Chapter 5: Assessing Model Accuracy

One of the key aims of this course is to introduce you to a wide range of statistical learning techniques. Why so many? Why not just the "best one"?

There is no BEST one for every situation! Lounless you know the model the dota came from (which you won't).

Hence, it's important to decide for any given set of data which method produces the best results.





https://xkcd.com/1838/

1 Measuring Quality of Fit

With linear regression we talked about some ways to measure fit of the model

R², Residual standard error.

In general, we need a way to measure fit and compare *across models*.

La not just linear repression.

One way could be to measure how well its predictions match the observed data. In a regression session, the most commonly used measure is the mean-squared error (MSE)



We don't really care how well our methods work on the training data.

Instead, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen data. Why?

test data

We already know the response values for fraining data.
Suppose we fit our modul to training data
$$\{(\underline{z}_i, y_i), \dots, (\underline{z}_n, y_n)\}$$
 and obtain \hat{f} .
We can support $\hat{f}(\underline{z}_i), \dots, \hat{f}(\underline{z}_n)$ if these are does to $y_{i_1, \dots, y_n} \Longrightarrow$ small training MSE.
But we care about
 $\hat{f}(\underline{z}_0) \simeq y_0$ for (\underline{z}_0, y_0) unseen data not used to fit the model
Want to discuss the model that gives locast test MSE
Average $[(y_0 - \hat{f}(\underline{z}_0))^2]$
for a large $\#$ of fest observations (\underline{z}_0, y_0) .

So how do we select a method that minimizes the test MSE?

Sometimes we have a test data cet available to us based on the scientific problem is access the set of dis. not used to fit the model.

But what if we don't have a test set available?

Maybe we just minimize train MSE?

30

25

>²⁰

15

"Problem : there is no guarantee bowerby traching MSE bowers test MSE!

because many stat learning methods estimate coef's to lover fraining MSE





3

frue touship

kinear regression

df = 2

1.1 Classification Setting

So far, we have talked about assessing model accuracy in the regression setting, but we also need a way to assess the accuracy of classification models.

Suppose we see to estimate f on the basis of training observations where now the response is categorical. The most common approach for quantifying the accuracy is the training error rate.

 $\frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(y_i \neq \hat{y}_i) \quad \text{where} \quad \mathbb{I}(y_i \neq \hat{y}_i) = \begin{cases} 1 & i \neq y_i \neq \hat{y}_i \\ 0 & o.w. \end{cases}$ $Pabel \text{ for } i^h \quad \text{for } i^h \text{ observation}$ $Pabel \text{ for } i^h \quad \text{for } i^h \text{ observation}$

This is called the *training error rate* because it is based on the data that was used to train the classifier.

As with the regression setting, we are mode interested in error rates for data *not* in our training data, i.e. Lest lat_{α} ($\mathcal{I}_{o}, \gamma_{o}$)

The fast error value is
Average
$$(I(y_0 \neq \hat{y}_0))$$

T
prelicited dass for fast dos. $v/$ predict to

A good classifier is one for which the fost error rate is small.

1.2 Bias-Variance Trade-off

The U-shape in the test MSE curve compared with flexibility is the result of two competing properties of statistical learning methods. It is possible to show that the expected test MSE, for a given test value x_0 , can be decomposed include

overage test MSE we would obtain if he repeatedly est. f at many training data Gets and predict Zo.

 $\rightarrow E\left[\left(y_{0}-\hat{S}(\underline{x}_{0})\right)^{a}\right] = Var\left(\hat{f}(x_{0})\right) + Bias\left(\hat{f}(x_{0})\right)^{a} + Var\left(\varepsilon\right)^{er}$ Leta ξ_{0} .

This tells us in order to minimize the expected test error, we need to select a statistical learning method that siulatenously achieves *low variance* and *low bias*.

Variance - the amount by which f would change if we estimated it u/ diffect thanking data. In gareal, more flocible methods have higher variance because they fit the transing data solocity ⇒ new data mean big changes in f.
Bias - the error that is introduced by approximating a real hife problem by a much simpler model:
Ox. linear regression assumes a linear form. It is unlikely that any real world data are actually linear ⇒ there will be some bias.

In general: 1 flexibility => 1 bies + 1 variance.

how much these charge determine fast MSE.

Similar ideas hold for dassification setting and tost error rule.

2 Cross-Validation

As we have seen, the test error can be easily calculated when there is a test data set available.

```
Unfortunately this is not always be case.
```

In contrast, the training error can be easily calculated.

```
But trading error can wildly under estimate test error rate.
```

In the absense of a very large designated test set that can be used to estimate the test error rate, what to do?

```
Split up data you already have (transvzdata).
```

For now we will assume we are in the regression setting (quantitative response), but concepts are the same for classification.

2.1 Validation Set

Suppose we would like to estimate the test error rate for a particular statistical learning method on a set of observations. What is the easiest thing we can think to do?

We could randomly divide the available data noto two ports : training + validation.





Let's do this using the mpg data set. Recall we found a non-linear relationship between displ and hwy mpg.



We fit the model with a squared term displ², but we might be wondering if we can get better predictive performance by including higher power terms!

displ³, displ⁴

```
## get index of training observations
              # take 60% of observations as training and 40% for validation
              mpg val <- validation_split(mpg, prop = 0.6)</pre>
                                                          × splits data in 2 parts randomly.
              ## models
              lm spec <- linear reg()</pre>
              linear recipe <- recipe(hwy ~ displ, data = mpg)</pre>
              quad recipe <- linear recipe |> step mutate(displ2 = displ^2)
              cubic recipe <- quad recipe |> step mutate(displ3 = displ^3)
              quart recipe <- cubic recipe |> step mutate(displ4 = displ^4)
              m0 <- workflow() |> add model(lm spec) |> add recipe(linear recipe) |>
                       fit resamples(resamples = mpg val)
              m1 <- workflow() |> add model(lm spec) |> add recipe(quad recipe) |>
                       fit_resamples(resamples = mpg_val)
pulid somethis test the
              m2 <- workflow() |> add model(lm spec) |> add recipe(cubic recipe) |>
                       fit resamples(resamples = mpg val)
              m3 <- workflow() |> add_model(lm_spec) |> add_recipe(quart_recipe) |>
                       fit_resamples(resamples = mpg_val)
              ## estimate test MSE
              collect metrics(m0) |> mutate(model = "linear") |>
                bind rows(collect metrics(m1) |> mutate(model = "quadratic")) |>
                bind_rows(collect_metrics(m2) |> mutate(model = "cubic")) |>
                bind_rows(collect_metrics(m3) |> mutate(model = "quartic")) |>
                select(model, .metric, mean) |>
                pivot wider(names_from = .metric, values_from = mean) |>
                select(-rsq) |>
                                                       5 root MSE
                kable()
                                          model
                                                      rmse
                                          linear
                                                  4.318968
                                          quadratic 3.882112
```

cubic

quartic

3.866194

3.860612 <- lodes like best model.

8

linear would

specifict



- The validation estimate of the test error is highly variable! Depends on which dis. we hold out. - Only a subset used to fit the model. Artificially reducing sample size.

cross-validation is a method to address frese mathemasses ...

2.2 Leave-One-Out Cross Validation

Leave-one-out cross-validation (LOOCV) is closely related to the validation set approach, but it attempts to address the method's drawbacks.



The LOOCV estimate for the test MSE is

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^{n} MSE_{i} = \frac{1}{n} \sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}$$

LOOCV has a couple major advantages and a few disadvantages. over the validation wethod,

- Advantages - since we fit using n-1 observations (instead of $\frac{n}{2}$ for the volidation approach), => LOOCV does not overestimate the true test error as much as volidation approach.
- No randomness in the approach. => will got be some answer every time.

Dîsadvantage

- Sometimes stat learning models can be expensive the fit (i.e. on order of days) LOOCV requires us the fit the model in times. => could be glow.



The k-fold CV estimate is computed by averaging

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^{k} MSE_{i} = \frac{1}{k} \sum_{i=1}^{k} \frac{1}{|F_{i}|} \sum_{j \in F_{i}} (\gamma_{i} - \hat{\gamma}_{j})^{2}$$

estante MSE based on its fold.

Usually use k=5 or k=10. Why k-fold over LOOCV? LOOCV is a special case of k-fold inchich k=n.

Computational advantage! Now have to fit model k times (not n).

Another less obvieus advantage due to bias-voriance trade-off (ome back to later).

```
## perform k-fold on the mpg dataset
                                Change number of filds!
mpg_10foldcv <- vfold_cv(mpg, v = 10)</pre>
## models
m0 <- workflow() |> add_model(lm_spec) |> add_recipe(linear_recipe) |>
        fit_resamples(resamples = mpg_10foldcv)
m1 <- workflow() |> add_model(lm_spec) |> add_recipe(quad_recipe) |>
        fit resamples(resamples = mpg 10foldcv)
m2 <- workflow() |> add_model(lm_spec) |> add_recipe(cubic_recipe) |>
         fit_resamples(resamples = mpg_10foldcv)
m3 <- workflow() |> add_model(lm_spec) |> add_recipe(quart_recipe) |>
        fit_resamples(resamples = mpg_10foldcv)
## estimate test MSE
collect_metrics(m0) |> mutate(model = "linear") |>
  bind rows(collect metrics(m1) |> mutate(model = "quadratic")) |>
  bind rows(collect metrics(m2) |> mutate(model = "cubic")) |>
  bind rows(collect metrics(m3) |> mutate(model = "quartic")) |>
  select(model, .metric, mean) |>
  pivot wider(names from = .metric, values from = mean) |>
  select(-rsq) |>
  kable()
```

	model	rmse	
	linear	3.805566	
	quadratic	3.432052	
	cubic	3.409391	K dos
	quartic	3.408420	
-			



When we perform CV we can be interested in estimating test error. Most after we use it the field minimum estimated test error the help us Eluoose a model (or model parameter)

2.4 Bias-Variance Trade-off for k-Fold Cross Validation

k-Fold CV with k < n has a computational advantage to LOOCV.

There is a less obvious advantage (but potentially more important),

- K-fold often gives more accurate estimates of test error than LOGCV!

We know the validation approach can overestimate the test error because we use only half of the data to fit the statistical learning method.

By this logic, LOOCV gives approximately unbiased estimates of the test error Cuses n-1 x n points to fit).
k-fold gives notomediate level of bias (uses (k-1) n obs to fit))
⇒ Loo cv gives lowest bias.

But we know that bias is only half the story! We also need to consider the procedure's variance. $\Gamma_{\chi, \varsigma} \times \chi$

LOO CV has higher variance then k-fold CV when k < n. War $\left(\frac{x+y}{2}\right) = \frac{1}{4} \left[Var x + Var y + 2 cou (x, y) \right]$ Why? LOOCV fit is modules on almost identical desta points \Longrightarrow averages outputs highly concluded \neg each other! k-fold averages k outputs w more different discriminants (or clap is smaller).

mean of highly concluded quantities has higher variance than mean of less correlated quantities! >> LOOCV has higher variance than K-fild CV!

To summarise, there is a bias-variance trade-off associated with the choice of k in k-fold CV. Typically we use k = 5 or k = 10 because these have been shown empirically to yield test error rates closest to the truth.

in numerical experiments.

2.5 Cross-Validation for Classification Problems

So far we have talked only about CV for regression problems.

Use MSE to quantify test ever.

But CV can also be very useful for classification problems! For example, the LOOCV error rate for classification problems takes the form

$$CY_{n} = \frac{1}{n} \sum_{i=1}^{n} Err_{i}^{i}$$

where $Err_{i}^{i} = \frac{1}{n} (y_{i} \neq \hat{y}_{i}^{i}) = \begin{cases} 1 & i \neq y_{i} \neq \hat{y}_{i}^{i} \\ 0 & 0.0. \end{cases}$

K-fold and validition errors estimated accordingly.



2.5 Cross-Validation for Classifi...



Minimum CV error of 0.23 found at K = 7.

So we might choose k=7. And fit KNN v/ k=7 on the entire training data set.

17