

#1

a)

Test-error criterion can be used to indirectly estimate training error by making an adjustment to the training error to account for the bias due to overfitting.

Advantages: works well for simple models such as linear regression; more interpretability w/ the adjustment on training error; closed form criterion, can be computed efficiently and easy to apply.

Disadvantages: sometimes perform poorly if the adjustment made for training error is far away from the true data set or data distribution.

CV can be used to directly estimate the test error.

Advantages: works for a wider range of models such as complicated non-linear models; provides a direct estimate; no need to estimate irreducible noise variance σ^2 .

Disadvantages: results may be highly variable because it depends on which part is used for training; computationally expensive

Therefore, when we are working w/ a simple model (e.g. linear regression), when computational efficiency is a priority, or when dataset is small, we could use test-error criterion. When we have a more complex / non-parametric model, when dataset is large, or when a more accurate out-of-sample error is required, CV might be a better option.

b)

$$AIC = \frac{1}{n} (RSS + 2d\hat{\sigma}^2) \quad > \quad \begin{aligned} \text{where } n &= \# \text{ observations} \\ d &= \# \text{ predictors} \\ \hat{\sigma}^2 &= \text{Variance of the } \epsilon \text{ associated w/} \\ &\text{each response.} \end{aligned}$$

$$BIC = \frac{1}{n} (RSS + \log(n)d\hat{\sigma}^2)$$

Formula difference: BIC replaces 2 w/ log(n).

$\log(n) \geq 2$ when $n > 7$, thus BIC places a heavier penalty on model when there are more than 7 observations, resulting in a model w/ less predictors.

AIC would be preferable when prediction accuracy is the goal, and BIC should be chosen if we want a simpler and more interpretable model.

c)

penalizing more parameters (d): adding more parameters can make a model fit the training data more closely but it may also capture noise (overfitting). The penalty term discourages very complex models by adding a cost for each additional parameter.

penalizing larger estimated irreducible noise variance ($\hat{\sigma}^2$): a larger $\hat{\sigma}^2$ indicates that the model is not capturing the underlying pattern in the data well, suggesting the noise is high. Penalizing based on $\hat{\sigma}^2$ discourages models that do not significantly reduce the unexplained variability in the data.

d) and lower bias

$\hat{f}_1(AIC)$ is likely to have higher variance than $\hat{f}_1(BIC)$ because BIC tends to favor simpler models ^{than AIC} due to its heavier penalty for additional parameters ($BIC > AIC$ when $n > 7$), and we know that a model with too many parameters may have low bias but high variance (overfitting), while a model with too few parameters may have high bias but low variance (underfitting).

2

a)

They each need 3, 5, 7, and 9 parameters, respectively. For a polynomial of degree M , we need $2M+1$ parameters.

b)

| | Bias | Variance |
|----------|-----------|-----------|
| Degree 1 | large | small |
| 2 | | |
| 3 | ↓ similar | ↓ similar |
| 4 | ↓ smaller | ↑ larger |

The bias is bigger with model of degree 1 and decrease with degree 2 and degree 3 (degree 4 looks very similar to degree 3, but bias should be getting smaller) as the fitted line moves closer to the Bayes optimal.

The variance is smaller with degree 1 model and increases as more parameters are added (degree 2 \rightarrow degree 3 \rightarrow degree 4) because model complexity brings more variance to the fit.

c)

True optimal model : order = 4

Est. optimal model : order = 3

The black curve is generally better than the orange curve because the 10-fold CV overestimated the test-error because splitting data into $k=10$ folds and using 90% of the data for training introduces bias that lead to an overestimation of the test error. Instead of 100%

d)

Way #1 : # of parameters. If we increase the number of predictors, the model complexity will increase.

Way #2 : Interaction terms. If we add terms of interactions b/w 2 variables that account for situations where the effect of one variable on the outcome is affected by the level of another variable, thus adding more complexity to the model.

In conclusion, both ways — increasing the # of parameters and including interaction terms — allows the model to have more degrees of freedom and can fit the training data more closely, though there is the concern of overfitting.

#4

a)

The model chosen by the best subset method will have the smaller training RSS because it considers all possible models (2^k of them) and chooses the one w/ smallest RSS.

b)

It's not possible to definitively say which method (best subset, forward stepwise, or backward stepwise) will provide the model w/ the smallest test RSS for k predictors w/o additional info or testing. This is because the model w/ the smallest test RSS is contingent upon the true underlying data-generating process, which might favor one method over the others in certain contexts. CV or other model validation techniques would typically be used to estimate test error rates and select the best model.

c)

- i. True. In forward stepwise selection, we start w/ a model w/ no predictors and sequentially add predictors that provide the best improvement to the fit.
- ii. True. In backward stepwise selection, we start w/ a model containing all predictors and sequentially remove the predictor that has the least impact on the fit.
- iii. False. Forward and backward stepwise selection do not necessarily yield the same predictors because they add/remove predictors based on different criteria at each step.
- iv. False. Explanation same as iii.
- v. False. Best subset selection chooses the best model at each level k based on some criterion (like RSS/R²) w/o considering the models at other levels.