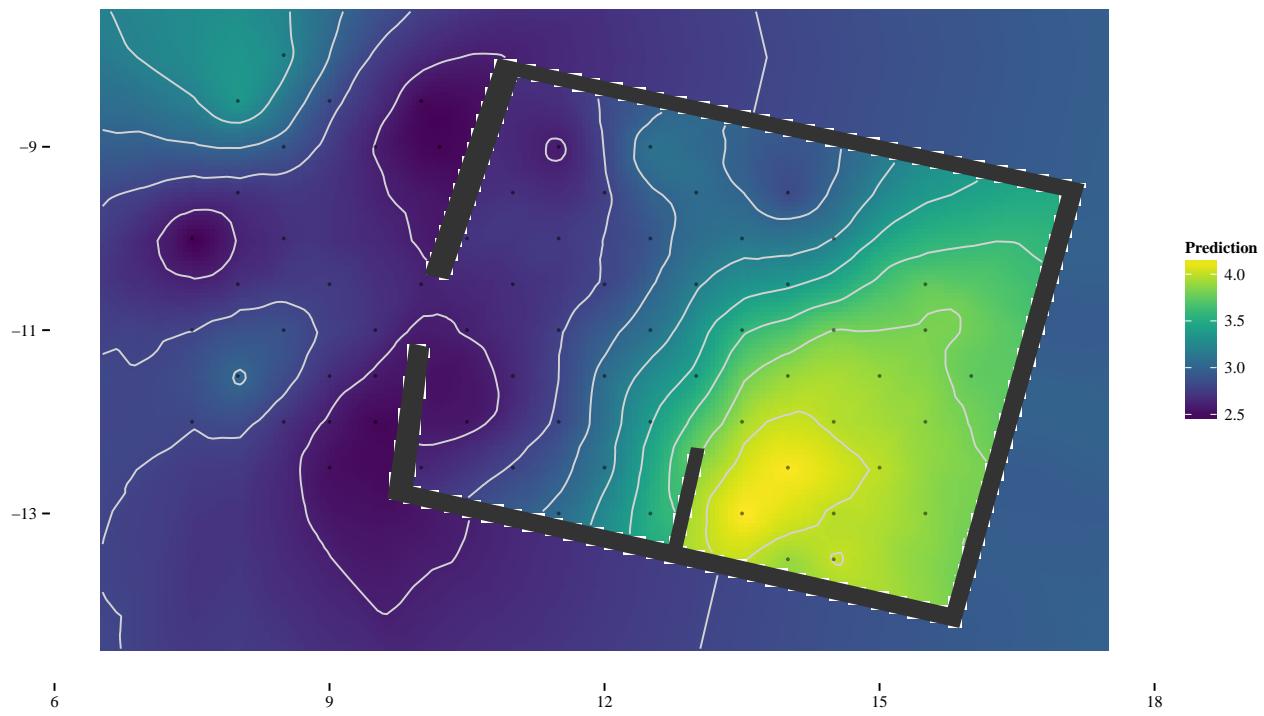


Figure 8: Cost-based kriging prediction

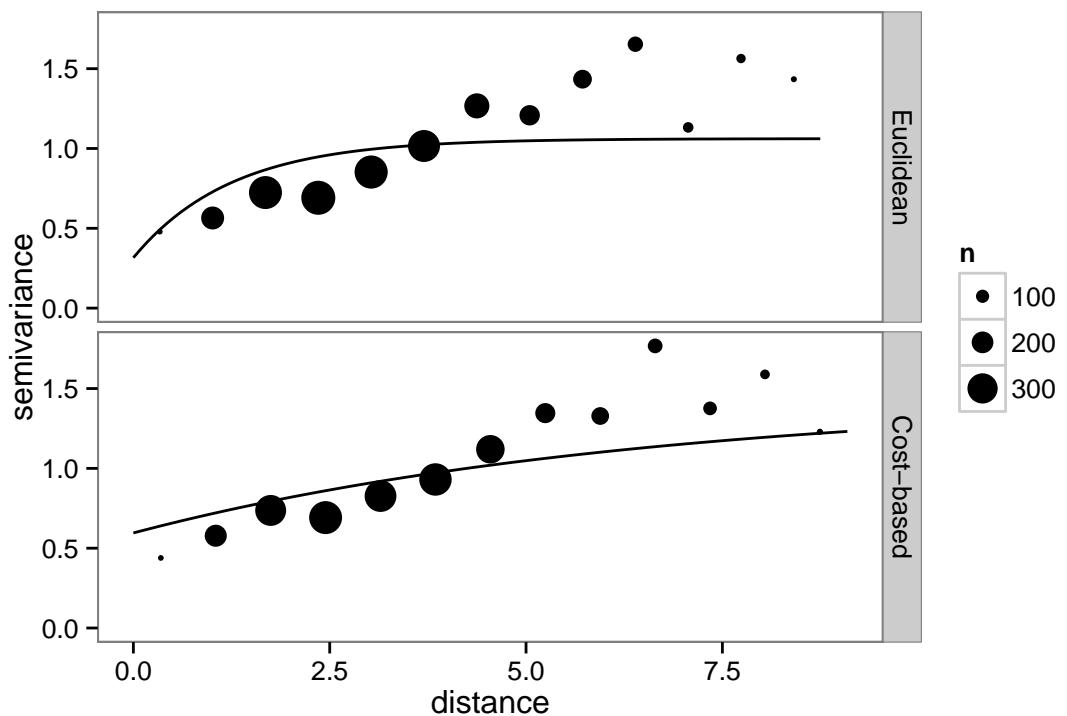


Figure 9: Empirical variogram and fitted models by method for Calcium.

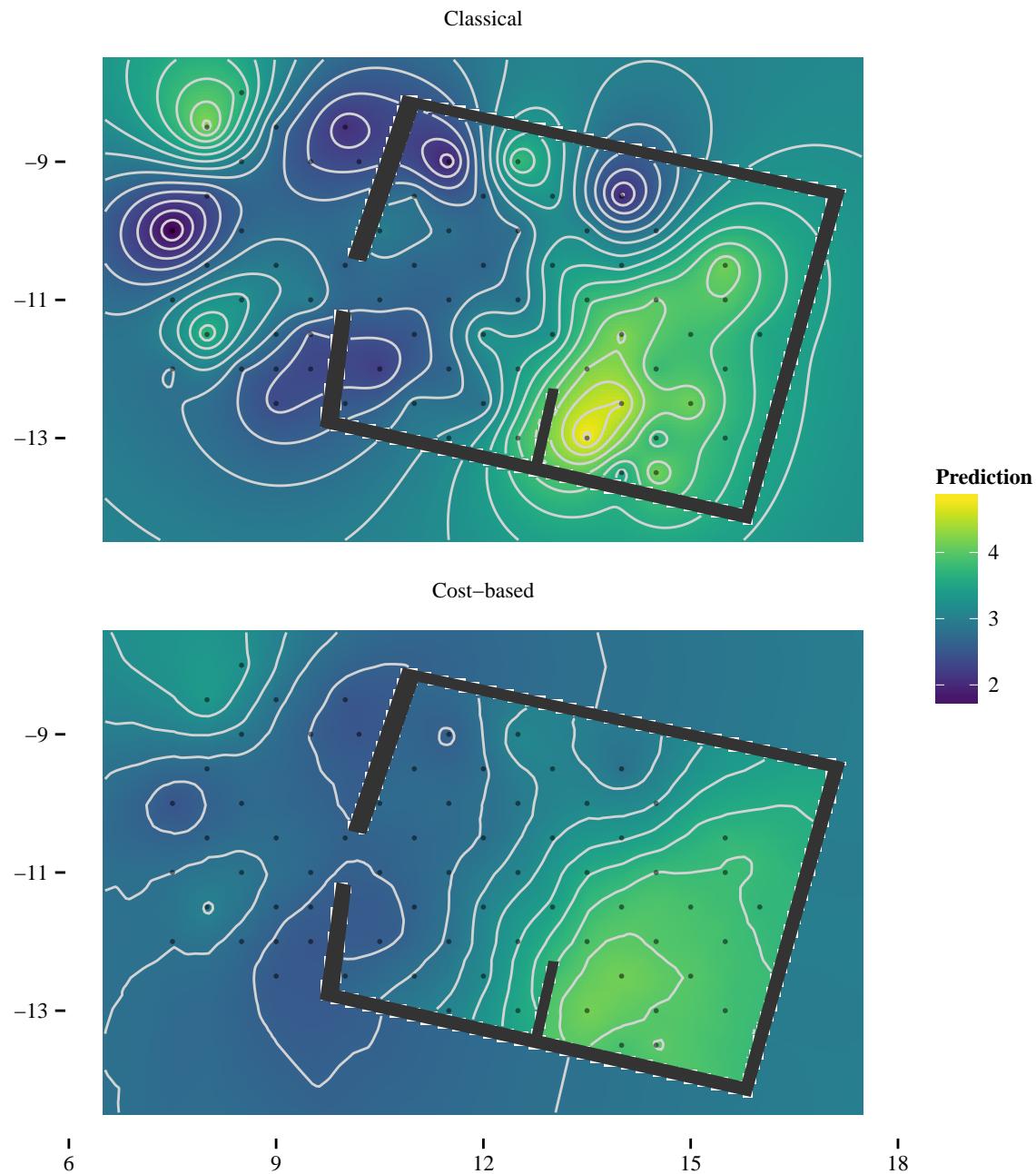


Figure 10: Comparison of Kriging estimates.

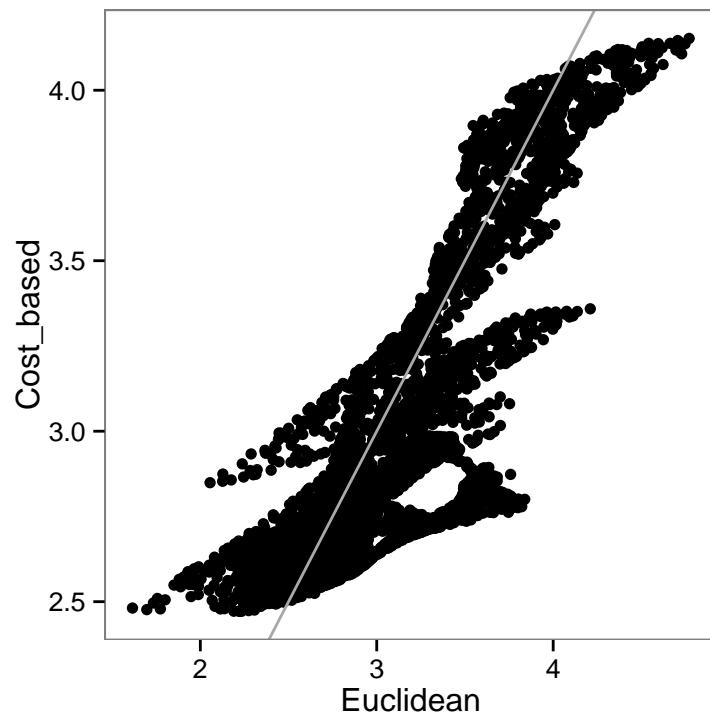


Figure 11: Pointwise comparison of predictions by method.

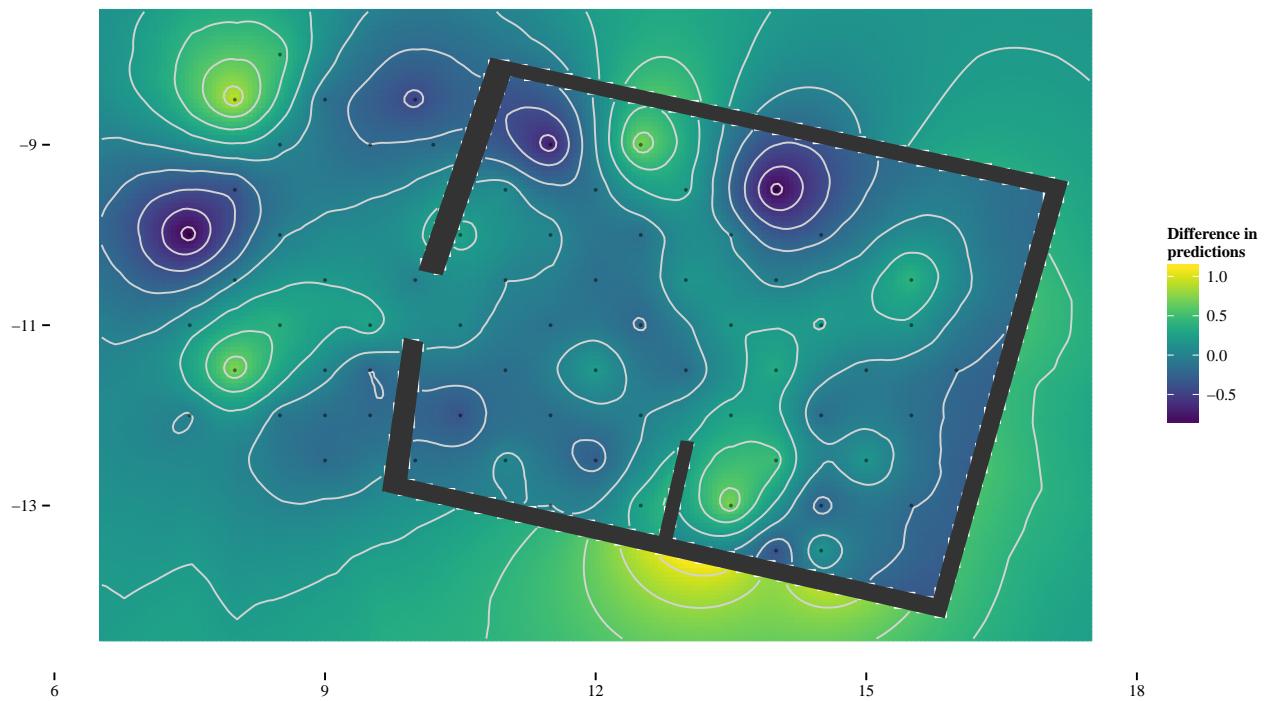


Figure 12: Difference between the Euclidean and the cost-based predictions.

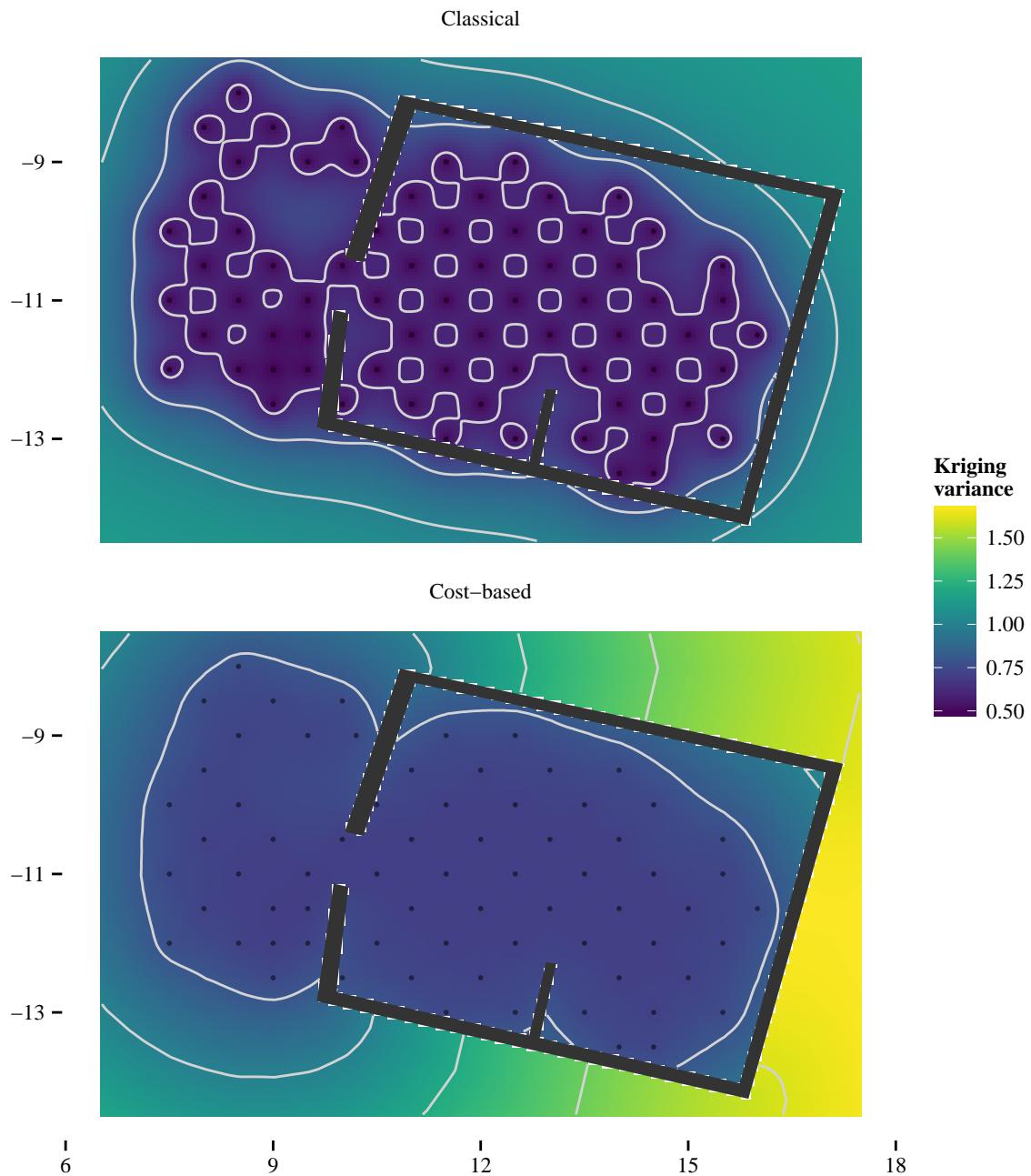


Figure 13: Comparison of prediction error by method.

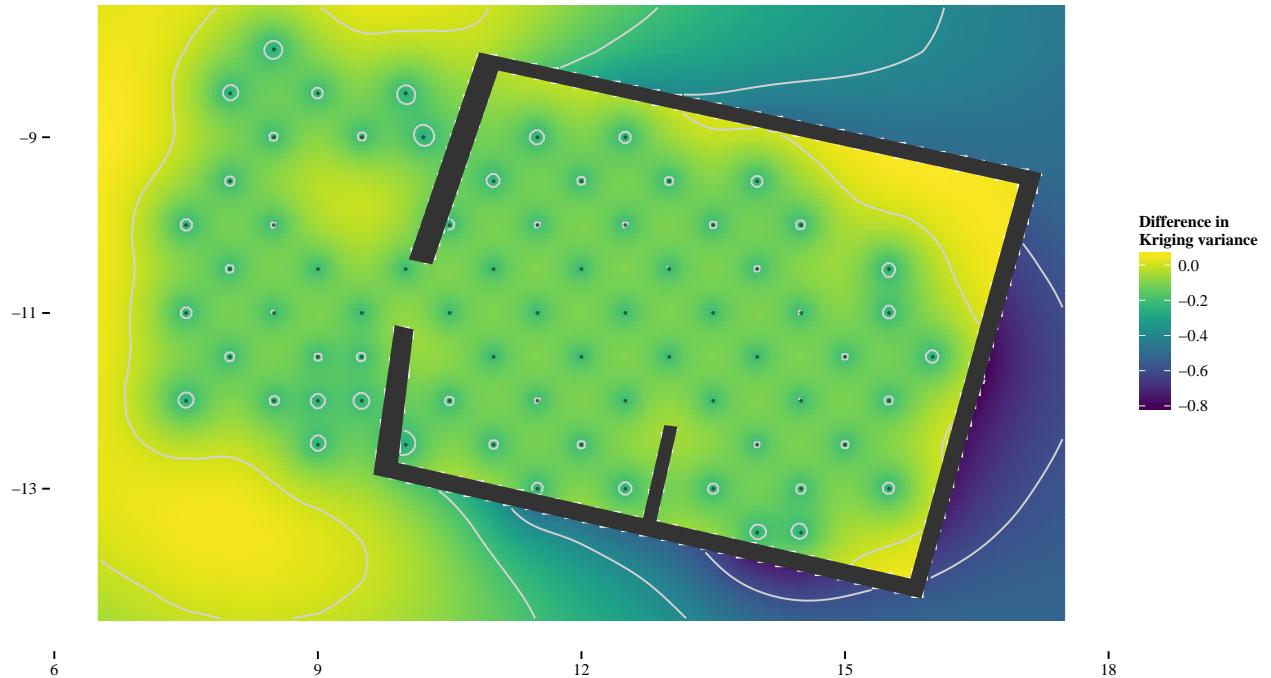


Figure 14: Difference between the Euclidean and the cost-based prediction errors

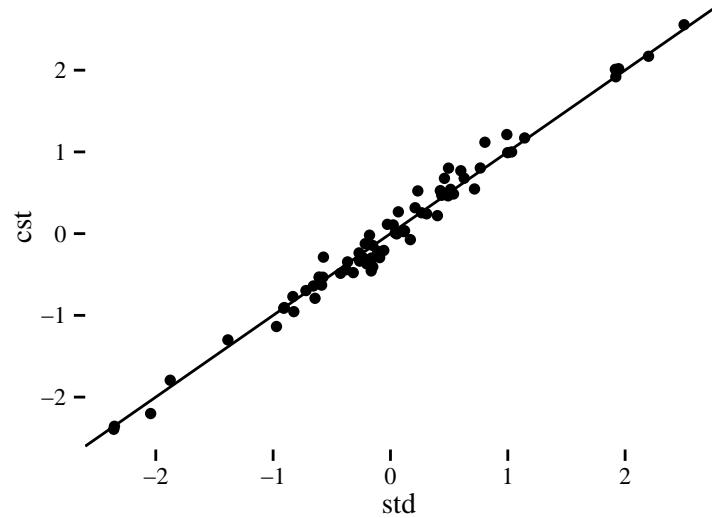


Figure 15: Pointwise leave-one-out prediction error by method.

4 Analysis of Copper

4.1 Euclidean kriging

The variogram model is Exponential. We choose to estimate the nugget effect, which may account for measurement error, for example.

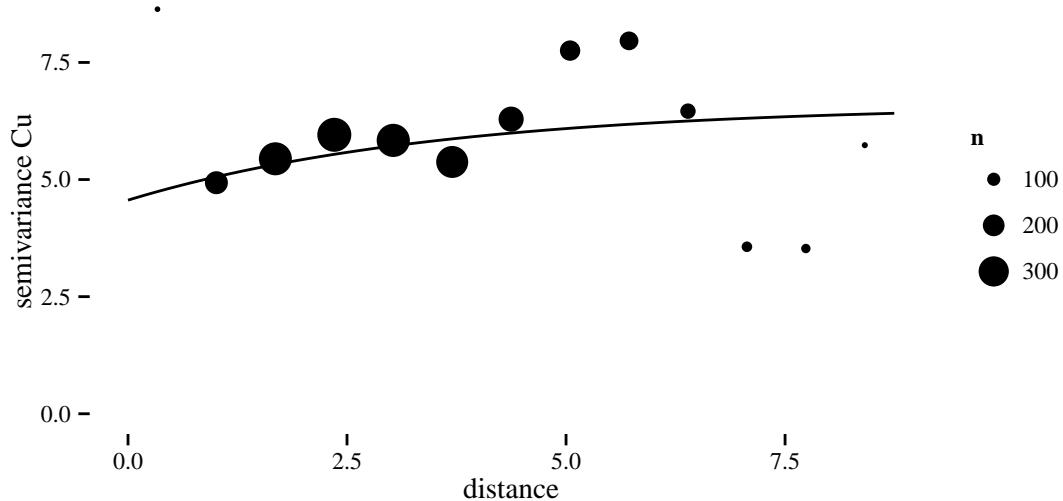


Figure 16: Empirical variogram and fitted model for Copper.

4.2 Cost-based kriging

4.3 Comparison of method outcomes

	Euclidean	Cost_based
Intercept	12.13	12.14
Nugget	4.56	4.58
Partial sill	2.03	1.98
kappa	0.51	0.51
phi	3.58	3.81
Pract. range	10.73	11.43

In the scatter plot, the horizontal patterns correspond to predictions on observed values. Otherwise, the differences are negligible.

Near the observations, the cost-based approach has a larger prediction error due to its increased estimation of the nugget (i.e. short-range variance). In the main area, the prediction errors are practically the same with both approaches. Behind the walls, the Euclidean prediction error is unrealistically low.

4.4 Leave-one-out Cross Validation (LOOCV)

method	rmse(error)
std	2.42

method	rmse(error)
cst	2.41

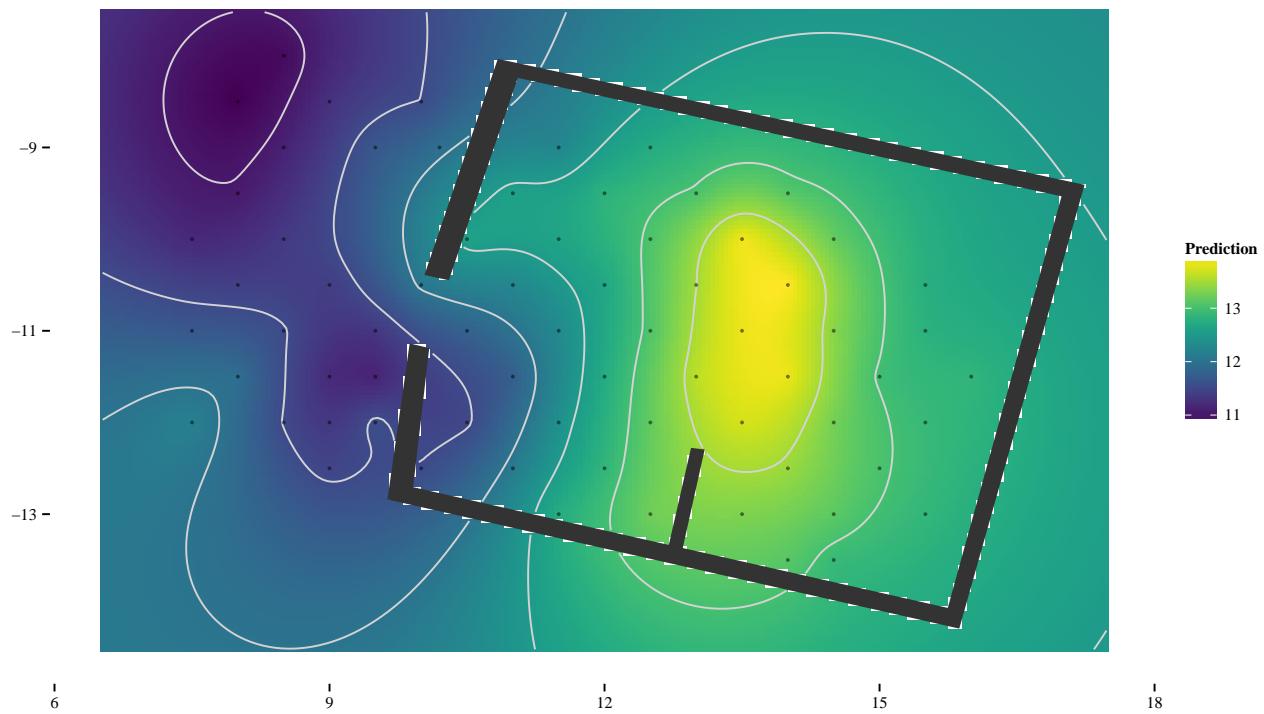


Figure 17: Euclidean kriging prediction for Copper.

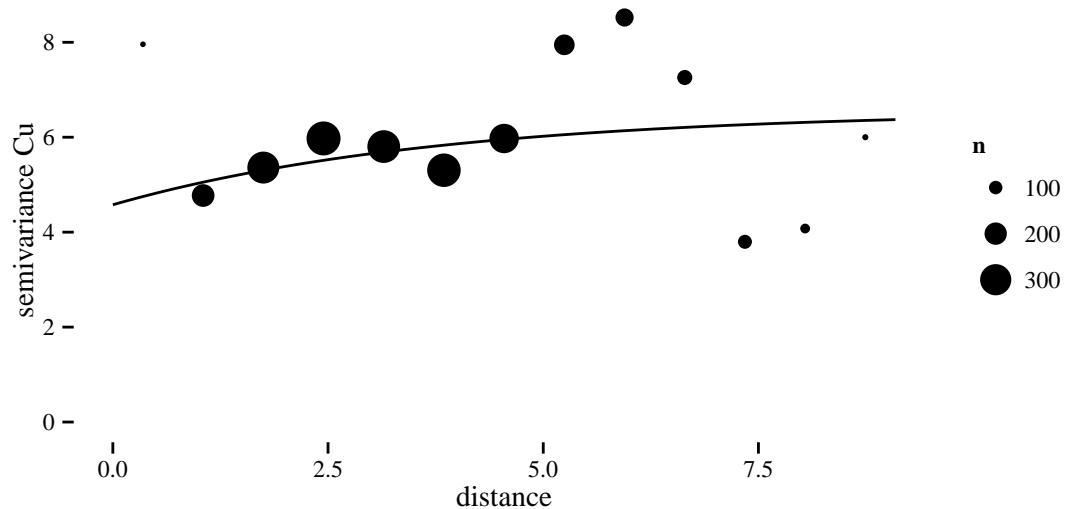


Figure 18: Empirical cost-based variogram and fitted model.

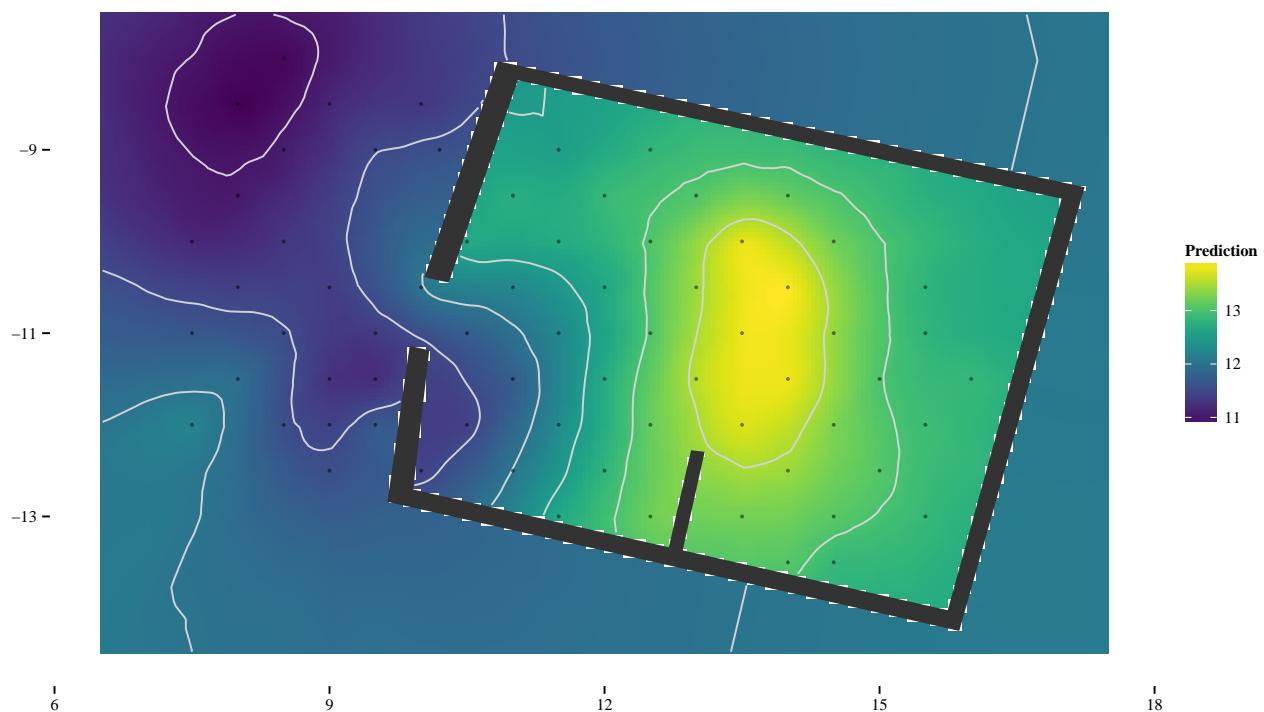


Figure 19: Cost-based kriging prediction

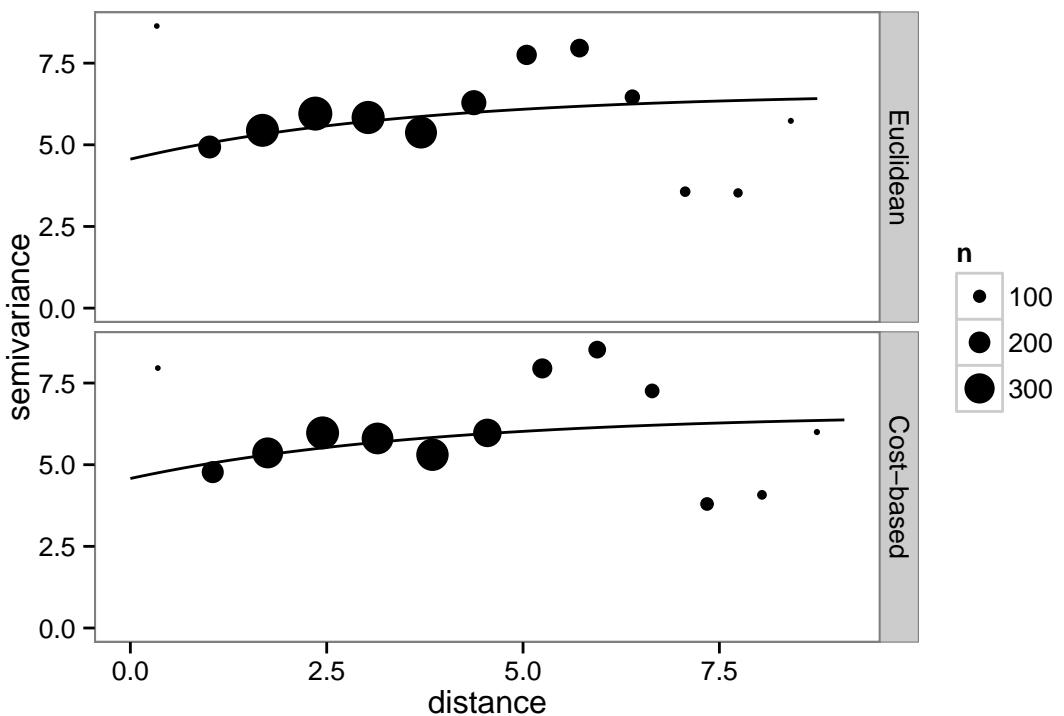


Figure 20: Empirical variogram and fitted models by method for Copper.

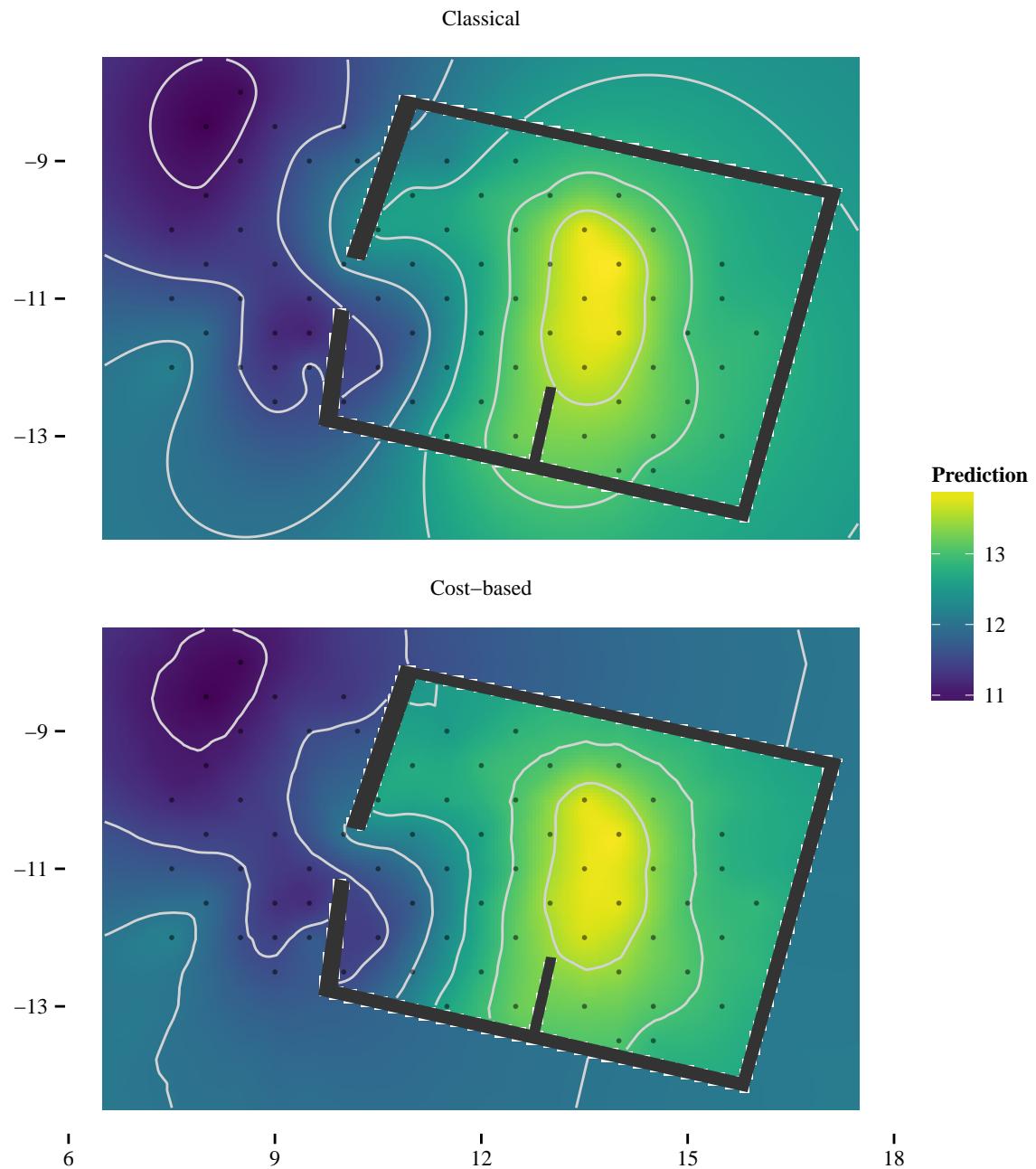


Figure 21: Comparison of Kriging estimates.

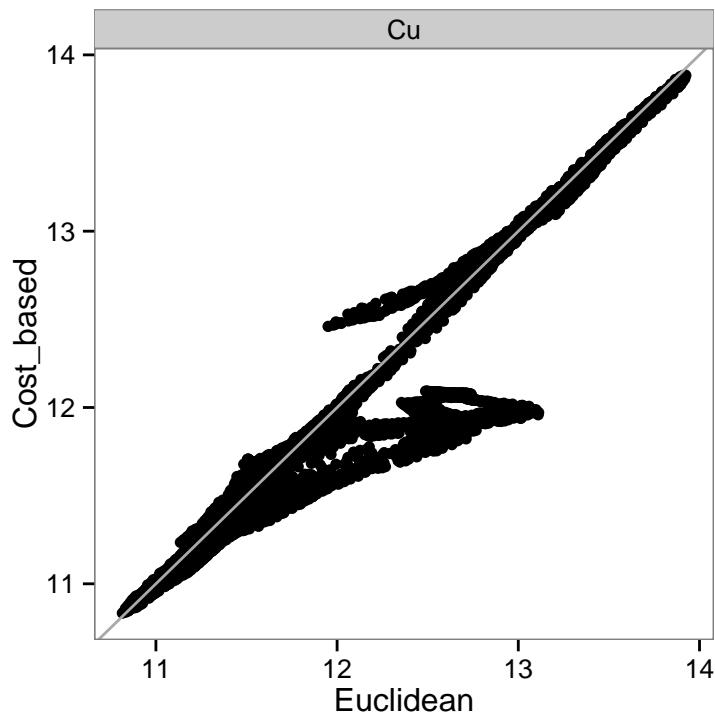


Figure 22: Pointwise comparison of predictions by method.

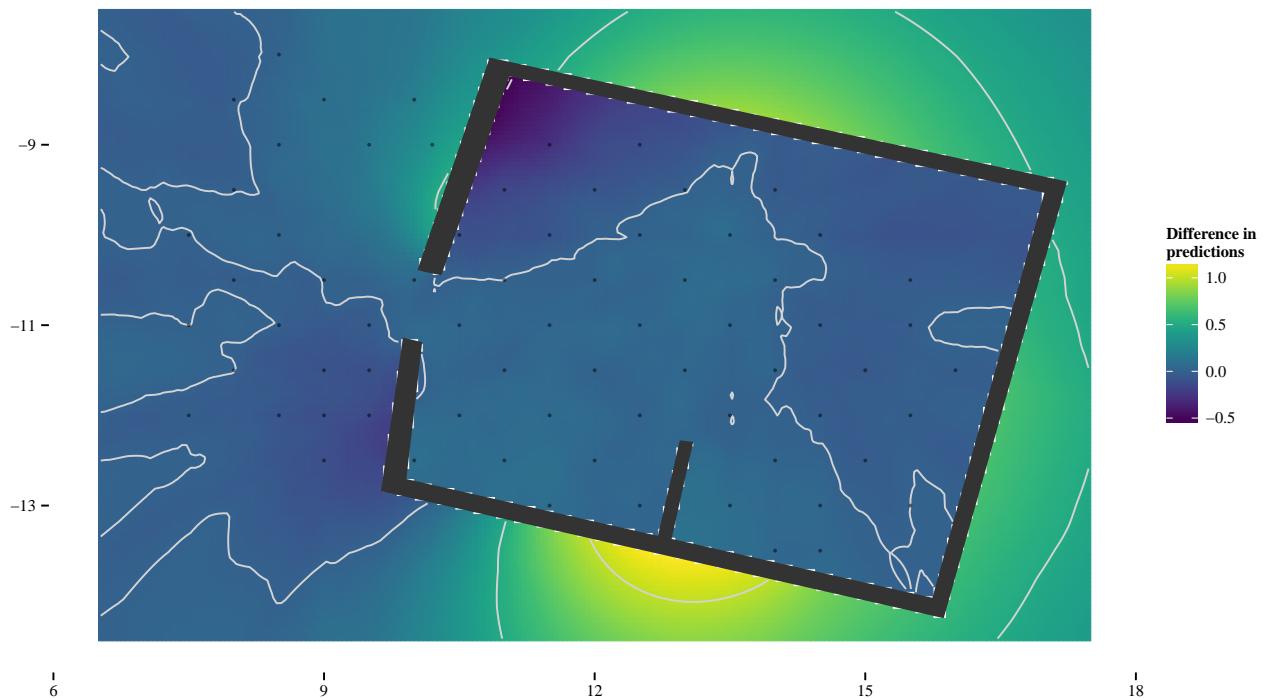


Figure 23: Difference between the Euclidean and the cost-based predictions.

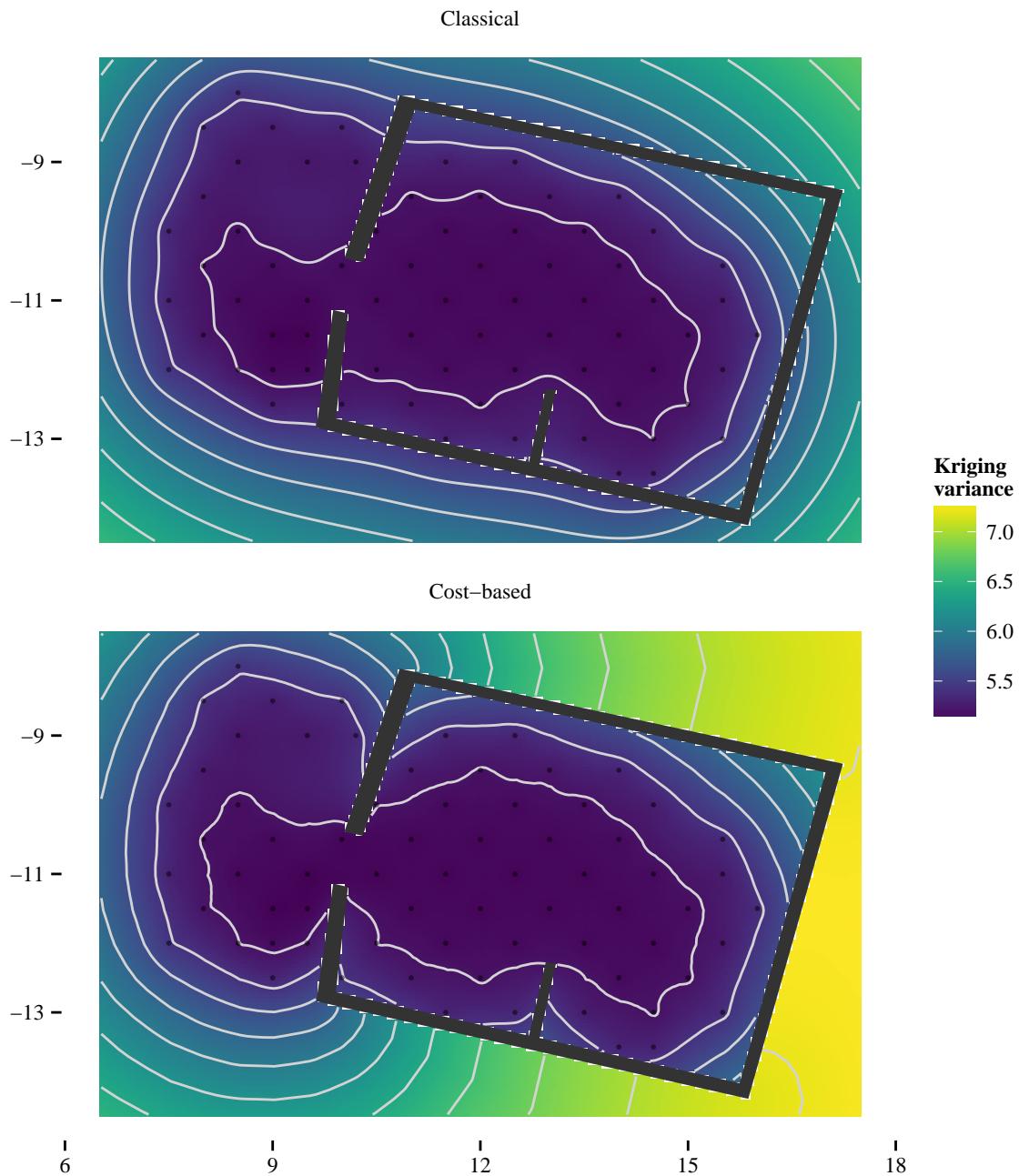


Figure 24: Comparison of prediction error by method.

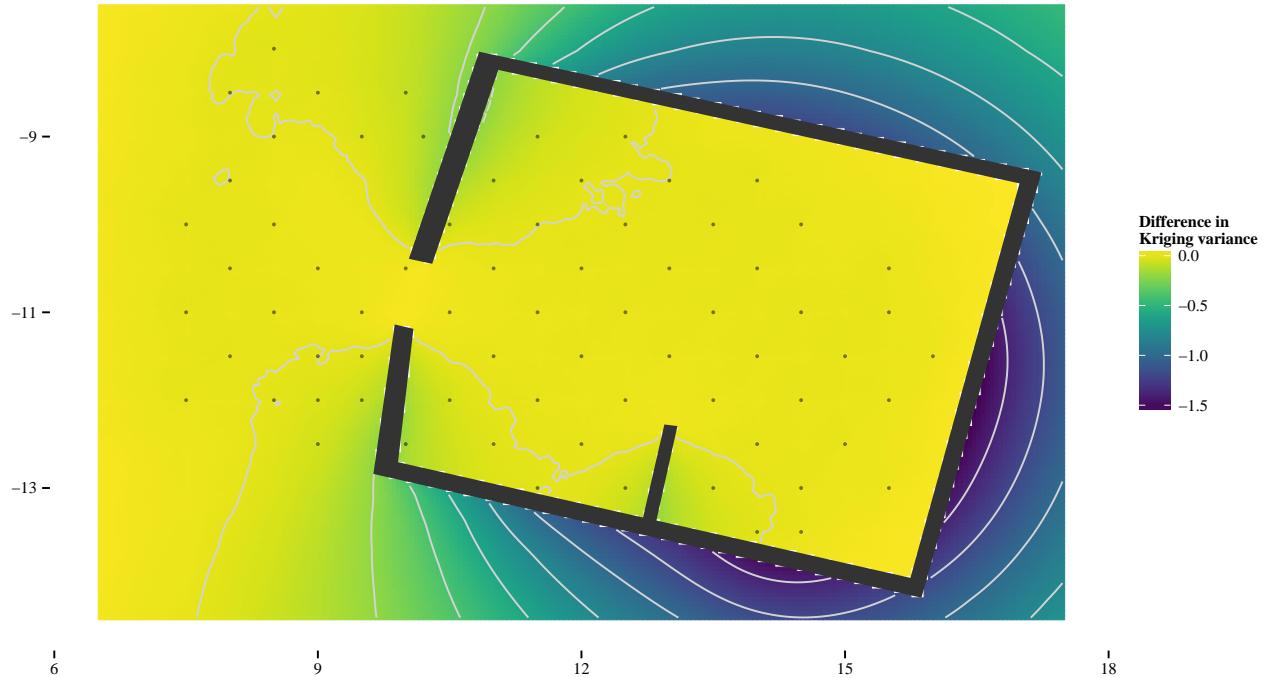


Figure 25: Difference between the Euclidean and the cost-based prediction errors

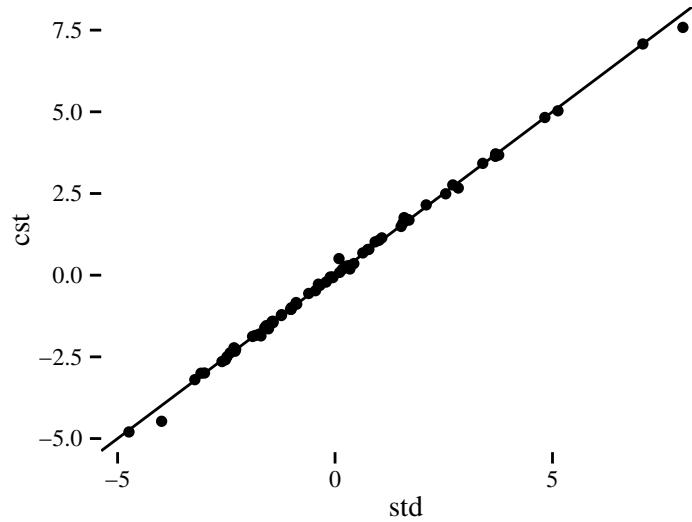


Figure 26: Pointwise leave-one-out prediction error by method.

5 Analysis of Ferrum

5.1 Euclidean kriging

The variogram model is Exponential. We choose to estimate the nugget effect, which may account for measurement error, for example.

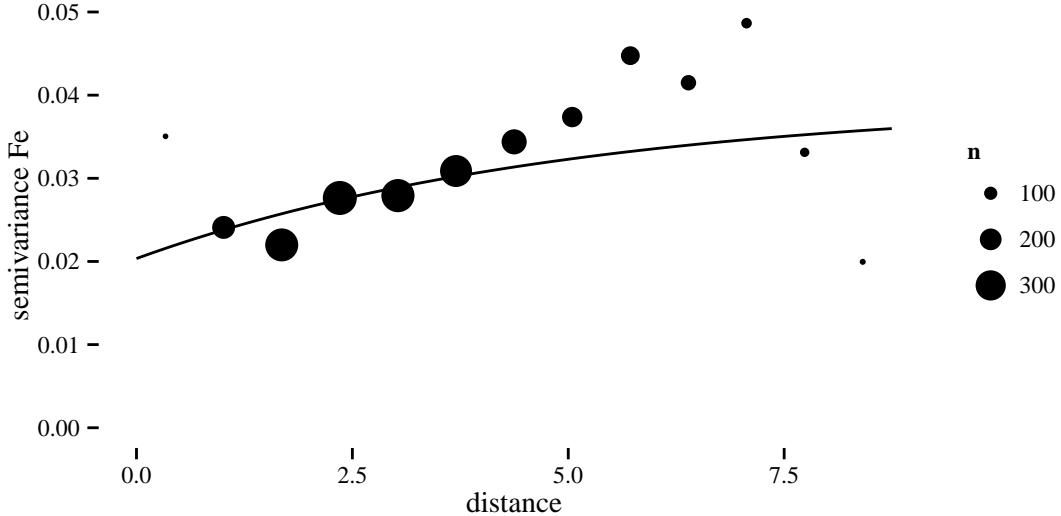


Figure 27: Empirical variogram and fitted model.

5.2 Cost-based kriging

5.3 Comparison of method outcomes

	Euclidean	Cost_based
Intercept	1.23	1.23
Nugget	0.02	0.02
Partial sill	0.02	0.02
phi	5.04	4.23
Pract. range	15.09	12.67
Log-likelihood	28.20	28.02

In the scatter plot, the horizontal patterns correspond to predictions on observed values. Otherwise, the differences are negligible.

Near the observations, the cost-based approach has a larger prediction error due to its increased estimation of the nugget (i.e. short-range variance). In the main area, the prediction errors are practically the same with both approaches. Behind the walls, the Euclidean prediction error is unrealistically low.

5.4 Leave-one-out Cross Validation (LOOCV)

method	rmse(error)
std	0.16

method	rmse(error)
cst	0.16

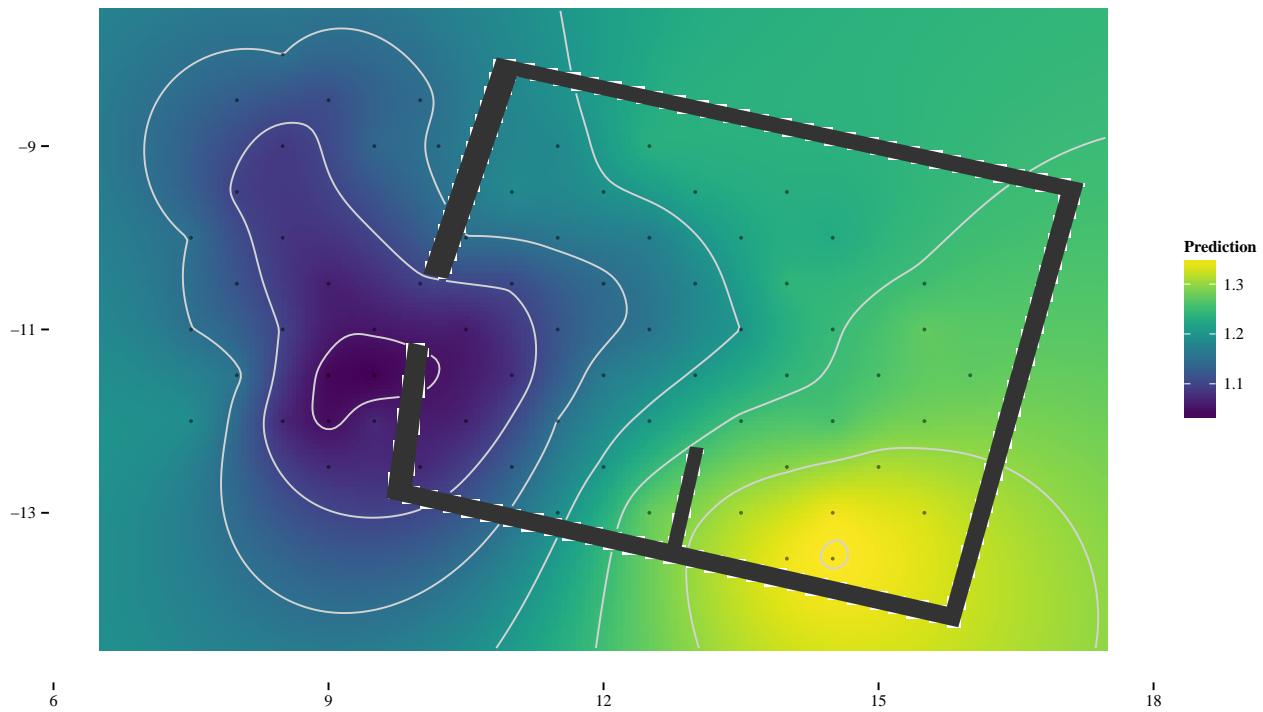


Figure 28: Euclidean kriging prediction

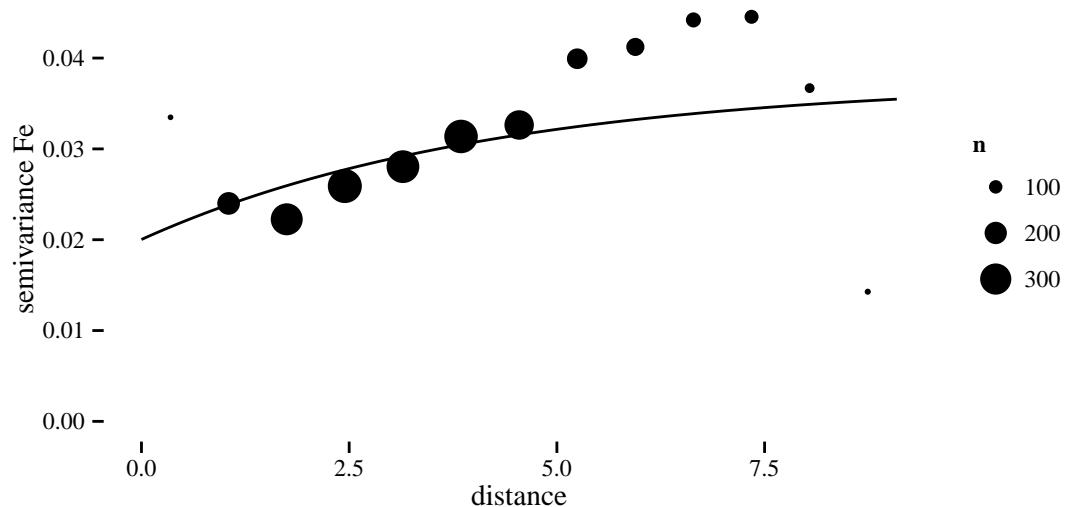


Figure 29: Empirical cost-based variogram and fitted model.

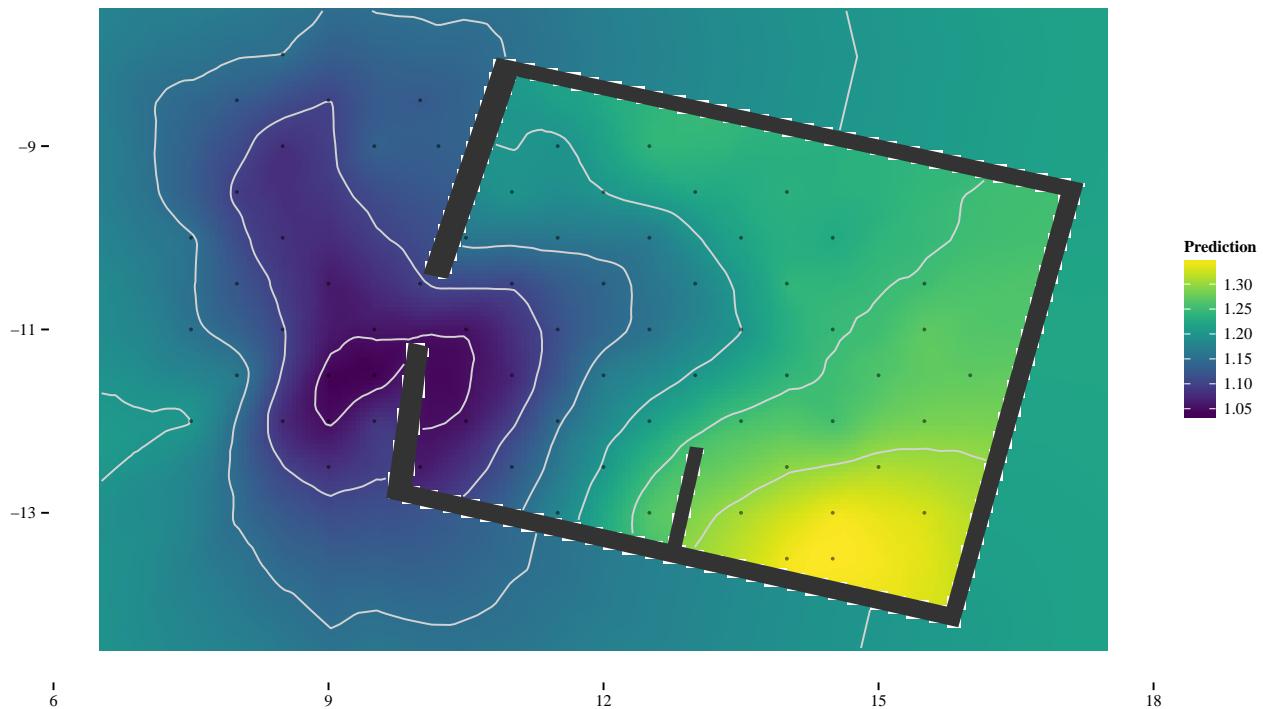


Figure 30: Cost-based kriging prediction

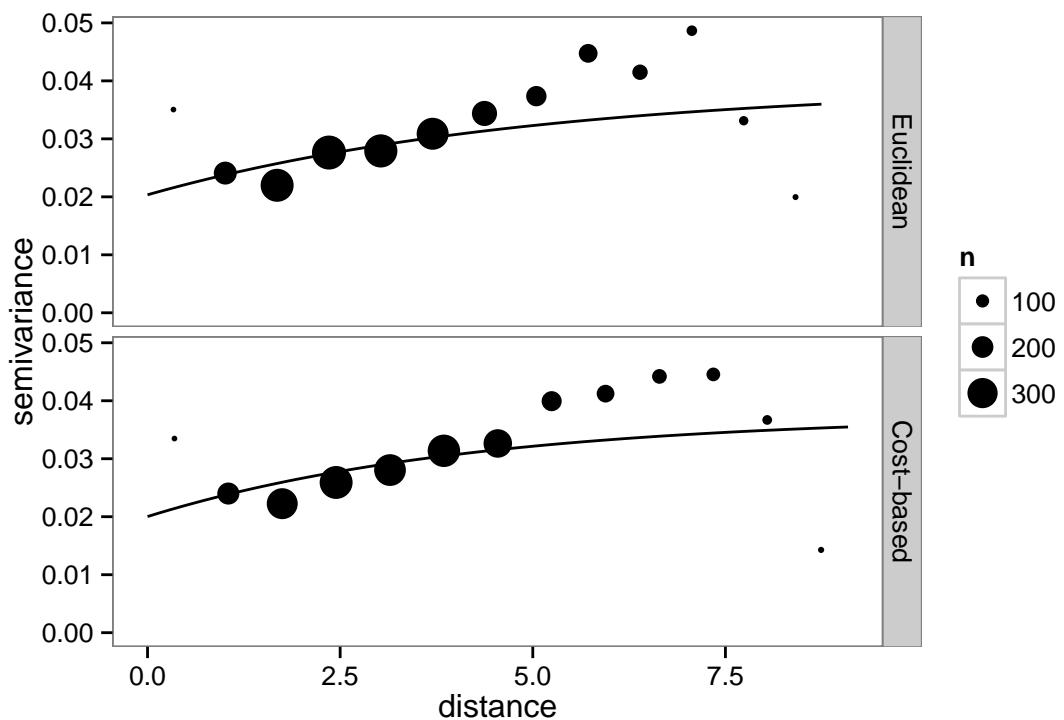


Figure 31: Empirical variogram and fitted models by method for Ferrum.

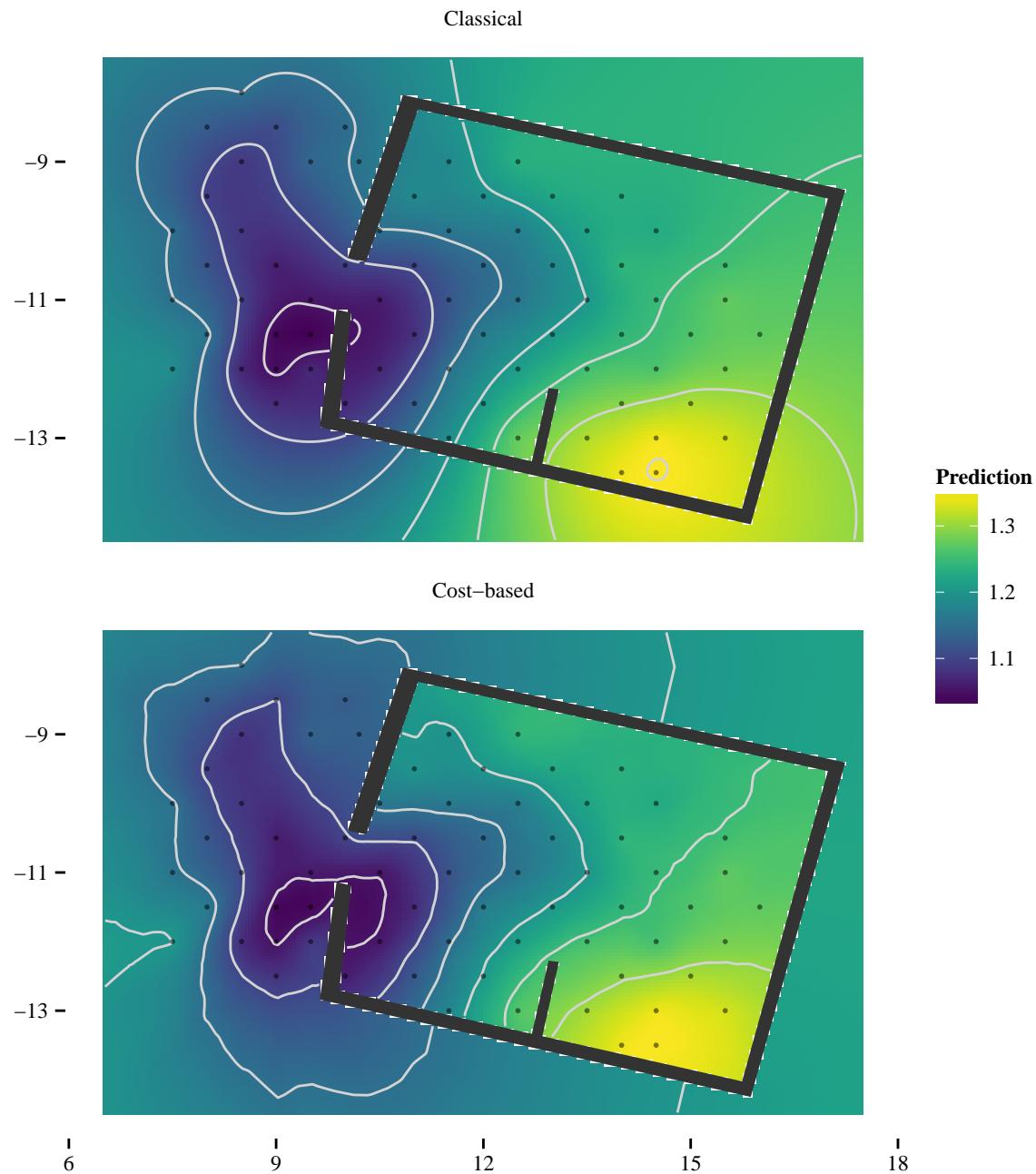


Figure 32: Comparison of Kriging estimates.

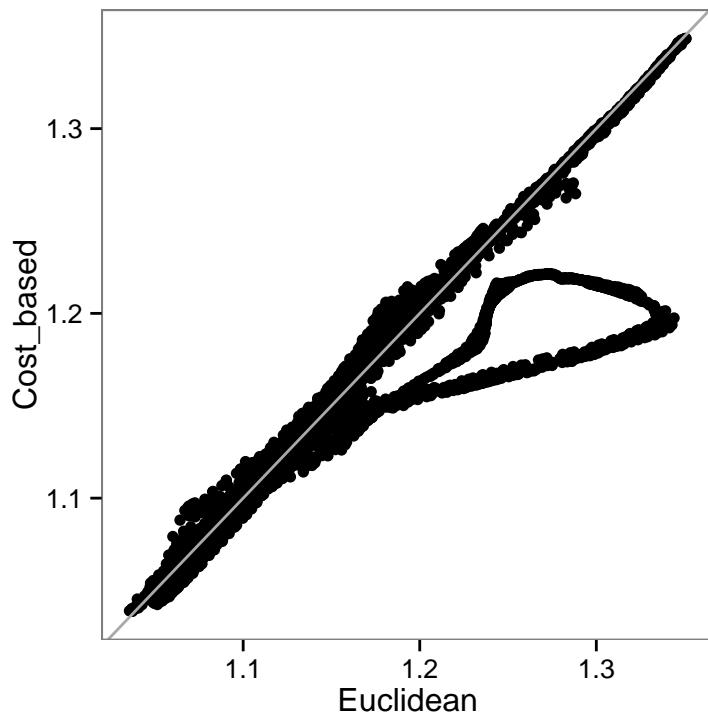


Figure 33: Pointwise comparison of predictions by method.

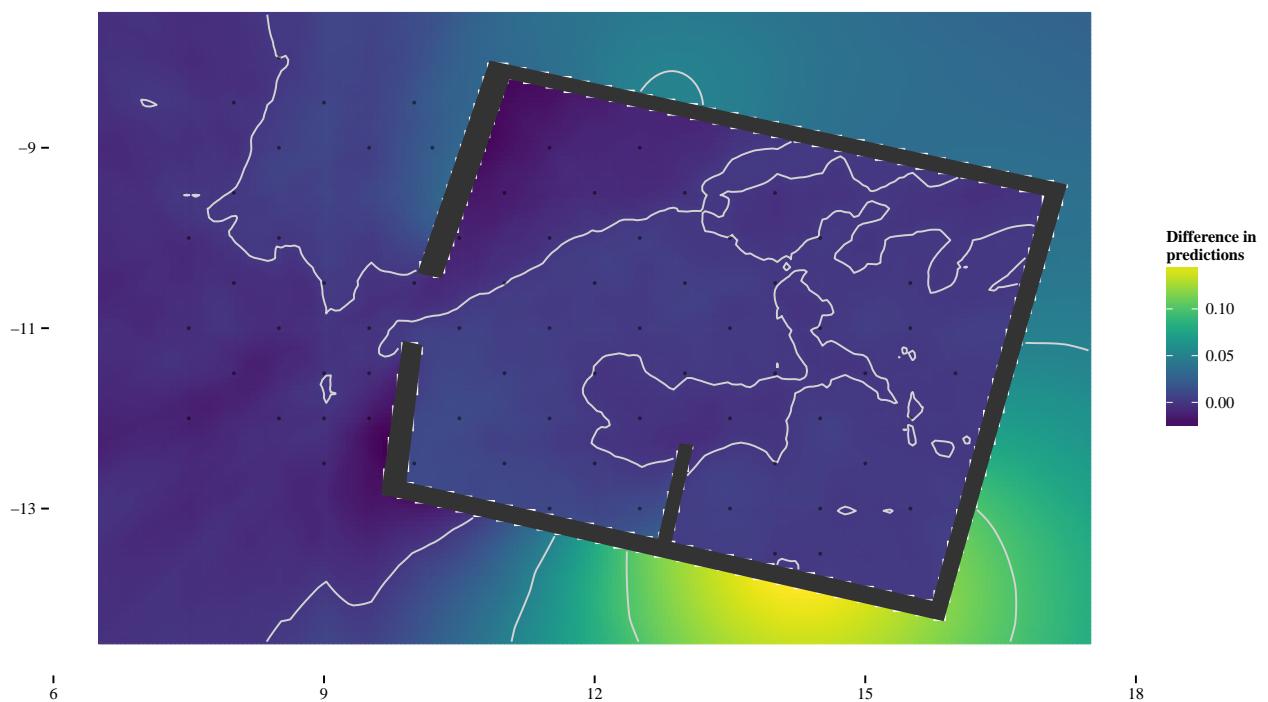


Figure 34: Difference between the Euclidean and the cost-based predictions.

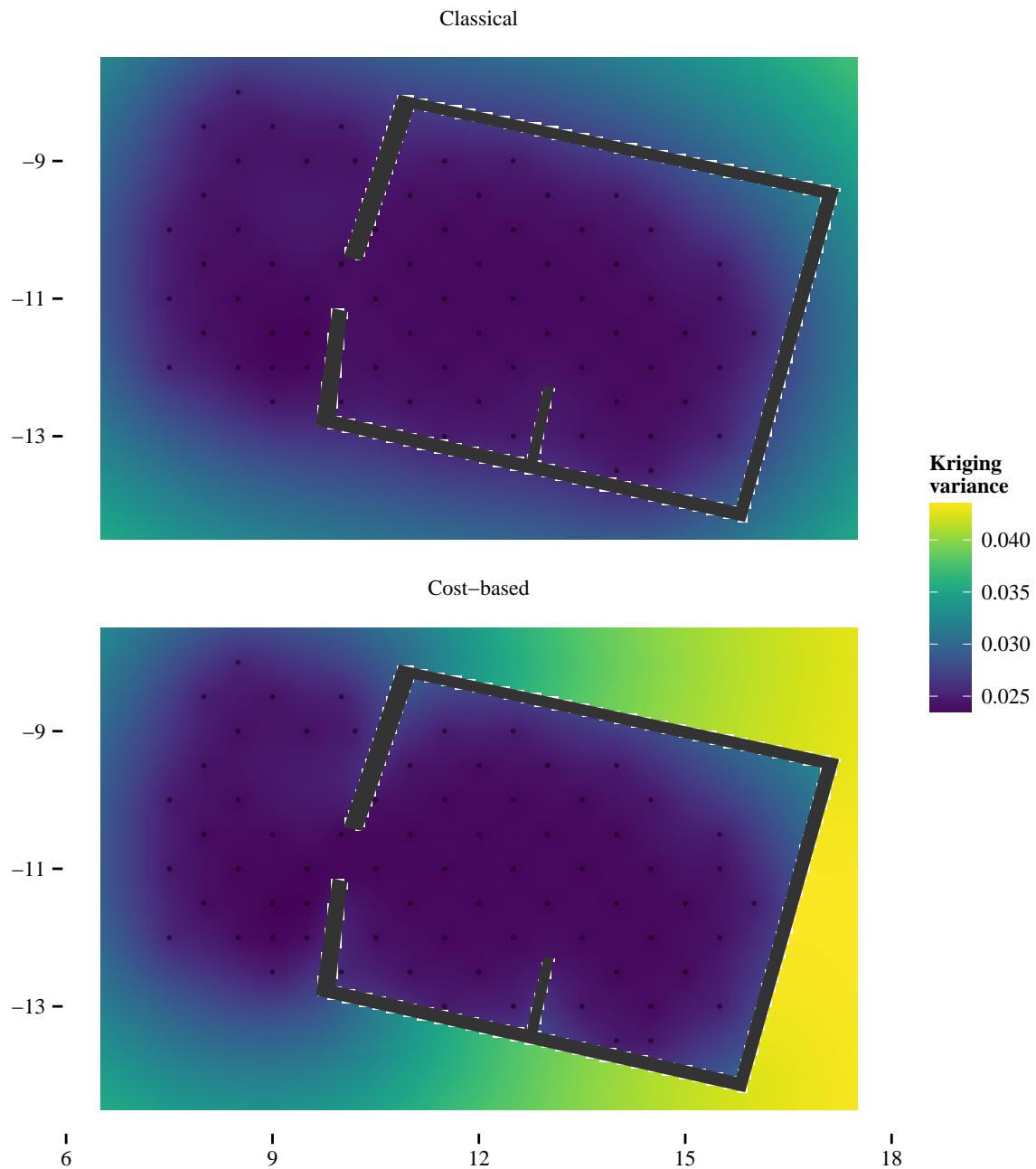


Figure 35: Comparison of prediction error by method.

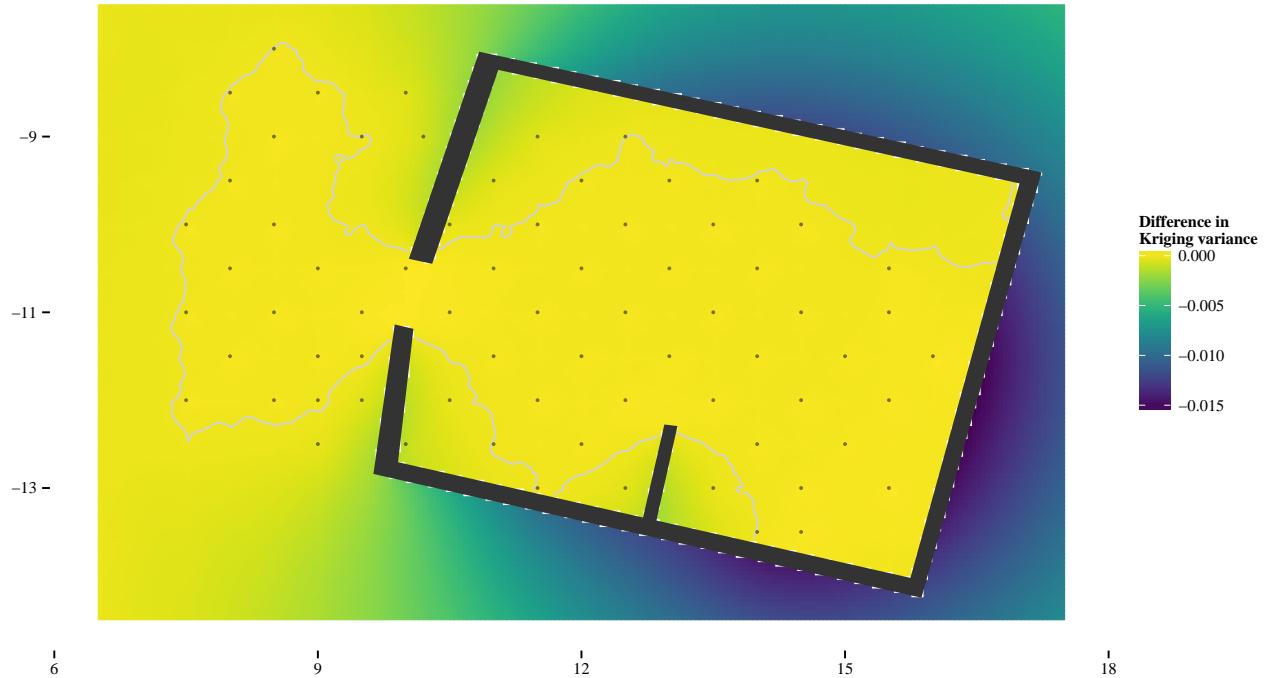


Figure 36: Difference between the Euclidean and the cost-based prediction errors

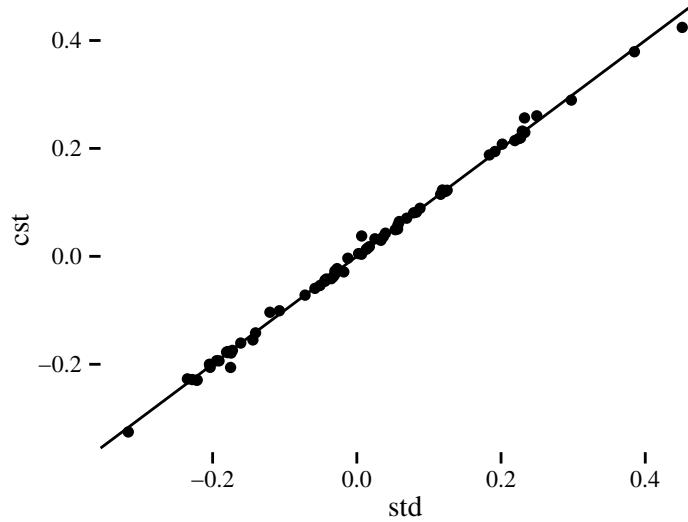


Figure 37: Pointwise leave-one-out prediction error by method.

6 Analysis of Potassium

6.1 Euclidean kriging

The variogram model is Exponential. We choose to estimate the nugget effect, which may account for measurement error, for example.

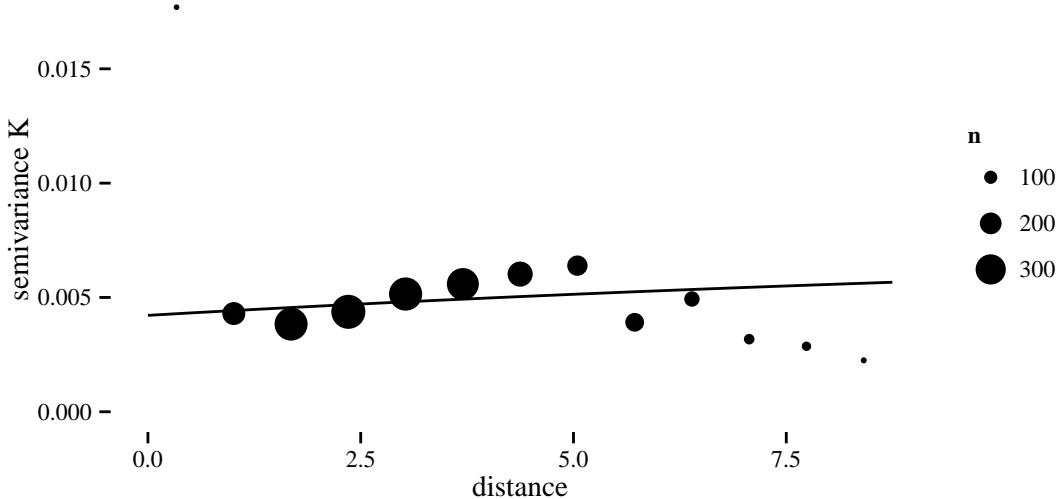


Figure 38: Empirical variogram and fitted model.

6.2 Cost-based kriging

6.3 Comparison of method outcomes

	Euclidean	Cost_based
Intercept	0.25	0.25
Nugget	0.00	0.00
Partial sill	0.00	0.00
phi	16.23	15.13
Pract. range	48.61	45.32
Log-likelihood	87.28	86.98

In the scatter plot, the horizontal patterns correspond to predictions on observed values. Otherwise, the differences are negligible.

Near the observations, the cost-based approach has a larger prediction error due to its increased estimation of the nugget (i.e. short-range variance). In the main area, the prediction errors are practically the same with both approaches. Behind the walls, the Euclidean prediction error is unrealistically low.

6.4 Leave-one-out Cross Validation (LOOCV)

method	rmse(error)
std	0.07

method	rmse(error)
cst	0.07

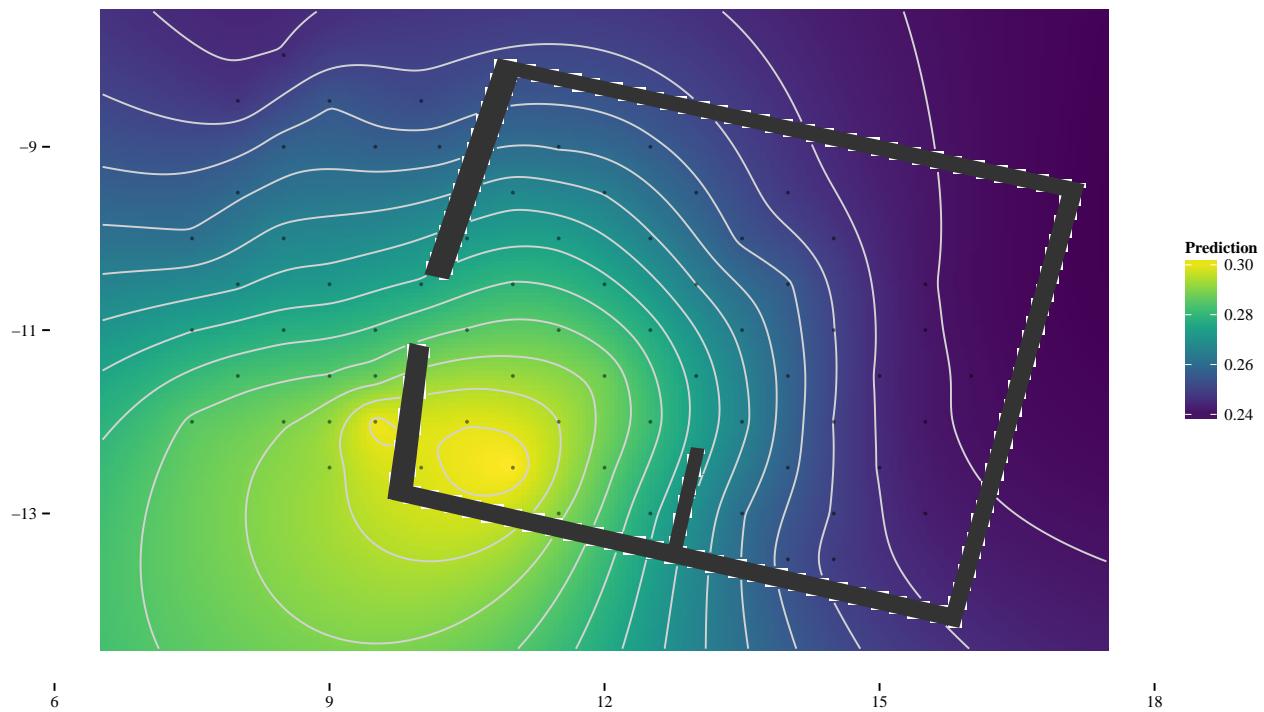


Figure 39: Euclidean kriging prediction

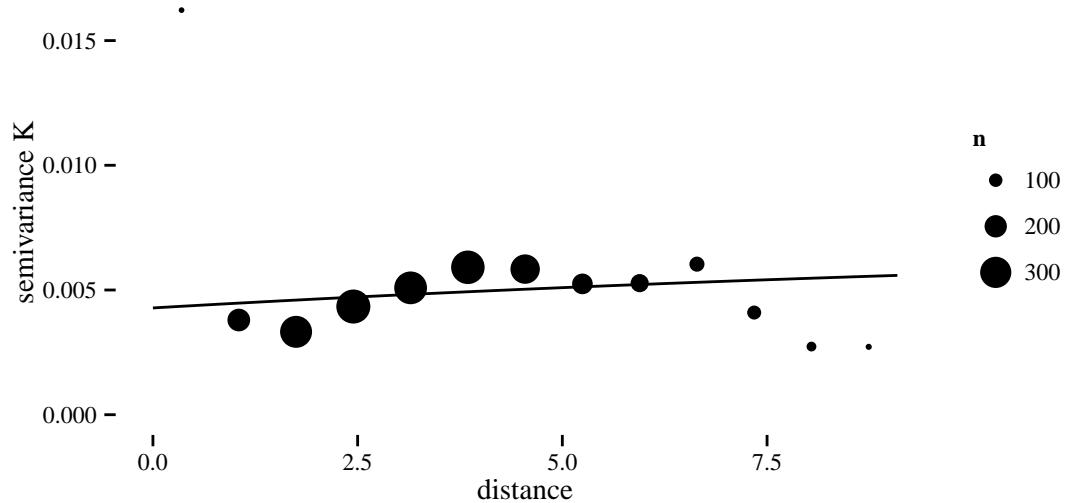


Figure 40: Empirical cost-based variogram and fitted model.

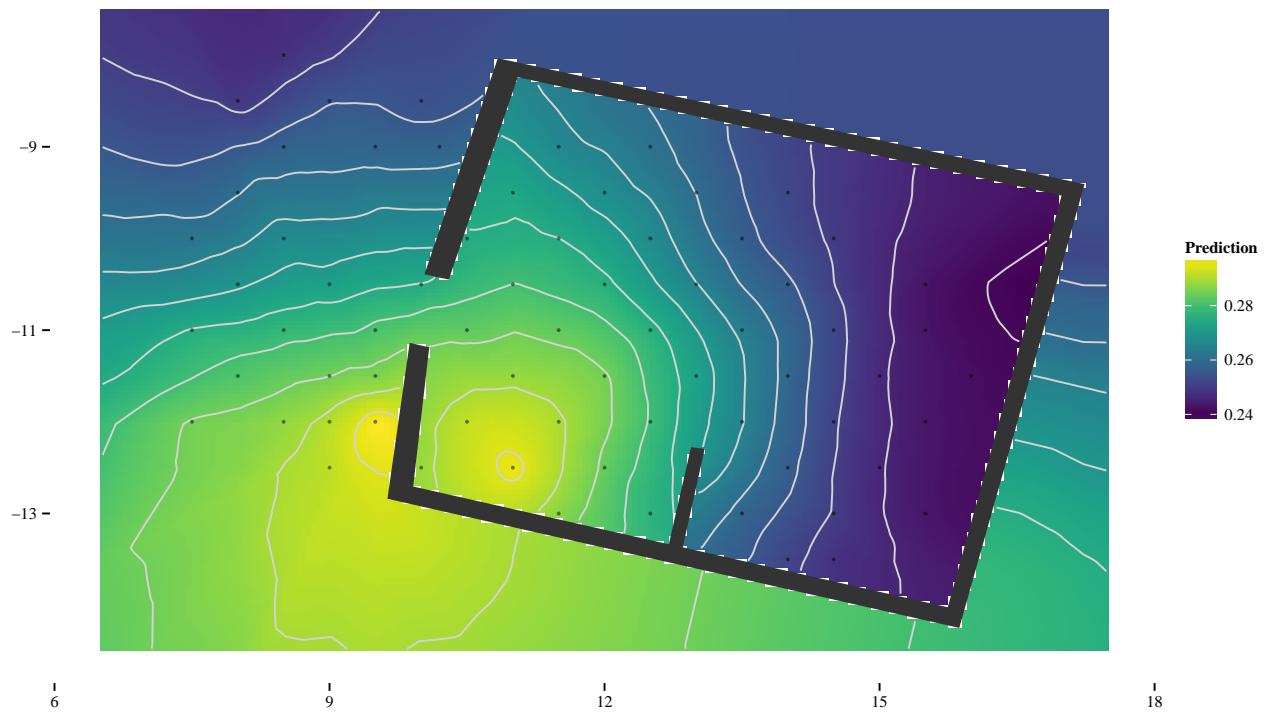


Figure 41: Cost-based kriging prediction

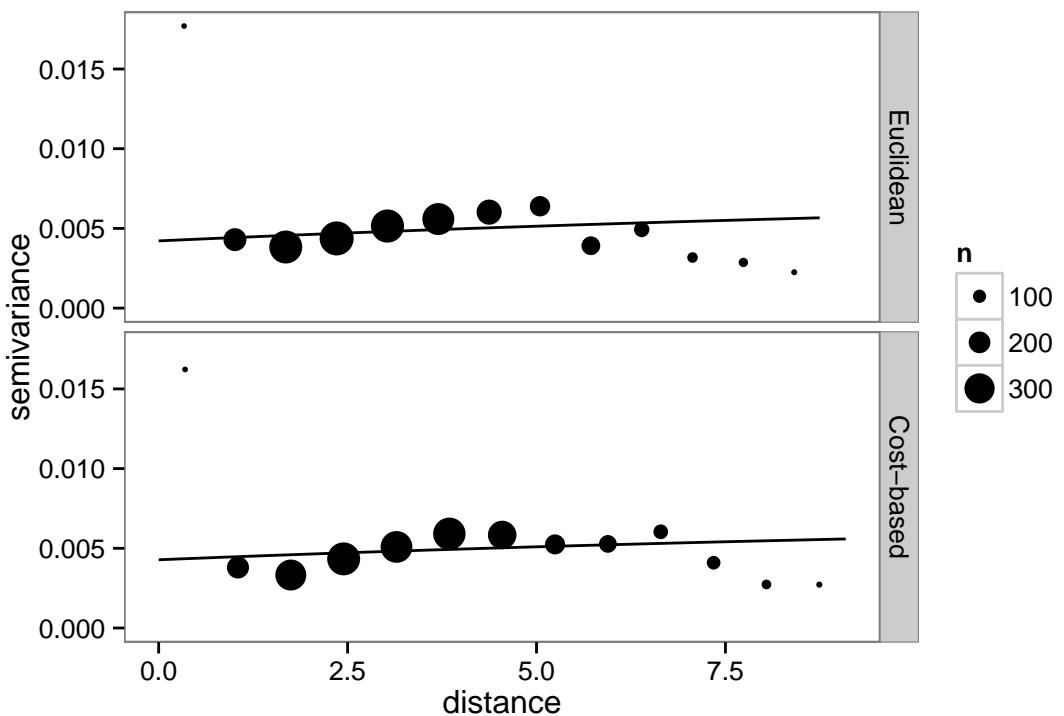


Figure 42: Empirical variogram and fitted models by method for Potassium.

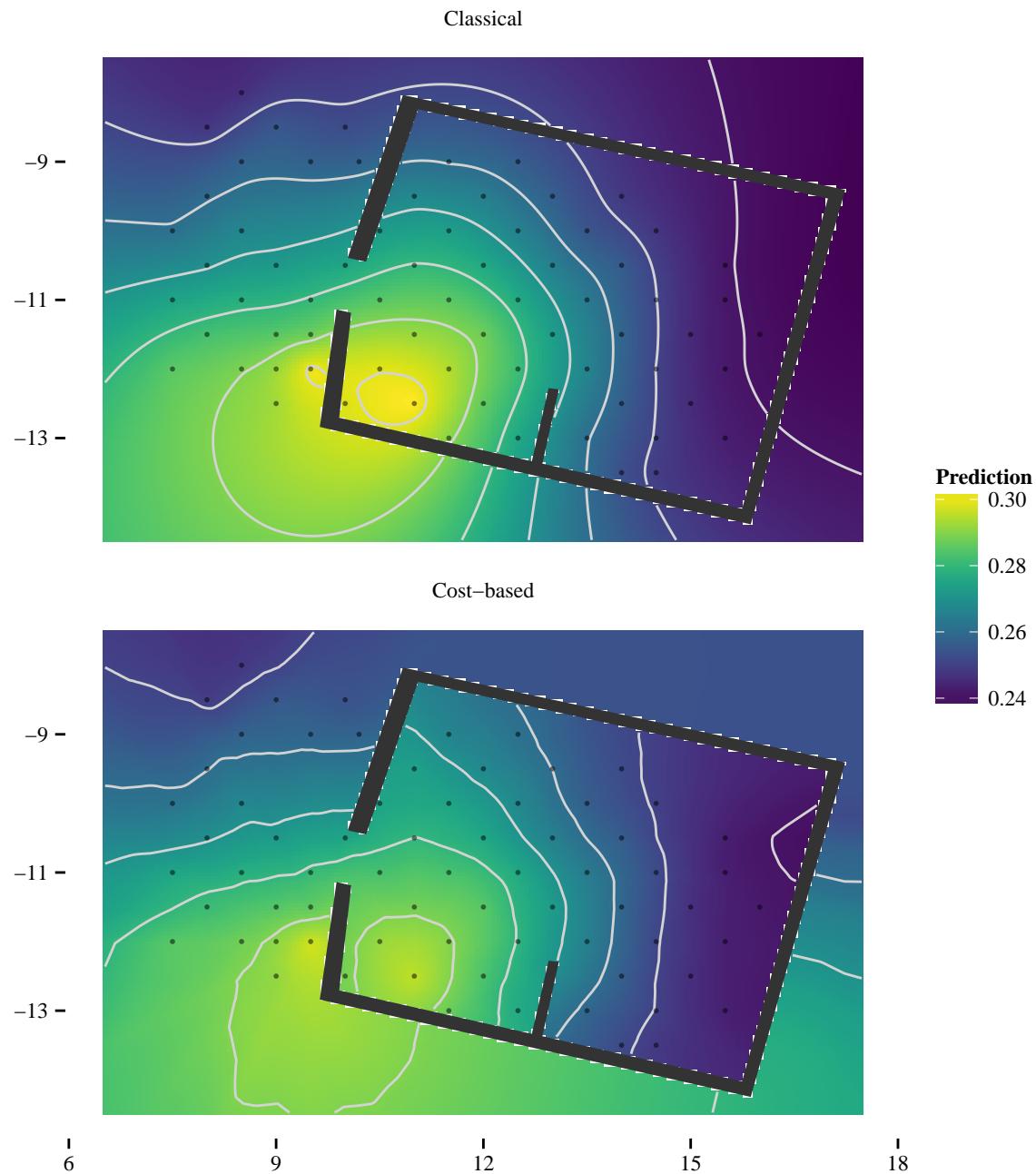


Figure 43: Comparison of Kriging estimates.

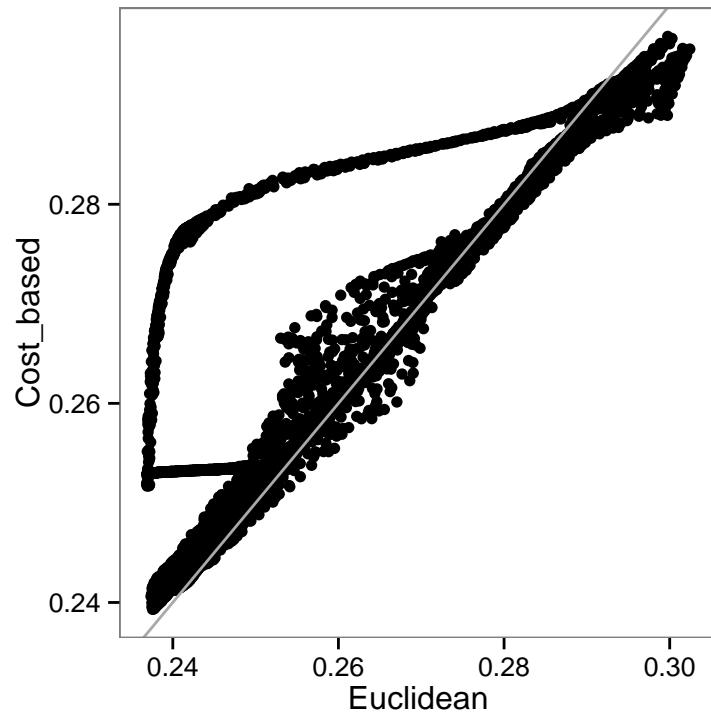


Figure 44: Pointwise comparison of predictions by method.

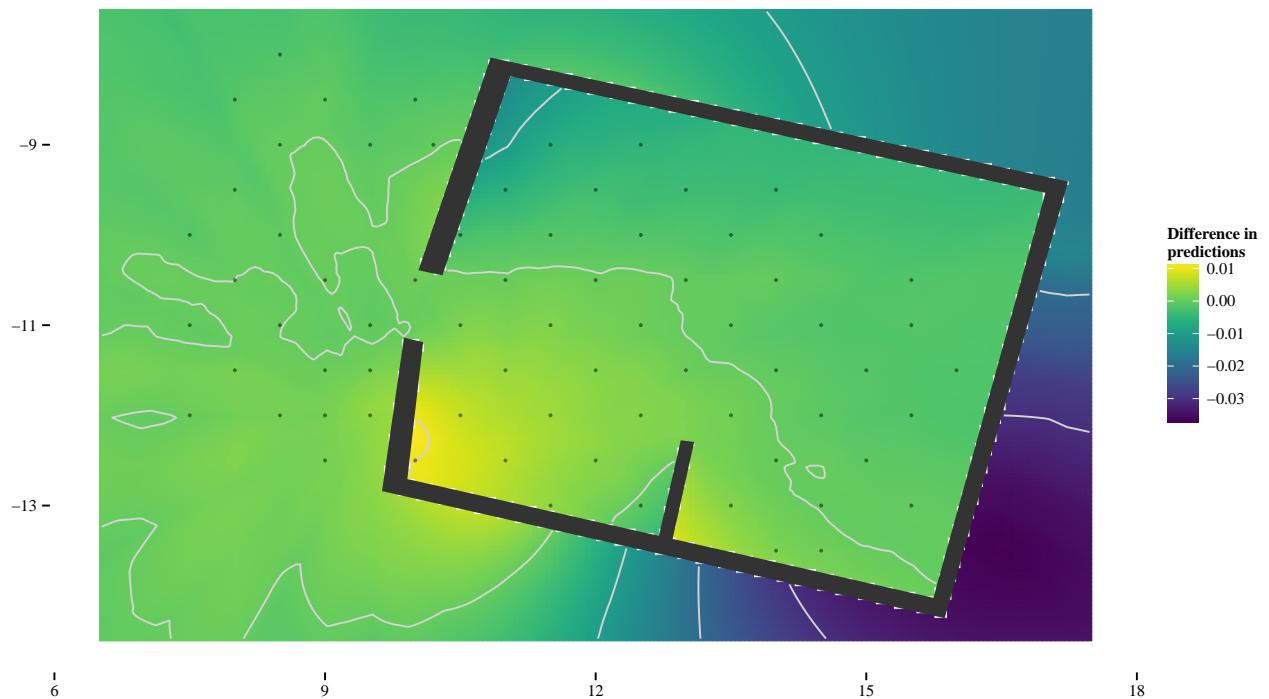


Figure 45: Difference between the Euclidean and the cost-based predictions.

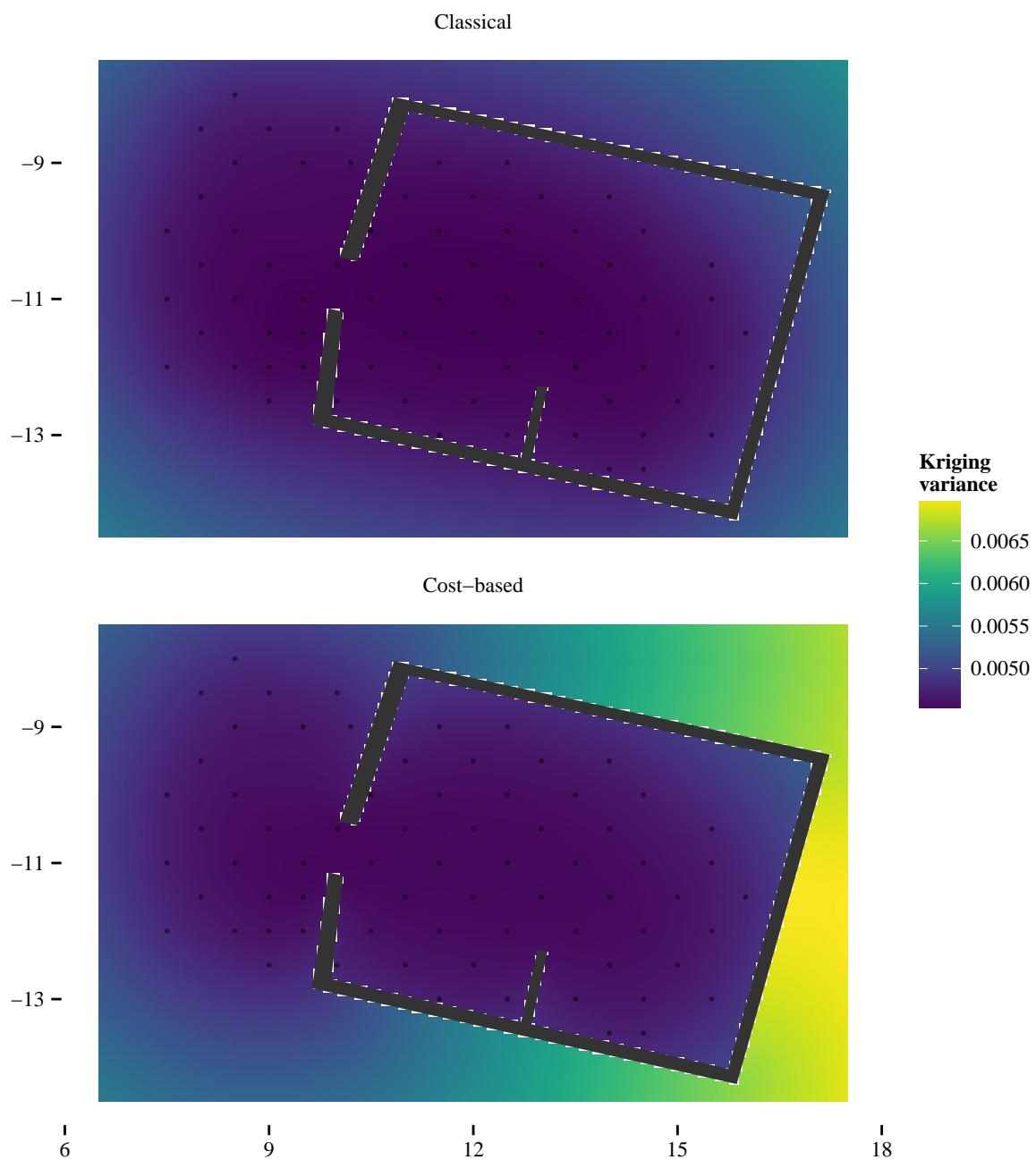


Figure 46: Comparison of prediction error by method.

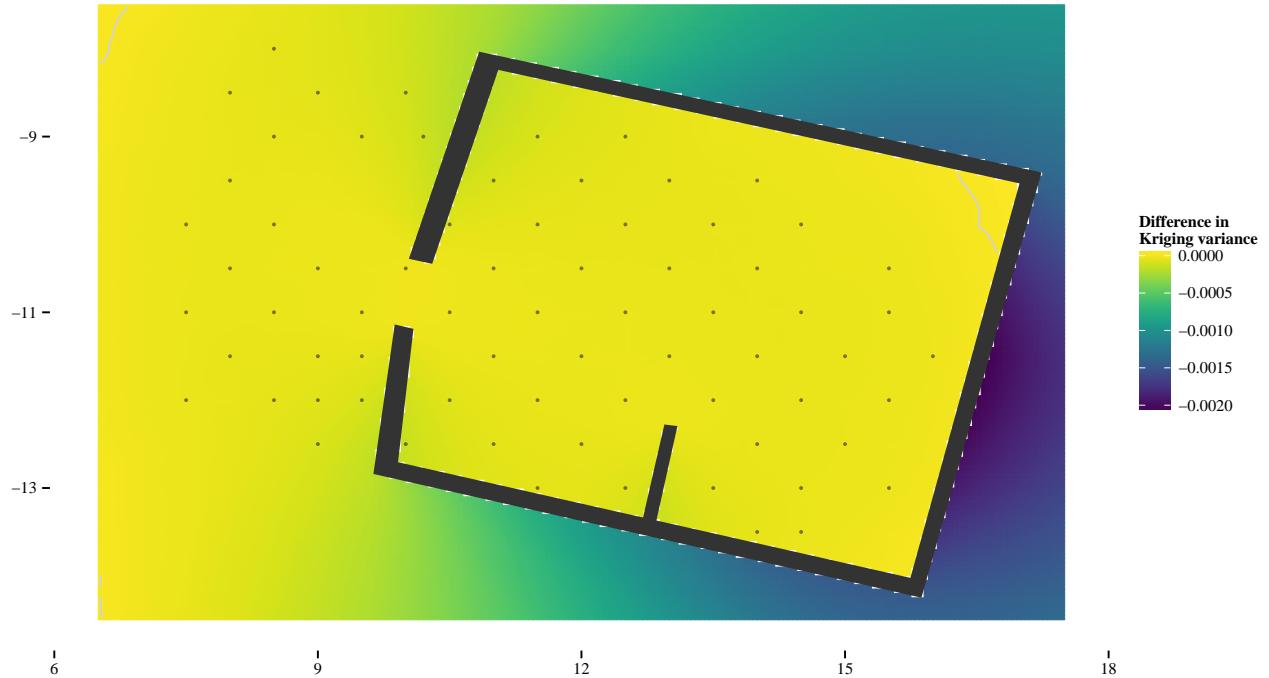


Figure 47: Difference between the Euclidean and the cost-based prediction errors

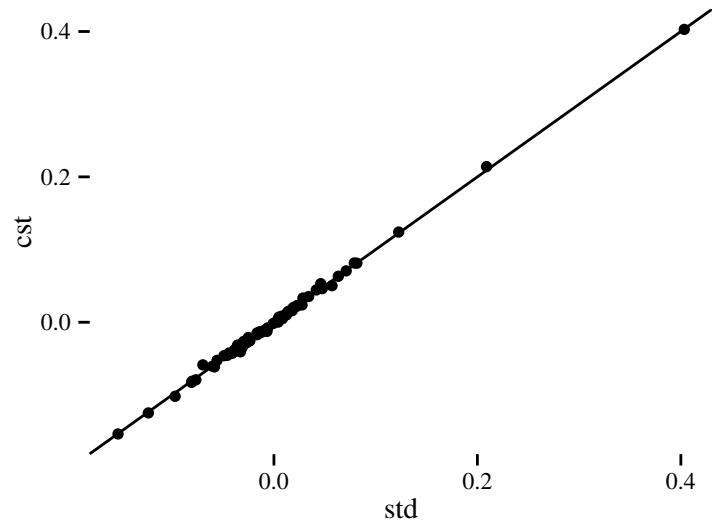


Figure 48: Pointwise leave-one-out prediction error by method.

7 Analysis of Magnesium

7.1 Euclidean kriging

The variogram model is Exponential. We choose to estimate the nugget effect, which may account for measurement error, for example.

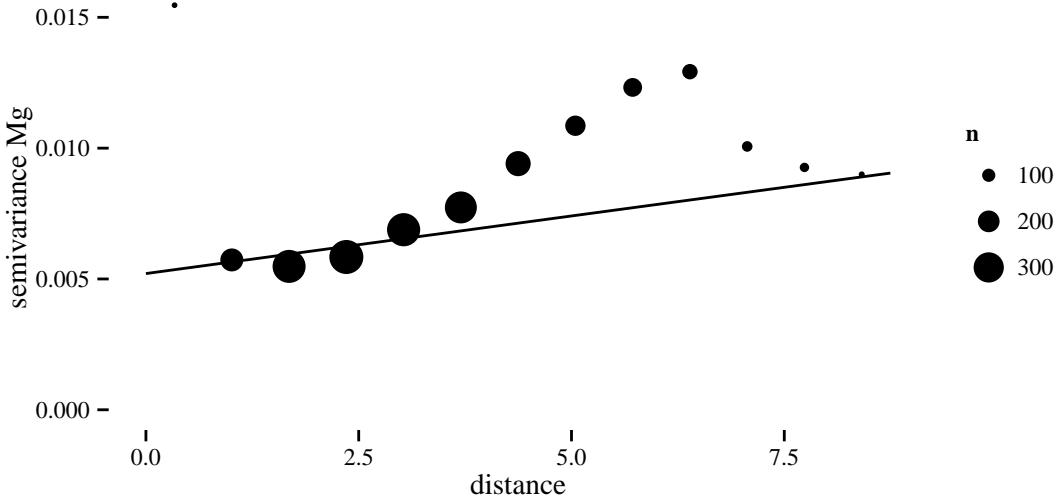


Figure 49: Empirical variogram and fitted model.

7.2 Cost-based kriging

7.3 Comparison of method outcomes

	Euclidean	Cost_based
Intercept	0.45	0.45
Nugget	0.01	0.01
Partial sill	0.30	0.28
phi	677.72	706.62
Pract. range	2030.28	2116.83
Log-likelihood	78.19	78.02

In the scatter plot, the horizontal patterns correspond to predictions on observed values. Otherwise, the differences are negligible.

Near the observations, the cost-based approach has a larger prediction error due to its increased estimation of the nugget (i.e. short-range variance). In the main area, the prediction errors are practically the same with both approaches. Behind the walls, the Euclidean prediction error is unrealistically low.

7.4 Leave-one-out Cross Validation (LOOCV)

method	rmse(error)
std	0.08

method	rmse(error)
cst	0.08

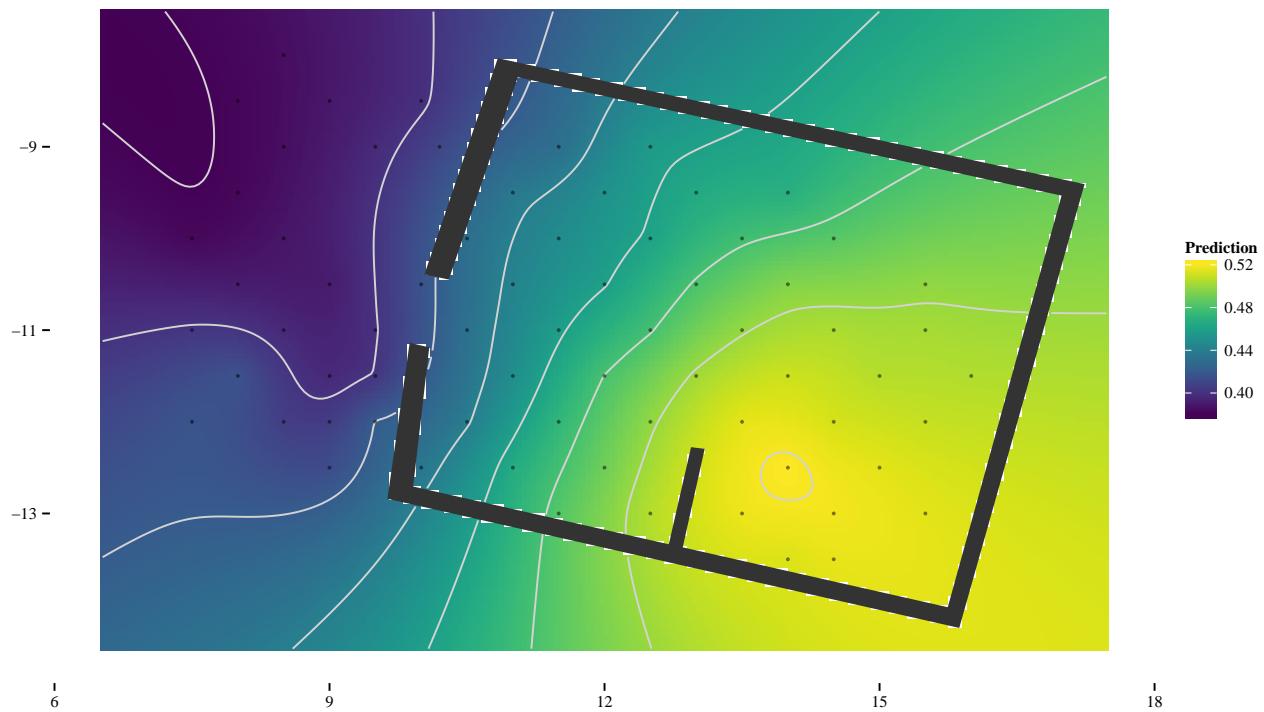


Figure 50: Euclidean kriging prediction

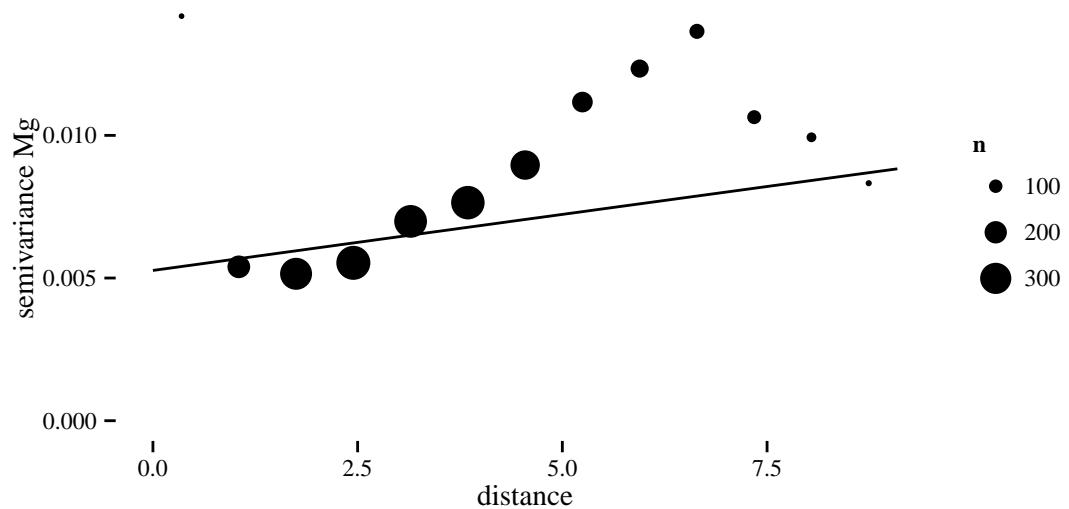


Figure 51: Empirical cost-based variogram and fitted model.

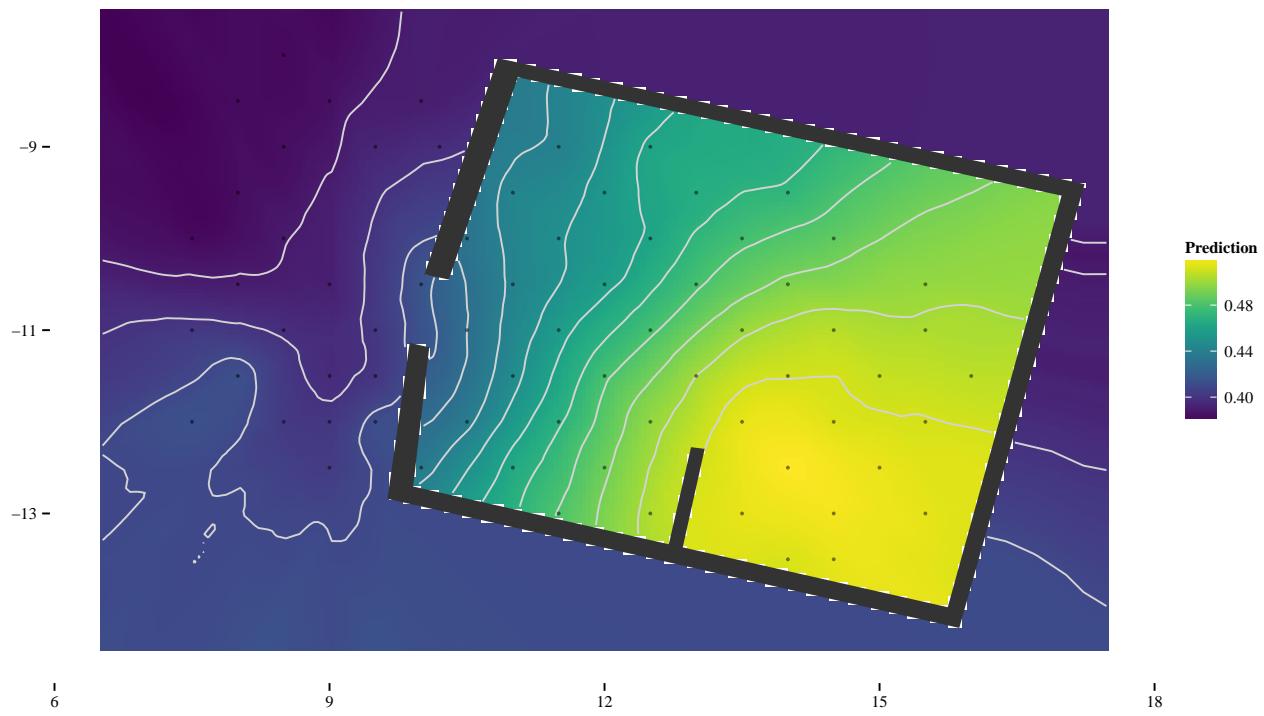


Figure 52: Cost-based kriging prediction

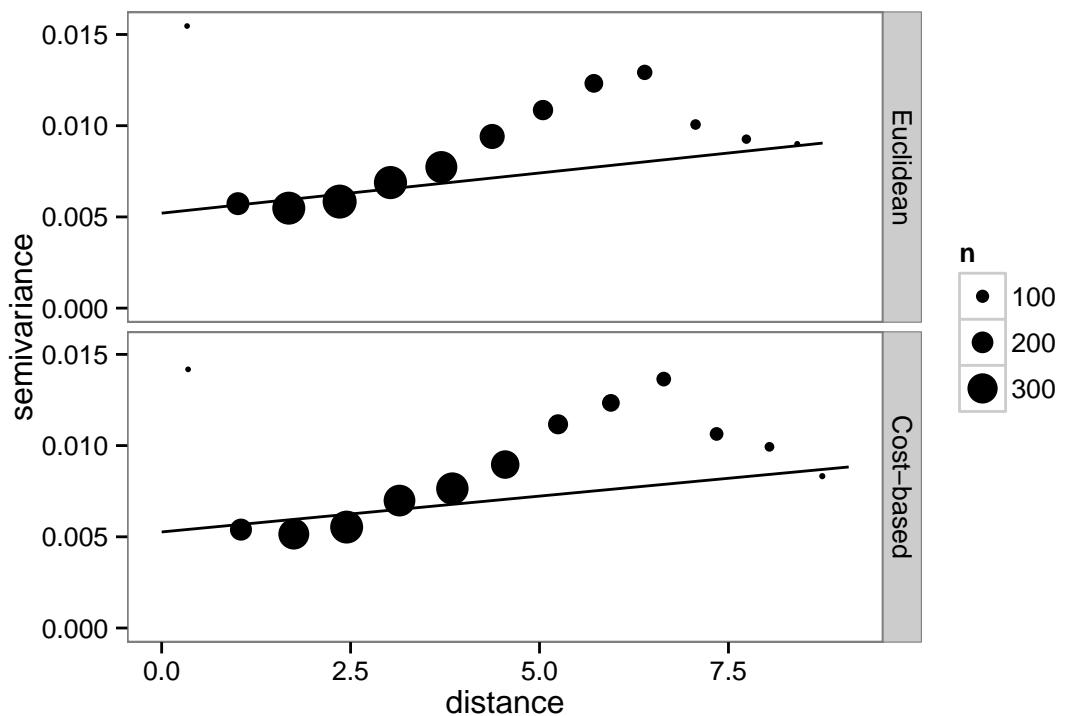


Figure 53: Empirical variogram and fitted models by method for Magnesium.

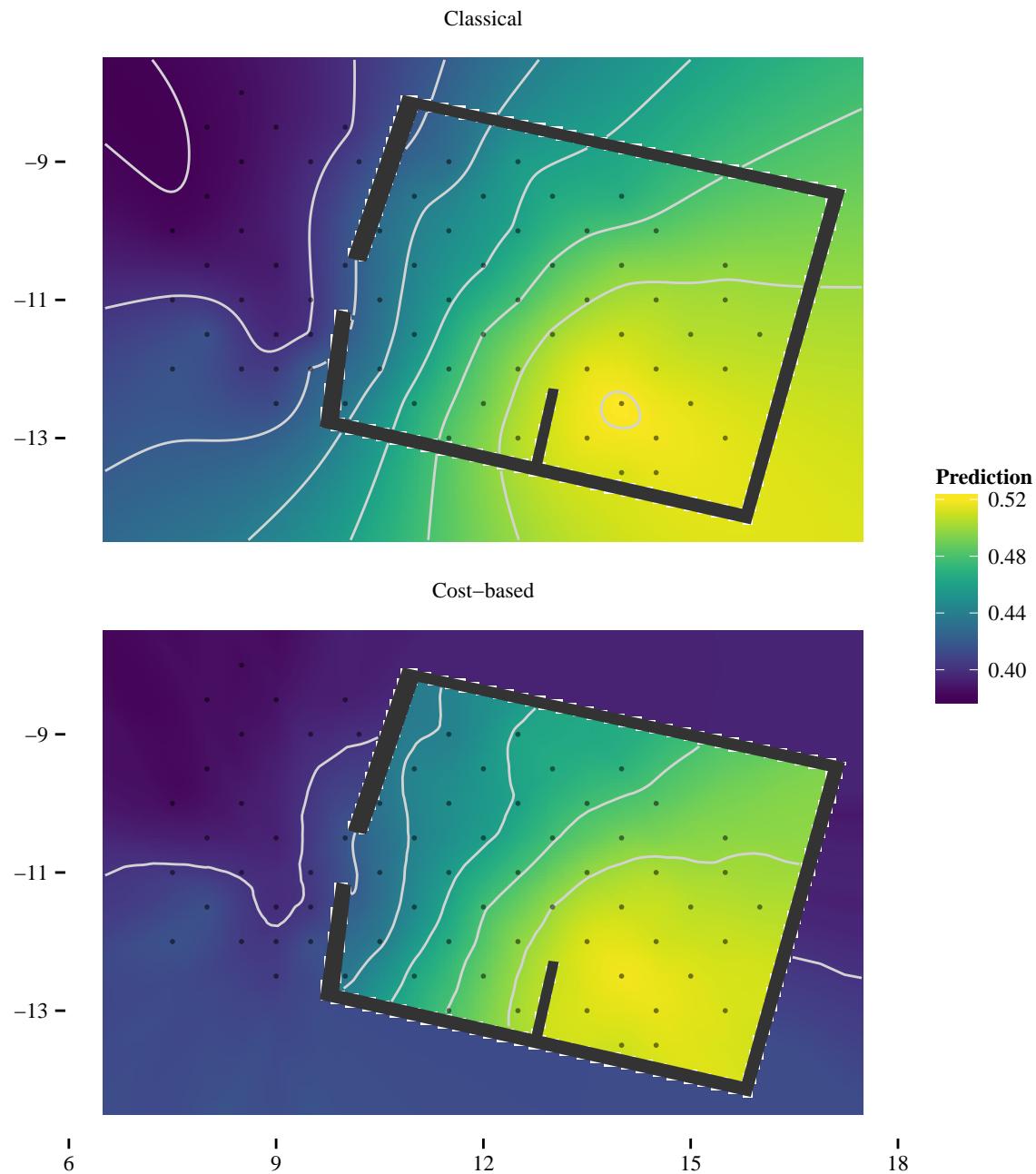


Figure 54: Comparison of Kriging estimates.

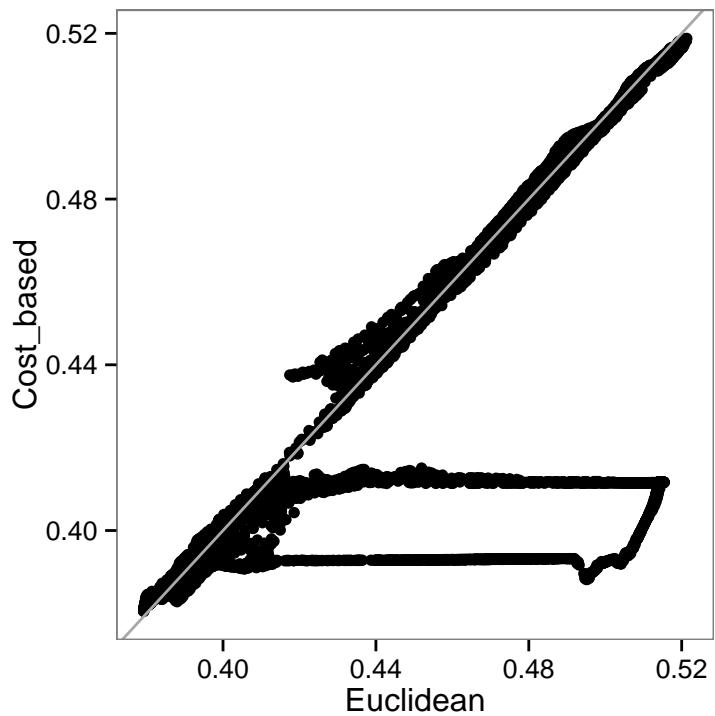


Figure 55: Pointwise comparison of predictions by method.

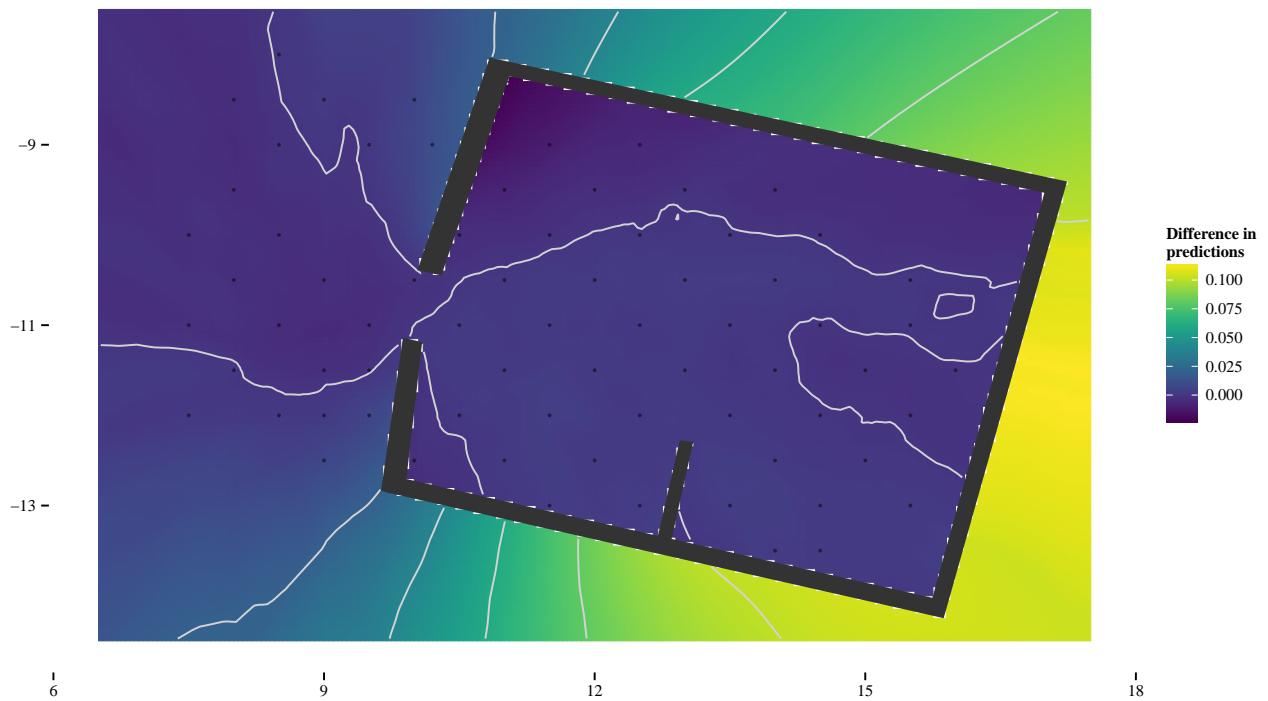


Figure 56: Difference between the Euclidean and the cost-based predictions.

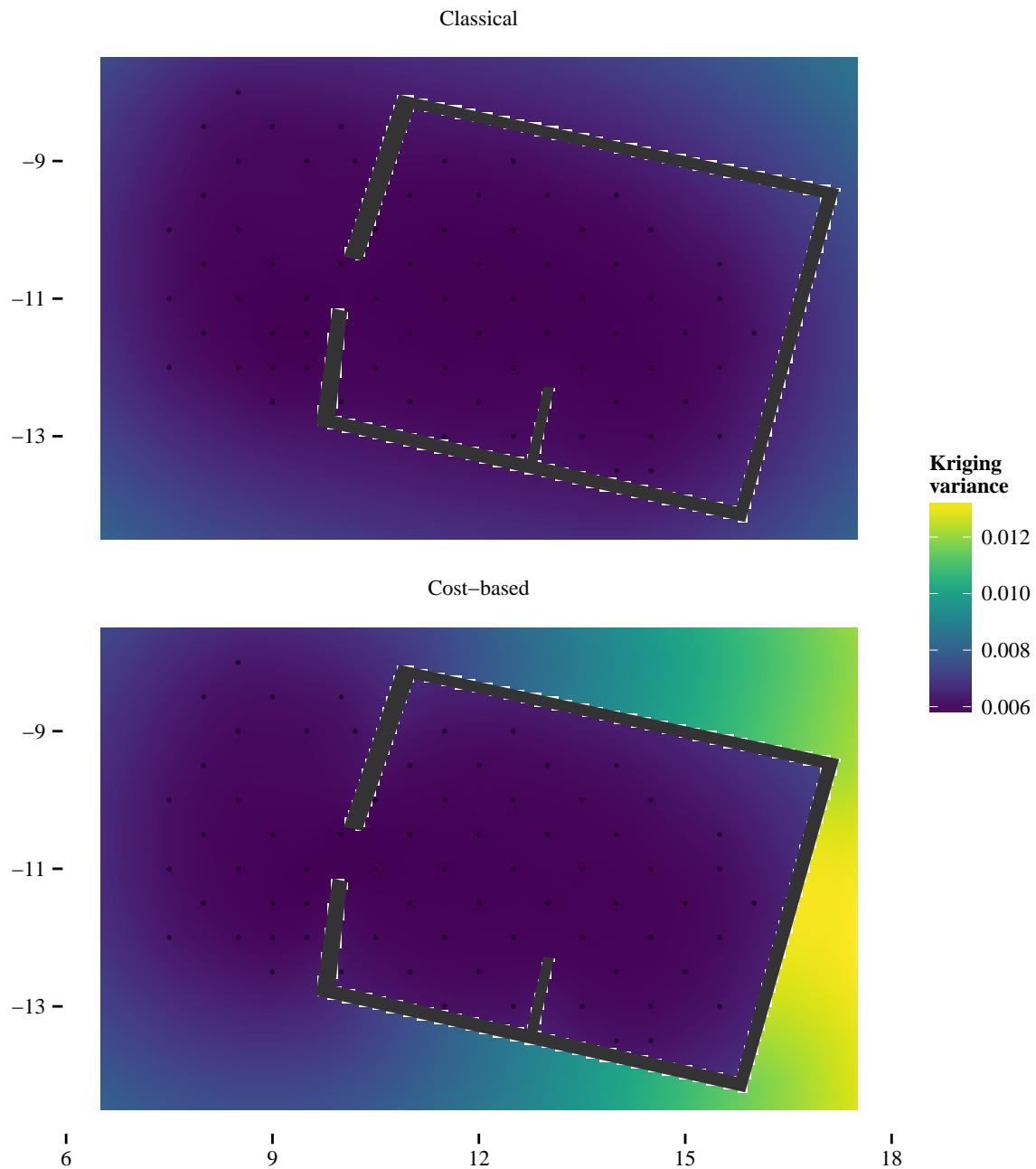


Figure 57: Comparison of prediction error by method.

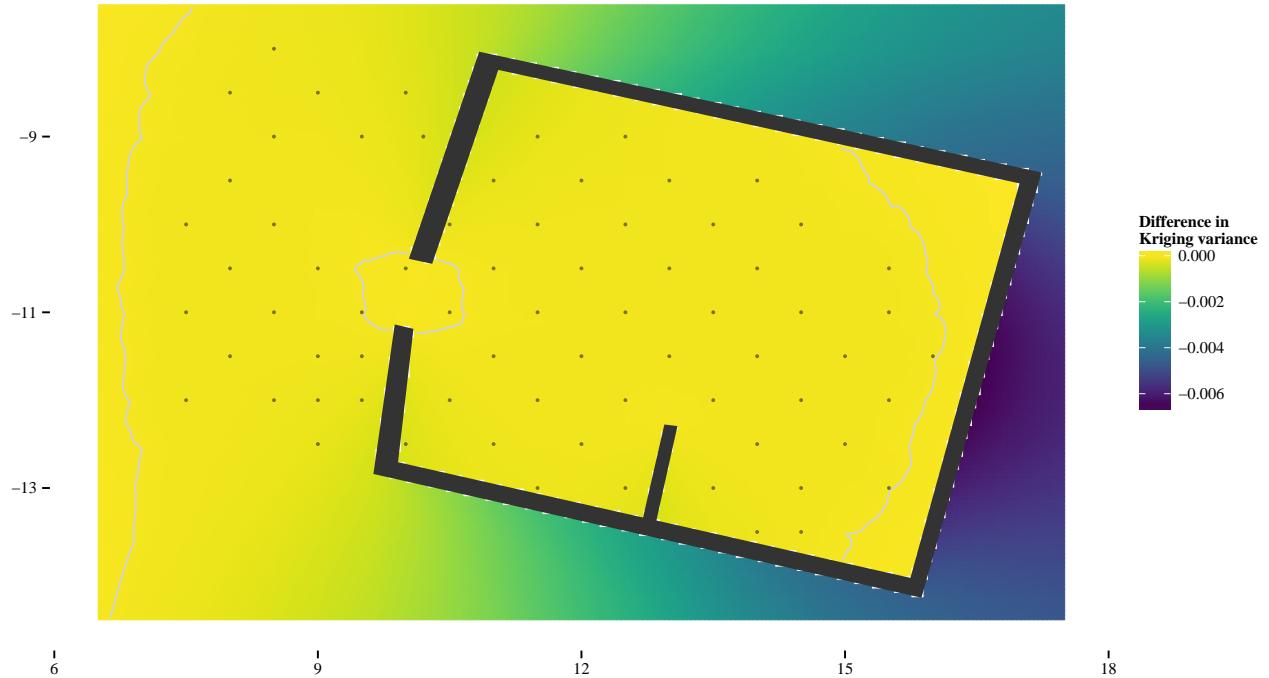


Figure 58: Difference between the Euclidean and the cost-based prediction errors

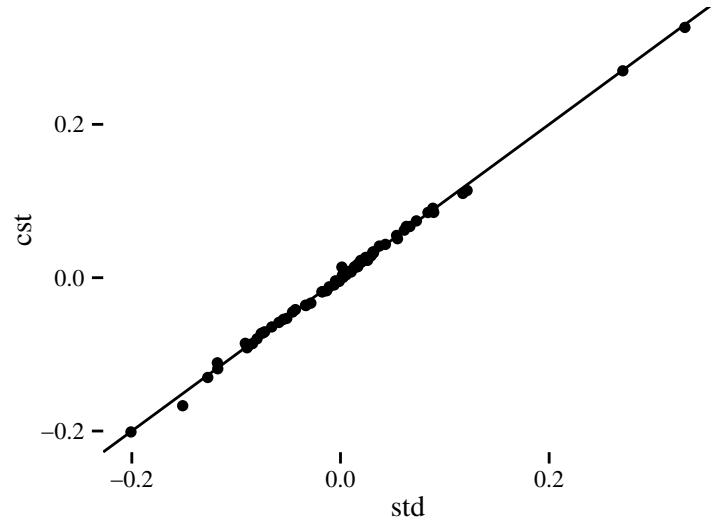


Figure 59: Pointwise leave-one-out prediction error by method.

8 Analysis of Zinc

8.1 Euclidean kriging

The variogram model is Exponential. We choose to estimate the nugget effect, which may account for measurement error, for example.

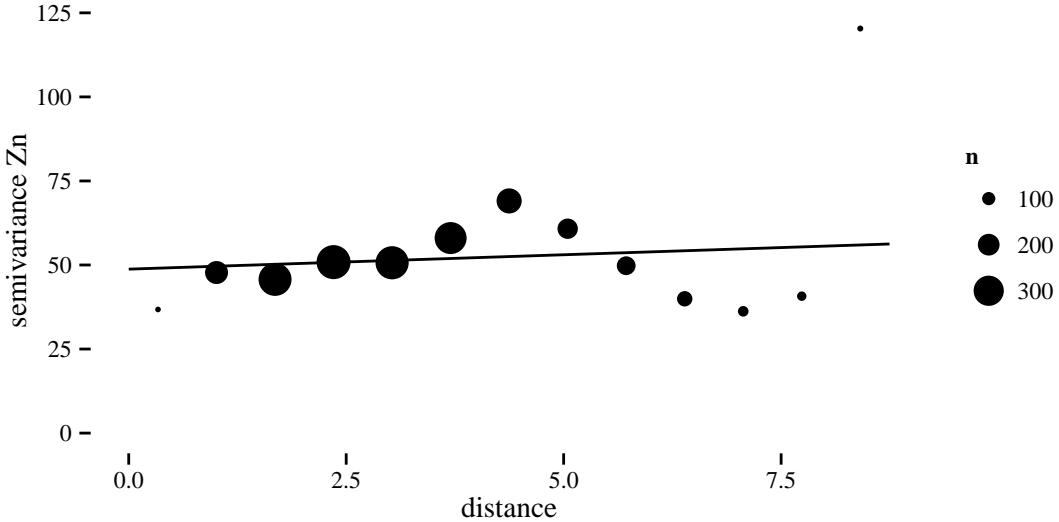


Figure 60: Empirical variogram and fitted model.

8.2 Cost-based kriging

8.3 Comparison of method outcomes

	Euclidean	Cost_based
Intercept	31.71	31.63
Nugget	48.76	48.23
Partial sill	1858.03	2040.56
phi	2167.78	2205.17
Pract. range	6494.08	6606.10
Log-likelihood	-233.61	-233.43

In the scatter plot, the horizontal patterns correspond to predictions on observed values. Otherwise, the differences are negligible.

Near the observations, the cost-based approach has a larger prediction error due to its increased estimation of the nugget (i.e. short-range variance). In the main area, the prediction errors are practically the same with both approaches. Behind the walls, the Euclidean prediction error is unrealistically low.

8.4 Leave-one-out Cross Validation (LOOCV)

method	rmse(error)
std	7.23

method	rmse(error)
cst	7.23

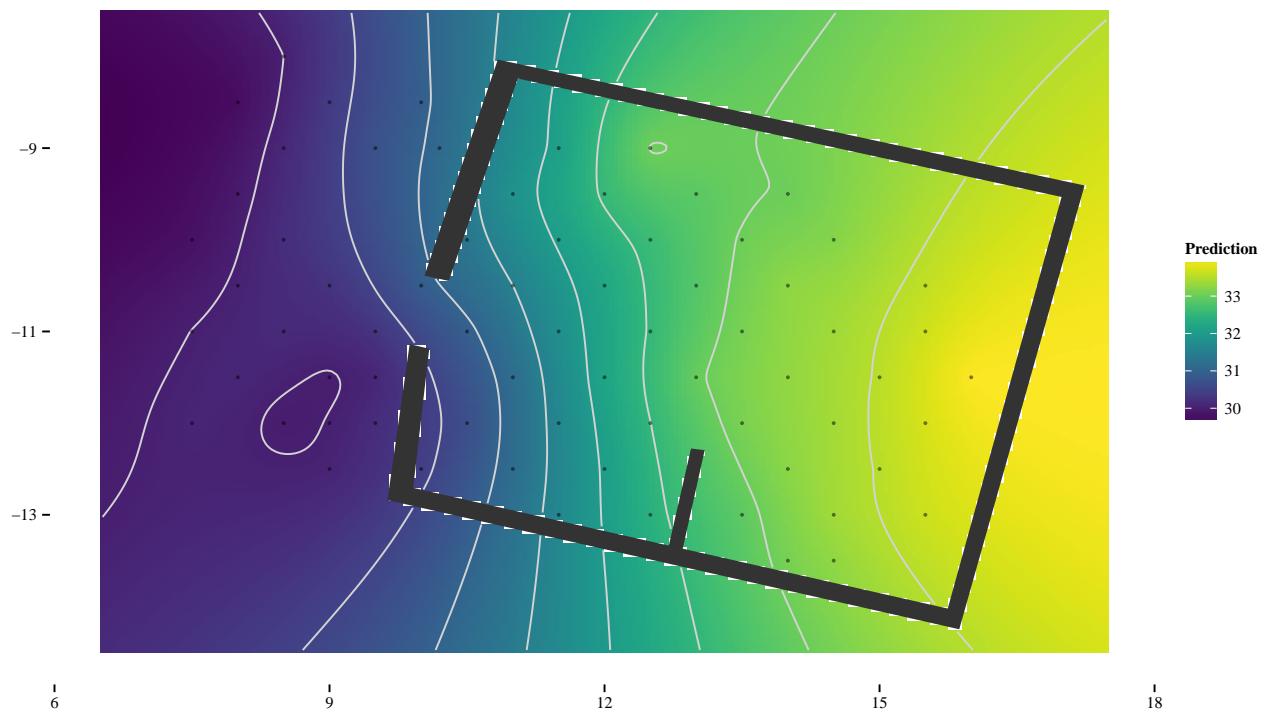


Figure 61: Euclidean kriging prediction

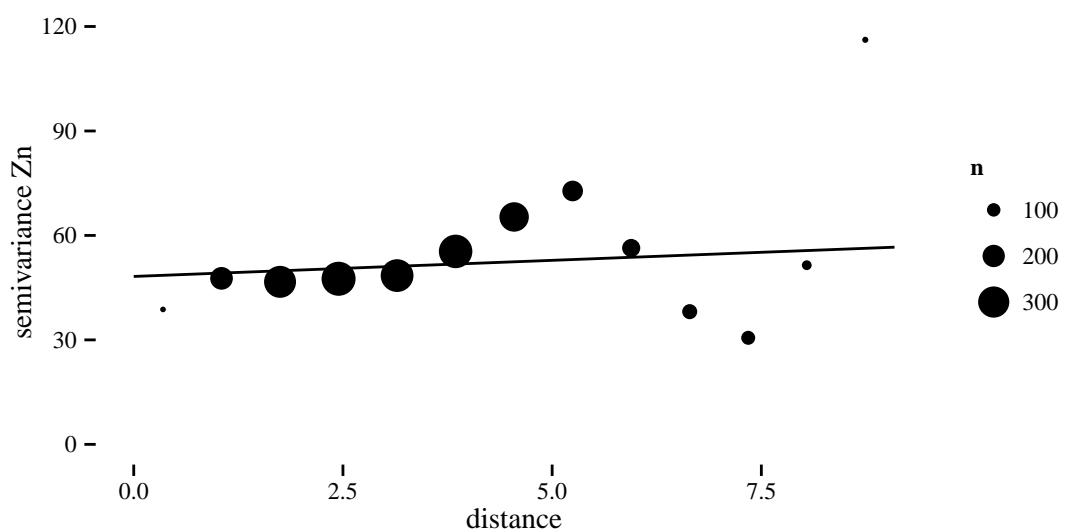


Figure 62: Empirical cost-based variogram and fitted model.

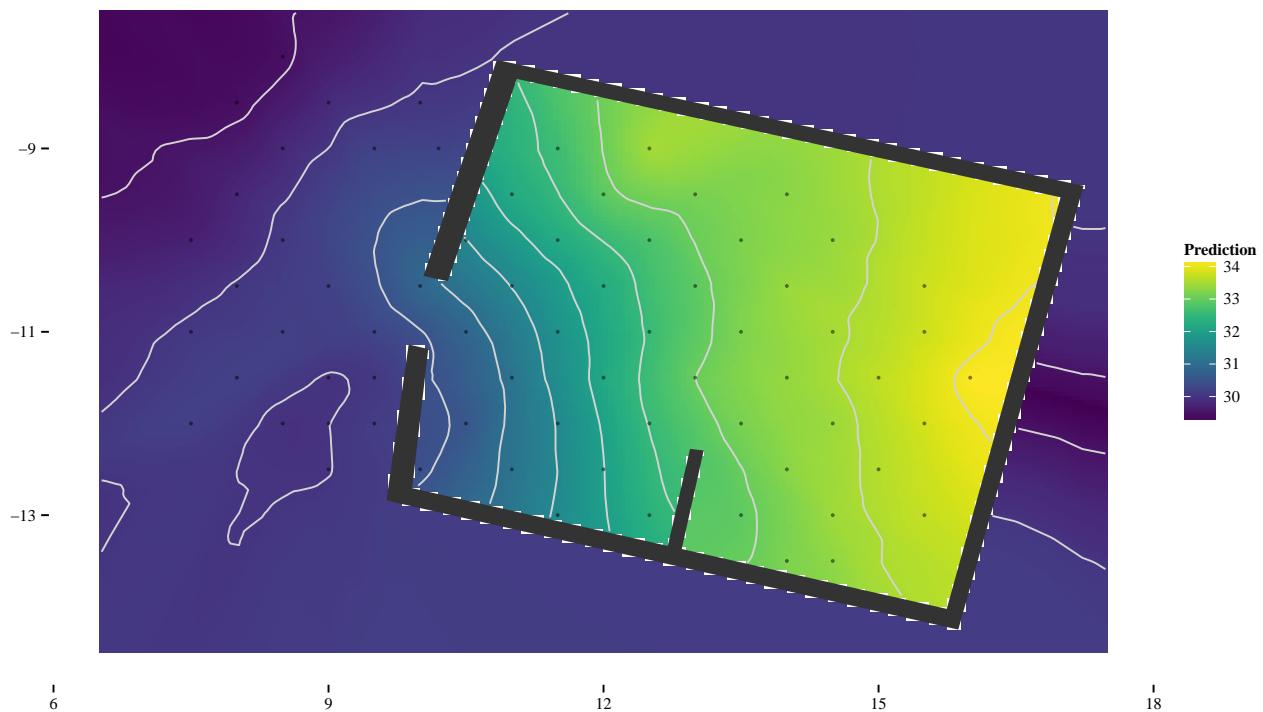


Figure 63: Cost-based kriging prediction

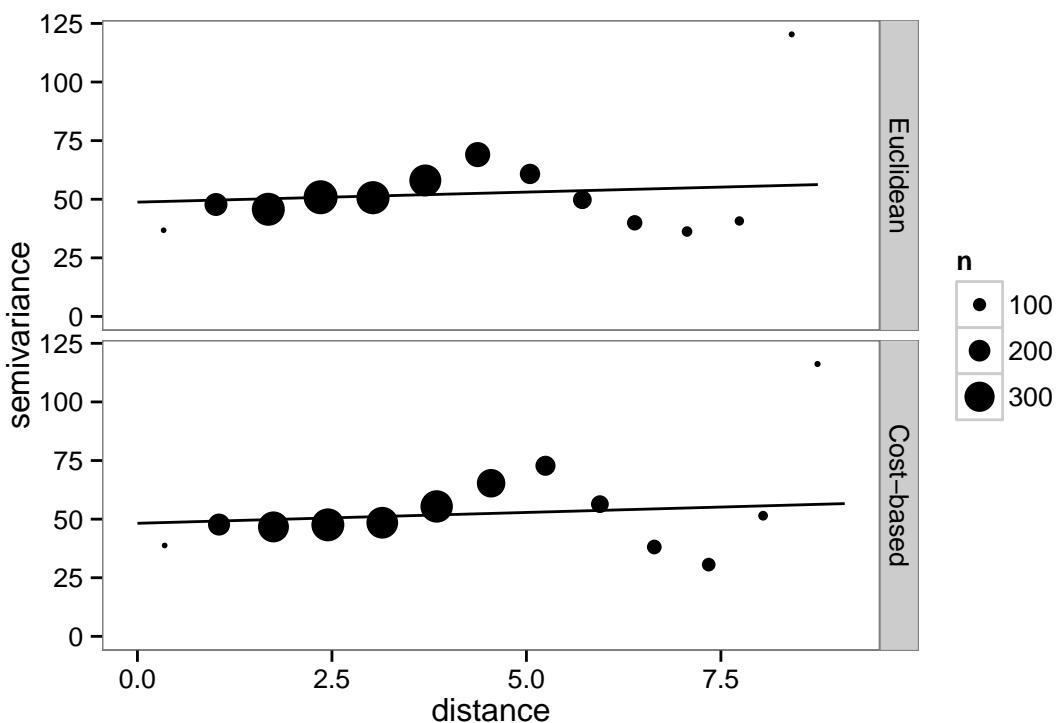


Figure 64: Empirical variogram and fitted models by method for Zinc.

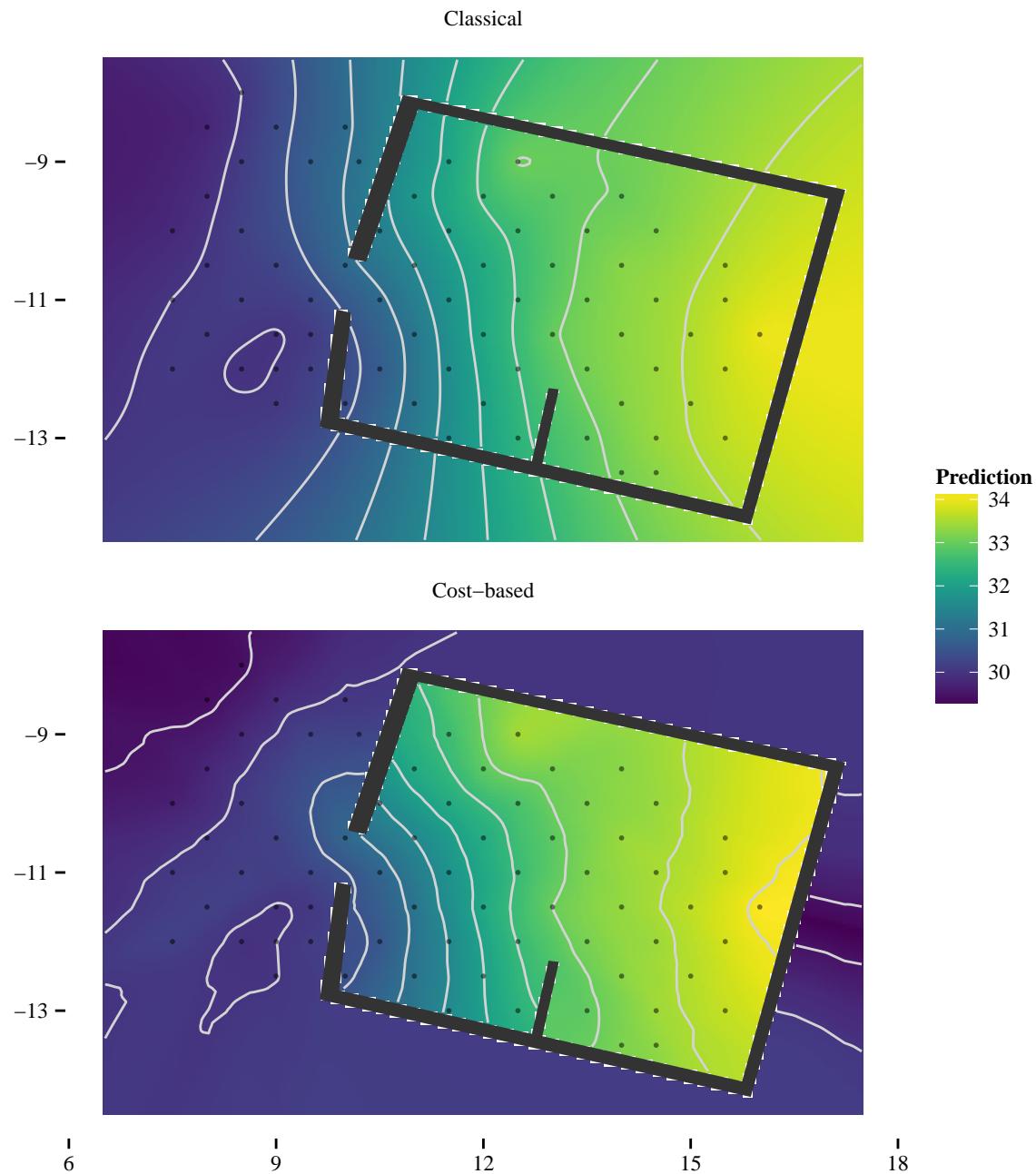


Figure 65: Comparison of Kriging estimates.

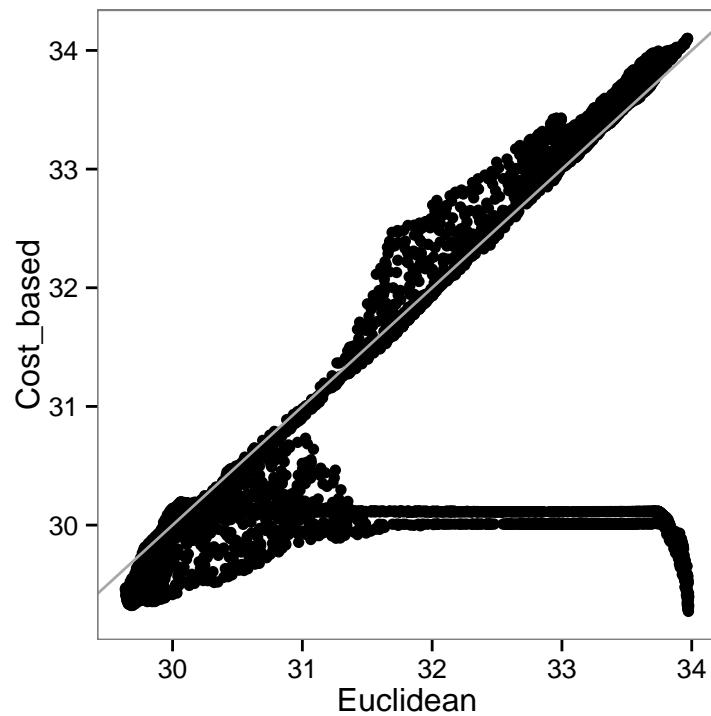


Figure 66: Pointwise comparison of predictions by method.

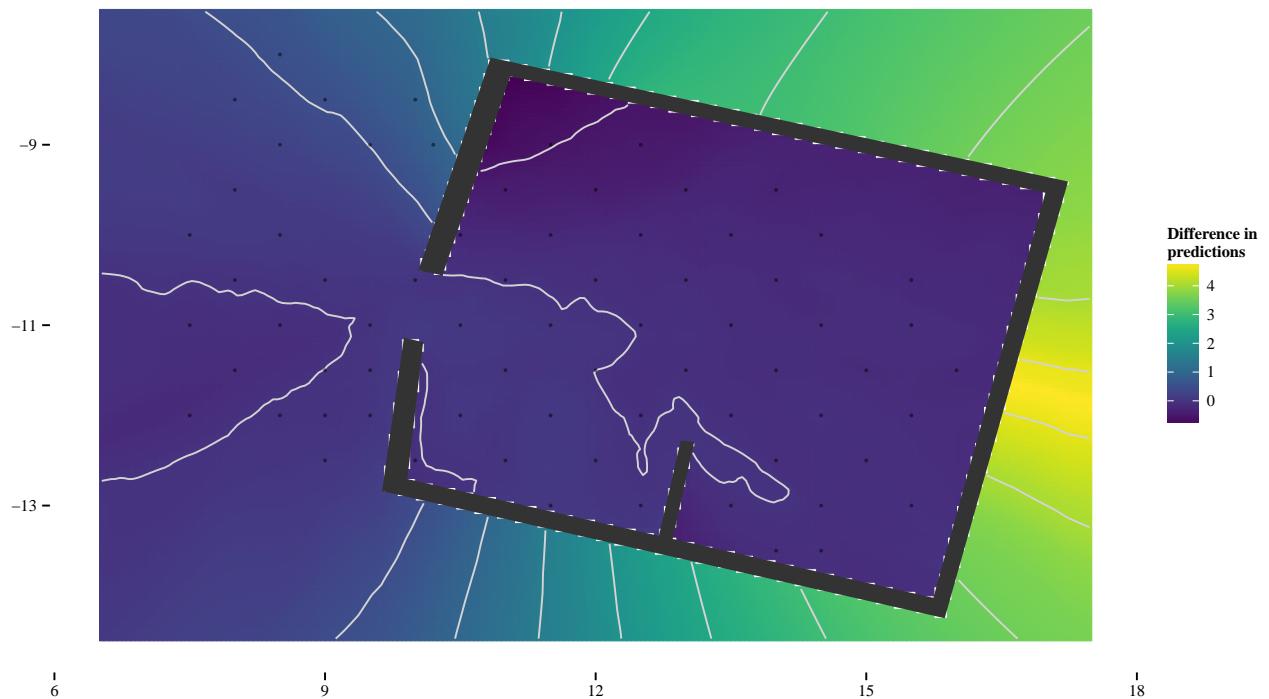


Figure 67: Difference between the Euclidean and the cost-based predictions.

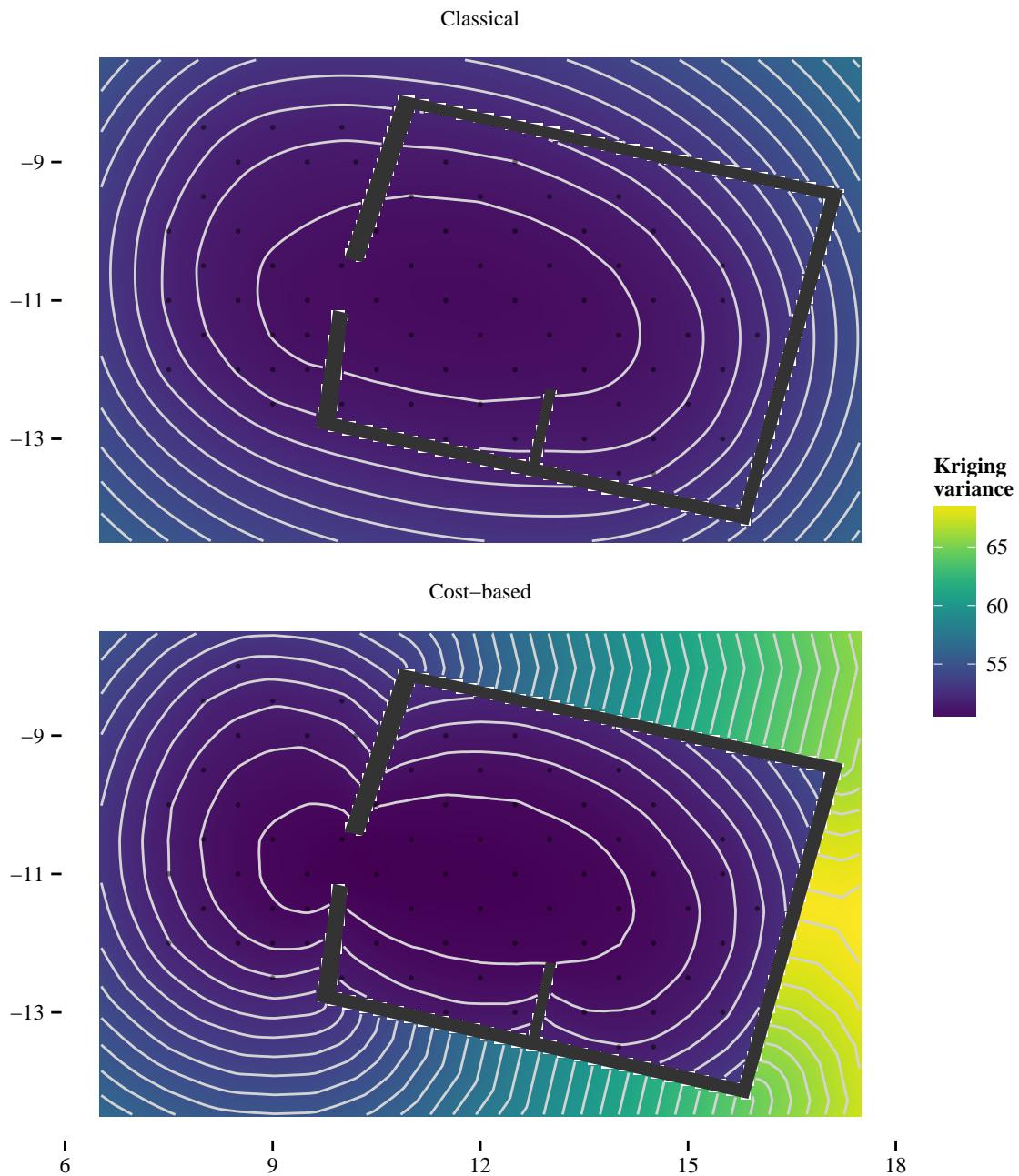


Figure 68: Comparison of prediction error by method.

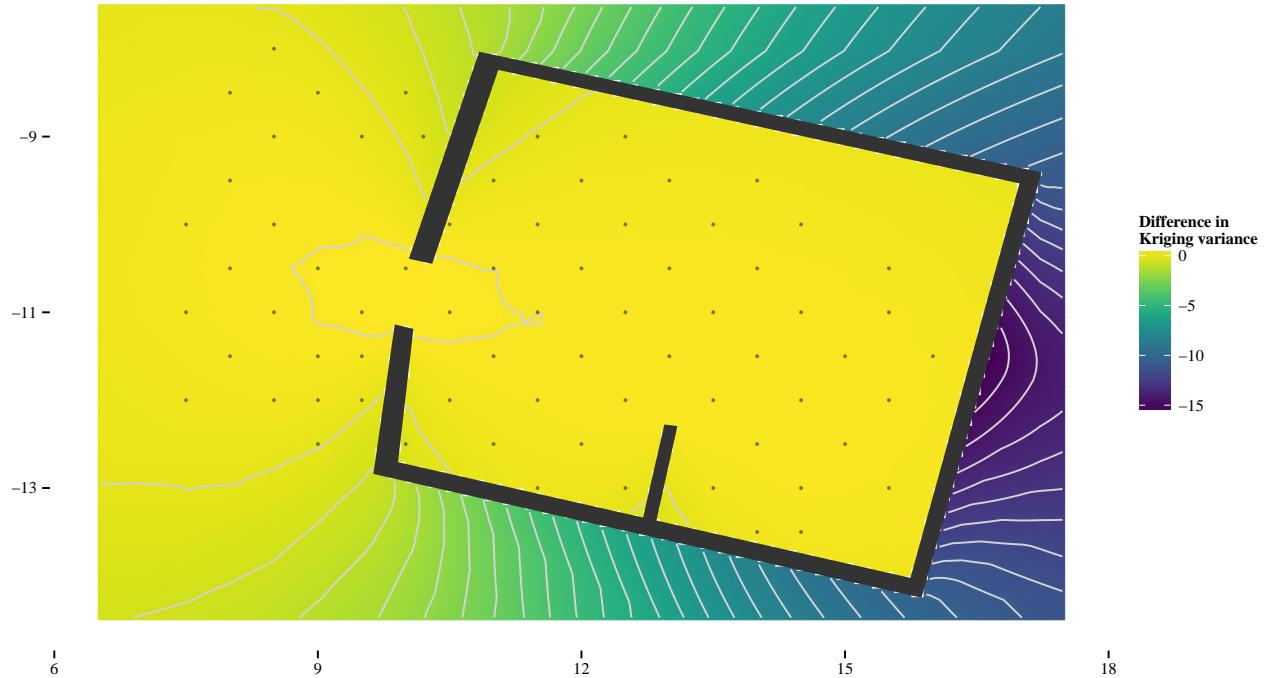


Figure 69: Difference between the Euclidean and the cost-based prediction errors

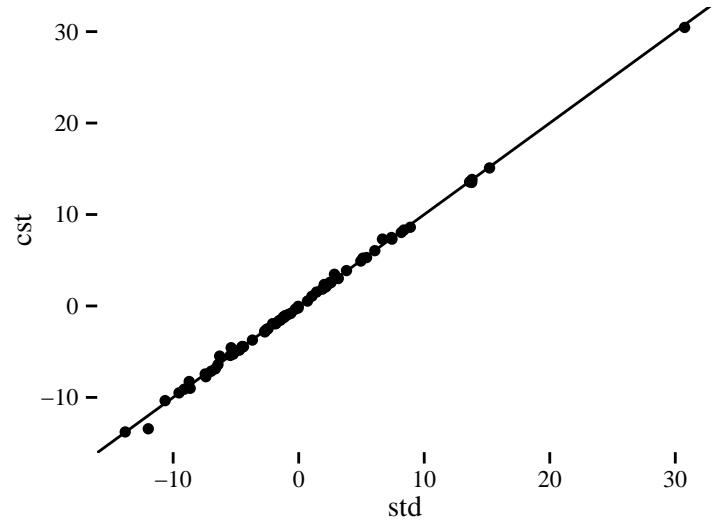


Figure 70: Pointwise leave-one-out prediction error by method.

9 Conclusions

- Actually, the kriging model for Calcium is not adjusting very well the tails of the data, which are heavier than expected. This happens both for the Euclidean and cost-based models. This means that none of both approaches will be really good predictors anyways.

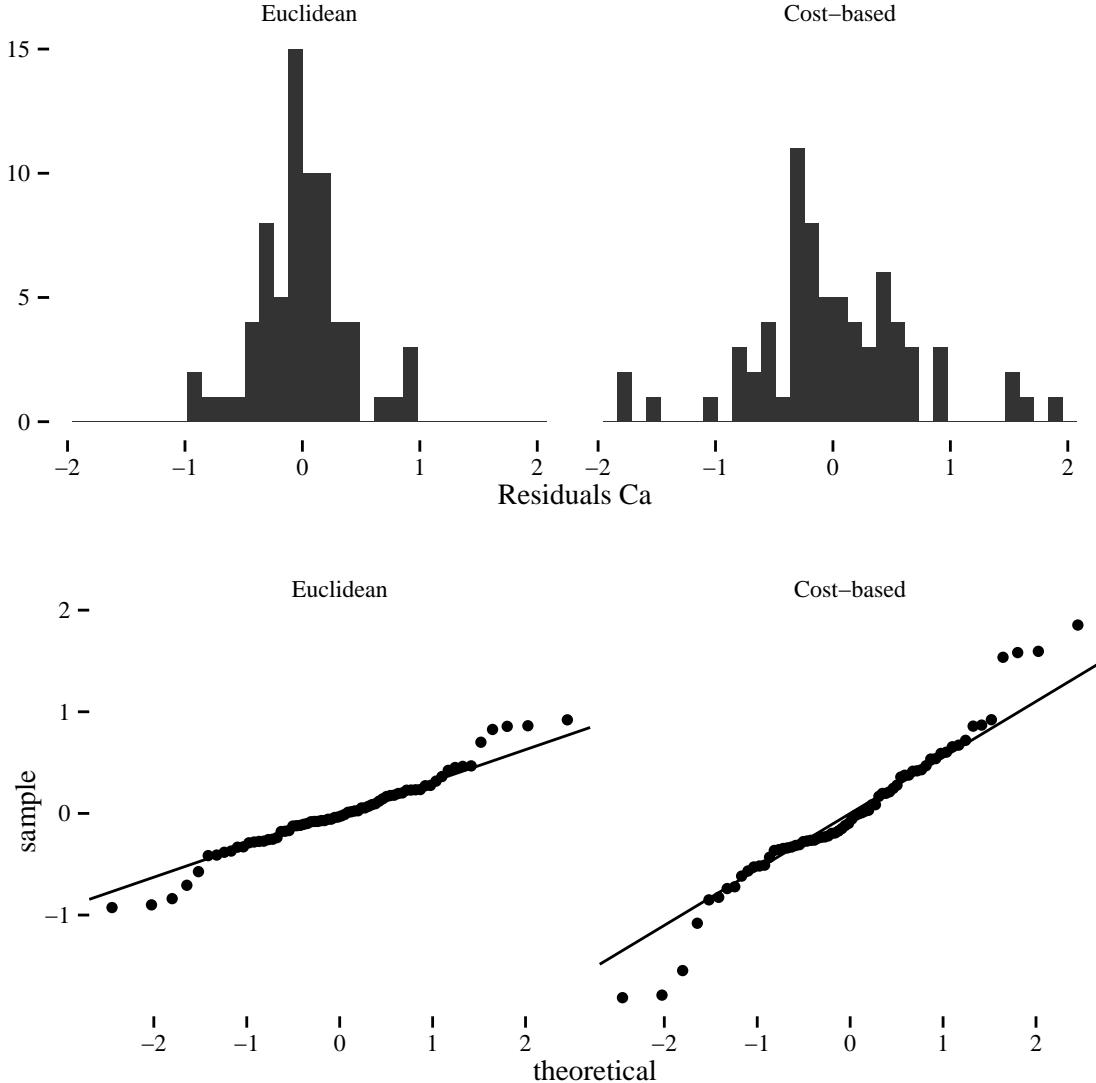


Figure 71: Diagnostics of residuals for Calcium.

- The rest of the variables do not show a clear discontinuity as a consequence of the walls. The Euclidean and cost-based predictions are very similar. Furthermore, the Euclidean and cost-based empirical variograms display practically the same shape in all cases. This suggests that the solid structures are not really affecting the spatial distribution of chemicals.
- Many variables display an initial drop in the semivariance. Even this is based on only 10 or 12 pairs of observations in the first lag, the drop is dramatic for several variables like Copper, Potassium or Magnesium. This may be due to the presence of some extreme values which contrast heavily with neighbouring values. The impact in higher lags is absorbed by the high number of pairs.