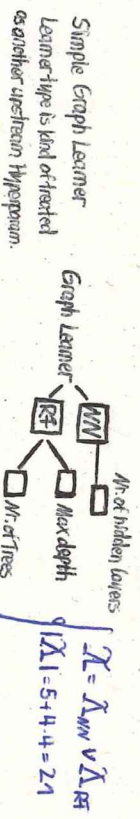


Exercise 1: Recap Nested Resampling

Assume we have a dataset $\mathcal{D} = ((x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}))$ with n observations of a continuous target variable y and p features x_1, \dots, x_p . We want to build a prediction model that can be deployed and we want to estimate the corresponding generalization error. For this, we build a graph learner that consists of a neural network in one arm and a random forest in the other arm. The neural network shall have one hyperparameter, the number of hidden layers; assume the number of nodes per hidden layer and all other possible hyperparameters are fixed. The random forest shall have two hyperparameters, the maximal depth and the number of trees; assume that all other possible hyperparameters are fixed. In total, we pursue three goals (not necessarily in this order):

- Train a final model \hat{f} that can be deployed.
- Tune the graph learner.
- Estimate the generalization error.



Answer the following questions:

1) For each goal:

- Do we need resampling, nested resampling, or no resampling?
- Which fraction of the available dataset can be used?

Final Tuning on full \mathcal{D} : 85 model fts, last one is final \hat{f} on \mathcal{D} !
Overall we need 256+85=340 total model fts
Make computation fast with Parallelization

4) Per Outer Fold: $|\mathcal{A}| \cdot 4 + 1 = 21 \cdot 4 + 1 = 85$ model fts
across three outer layers: 3 · 85 = 255

- Write down a pseudo-algorithm for carrying out all three steps (in a sensible order as derived in 2))
- Assume the number of hidden layers is $\in \{1, 2, 3, 4, 5\}$, the number of trees is $\in \{10, 50, 100, 200\}$ and the maximal depth is $\in \{2, 3, 4, 5\}$. Use 3-fold cross-validation as outer resampling and 4-fold cross-validation as inner resampling. Use grid search and consider all possible hyperparameter combinations. Compute the total number of model trainings carried out in 3).

A) Train final learner \hat{f}

a) No resampling, use full dataset \mathcal{D} for final fit

B) Tune the graph learner

a) Use resampling to estimate \hat{G}_E of different HPs and choose best λ^*

b) Use all data, but with repeated train-test splits (here potentially other splits for each graph branch)

C) Estimate \hat{G}_E of the whole Tuned Learning Algorithm Process

a) Nested Resampling for unbiased \hat{G}_E estimation, outer loop for eval to measure \hat{G}_E (average across all λ^*)

inner loop is for Tuning and refits to B

b) Use all data, but two-level nested train-test splits, one outer k and one inner resampling

3) Pseudo Code | Import data \mathcal{D} , indices $\mathcal{I} \subset \mathcal{D} \times \mathcal{A} \rightarrow \mathcal{H}$ of Perf. measure ρ , HP Search space \mathcal{A}
Split \mathcal{D} into $(\mathcal{D}_{train}^{(i)}, \mathcal{D}_{test}^{(i)}) \leftarrow \frac{1}{3}, \frac{2}{3}$ times (3-fold CV in outer resampling)
for \mathcal{I} :
split $\mathcal{D}_{train}^{(i)}$ into $(\mathcal{D}_{train}^{(i,k)}, \mathcal{D}_{valid}^{(i,k)}) \leftarrow \frac{1}{4}, \frac{3}{4}$ times (4-fold CV in inner resampling)
for each $\lambda \in \mathcal{A}$: Train $\pm (\mathcal{D}_{train}^{(i,k)}, \lambda)$ and measure inner performance on $\mathcal{D}_{valid}^{(i,k)}$
Grid search: [for each $k \in \mathcal{K}$: Calc. $\hat{G}_E = \frac{1}{|\mathcal{A}|} \sum_{\lambda \in \mathcal{A}} \rho(\mathcal{D}_{valid}^{(i,k)}, \lambda)$, $\hat{G}_E^{(i,k)} = \frac{1}{|\mathcal{A}|} \sum_{\lambda \in \mathcal{A}} \rho(\mathcal{D}_{valid}^{(i,k)}, \lambda)$]
End for loop: select $\lambda^* = \argmin_{\lambda} \hat{G}_E^{(i,k)}$
• Refit λ^* on all $\mathcal{D}_{train}^{(i)} \Rightarrow \hat{f}_i = \pm (\mathcal{D}_{train}^{(i)}, \lambda^*)$
Do outer-test evaluation $\hat{G}_E^{(i)} = \rho(\mathcal{D}_{test}^{(i)}, \hat{f}_i)$
End for loop, Estimate overall Generalization Error $\hat{G}_E = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \hat{G}_E^{(i)}$

* Tune on full dataset (same grid search, 4-fold CV on \mathcal{D})
yielding us minimal $\lambda^* \in \mathcal{A}$ (Goal B)
Final fit: Train once on 100% of \mathcal{D} !
 $\hat{f} = \pm (\mathcal{D}, \lambda^*) \leftarrow \text{Goal A}$

