# Cross Validation and Sample Splitting

C.Conlon

February 14, 2023

Applied Econometrics II

# Cross Validation

Cross Validation appears superficially similar to bootstrap but asks a different question.

- ▶ Bootstrap tries to construct an empirical analogue to the sampling distribution of $\hat{\theta}$.

- ▶ CV tries to measure what the expected out of sample (OOS or EPE) prediction error of a new never seen before dataset.

- ▶ The main consideration is to prevent overfitting.
  - In sample fit is always going to be maximized by the most complicated model.
  - OOS fit might be a different story.
  - 1-NN might do really well in-sample, but with a new sample might perform badly.

## Sample Splitting/Holdout Method and CV

Cross Validation is actually a more complicated version of sample splitting that is one of the organizing principles in machine learning literature.

**Training Set** This is where you estimate parameter values.

**Validation Set** This is where you choose a model- a bandwidth $h$ or tuning parameter $\lambda$ by computing the error.

**Test Set** You are only allowed to look at this after you have chosen a model. Only Test Once: compute the error again on fresh data.

▶ Conventional approach is to allocate 50-80% to training and 10-20% to Validation and Test.

▶ Sometimes we don't have enough data to do this reliably.

**FIGURE 5.1.** *A schematic display of the validation set approach. A set of n observations are randomly split into a training set (shown in blue, containing observations 7, 22, and 13, among others) and a validation set (shown in beige, and containing observation 91, among others). The statistical learning method is fit on the training set, and its performance is evaluated on the validation set.*
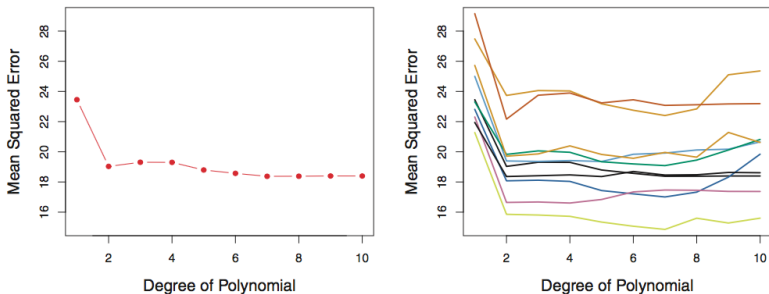
## Challenge with Sample Splitting



**FIGURE 5.2.** *The validation set approach was used on the* `Auto` *data set in order to estimate the test error that results from predicting* `mpg` *using polynomial functions of* `horsepower`. *Left: Validation error estimates for a single split into training and validation data sets. Right: The validation method was repeated ten times, each time using a different random split of the observations into a training set and a validation set. This illustrates the variability in the estimated test MSE that results from this approach.*
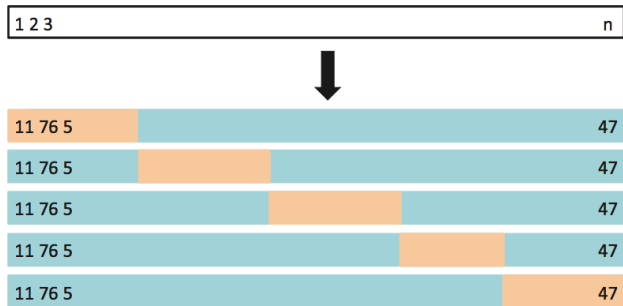
**FIGURE 5.5.** *A schematic display of 5-fold CV. A set of n observations is randomly split into five non-overlapping groups. Each of these fifths acts as a validation set (shown in beige), and the remainder as a training set (shown in blue). The test error is estimated by averaging the five resulting MSE estimates.*

# k-fold Cross Validation

▶ Break the dataset into $k$ equally sized "folds" (at random).

▶ Withhold $i = 1$ fold

- Estimate the model parameters $\hat{\theta}^{(-i)}$ on the remaining $k - 1$ folds
- Predict $\hat{y}^{(-i)}$ using $\hat{\theta}^{(-i)}$ estimates for the $i$th fold (withheld data).
- Compute $MSE_i = \frac{1}{k \cdot N} \sum_j (y_j^{(-i)} - \hat{y}_j^{(-i)})^2$.
- Repeat for $i = 1, \ldots, k$.

▶ Construct $\widehat{MSE}_{k,CV} = \frac{1}{k} \sum_i MSE_i$

**Leave One Out Cross Validation (LOOCV)**

Same as $k$-fold but with $k = N$.

- ▶ Withhold a single observation $i$
- ▶ Estimate $\hat{\theta}_{(-i)}$.
- ▶ Predict $\hat{y}_i$ using $\hat{\theta}^{(-i)}$ estimates
- ▶ Compute $MSE_i = \frac{1}{N} \sum_j (y_i - \hat{y}_i(\hat{\theta}^{(-i)}))^2$.

Note: this requires estimating the model $N$ times which can be costly.
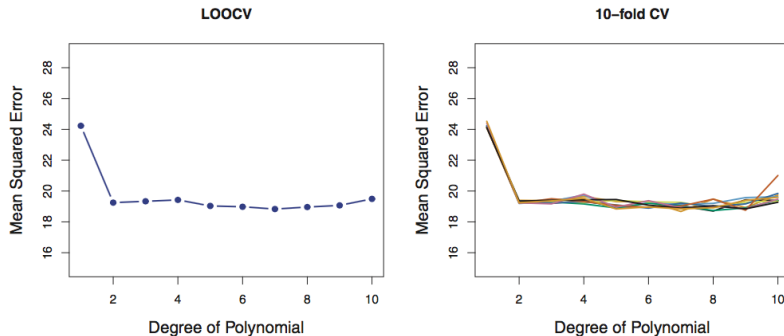
# Cross Validation



**FIGURE 5.4.** *Cross-validation was used on the* `Auto` *data set in order to estimate the test error that results from predicting* `mpg` *using polynomial functions of* `horsepower`. Left: *The LOOCV error curve.* Right: *10-fold CV was run nine separate times, each with a different random split of the data into ten parts. The figure shows the nine slightly different CV error curves.*

## Cross Validation

- Main advantage of cross validation is that we use all of the data in both estimation and in validation.
    - For our purposes validation is mostly about choosing the right bandwidth or tuning parameter.
- We have much lower variance in our estimate of the OOS mean squared error.
    - Hopefully our bandwidth choice doesn't depend on randomness of splitting sample.

## Test Data

- In Statistics/Machine learning there is a tradition to withhold 10% of the data as Test Data.
- This is completely new data that was not used in the CV procedure.
- The idea is to report the results using this test data because it most accurately simulates true OOS performance.
- We don't do much of this in economics.
  (Should we do more?)