

Regularized regression I

Stepwise Variable Selection

Stepwise Variable Selection

Algorithm 1: Classical forward selection

```
1 Set  $p$ -value threshold  $\tau$ ;  
2 Initialize empty model;  
3 repeat  
4   for each predictor not in the model do  
5     | Add the predictor to the current model;  
6     | Estimate the statistical significance of the new term;  
7   end  
8   if the smallest  $p$  is less than  $\tau$  then  
9     | Include the corresponding predictor in the model;  
10  end  
11 until no significant predictor remains outside the model;
```

Stepwise Variable Selection

There are a number of **serious** problems here!

- ① Multiple testing issue
- ② Objective function does not focus on prediction accuracy
- ③ Prone to performance evaluation bias

Adjusting stepwise selection approaches

- ① Usage of performance measures instead of p -values
- ② Implement feature selection in a (proper) **resampling setting**
- ③ Interweave feature selection in the model-building process

Stepwise Variable Selection

Algorithm 2: Forward selection with resampling

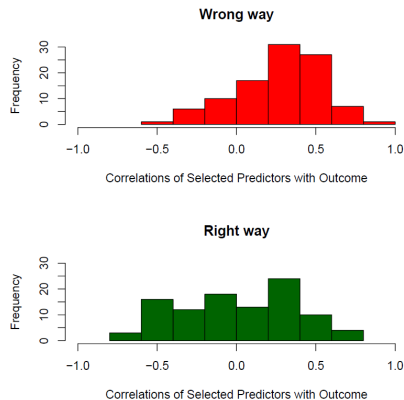
- 1 Set the number of resampling iterations;
 - 2 Set the number of features p ;
 - 3 Initialize empty model;
 - 4 **for** *each resampling iteration* **do**
 - 5 Partition data into training and hold-out set;
 - 6 **for** $k = 0, \dots, p - 1$ **do**
 - 7 Consider all $p - k$ models that add an additional predictor to the current model;
 - 8 Choose the best among these models in terms of loss in the training data;
 - 9 Evaluate the chosen model in the hold-out set;
 - 10 **end**
 - 11 **end**
 - 12 Determine the best number of predictors over all hold-out sets;
-

Stepwise Variable Selection

Cross-Validation done wrong

- **Never** separate feature selection and CV
 - CV *after* selection on full data *biases* performance measures
 - Hold-out samples are no longer independent test sets
- Include feature selection within the CV loop
- Unsupervised screening on full data is valid

Figure: Correlations of y with unrelated x 's with incorrect and correct CV



Hastie et al. (2009)