

02/01/2012

# RANDOM SEARCH FOR HYPER-PARAMETER OPTIMIZATION

So Lots of people use manual & grid search when looking to find optimal hyper-parameter values.

They start by explaining the optimization problem & saying that it suffers from the curse of dimensionality when you are trying to solve it using grid search.

But they claim they will show how has the practical advantages of Random Search (conceptual simplicity, easy to implement, easy parallelism) & is more efficient in higher dimensional space.



Take:  $f(x, y) = g(x) + h(y) \approx g(x)$

the function has low effective Dimensionality

In this case Random search will work better.

Why Because with grid search we are

searching every Dimension even those that are not important. Due to the randomness

It searches those Relevant Subspaces Just

as if it knew to Just search those.

And we cannot know which ones are relevant & which are not since it changes for each dataset (they found).

So they say they will Redo the results of "Larochelle (2007)" where there they did grid search & here they will do random search.

In (2007) it was grid search but here they define a distribution to mimic the space of (2007) & then they sample from it drawing 256 samples



Now for the results.

They coin this Central trade off.  
Random search. exploration v.s  
exploitation.

In a restricted space less trials "models  
to sample + train" are needed to  
outperform grid search.

But better results can be achieved  
with a less Restrictive sample space  
but more trials are needed.

ex. of exploration v.s. exploitation.

They only needed to do 8 trials  
to outperform grid search with  
no preprocessing.

but when they used preprocessing  
(PCA and Normalization) they needed

32 trials, this indicates that there are  
many bad ways to preprocess.

But when they did 64 trials the  
results were better than those found  
in the restricted smaller space (NO  
preprocessing) but those were found  
faster obviously.



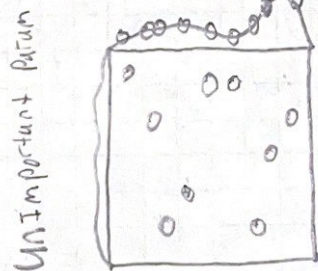
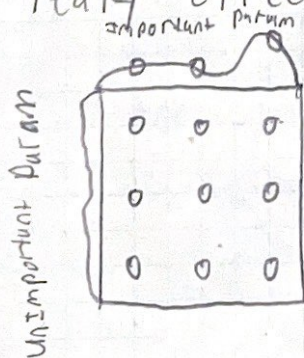
They also found that there are different shapes to the hyper-parameter optimization function  $\Psi$  for different datasets.

ex. In the classic MNIST dataset it seems that after 4-8 trials the best model in each of the 4-8 trials (1 experiment) converges to the same results. Meaning the good region of  $\Psi$  is prob  $1/4$  to  $1/8$  of the entire region.

But for a problem like CIFAR-10 or MNIST Rotated with background images even experiments each with 16 to 32 trials had large variation in results from the best model. This means that here  $\Psi$  has smaller region of good performance.

Now we start with saying why this Random Sampling works well. The claim is this is because  $\Psi$  has a low effective dimension. This is where one or multiple dimensions do not really effect the output.

$$f(x, y) = g(x) + h(y) \approx g(x)$$





Now they will go on to show why  $\mathbb{I}$  often has low effective dimensionality.

LOOK AT Figure #7. They show a small fraction of Hyperparameters matter for any one dataset, but it differs from dataset to dataset which ones matter.

---

Then they note that there might be a better way to sample points than pure Random and they do it with low-discrepancy Sets Sampling methods. Sobol seems to do a little better than random.

---

When it comes to Deep Belief Networks (DBN) Random Search can be competitive but is not consistently outperforming

---

In conclusion Random Search is more efficient than grid Search not because all Hyper-params are as important but precisely because some are not. This is why grid Search gives too many trials to the unimportant dimensions.