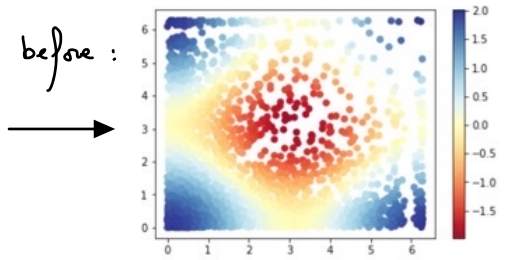
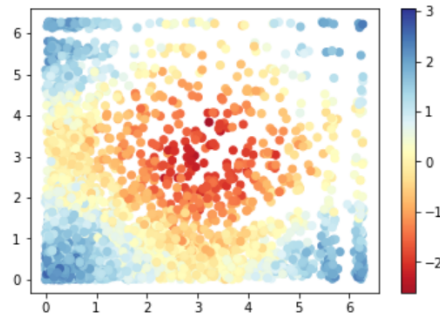
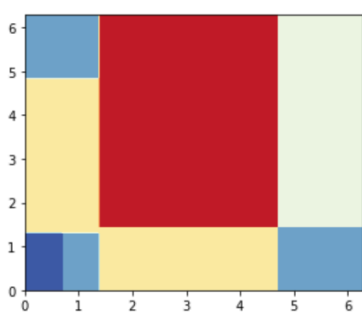


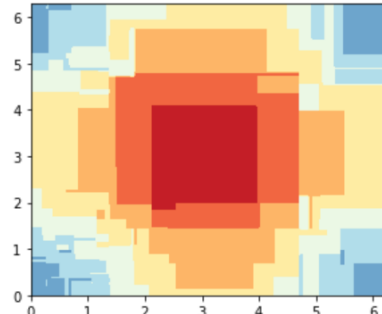
The new data distribution:  
the output (that I try to predict) is more noisy!



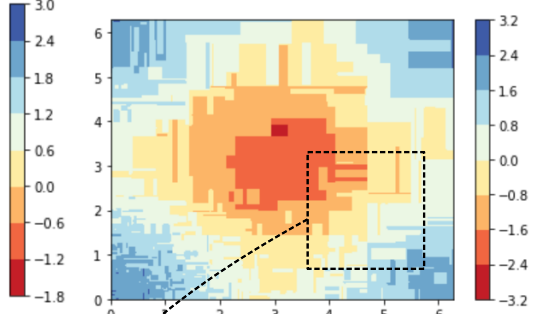
Decision trees:



max\_depth = 3



max\_depth = 8



max\_depth = 14

we fit poorly the train data  
→ "under fitting"

fit is on! ✓  
predictor does not seem "too complicated" ✓  
→ on!

the trained tree fits too much the learning data, and varies too much:



→ over fitting

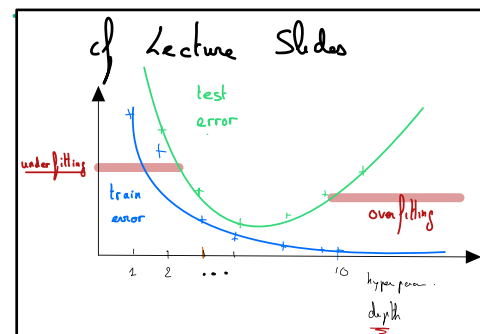
We can check on the cross-validation score (test . score) of a tree over multiple train-test splits:

depth	score
1	0.3059732077354048
2	0.5798514521408162
3	0.7397198103247626
4	0.8207684060499227
5	0.8553084624998055
6	0.875657408869172
7	0.8845319514551624
8	0.8825718728674463
9	0.8790144221038204
10	0.8699275541790745
11	0.858754696223517
12	0.8473440408526696
13	0.8388846695486368
14	0.8322986395600955

underfitting

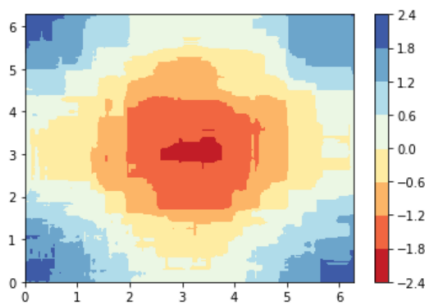
→ optimal depth for a single tree

overfitting

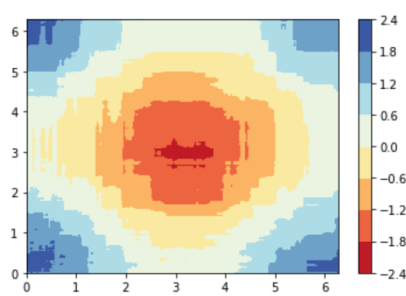


---

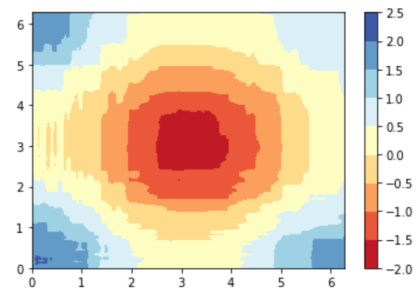
Still, the prediction does not look great - let us use random forests!



with 200 trees



with some parameters  
selected w. RandomizedSearchCV  
(tested 20 values of the 8000 possible  
set of hyperparameters)



with 2000 trees,  
 $\text{min\_samples\_leaf} = 6$ .