

# Nonparametrics and Local Methods

Charlie Murry & Richard L. Sweeney

based on slides by Chris Conlon

Empirical Methods  
Spring 2020

## ① Nonparametric Density Estimation

## ② Cross-Validation

### Why nonparametrics?

Not always just interested in the mean of a (conditional) distribution.

- sometimes just interested in the distribution
- sometimes this is the first stage and we want to integrate
- sometimes want to do something semiparametric

In this section, we are interested in estimating the **density**  $f(x)$  under minimal assumptions.

## Examples from my own research

1. Estimation involves matching whole distribution (not just mean/variance).
  - Ciliberto, Tamer, and Murry – Estimating static full-info games.
  - Histogram
2. First step to recover distributions as input into structural estimation.
  - Gaurab, Murry, and Williams – Dynamic price discrimination in airline industry.
  - Kernel and NN.
3. Distribution is outcome of estimation.
  - Murry and Schurter – Estimate w.t.p. of used car purchasers. [like an auction]
  - semi-parametric regression.

## Let's start with the histogram

One of the more successful and popular uses of nonparametric methods is estimating the density or distribution function  $f(x)$  or  $F(x)$ .

$$\hat{f}_{HIST}(x_0) = \frac{1}{N} \sum_{i=1}^N \frac{\mathbf{1}(x_0 - h < x_i < x_0 + h)}{2h}$$

- Divide the dataset into bins, count up fraction of observations in each bins.
- $2h$  is the length of a bin.

Let's rewrite the histogram estimator

$$\hat{f}_{HIST}(x_0) = \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x_i - x_0}{h}\right)$$

Where  $K(z) = \frac{1}{2} \cdot \mathbf{1}(|z| < 1)$

We can think of more general forms of  $K(z; h)$ .

## Density estimator interpretation

- for each observation, there is probability mass 1 to spread around
- use the function  $K(\cdot)$  and smoothing parameter  $h$  to choose how to allocate this mass
- then, for any given  $x_0$ , sum over these functions that spread out mass, and normalize by dividing by  $N$

## Smooth Kernels

We call  $K(\cdot)$  a **Kernel function** and  $h$  the **bandwidth**. We usually assume

- i  $K(z)$  is symmetric about 0 and continuous.
- ii  $\int K(z)dz = 1$ ,  $\int zK(z)dz = 0$ ,  $\int |K(z)|dz < \infty$ .
- iii Either (a)  $K(z) = 0$  if  $|z| \geq z_0$  for some  $z_0$  or  
(b)  $|z|K(z) \rightarrow 0$  as  $|z| \rightarrow \infty$ .
- iv  $\int z^K(z)dz = \kappa$  where  $\kappa$  is a constant.



## Smooth Kernels

We call  $K(\cdot)$  a **Kernel function** and  $h$  the **bandwidth**. We usually assume

- i  $K(z)$  is symmetric about 0 and continuous.
- ii  $\int K(z)dz = 1, \int zK(z)dz = 0, \int |K(z)|dz < \infty.$
- iii Either (a)  $K(z) = 0$  if  $|z| \geq z_0$  for some  $z_0$  or  
(b)  $|z|K(z) \rightarrow 0$  as  $|z| \rightarrow \infty.$
- iv  $\int z^K(z)dz = \kappa$  where  $\kappa$  is a constant.

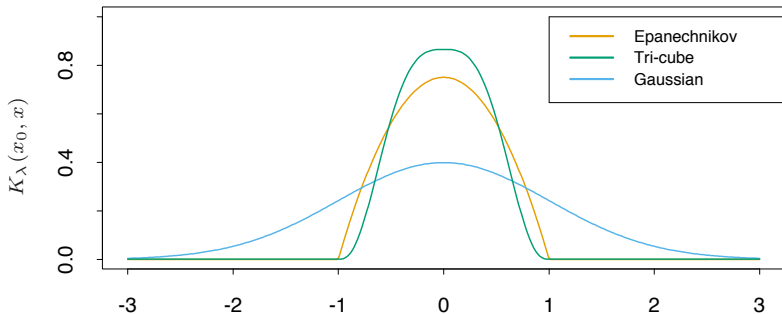
Usually we choose a smooth, symmetric  $K$ . But a common nonsmooth choice:  $K(x) = (|x| < 1/2)$  gives the *histogram* estimate.

## Some Common Kernels

**Table 9.1.** *Kernel Functions: Commonly Used Examples<sup>a</sup>*

Kernel	Kernel Function $K(z)$	$\delta$
Uniform (or box or rectangular)	$\frac{1}{2} \times \mathbf{1}( z  < 1)$	1.3510
Triangular (or triangle)	$(1 -  z ) \times \mathbf{1}( z  < 1)$	—
Epanechnikov (or quadratic)	$\frac{3}{4}(1 - z^2) \times \mathbf{1}( z  < 1)$	1.7188
Quartic (or biweight)	$\frac{15}{16}(1 - z^2)^2 \times \mathbf{1}( z  < 1)$	2.0362
Triweight	$\frac{35}{32}(1 - z^2)^3 \times \mathbf{1}( z  < 1)$	2.3122
Tricubic	$\frac{70}{81}(1 -  z ^3)^3 \times \mathbf{1}( z  < 1)$	—
Gaussian (or normal)	$(2\pi)^{-1/2} \exp(-z^2/2)$	0.7764
Fourth-order Gaussian	$\frac{1}{2}(3 - z^2)(2\pi)^{-1/2} \exp(-z^2/2)$	—
Fourth-order quartic	$\frac{15}{32}(3 - 10z^2 + 7z^4) \times \mathbf{1}( z  < 1)$	—

## Kernel Comparison



**FIGURE 6.2.** A comparison of three popular kernels for local smoothing. Each has been calibrated to integrate to 1. The tri-cube kernel is compact and has two continuous derivatives at the boundary of its support, while the Epanechnikov kernel has none. The Gaussian kernel is continuously differentiable, but has infinite support.

Mean and Variance of  $\hat{f}(x_0)$ 

Assume that the derivative of  $f(x)$  exists and is bounded, and  $\int zK(z)dz = 0$

Then the estimator has **bias**

$$b(x_0) = E \left[ \hat{f}(x_0) \right] - f(x_0) = \frac{1}{2}h^2 f''(x_0) \int z^2 K(z) dz$$

The **variance** of the estimator is

$$V \left[ \hat{f}(x_0) \right] = \frac{1}{Nh} f(x_0) \int K(z)^2 dz \left\{ +o\left(\frac{1}{Nh}\right) \right\}$$

So, unsurprisingly, the bias is *increasing* in  $h$ , and the variance is *decreasing* in  $h$ .

## How to Choose $h$

- We want both bias and variance to be as small as possible, as usual.
- In parametric estimation, it is not a problem: they both go to zero as sample size increases.
- In nonparametric estimation reducing  $h$  reduces bias, but increases variance; how are we to make his trade off?
- Note that how we set  $h$  is going to be much more important than the choice of  $K(\cdot)$

## Mean Integrated Square Error

- Start with the *local* performance at  $x_0$

$$MSE \left[ \hat{f}(x_0) \right] = E \left[ \left( \hat{f}(x_0) - f(x_0) \right)^2 \right]$$

- Calculate the *integrated* (as opposed to expected) squared error

$$\int \left( \hat{f}(x) - f(x) \right)^2 dx = \int \text{bias}^2 \left( \hat{f}(x) \right) + \text{var} \left( \hat{f}(x) \right) dx$$

- Simple approximate expression (symmetric order 2 kernels):

$$(\text{bias})^2 + \text{variance} = Ah^4 + B/nh$$

$$\text{with } A = \int (f''(x))^2 \left( \int u^2 K \right)^2 / 4 \text{ and } B = \int f(x) \int K^2$$

## Optimal bandwidth

- The AMISE is

$$Ah^4 + B/nh$$

- Minimize by taking the FOC

$$h_n^* = \left( \frac{B}{4An} \right)^{1/5}$$

- bias and standard error are *both* in  $n^{-2/5}$
- and the AMISE is  $n^{-4/5}$ —**not**  $1/n$  as it is in parametric models.
- But:  $A$  and  $B$  both depend on  $K$  (known) and  $f(y)$  (unknown), and especially “wiggleness”  $\int (f'')^2$  (unknown, not easily estimated). Where do we go from here?

## Optimal bandwidth

Can be shown that the optimal bandwidth is

$$h^* = \delta \left( \int f''(x_0)^2 dx_0 \right)^{-0.2} N^{-0.2}$$

where  $\delta$  depends on the kernel used (Silverman 1986) [these  $\delta$ 's are given in the [kernel table](#)]

Note the "optimal" kernel is Epanechnikov, although the difference is small.



## Silverman's Rule of Thumb

- If  $f$  is normal with variance  $\sigma^2$  (may not be a very appropriate benchmark!), the optimal bandwidth is

$$h_n^* = 1.06\sigma n^{-1/5}$$

- In practice, typically use **Silverman's plug-in estimate**:

$$h_n^* = 0.9 * \min(s, IQ/1.34) * n^{-1/5}$$

where IQ=interquartile distance

- Investigate changing it by a reasonable multiple.

## Silverman's Rule of Thumb

- If  $f$  is normal with variance  $\sigma^2$  (may not be a very appropriate benchmark!), the optimal bandwidth is

$$h_n^* = 1.06\sigma n^{-1/5}$$

- In practice, typically use **Silverman's plug-in estimate**:

$$h_n^* = 0.9 * \min(s, IQ/1.34) * n^{-1/5}$$

where IQ=interquartile distance

- Investigate changing it by a reasonable multiple.

This tends to work pretty well. But can we do better?

## Why not search for optimal $h$ in our data?

- Know we want to minimize MISE.
- One option is to find the  $h$  that minimizes it *in sample*
  - Loop through increments of  $h$
  - Calculate MISE
- Example: Old Faithful R data
  - Waiting time between eruptions and the duration of the eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA.
  - See R code in this folder.

## Cross-validation

- General concept in the whole of nonparametrics: choose  $h$  to minimize a criterion  $CV(h)$  that approximates

$$AMISE(h) = \int E(\hat{f}_n(x) - f(x))^2 dx.$$

- Usually programmed in metrics software. *If you can do it, do it on a subsample, and rescale.*
- 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.
  - ie 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.

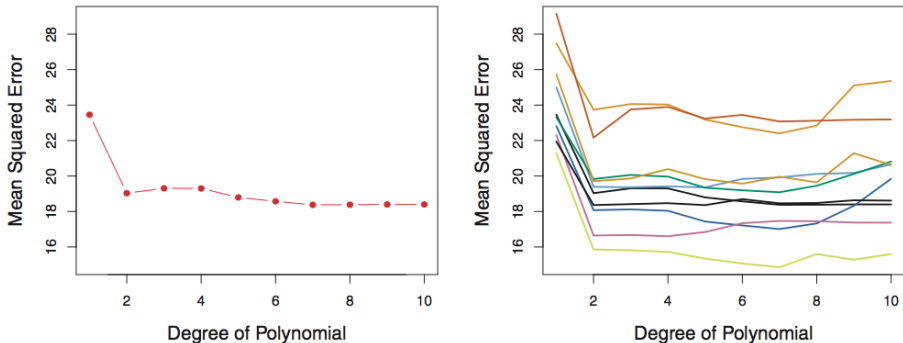
- Conventional approach is to allocate 50-80% to training and 10-20% to Validation and Test.
- Sometimes we don't have enough data to do this reliably.

## Sample Splitting/Holdout Method



**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

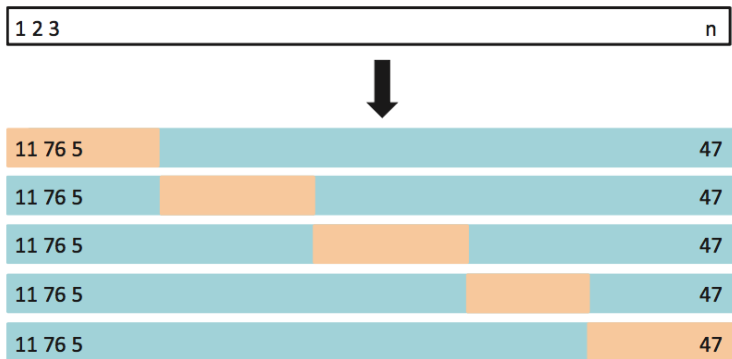


**FIGURE 5.2.** The validation set approach was used on the **Auto** data set in order to estimate the test error that results from predicting **mpg** using polynomial functions of **horsepower**. Left: Validation error estimates for a single split into training and validation data sets. Right: The validation 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.

## $k$ -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  $i = 1, \dots, k$ .
- Construct  $\widehat{MSE}_{k,CV} = \frac{1}{k} \sum_i MSE_i$



$k$ -fold Cross Validation

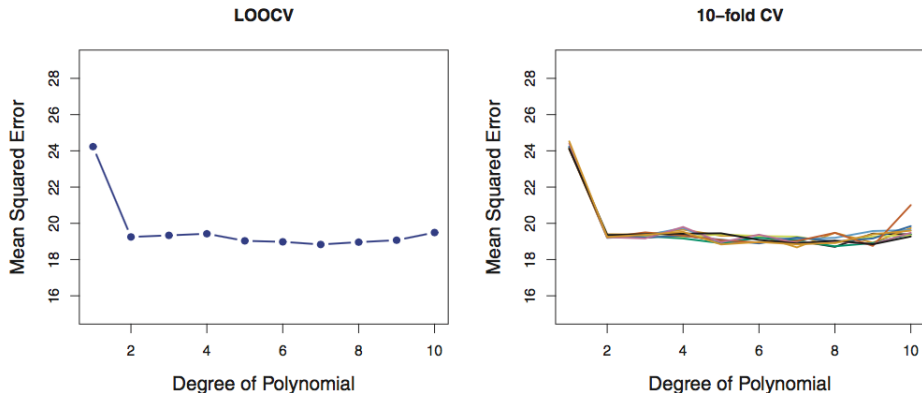
**FIGURE 5.5.** A schematic display of 5-fold CV. A set of  $n$  observations is randomly split into five non-overlapping groups. Each of these fifths acts as a validation set (shown in beige), and the remainder as a training set (shown in blue). The test error is estimated by averaging the five resulting MSE estimates.

# 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.

LOOCV vs  $k$ -fold CV

**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 25 / 27

## 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?)